



# OmniDepth: Dense Depth Estimation for Indoors Spherical Panoramas

Nikolaos Zioulis<sup>(✉)</sup>, Antonis Karakottas, Dimitrios Zarpalas,  
and Petros Daras

Centre for Research and Technology Hellas (CERTH) - Information Technologies  
Institute (ITI) - Visual Computing Lab (VCL), Thessaloniki, Greece  
{nzioulis,ankarako,zarpalas,daras}@iti.gr  
<http://vcl.iti.gr>

**Abstract.** Recent work on depth estimation up to now has only focused on projective images ignoring 360° content which is now increasingly and more easily produced. We show that monocular depth estimation models trained on traditional images produce sub-optimal results on omnidirectional images, showcasing the need for training directly on 360° datasets, which however, are hard to acquire. In this work, we circumvent the challenges associated with acquiring high quality 360° datasets with ground truth depth annotations, by re-using recently released large scale 3D datasets and re-purposing them to 360° via rendering. This dataset, which is considerably larger than similar projective datasets, is publicly offered to the community to enable future research in this direction. We use this dataset to learn in an end-to-end fashion the task of depth estimation from 360° images. We show promising results in our synthesized data as well as in unseen realistic images.

**Keywords:** Omnidirectional media · 360° · Spherical panorama  
Scene understanding · Depth estimation · Synthetic dataset  
Learning with virtual data

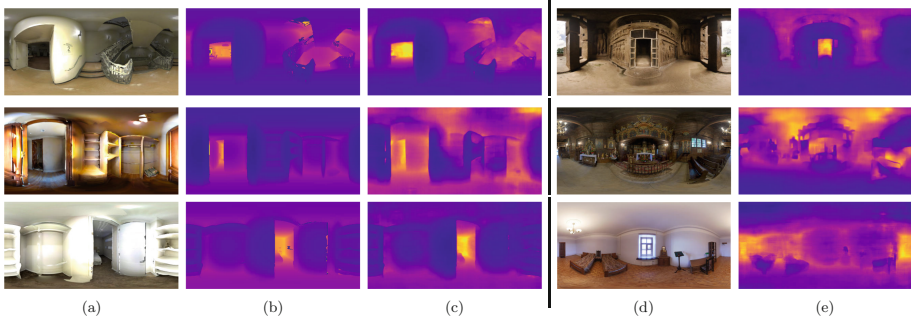
## 1 Introduction

One of the fundamental challenges in computer and 3D vision is the estimation of a scene's depth. Depth estimation leads to a three-dimensional understanding of the world which is very important to numerous applications. These vary from creating 3D maps [1] and allowing navigation in real-world environments [2], to enabling stereoscopic rendering [3], synthesizing novel views out of pre-captured content [4] and even compositing 3D objects into it [5]. Depth information has been shown to boost the effectiveness of many vision tasks related to scene understanding when utilized jointly with color information [6, 7].

---

N. Zioulis and A. Karakottas—Equal contribution.

**Electronic supplementary material** The online version of this chapter ([https://doi.org/10.1007/978-3-030-01231-1\\_28](https://doi.org/10.1007/978-3-030-01231-1_28)) contains supplementary material, which is available to authorized users.



**Fig. 1.** We learn the task of predicting depth directly from omnidirectional indoor scene images. Results from our RectNet model are presented (left to right): (a) 360° image samples from our test set, (b) corresponding ground truth depth, (c) predicted depth maps of the test image samples, (d) 360° unseen image samples from the Sun360 dataset, (e) predicted depth maps of the Sun360 image samples.

Similar to how babies start to perceive depth from two viewpoints and then by ego-motion and observation of objects’ motions, researchers have tackled the problem of estimating depth via methods built on multi-view consistency [8,9] and structure-from-motion (SfM) [10]. But humans are also driven by past experiences and contextual similarities and apply this collective knowledge when presented with new scenes. Likewise, with the advent of more effective machine learning techniques, recent research focuses on learning to predict depth and has led to impressive results even with completely unsupervised learning approaches.

However, learning based approaches have only focused on traditional 2D content captured by typical pinhole projection model based cameras. With the emergence of efficient spherical cameras and rigs, omnidirectional (360°) content is now more easily and consistently produced and is witnessing increased adoption in entertainment and marketing productions, robotics and vehicular applications as well as coverage of events and even journalism. Consumers can now experience 360° content in mobile phones, desktops and, more importantly, the new arising medium – Virtual Reality (VR) – headsets.

Depth and/or geometry extraction from omnidirectional content has been approached similar to traditional 2D content via omnidirectional stereo [11–14] and SfM [4] analytical solutions. There are inherent problems though to applying learning based methods to 360° content as a result of its acquisition process that inhibits the creation of high quality datasets. Coupling them with 360° LIDARs would produce low resolution depths and would also insert the depth sensor into the content itself, a drawback that also exists when aiming to acquire stereo datasets. One alternative would be to manually re-position the camera but that would be tedious and error prone as a consistent baseline would not be possible.

In this work, we train a CNN to learn to estimate a scene’s depth given an omnidirectional (equirectangular) image as input<sup>1</sup> (Fig. 1). To circumvent the lack of available training data we resort to re-using existing 3D datasets and repurposing them for use within a 360° context. This is accomplished by generating diverse 360° views via rendering. We use this dataset for learning to infer depth from omnidirectional content. In summary, our contributions are the following:

1. We present the first, to the authors’ knowledge, learning based dense depth estimation method that was trained with and operates directly on omnidirectional content in the form of equirectangular images.
2. We offer a dataset consisting of 360° color images paired with ground truth 360° depth maps in equirectangular format. The dataset is available online<sup>2</sup>.
3. We propose and validate, a CNN auto-encoder architecture specifically designed for estimating depth directly on equirectangular images.
4. We show how monocular depth estimation methods trained on traditional 2D images fall short or produce low quality results when applied to equirectangular inputs, highlighting the need for learning directly on the 360° domain.

## 2 Related Work

Since this work aims to learn the task of omnidirectional dense depth estimation, and given that - to the authors’ knowledge - no other similar work exists, we first review non-learning based methods for geometric scene understanding based on 360° images. We then examine learning based approaches for spherical content and, finally, present recent monocular dense depth estimation methods.

### 2.1 Geometric Understanding on 360° Images

Similar to pinhole projection model cameras, the same multi-view geometry [8] principles apply to 360° images. By observing the scene from multiple viewpoints and establishing correspondences between them, the underlying geometrical structure can be estimated. For 360° cameras the conventional binocular or multi-view stereo [9] problem is reformulated to binocular or multi-view spherical stereo [11] respectively, by taking into account the different projection model and after defining the disparity as angular displacements. By estimating the disparity (i.e. depth), complete scenes can be 3D reconstructed from multiple [14, 15] or even just two [12, 13] spherical viewpoints. However, all these approaches require multiple 360° images to estimate the scene’s geometry. Recently it was shown that 360° videos acquired with a moving camera can be used to 3D reconstruct a scene’s geometry via SfM [4] and enable 6 DOF viewing in VR headsets.

