

# Monocular Fisheye Camera Depth Estimation Using Sparse LiDAR Supervision

Varun Ravi Kumar<sup>1</sup>, Stefan Milz<sup>1</sup>, Christian Witt<sup>1</sup>, Martin Simon<sup>1</sup>, Karl Amende<sup>1</sup>, Johannes Petzold<sup>1</sup>,  
Senthil Yogamani<sup>2</sup> and Timo Pech<sup>3</sup>

**Abstract**—Near-field depth estimation around a self-driving car is an important function that can be achieved by four wide-angle fisheye cameras having a field of view of over 180°. Depth estimation based on convolutional neural networks (CNNs) produce state of the art results, but progress is hindered because depth annotation cannot be obtained manually. Synthetic datasets are commonly used but they have limitations. For instance, they do not capture the extensive variability in the appearance of objects like vehicles present in real datasets. There is also a domain shift while performing inference on natural images illustrated by many attempts to handle the domain adaptation explicitly. In this work, we explore an alternate approach of training using sparse LiDAR data as ground truth for depth estimation for fisheye camera. We built our own dataset using our self-driving car setup which has a 64-beam Velodyne LiDAR and four wide angle fisheye cameras. To handle the difference in view-points of LiDAR and fisheye camera, an occlusion resolution mechanism was implemented. We started with Eigen’s multiscale convolutional network architecture [1] and improved by modifying activation function and optimizer. We obtained promising results on our dataset with RMSE errors comparable to the state-of-the-art results obtained on KITTI.

## I. INTRODUCTION

Depth estimation from single camera images is an important basic task for self driving cars such as driver assistance systems to solve localization and perception problems. Predominantly, the challenge is an arduous process and it cannot be decoded directly from bottom-up geometric cues. A single captured image scene may be congruous with infinite real world scenarios [2]. Successful approaches have relied on structure from motion, shape-from-X, binocular and multi-view stereo. These techniques hinge on the assumption of prior knowledge about the characteristic appearance and multiple observations of the scene of interest that are available. The aforementioned can occur via multiple viewpoints, layout and size of object needs, cues such as shading, or observations of the scene under different lighting conditions. To overcome this limitation, there has recently been a rise in the number of works that pose the task of single image depth estimation as a supervised learning problem [1], [2], [3]. These methods seek to directly predict the depth from a single RGB image for each pixel through deep learning

models that have been modeled on large collections of ground truth depth data.

Humans excel at monocular depth estimation by exploiting cues such as motion parallax, linear perspective, shape from shading, relative size and occlusion [4]. Full scene understanding with our capability to precisely estimate depth appears to bolster from the combination of both top-down and bottom-up cues [5].

For supervised deep learning a large amount of training data is required in order to achieve high accuracy and to generalize on new scenes. In indoor environments, RGB-D cameras are used to generate ground truth depth data for this task. However, strong sunlight has an adverse effect on infrared interference and make depth information of those sensing devices extremely noisy. In outdoor applications, especially in the domain of self driving cars, LiDAR or other laser scanners-are used to capture ground truth data. Since measurements from 3D lasers have usually a sparse nature, the depth variations are captured with less details than visible in the image.

Additional to the use of real data, synthetic rendering of depth maps are used to generate ground truth data. Rendered images do not unveil the scene and fail to implement real image noise characteristics-which are the two drawbacks of this method [6]. Also, there is an inefficiency to generalize on new scenes by the model trained on this approach.

The motivation of this paper is to provide a baseline for single frame depth estimation based on sparse Velodyne data as ground truth for training. This paper builds upon the authors’ previous work published in a short paper [7] and the contributions of this paper include:

- 1) Demonstration of a working prototype purely trained on sparse Velodyne LiDAR data.
- 2) Demonstration of fisheye camera depth estimation using CNN.
- 3) Adapting training data to handle occlusion due to difference in camera and Velodyne LiDAR viewpoint.
- 4) Tailoring the loss function and training algorithm to handle sparse depth data.

The rest of the paper is structured as follows. Section II provides a survey of convolution neural networks (CNN) based depth estimation. Section III discusses the details of the network architecture, loss function tailoring and training algorithms. Section IV summarizes results on our internal fisheye camera dataset and provide a comparison with publicly available KITTI results. Finally, Section V concludes the paper and provides potential future directions.