There are also approaches that require only a single image to understand indoors scenes and estimate their layout. PanoContext [16], generates a 3D room

<sup>1</sup> We use the terms omnidirectional image, 360° image, spherical panorama and equirectangular image interchangeably in this document.

<sup>2</sup> <http://vcl.itl.gr/360-dataset/>.

layout hypothesis given an indoor  $360^\circ$  image in equirectangular format. With the estimations being bounding boxes, the inferred geometry is only a coarse approximation of the scene. Similar in spirit, the work of Yang et al. [17] generates complete room layouts from panoramic indoor images by combining super-pixel information, vanishing points estimation and a geometric context prior under a Manhattan world assumption. However, focusing on room layout estimation, it is unable to recover finer details and structures of the scene. Another similar approach [18] addresses the problem of geometric scene understanding from another perspective. Under a maximum a posteriori estimation it unifies semantic, pose and location cues to generate CAD models of the observed scenes. Finally, in [19] a spherical stereo pair is used to estimate both the room layout but also object and material attributes. After retrieving the scene's depth by stereo matching and subsequently calculating the normals, the equirectangular image is projected to the faces of a cube that are then fed to a CNN whose object predictions are fused into the  $360^\circ$  image to finally reconstruct the 3D layout.

## 2.2 Learning for $360^\circ$ Images

One of the first approaches to estimate distances purely from omnidirectional input [20] under a machine learning setting utilized Gaussian processes. Instead of estimating the distance of each pixel, a range value per image column was predicted to drive robotic navigation. Nowadays, with the establishment of CNNs, there are two straightforward ways to apply current CNN processing pipelines to spherical input. Either directly on a projected (typically equirectangular) image, or by projecting the spherical content to the faces of a cube (cubemap) and running the CNN predictions on them, which are then merged by back-projecting them to the spherical domain. The latter approach was selected by an artistic style transfer work [21], where each face was re-styled separately and then the cubemap was re-mapped back to the equirectangular domain. Likewise, in SalNet360 [22], saliency predictions on the cube's faces are refined using their spherical coordinates and then merged back to  $360^\circ$ . The former approach, applying a CNN directly to the equirectangular image, was opted for in [23] to increase the dynamic range of outdoor panoramas.

More recently, new techniques for applying CNNs to omnidirectional input were presented. Given the difficulty to model the projection's distortion directly in typical CNNs as well as achieve invariance to the viewpoint's rotation, the alternative pursued by [24] is based on graph-based deep learning. Specifically they model distortion directly into the graph's structure and apply it to a classification task. A novel approach taken in [25] is to learn appropriate convolution weights for equirectangular projected spherical images by transferring them from an existing network trained on traditional 2D images. This conversion from the 2D to the  $360^\circ$  domain is accomplished by enforcing consistency between the predictions of the 2D projected views and those in the  $360^\circ$  image. Moreover, recent work on convolutions [26, 27] that in addition to learning their weights also learn their shape, are very well suited for learning the distortion model of

spherical images, even though they have only been applied to fisheye lenses up to now [28]. Finally, very recently, Spherical CNNs were proposed in [29,30] that are based in a rotation-equivariant definition of spherical cross-correlation. However these were only demonstrated in classification and single variable regression problems. In addition, they are also applied in the spectral domain while we formulate our network design for the spatial image domain.

### 2.3 Monocular Depth Estimation

Depth estimation from monocular input has attracted lots of interest lately. While there are some impressive non learning based approaches [31–33], they come with their limitations, namely reliance on optical flow and relevance of the training dataset. Still, most recent research has focused on machine learning to address the ill-posed depth estimation problem. Initially, the work of Eigen et al. [34] trained a CNN in a coarse-to-fine scheme using direct depth supervision from RGB-D images. In a subsequent continuation of their work [6], they trained a multi-task network that among predicting semantic labels and normals, also estimated a scene’s depth. Their results showed that jointly learning the tasks achieved higher performance due to their complementarity. In a recent similar work [35], a multi-task network that among other modalities also estimated depth, was trained using synthetic data and a domain adaptation loss based on adversarial learning, to increase its robustness when running on real scenes. Laina et al. [36] designed a directly supervised fully convolutional residual network (FCRN) with novel up-projection blocks that achieved impressive results for indoor scenes and was also used in a SLAM pipeline [1].

Another body of work focused on applying Conditional Random Fields (CRFs) to the depth estimation problem. Initially, the output of a deep network was refined using a hierarchical CRF [37], with Liu et al. [38] further exploring the interplay between CNNs and CRFs for depth estimation in their work. Recently, multi-scale CRFs were used and trained in an end-to-end manner along with the CNN [39]. Dense depth estimation has also been addressed as a classification problem. Since perfect regression is usually impossible, dense probabilities were estimated in [40] and then optimized to estimate the final depth map. Similarly, in [41] and [42] depth values were discretized in bins and densely classified, to be afterwards refined either via a hierarchical fusion scheme or through the use of a CRF respectively. Taking a step further, a regression-classification cascaded network was proposed in [43] where a low spatial resolution depth map was regressed and then refined by a classification branch.

The concurrent works of Garg et al. [44] and Godard et al. [45] showed that unsupervised learning of the depth estimation task is possible. This is accomplished by an intermediate task, view synthesis, and allowed training by only using stereo pair input with known baselines. In a similar fashion, using view synthesis as the main supervisory signal, learning to estimate depth was also achieved by training with pure video sequences in a completely unsupervised manner [46–50]. Another novel unsupervised depth estimation method relies on

aperture supervision [51] by simply acquiring training data in various focus levels. Finally, in [52] it was shown that a CNN can be trained to estimate depth from monocular input with only relative depth annotations.

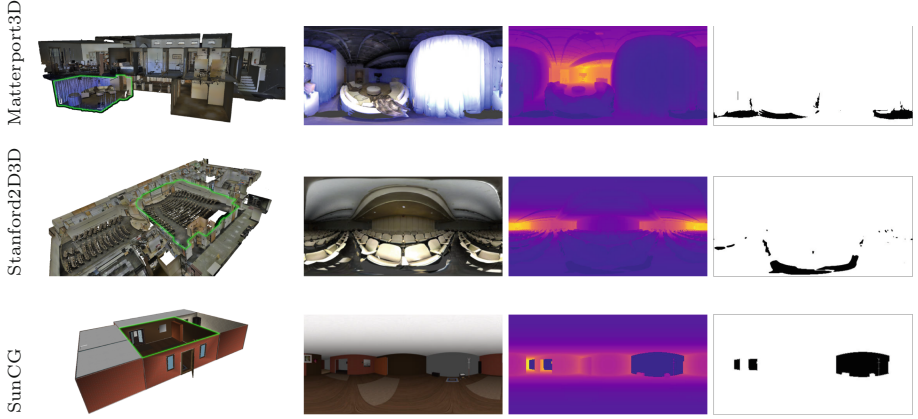
### 3 Synthesizing Data

End-to-end training of deep networks requires a large amount of annotated ground truth data. While for typical pinhole camera datasets this was partly addressed by using depth sensors [53] or laser scanners [54] such an approach is impractical for spherical images due to a larger diversity in resolution for 360° cameras and laser scanners, and because each 360° sensor would be visible from the other one. As much as approaches like the one employed in [55] could be used to in-paint the sensor regions, these would still be the result of an algorithmic process and not the acquisition process itself, potentially introducing errors and artifacts that would reduce the quality of the data. This also applies to unsupervised stereo approaches that require the simultaneous capture of the scene from two viewpoints. Although one could re-position the same sensor to acquire clean panoramas, a consistent baseline would not be possible. More recently, unsupervised approaches for inferring a scene’s depth have emerged that are trained with video sequences. However, they assume a moving camera as they rely on view synthesis as the supervisory signal which is not a typical setting for indoors 360° captures, but for action camera like recordings.

**360D Dataset:** Instead, we rely on generating a dataset with ground truth depth by synthesizing both the color and the depth image via rendering. To accomplish that we leverage the latest efforts in creating publicly available textured 3D datasets of indoors scenes. Specifically, we use two computer generated (CG) datasets, SunCG [56] and SceneNet [57], and two realistic ones acquired by scanning indoor buildings, Stanford2D3D [58, 59] and Matterport3D [60]. We use a path tracing renderer<sup>3</sup> to render our dataset by placing a spherical camera and a uniform point light at the same position  $\mathbf{c} \in \mathbb{R}^3$  in the scene. We then acquire the rendered image  $I(\mathbf{p}) \in \mathbb{R}$ ,  $\mathbf{p} = (u, v) \in \mathbb{N}^2$ , as well as the underlying  $z$ -buffer that was generated as a result of the graphics rendering process, that serves as the ground truth depth  $D(\mathbf{p}) \in \mathbb{R}$ . It should be noted that unlike pinhole camera model images, the  $z$ -buffer in this case does not contain the  $z$  coordinate value of the 3D point  $\mathbf{v}(\mathbf{p}) \in \mathbb{R}^3$ , corresponding to pixel  $\mathbf{p}$ , but instead the 3D point’s radius  $r = \|\mathbf{v} - \mathbf{c}\|$ , in the camera’s spherical coordinate system.

For the two CG indoors datasets we place the camera and the light at the center of each house, while for the two scanned indoors datasets we use the pose information available (estimated during the scanning process) and thus, for each building we generate multiple 360° data samples. Given that the latter two datasets were scanned, their geometries contain holes or inaccurate/coarse estimations, and also have lighting information baked into the models. On the other hand, the CG datasets contain perfect per pixel depth but lack the realism

<sup>3</sup> <https://www.cycles-renderer.org>.



**Fig. 2.** Example renders from our dataset, from left to right: the 3D building with a green highlight denoting the rendered scene, color output, corresponding depth map, and the binary mask depicting the missing regions in black. (Color figure online)

of the scanned datasets, creating a complementary mix. However, as no scanning poses are available, the centered poses may sometimes be placed within or on top of objects and we also observed missing information in some scenes (e.g. walls/ceilings) that, given SunCG’s size, are impractical to manually correct.

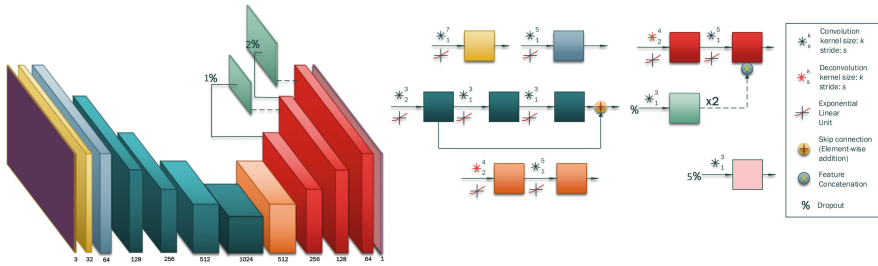
For each pose, we augment the dataset by rotating the camera in  $90^\circ$  resulting in 4 distinct viewpoints per pose sample. Given the size of SunCG, we only utilize a subset of it and end up using **11118** houses, resulting in a 24.36% utilization. The remaining three datasets are completely rendered. This results in a total of **88384** renders and **22096** unique viewpoints. Our generated *360D* dataset contains a mix of synthetic and realistic  $360^\circ$  color *I* and depth *D* image data in a variety of indoors contexts (houses, offices, educational spaces, different room layouts) and is publicly available at <http://vcl.iti.gr/360-dataset/>.

## 4 Omnidirectional Depth Estimation

The majority of recent CNN architectures for dense estimation follow the autoencoder structure, in which an encoder encodes the input, by progressively decreasing its spatial dimensions, to a representation of much smaller size, and a decoder that regresses to the desired output by upscaling this representation.

We use two encoder-decoder network architectures that are structured differently. The first resembles those found in similar works in the literature [36, 45], while the second is designed from scratch to be more suitable for learning with  $360^\circ$  images. Both networks are fully convolutional [61] and predict an equirectangular depth map with the only input being a  $360^\circ$  color image in equirectangular format. We use ELUs [62] as the activation function which also remove the need for batch normalization [63] and its added computational complexity.





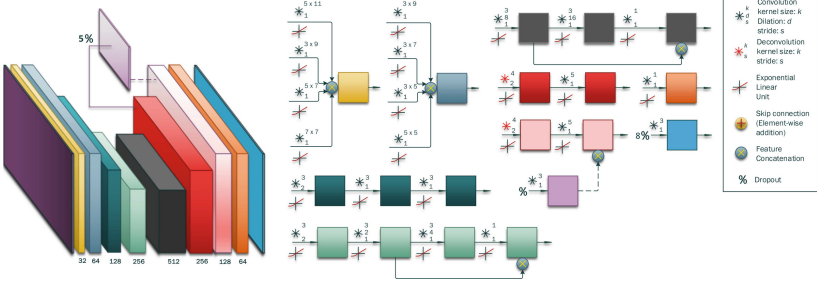
**Fig. 3.** UResNet Architecture: the encoder consists of two input preprocessing blocks, and four down-scaling blocks (dark green). The former are single convolutional (conv) layers while the latter consist of a strided conv and two more regular convs with a skip/residual connection. The decoder contains one upscaling block (orange) and three up-prediction blocks (red), followed by the prediction layer (pink). Up-scaling is achieved with a strided deconv followed by a conv, and similarly, up-predictions additionally branch out to estimate a depth prediction at the corresponding scale with an extra conv that is concatenated with the features of the next block’s last layer. (Color figure online)

**UResNet:** In this unbalanced ResNet, the encoding and decoding parts are not symmetrical, with the decoder being shallower. The encoder is built with skip connections [64], a technique that helps when training deeper architectures by preventing gradient degradation, allowing for larger receptive fields. More detailed architectural information is presented in Fig. 3 where the network is decomposed into processing blocks.