<sup>1</sup>Valeo und Schalter und Sensoren GmbH, Driving Assistance Advanced Research, Kronach [varun-ravi.kumar@valeo.com](mailto:varun-ravi.kumar@valeo.com)

<sup>2</sup>Senthil Yogamani is with Valeo Vision Systems, Ireland [senthil.yogamani@valeo.com](mailto:senthil.yogamani@valeo.com)

<sup>3</sup>Timo Pech is with Technische University, Chemnitz, Germany [timo.pech@etit.tu-chemnitz.de](mailto:timo.pech@etit.tu-chemnitz.de)

## II. RELATED WORK

It has been noted that in recent years, several deep learning based approaches to monocular depth estimation are trained in a supervised way - which requires a single input image - with no assumptions about the scene geometry or types of objects which are present. In monocular depth estimation only single images are used at the inference time. Saxena et al. [8] pioneered the supervised-learning based approach called Make3D patch-based model. The input images are initially over-segmented into patches and the 3D location and orientation of local planes are estimated which illustrates each patch. Markov Random Fields are used to combine the monocular cues with the stereo correspondences. The drawback of planar based approximations including [9] is realistic outputs can not be generated as they lack global context since the estimates are made locally. They can be hindered when it comes to modeling of thin structures.

Liu et al. [3] formulated an approach for depth estimation as a deep continuous Conditional Random Fields (CRF) learning problem. Instead of hand-crafted features such as unary and pairwise terms, Liu used deep convolutional neural fields that permitted the CNN features of unary and pairwise potentials end-to-end for training by utilizing continuous depth and Gaussian assumptions on the pairwise potentials.

Ladicky et al. [10] improved the per pixel depth estimation to a lucid classifier estimating only the probability of a pixel present at an arbitrarily fixed canonical depth. After appropriate image transformations, the probability of any other depths can be achieved by implementing the same classifier. The vulnerability of independent approaches of depth estimation and semantic segmentation are aimed directly by improving and generalizing the overall approach.

Karsch et al. [11] recommended a k-Nearest-Neighbor (kNN) transfer mechanism which can achieve better alignment which hinges on SIFT Flow [12] to estimate depths from single images of static backgrounds. They accomplished better estimation with the scene of interest in videos with dynamic foreground coupled with augmentation of the latter with motion information. A major drawback of this approach is a requirement of a complete training dataset to be available at inference time.

In the last few years, it has been observed that object classification and recognition [13], [14], [15] reap great success with the application of Convolutional Neural Networks. CNNs perform classification of a single or multiple object label for a complete input image and apply bounding boxes on a few objects in each scene of an image. In addition to this, a variety of tasks like pose estimation [16], stereo depth [17] and instance segmentation [18] incorporate CNNs. Most of these models use CNNs to find only local features, or generate descriptors of discrete proposal regions; in contrast, Eigen's network uses both local and global views to predict a variety of output types.

Laina [19] illustrated that dense depth maps can be produced by using ResNet-based encoder-decoder architecture. Their approach is demonstrated to predict dense depth maps

in indoor scenes using RGB-images for training. Through example images [20], [21] it is found that the idea of depth transfer can be used to predict depth map or integrate depth map prediction with semantic segmentation [1], [10], [22] in supervised training.

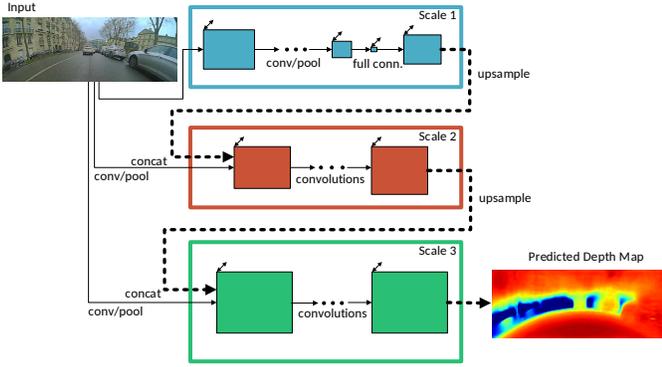
Single-image based depth estimation has various hardware-based solutions like performing depth from defocus using a modified camera aperture proposed by Levin et al. [23] and the Kinect v2 uses time-of-flight and active stereo to record depth.

We have incorporated Eigen's [1] core multi-scale architecture to adapt to a single task of estimating depth with an output resolution twice the original. We could achieve similar qualitative results with a sparse dataset, obtained from Velodyne HDL-64L rotating 3D laser scanner with valid depth points ranging from 3k-25k after occlusion removal.