**RectNet:** Omnidirectional images differ from traditional images in the sense that they capture global (full  $360^\circ$ ) visual information and, when in equirect-angular format, suffer from high distortions along their  $y$  (i.e. latitude) axis. Therefore, the second architecture’s design aims to exploit and address these properties of spherical panoramas while keeping some of the desirable properties of UResNet like skip connections. Capturing the  $360^\circ$  image’s global context is achieved by increasing the effective receptive field (RF) of each neuron by utilizing dilated convolutions [65]. Instead of progressive downscaling as in most depth estimation networks and similarly UResNet, we only drop the spatial dimensions by a factor of 4. Then, inspired by [66], we use progressively increasing dilations to increase the RF to about half the input’s spatial dimensions and increase global scene understanding. In addition, within each dilation block we utilize  $1 \times 1$  convolutions to reduce the spatial correlations of the feature maps.

The distortion factor of spherical panoramas increases towards the sphere’s poles and is therefore different for every image row. This means that information is scattered horizontally, as we vertically approach the two poles. In order to account for this varying distortion we alter our input blocks, as their features are closer to natural image ones (e.g. edges). Following [25], where 2D CNN filters are transferred into distorted (practically rectangular) row-wise versions to increase performance when applied to the  $360^\circ$  domain, we use rectangle filters





**Fig. 4.** RectNet Architecture: the encoder consists of two preprocessing blocks (yellow and blue) and a downscaling block (dark green), followed by two increasing dilation blocks (light green and black). The preprocessing blocks concatenate features produced by convolutions (convs) with different filter sizes, accounting for the equirectangular projection’s varying distortion factor. The down-scaling block comprises a strided and two regular convs. (Color figure online)

along with traditional square filters and vary the resolution of the rectangle ones to account for different distortion levels. However, this variation is done while also preserving the area of the filter to be as close as possible to the original square filter’s. The outputs of the rectangle and square filters are concatenated while preserving the overall output feature count. The detailed architecture is presented in Fig. 4.

**Training Loss:** Given that we synthesize perfect ground truth depth annotations  $D_{gt}$ , as presented in Sect. 3, we take a completely supervised approach. Even though most approaches using synthetic data fail to generalize to realistic input, our dataset contains an interesting mix of synthetic (CAD) renders as well as realistic ones. The scanned data are acquired from real environments and, as a result, their renders are very realistic. Following previous work, we predict depth  $D_{pred}^s$  against downsampled versions of the ground truth data  $D_{gt}^s$  at multiple scales (with  $s$  being the downscaling factor) and upsample these predictions using nearest neighbor interpolation to later concatenate them with the subsequent higher spatial dimension feature maps. We also use the dropout technique [67] in those layers used to produce each prediction. Further, we use L2 loss for regressing the dense depth output  $E_{depth}(\mathbf{p}) = \|D_{gt}(\mathbf{p}) - D_{pred}(\mathbf{p})\|^2$  and additionally add a smoothness term  $E_{smooth}(\mathbf{p}) = \|\nabla D(\mathbf{p})\|^2$  for the predicted depth map by minimizing its gradient.

Although our rendered depth maps are accurate in terms of depth, in practice there are missing regions in the rendered output. These are either because of missing information in the CAD models (e.g. walls/ceilings) or the imperfect process of large scale 3D scanning, with visual examples illustrated in Fig. 2. These missing regions/holes manifest as a specific color (“clear color”), selected during rendering, in the rendered image and as infinity (“far”) values in the rendered depth map. As these outlier values will destabilize the training process,

**Table 1.** Quantitative results of our networks for 360° dense depth estimation.

Network	Tested on	Abs Rel ↓	Sq Rel ↓	RMS ↓	RMSlog ↓	$\delta < 1.25$ ↑	$\delta < 1.2^2$ ↑	$\delta < 1.25^3$ ↑
UResNet	Test set	0.0835	0.0416	0.3374	0.1204	0.9319	0.9889	0.9968
RectNet	Test set	<b>0.0702</b>	<b>0.0297</b>	<b>0.2911</b>	<b>0.1017</b>	<b>0.9574</b>	<b>0.9933</b>	<b>0.9979</b>
UResNet	SceneNet	0.1218	0.0727	0.4066	0.1538	0.8598	0.9815	0.9962
RectNet	SceneNet	0.1077	0.699	0.3572	0.1386	0.8965	0.9879	0.9971
UResNet -S2R	Stanford	0.1226	0.0768	0.489	0.1667	0.8593	0.9756	0.9942
RectNet -S2R	Stanford	0.0824	0.0457	0.3998	0.1229	0.928	0.9879	0.9971
UResNet -S2R	SceneNet	0.1448	0.0991	0.517	0.1792	0.7898	0.9761	0.9935
RectNet -S2R	SceneNet	0.1079	0.0644	0.3778	0.1404	0.8966	0.9866	0.996

we ignore them during backpropagation by using a per pixel  $\mathbf{p}$  binary mask  $M(\mathbf{p})$  that is zero in these missing regions. This allows us to train the network even with incomplete or slightly inaccurate/erroneous 3D models. Thus, our final loss function is:

$$E_{loss}(\mathbf{p}) = \sum_s \alpha_s M(\mathbf{p}) E_{depth}(\mathbf{p}) + \sum_s \beta_s M(\mathbf{p}) E_{smooth}(\mathbf{p}), \quad (1)$$

where  $\alpha_s, \beta_s$  are the weights for each scale of the depth and smoothing term.

## 5 Results

We evaluate the performance of our two 360° depth estimation networks by first conducting an intra assessment of the two models and then offering quantitative comparisons with other depth estimation methods. Finally, we present comparative qualitative results in unseen, realistic data of everyday scenes.

**Training Details:** Our networks are trained using Caffe [68] on a single NVIDIA Titan X. We use Xavier weight initialization [69] and ADAM [70] as the optimizer with its default parameters  $[\beta_1, \beta_2, \epsilon] = [0.9, 0.999, 10^{-8}]$  and an initial learning rate of 0.0002. Our input dimensions are  $512 \times 256$  and are given in equirectangular format, with our depth predictions being equal sized.

We split our dataset into corresponding train and tests sets as follows: (i) Initially we remove 1 complete area from Stanford2D3D, 3 complete buildings from Matterport3D and 3 CAD scenes from SunCG for our test set totaling 1,298 samples. (ii) We skip SceneNet entirely and use it as our validation set. (iii) Then, from the remaining SunCG, Stanford2D3D and Matterport3D samples we automatically remove scenes which contain regions with very large or small depth values ( $>5\%$  of total image area above  $20m$  or under  $0.5m$ ). Finally, we are left with a train-set that consists of 34,679<sup>4</sup> RGB 360° images along with their corresponding ground truth depth map annotations. Our loss weights for UResNet are  $[\alpha_1, \alpha_2, \alpha_4, \beta_1] = [0.445, 0.275, 0.13, 0.15]$ , and for RectNet they are

<sup>4</sup> Only a subset of SunCG was used by prioritizing larger scenes given the length of the rendering process. However, a larger subset is publicly available.

**Table 2.** Quantitative results against other monocular depth estimation models.

Network		Abs Rel↓	Sq Rel↓	RMS↓	RMS(log)↓	$\delta < 1.25$ ↑	$\delta < 1.25^2$ ↑	$\delta < 1.25^3$ ↑
UResNet		0.0835	0.0416	0.3374	0.1204	0.9319	0.9889	0.9968
RectNet		<b>0.0702</b>	<b>0.0297</b>	<b>0.2911</b>	<b>0.1017</b>	<b>0.9574</b>	<b>0.9933</b>	<b>0.9979</b>
Equirect.	Godard et al. [45]	0.4747	2.3783	7.2097	0.82	0.297	0.79	0.751
	Laina et al. [36]	0.3181	0.4469	0.941	0.376	0.4922	0.7792	0.915
	Liu et al. [38]	0.4202	0.7597	1.1596	0.44	0.3889	0.7044	0.8774
Cubemap	Godard et al. [45]	0.2552	0.9864	4.4524	0.5087	0.3096	0.5506	0.7202
	Laina et al. [36]	0.1423	0.2544	0.7751	0.2497	0.5198	0.8032	0.9175
	Liu et al. [38]	0.1869	0.4076	0.9243	0.2961	0.424	0.7148	0.8705

$[\alpha_1, \alpha_2, \beta_1, \beta_2] = [0.535, 0.272, 0.134, 0.068]$ . For quantitative evaluation we use the same error metrics as previous works [6, 34, 36, 38, 45] (arrows next to each metric in the tables denote the direction of better performance).

**Model Performance:** Table 1 presents the results of our two models in our test set, and in the unseen synthetic SceneNet generated data, after training for 10 epochs in all of our train set. We observe that RectNet – which was designed with  $360^\circ$  input in mind – performs better than the standard UResNet even with far fewer parameters ( $\sim 8.8M$  vs  $\sim 51.2M$ ). In order to assess their efficacy and generalization capabilities we perform a leave-one-out evaluation. We train both networks initially only in the synthetic SunCG generated data for 10 epochs, and then finetune them in the realistic Matterport3D generated data for another 10 epochs. This train is suffixed with “-S2R”. We then evaluate them in the entirety of the Stanford2D3D generated dataset, as well as in the SceneNet one. Comparable results to the previous train with all datasets are observed. Again, RectNet outperforms UResNet – albeit both perform slightly worse as expected due to being trained with less amount of data.

The increased performance of RectNet against UResNet in every error metric or accuracy, can be attributed to its larger RF, which for  $360^\circ$  images is very important as it allows the network to capture the global context more efficiently<sup>5</sup>. Despite the fact that UResNet is much deeper than RectNet and significantly drops the input’s spatial dimensions, RectNet still achieves a larger receptive field. Specifically, UResNet has a  $190 \times 190$  RF compared to that of RectNet which is  $266 \times 276$ . In addition, RectNet drops the input’s spatial dimensions only by a factor of 4, maintaining denser information in the extracted features.

**Comparison Against Other Methods:** Given that there are no other methods to perform dense depth estimation for  $360^\circ$  images, we assess its performance against the state of the art in monocular depth estimation models. Since the predictions of these methods are defined in different scales, we scale the estimated depth maps by a scalar  $\tilde{s}$ , which matches their median with our ground truth like [46], i.e.  $\tilde{s} = \text{median}(D_{gt}) / \text{median}(D_{pred})$ . Moreover, we evaluate the

<sup>5</sup> Varying RF experiments supporting this claim can be found in the supplement.

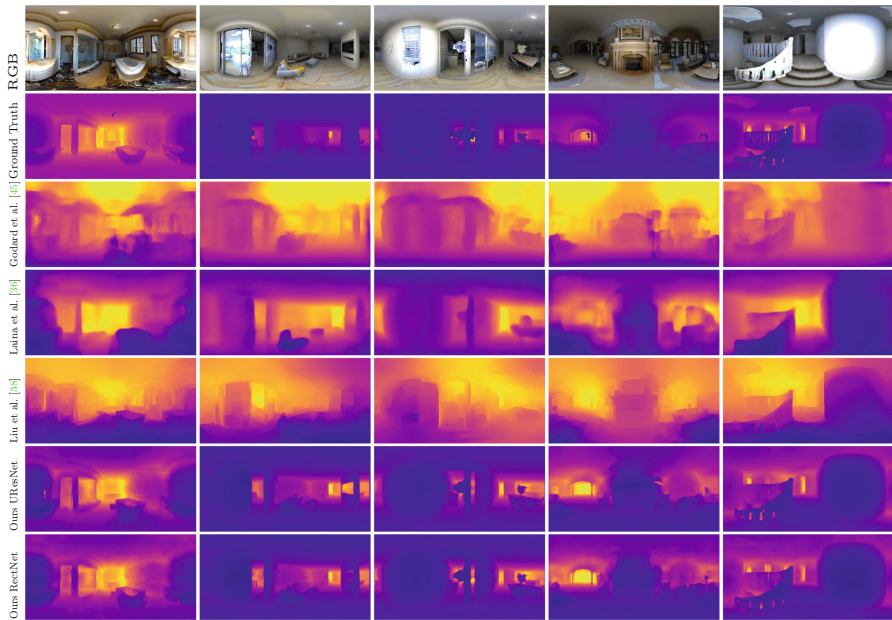


Fig. 5. Qualitative results on our test split.

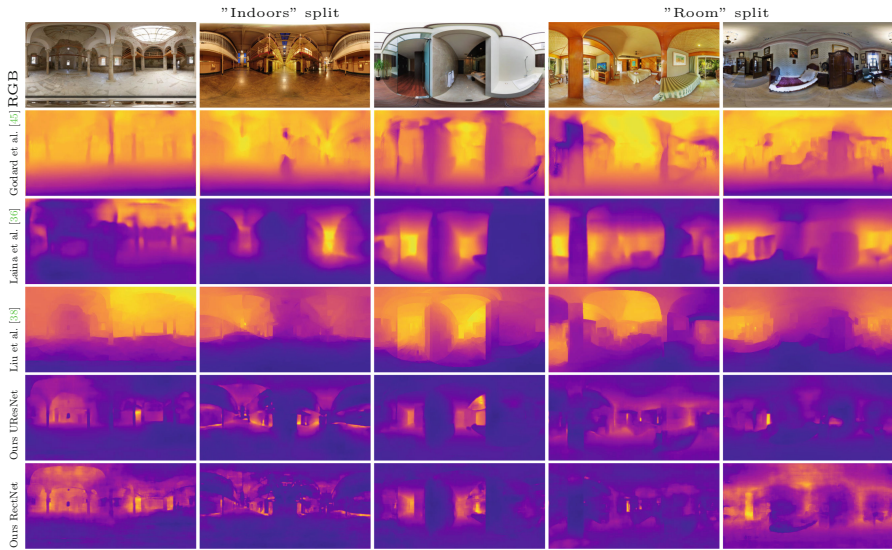
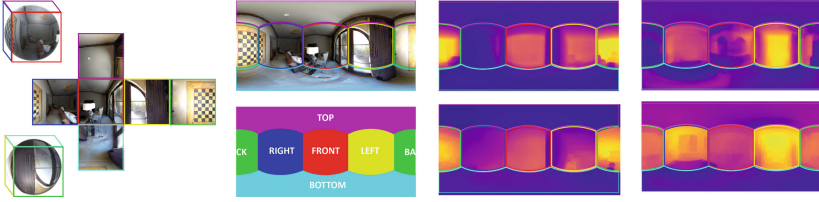


Fig. 6. Qualitative results on the "Room" and "Indoors" Sun360 splits.

**Table 3.** Per cube face quantitative results against other monocular models.

Network	AbsRel ↓	SqRel ↓	RMSE ↓	RMSElog ↓	$\delta < 1.25 \uparrow$	$\delta < 1.25^2 \uparrow$	$\delta N1.25^3 \uparrow$
UResNet	0.0097	0.0062	0.1289	0.041	0.9245	0.9853	0.9955
RectNet	<b>0.008</b>	<b>0.0042</b>	<b>0.1113</b>	<b>0.03504</b>	<b>0.9497</b>	<b>0.9907</b>	<b>0.9969</b>
Godard et al. [45]	0.0453	0.1743	1.6559	0.1958	0.4524	0.7023	0.8315
Laina et al. [36]	0.03	0.0549	0.3152	0.1033	0.6353	0.8616	0.9412
Liu et al. [38]	0.0312	0.0532	0.3048	0.107	0.603	0.8412	0.9338

**Fig. 7.** Cubemap projection (left) and merged monocular predictions (right).

masked depth maps as mentioned in Sect. 3 in order to ignore the missing values. Table 2 presents the results of state-of-the-art methods when applied directly on our test split in the equirectangular domain (with a sample of qualitative results presented in Fig. 5). We offer results for the model of Laina et al. [36], trained with direct depth supervision in indoor scenes, Godard et al. [45], trained in an unsupervised manner in outdoor driving scenes using calibrated stereo pairs, and the method of Liu et al. [38], which combines learning with CRFs and is trained in indoor scenes. As observed by the results, the performance of all the methods directly on equirectangular images is poor, and our main models outperform them. However, inferior performance is expected as these were not trained directly in the equirectangular domain but in perspective images. Nonetheless, Laina et al. [36] and Liu et al. [38] achieve much better results than Godard et al. [45]. This is also expected as the latter is trained in an outdoor setting, with very different statistics than our indoor dataset.

For a more fair comparison we use a cubemap projection (Fig. 7 (left)) of all spherical images and then run each model on the projected cube faces which are typical perspective images. After acquiring the predictions, we merge all cube faces' depth maps by projecting them back to the equirectangular domain to be evaluated. However, since the top and bottom cube face projections will be mostly planar, we ignore them during evaluation of all metrics. While monocular performance is improved compared to when applied directly to equirectangular images, their quantitative performance is still inferior to our models. Further, the runtime performance is also worse as multiple inferences need to run, one for each face, incurring a much higher computational cost. Moreover, another apparent issue is the lack of consistency between the predictions of each face. This is shown in Fig. 7 (right) where it is clear that the depth scales of each face are different. This is in line with the observations in [21], but is more pronounced

in the depth estimation case, than the style transfer one. Based on this observation, we evaluate each cube face separately against the ground truth values of that face alone which is also median scaled separately. The average values of the front, back, right and left faces for each monocular model against the obtained by our models on the same faces alone are presented in Table 3. Although the performance of the monocular models is further improved, our models still perform better. This can be attributed to various reasons besides training directly on equirectangular domain. One explanation is that  $360^\circ$  images capture global information which can better help reasoning about relative depth and overall increase inference performance. The other is that our generated dataset is considerably larger and more diverse than other indoor datasets. In addition, the cube faces are projected out of  $512 \times 256$  images and are thus, of lower quality/resolution than typical images these models were trained in.

**Qualitative Results:** To determine how well our models generalize, we examine their performance on completely unseen data found in the Sun360 dataset [71], where no ground truth depth is available. The Sun360 dataset comprises realistic environment captures and has also been used in the work of Yang et al. [17] for room layout estimation. We offer some qualitative results on a data split from [17], referred to as “Room”, as well as an additional split of indoor scenes that we select from the Sun360 dataset, referred to as “Indoors”. These are presented in Fig. 6 for our two models as well as the monocular ones that were quantitatively evaluated. Our models are able to estimate the scenes’ depth with the only monocular model to produce plausible results being the one of Laina et al. [36]. We also observe that UResNet offers smoother predictions than the better performing RectNet, unlike the results obtained on our test split. More qualitative results can be found in the supplementary material where comparison with the method of Yang et al. [17] is also offered.

## 6 Conclusions

We have presented a learning framework to estimate a scene’s depth from a single  $360^\circ$  image. Our models were trained in a completely supervised manner with ground truth depth. To accomplish this, we overcame the dataset unavailability and difficulty in acquisition for paired  $360^\circ$  color and depth image pairs. This was achieved by re-using 3D datasets with both synthetic and real-world scanned indoors scenes and synthesizing a  $360^\circ$  dataset via rendering.  $360^\circ$  depth information can be useful for a variety of tasks, like in adding automation in the composition of 3D elements within spherical content [72].

Since our approach is the first work for dense  $360^\circ$  depth estimation, there are many challenges that still need to be overcome. Our datasets only cover indoor cases, limiting the networks’ applicability to outdoor settings, and are generated with perfect camera vertical alignment with constant lighting and no stitching artifacts. This issue is further accentuated as the scanned datasets had lighting information baked into them during scanning. This can potentially



hamper robustness when applied in real world conditions that also contain a much higher dynamic range of luminosity.

For future work, we want to explore unsupervised learning approaches that are based on view synthesis as the supervisory signal. Furthermore, robustness to real world scenes can be achieved, either by utilizing GANs as generators of realistic content, or by using a discriminator to identify plausible/real images.

**Acknowledgements.** This work was supported and received funding from the European Union Horizon H2020 Framework Programme funded project Hyper360, under Grant Agreement no. 761934 (<http://www.hyper360.eu/>). We are also grateful and acknowledge the support of NVIDIA for a hardware donation.

## References

1. Tateno, K., Tombari, F., Laina, I., Navab, N.: CNN-SLAM: real-time dense monocular slam with learned depth prediction. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6565–6574, July 2017
2. Mo, K., Li, H., Lin, Z., Lee, J.Y.: The AdobeIndoorNav dataset: towards deep reinforcement learning based real-world indoor robot visual navigation (2018)
3. Hedman, P., Alsisan, S., Szeliski, R., Kopf, J.: Casual 3D photography. *ACM Trans. Graph. (TOG)* **36**(6), 234 (2017)
4. Huang, J., Chen, Z., Ceylan, D., Jin, H.: 6-DOF VR videos with a single 360-camera. In: 2017 IEEE Virtual Reality (VR), pp. 37–44. IEEE (2017)
5. Karsch, K.: Automatic scene inference for 3D object compositing. *ACM Trans. Graph. (TOG)* **33**(3), 32 (2014)
6. Eigen, D., Fergus, R.: Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 2650–2658 (2015)
7. Ren, X., Bo, L., Fox, D.: RGB-(D) scene labeling: features and algorithms. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2759–2766. IEEE (2012)
8. Hartley, R., Zisserman, A.: *Multiple View Geometry in Computer Vision*, 2nd edn. Cambridge University Press, Cambridge (2000)
9. Furukawa, Y., Hernández, C., et al.: Multi-view stereo: a tutorial. *Found. Trends® in Comput. Graph. Vis.* **9**(1–2), 1–148 (2015)
10. Özyeşil, O., Voroninski, V., Basri, R., Singer, A.: A survey of structure from motion\*. *Acta Numerica* **26**, 305–364 (2017)
11. Li, S.: Binocular spherical stereo. *IEEE Trans. Intell. Transp. Syst.* **9**(4), 589–600 (2008)
12. Ma, C., Shi, L., Huang, H., Yan, M.: 3D reconstruction from full-view fisheye camera. *arXiv preprint [arXiv:1506.06273](https://arxiv.org/abs/1506.06273)* (2015)
13. Pathak, S., Moro, A., Yamashita, A., Asama, H.: Dense 3D reconstruction from two spherical images via optical flow-based equirectangular epipolar rectification. In: 2016 IEEE International Conference on Imaging Systems and Techniques (IST), pp. 140–145. IEEE (2016)
14. Li, S., Fukumori, K.: Spherical stereo for the construction of immersive VR environment. In: Proceedings of Virtual Reality, VR 2005, pp. 217–222. IEEE (2005)
15. Kim, H., Hilton, A.: 3D scene reconstruction from multiple spherical stereo pairs. *Int. J. Comput. Vis.* **104**(1), 94–116 (2013)



16. Zhang, Y., Song, S., Tan, P., Xiao, J.: PanoContext: a whole-room 3D context model for panoramic scene understanding. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) ECCV 2014. LNCS, vol. 8694, pp. 668–686. Springer, Cham (2014). [https://doi.org/10.1007/978-3-319-10599-4\\_43](https://doi.org/10.1007/978-3-319-10599-4_43)
17. Yang, H., Zhang, H.: Efficient 3D room shape recovery from a single panorama. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5422–5430 (2016)
18. Xu, J., Stenger, B., Kerola, T., Tung, T.: Pano2CAD: room layout from a single panorama image. In: 2017 IEEE Winter Conference on Applications of Computer Vision (WACV), pp. 354–362. IEEE (2017)
19. Kim, H., de Campos, T., Hilton, A.: Room layout estimation with object and material attributes information using a spherical camera. In: 2016 Fourth International Conference on 3D Vision (3DV), pp. 519–527. IEEE (2016)
20. Plagemann, C., Stachniss, C., Hess, J., Endres, F., Franklin, N.: A nonparametric learning approach to range sensing from omnidirectional vision. *Robot. Auton. Syst.* **58**(6), 762–772 (2010)
21. Ruder, M., Dosovitskiy, A., Brox, T.: Artistic style transfer for videos and spherical images. *Int. J. Comput. Vis.* **126**(11), 1199–1219 (2018)
22. Monroy, R., Lutz, S., Chalasani, T., Smolic, A.: SalNet360: saliency maps for omnidirectional images with CNN. arXiv preprint [arXiv:1709.06505](https://arxiv.org/abs/1709.06505) (2017)
23. Zhang, J., Lalonde, J.F.: Learning high dynamic range from outdoor panoramas. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4519–4528 (2017)
24. Frossard, P., Khasanova, R.: Graph-based classification of omnidirectional images. In: 2017 IEEE International Conference on Computer Vision Workshop (ICCVW), pp. 860–869. IEEE (2017)
25. Su, Y.C., Grauman, K.: Learning spherical convolution for fast features from 360 imagery. In: Advances in Neural Information Processing Systems, pp. 529–539 (2017)
26. Jeon, Y., Kim, J.: Active convolution: learning the shape of convolution for image classification. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1846–1854. IEEE (2017)
27. Dai, J., et al.: Deformable convolutional networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 764–773 (2017)
28. Deng, L., Yang, M., Li, H., Li, T., Hu, B., Wang, C.: Restricted deformable convolution based road scene semantic segmentation using surround view cameras. arXiv preprint [arXiv:1801.00708](https://arxiv.org/abs/1801.00708) (2018)
29. Cohen, T., Geiger, M., Welling, M.: Convolutional networks for spherical signals. In: Principled Approaches to Deep Learning Workshop ICML 2017 (2017)
30. Cohen, T.S., Geiger, M., Köhler, J., Welling, M.: Spherical CNNs. In: International Conference on Learning Representations (ICLR) (2018)
31. Ranftl, R., Vineet, V., Chen, Q., Koltun, V.: Dense monocular depth estimation in complex dynamic scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4058–4066 (2016)
32. Liu, M., Salzmann, M., He, X.: Discrete-continuous depth estimation from a single image. In: 2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 716–723. IEEE (2014)
33. Karsch, K., Liu, C., Kang, S.B.: Depth transfer: depth extraction from videos using nonparametric sampling. In: Hassner, T., Liu, C. (eds.) Dense Image Correspondences for Computer Vision, pp. 173–205. Springer, Cham (2016). [https://doi.org/10.1007/978-3-319-23048-1\\_9](https://doi.org/10.1007/978-3-319-23048-1_9)

34. Eigen, D., Puhrsch, C., Fergus, R.: Depth map prediction from a single image using a multi-scale deep network. In: *Advances in Neural Information Processing Systems*, pp. 2366–2374 (2014)
35. Ren, Z., Lee, Y.J.: Cross-domain self-supervised multi-task feature learning using synthetic imagery. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2018)
36. Laina, I., Rupprecht, C., Belagiannis, V., Tombari, F., Navab, N.: Deeper depth prediction with fully convolutional residual networks. In: *2016 Fourth International Conference on 3D Vision (3DV)*, pp. 239–248. IEEE (2016)
37. Li, B., Shen, C., Dai, Y., van den Hengel, A., He, M.: Depth and surface normal estimation from monocular images using regression on deep features and hierarchical CRFs. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1119–1127 (2015)
38. Liu, F., Shen, C., Lin, G., Reid, I.: Learning depth from single monocular images using deep convolutional neural fields. *IEEE Trans. Pattern Anal. Mach. Intell.* **38**(10), 2024–2039 (2016)
39. Xu, D., Ricci, E., Ouyang, W., Wang, X., Sebe, N.: Multi-scale continuous CRFs as sequential deep networks for monocular depth estimation. In: *Proceedings of CVPR* (2017)
40. Chakrabarti, A., Shao, J., Shakhnarovich, G.: Depth from a single image by harmonizing overcomplete local network predictions. In: *Advances in Neural Information Processing Systems*, pp. 2658–2666 (2016)
41. Li, B., Dai, Y., He, M.: Monocular depth estimation with hierarchical fusion of dilated CNNs and soft-weighted-sum inference. *Pattern Recogn.* **83**, 328–339 (2018)
42. Cao, Y., Wu, Z., Shen, C.: Estimating depth from monocular images as classification using deep fully convolutional residual networks. *IEEE Trans. Circ. Syst. Video Technol.* (2017)
43. Fu, H., Gong, M., Wang, C., Tao, D.: A compromise principle in deep monocular depth estimation. *arXiv preprint [arXiv:1708.08267](https://arxiv.org/abs/1708.08267)* (2017)
44. Garg, R., B.G., V.K., Carneiro, G., Reid, I.: Unsupervised CNN for single view depth estimation: geometry to the rescue. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) *ECCV 2016. LNCS*, vol. 9912, pp. 740–756. Springer, Cham (2016). [https://doi.org/10.1007/978-3-319-46484-8\\_45](https://doi.org/10.1007/978-3-319-46484-8_45)
45. Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: *CVPR* (2017)
46. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: *CVPR*, vol. 2, p. 7 (2017)
47. Wang, C., Buenaposada, J.M., Zhu, R., Lucey, S.: Learning depth from monocular videos using direct methods. *IEEE Conference on Computer Vision and Pattern Recognition* (2018)
48. Mahjourian, R., Wicke, M., Angelova, A.: Unsupervised learning of depth and ego-motion from monocular video using 3D geometric constraints. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2018)
49. Yang, Z., Wang, P., Xu, W., Zhao, L., Nevatia, R.: Unsupervised learning of geometry with edge-aware depth-normal consistency. *arXiv preprint [arXiv:1711.03665](https://arxiv.org/abs/1711.03665)* (2017)
50. Yin, Z., Shi, J.: GeoNet: unsupervised learning of dense depth, optical flow and camera pose. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2018)

51. Srinivasan, P.P., Garg, R., Wadhwa, N., Ng, R., Barron, J.T.: Aperture supervision for monocular depth estimation (2017)
52. Chen, W., Fu, Z., Yang, D., Deng, J.: Single-image depth perception in the wild. In: *Advances in Neural Information Processing Systems*, pp. 730–738 (2016)
53. Silberman, N., Hoiem, D., Kohli, P., Fergus, R.: Indoor segmentation and support inference from RGBD images. In: Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., Schmid, C. (eds.) *ECCV 2012. LNCS*, vol. 7576, pp. 746–760. Springer, Heidelberg (2012). [https://doi.org/10.1007/978-3-642-33715-4\\_54](https://doi.org/10.1007/978-3-642-33715-4_54)
54. Saxena, A., Sun, M., Ng, A.Y.: Make3D: learning 3D scene structure from a single still image. *IEEE Trans. Pattern Anal. Mach. Intell.* **31**(5), 824–840 (2009)
55. Matzen, K., Cohen, M.F., Evans, B., Kopf, J., Szeliski, R.: Low-cost 360 stereo photography and video capture. *ACM Trans. Graph. (TOG)* **36**(4), 148 (2017)
56. Song, S., Yu, F., Zeng, A., Chang, A.X., Savva, M., Funkhouser, T.: Semantic scene completion from a single depth image. In: *IEEE Conference on Computer Vision and Pattern Recognition* (2017)
57. Handa, A., Pătrăucean, V., Stent, S., Cipolla, R.: SceneNet: an annotated model generator for indoor scene understanding. In: *2016 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5737–5743. IEEE (2016)
58. Armeni, I., Sax, S., Zamir, A.R., Savarese, S.: Joint 2D-3D-semantic data for indoor scene understanding. *arXiv preprint [arXiv:1702.01105](https://arxiv.org/abs/1702.01105)* (2017)
59. Armeni, I., et al.: 3D semantic parsing of large-scale indoor spaces. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1534–1543 (2016)
60. Chang, A., et al.: Matterport3D: learning from RGB-D data in indoor environments. In: *International Conference on 3D Vision (3DV)* (2017)
61. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3431–3440 (2015)
62. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (ELUs). *arXiv preprint [arXiv:1511.07289](https://arxiv.org/abs/1511.07289)* (2015)
63. Ioffe, S., Szegedy, C.: Batch normalization: accelerating deep network training by reducing internal covariate shift. In: *International Conference on Machine Learning*, pp. 448–456 (2015)
64. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778 (2016)
65. Yu, F., Koltun, V., Funkhouser, T.: Dilated residual networks. In: *Computer Vision and Pattern Recognition*, vol. 1 (2017)
66. van Noord, N., Postma, E.O.: Light-weight pixel context encoders for image inpainting. *CoRR abs/1801.05585* (2018)
67. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **15**(1), 1929–1958 (2014)
68. Jia, Y., et al.: Caffe: convolutional architecture for fast feature embedding. In: *Proceedings of the 22nd ACM International Conference on Multimedia, MM 2014*, pp. 675–678. ACM, New York (2014)
69. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, Proceedings of Machine Learning Research*, vol. 9. PMLR, pp. 249–256, 13–15 May 2010

70. Kingma, D.P., Ba, J.: Adam: a method for stochastic optimization. arXiv preprint [arXiv:1412.6980](https://arxiv.org/abs/1412.6980) (2014)
71. Xiao, J., Ehinger, K.A., Oliva, A., Torralba, A.: Recognizing scene viewpoint using panoramic place representation. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2695–2702. IEEE (2012)
72. Rhee, T., Petikam, L., Allen, B., Chalmers, A.: MR360: mixed reality rendering for 360 panoramic videos. *IEEE Trans. Visual. Comput. Graph.* **23**(4), 1379–1388 (2017)