## III. MODEL ARCHITECTURE

Our model offers several architectural improvements to [1] which is initially based on Eigen et al. [2]. We adopted a simple architecture for Scale 1 based on AlexNet [13] to achieve real time on an embedded platform Nvidia TX 2. However, the usage of new model architectures such as ResNet-50 [24] which have a bigger field of view could improve the results. These models take images of bigger dimensions as input and hence can provide a better global view of the image to the learning algorithm. Depending on the whole image area, a multi-scale deep neural network first predicts a coarse global output and refines it using finer-scale local networks. This scheme is described in Fig. 1. The model is deeper with more convolutional layers compared to [2]. Second, with the added third scale from [1] at higher resolution, bringing the final output resolution up to half the input, or  $284 \text{ px} \times 80 \text{ px}$  for our sparse LiDAR fisheye camera dataset. In addition, we use swish [25] as the activation function rather than the mostly preferred rectified linear unit (ReLU) [26]. Finally, we adopt Adam optimizer [27] which yields faster converging instead of the stochastic gradient descent (SGD) used by Eigen et al. [1], [2]. Multi channel feature maps were passed similarly to [1] avoiding the flow of output predictions from the coarse scale to the refine scale.

*a) Scale 1: Full-Image View:* The first scale of the neural network analyses the global structure of the image and extracts global features. Global understanding of the scene requires an effective use of depth cues like object locations, vanishing points and alignment of structures [1]. The local view of the image is inadequate to capture these features. Scale 1 is based on an ImageNet-trained AlexNet [13] with initialization of pre-trained AlexNet weights only on convolutional layers. The global understanding of the image is achieved by two fully connected layers at the end. A very large field of view is obtained as each spatial location in the output connects to all the image features. The neural network takes fisheye images of size  $576 \text{ px} \times 172 \text{ px}$  as input. The output of the scale is a 64-channel feature map with a resolution  $142 \text{ px} \times 40 \text{ px}$ .



	Layer	1.1	1.2	1.3	1.4	1.5	1.6	1.7	upsamp
Scale 1 (AlexNet)	Size	142x41	71x21	36x11	36x11	36x11	1x1	36x10	144x40
	#convs	1	1	1	1	1	-	-	-
	#chan	96	256	384	384	256	4096	64	64
	ker. sz	11x11	5x5	3x3	3x3	3x3	-	-	-
	Ratio	/8	/16	/16	/16	/32	-	-	/16
	stride	4	1	1	1	1	-	-	-
	Layer	2.1	2.2	2.3	2.4	2.5	upsamp		
Scale 2	Size	284x82	142x40	142x40	142x40	142x40	284x80		
	#chan	96+64	64	64	64	1	1		
	ker. sz	9x9	5x5	5x5	5x5	5x5	-		
	Ratio	/4	/4	/4	/4	/4	-		
	stride	2	1	1	1	1	/2		
	Layer	3.1	3.2	3.3	3.4	final			
Scale 3	Size	284x82	284x80	284x80	284x80	284x80			
	#chan	64	64	64	1	1			
	ker. sz	9x9	5x5	5x5	5x5	-			
	Ratio	/2	/2	/2	/2	-			
	stride	1	1	1	1	/2			

Fig. 1. Multi-scale architecture for depth prediction on raw fisheye images with a sparse velodyne (HD64L) ground truth. The input to the network is 576x172. Occlusion correction is essential, if velodyne points are mapped to the fisheye eye image plane, because of the different mounting positions of camera and LiDAR (see Section III-D).

*b) Scale 2: Predictions:* This scale incorporates a narrow view of the image and makes depth predictions at a resolution one-fourth of the input image [1]. While making predictions, the global scene information supplied by the Scale 1 is also considered by concatenation of feature maps. The input to this scale is the same RGB image which was given as input to Scale 1. Scale 2 corrects the coarse prediction it receives from Scale 1 to align with local details such as object and car edges, by concatenating the feature maps of the coarse network with those from a single layer of convolution and pooling. The output of the second scale is a 284 px × 80 px prediction for our sparse fisheye cameras dataset, with a single channel as a gray scale image.

*c) Scale 3: Higher Resolution:* Scale 3 refines the predictions made by Scale 2. It contains a set of convolutional operations with a small stride that can blend detailed structure of the image into the predictions. The alignment of output to higher-resolution details is further refined which produces detailed spatially coherent depth map predictions. The final linear layer of this scale predicts the depth map with a resolution of 284 px × 80 px.

#### A. Sparse ground-truth depth maps

A Velodyne HDL-64ES2 sensor can fire only 64 beams of lasers at different vertical angles with a vertical field of view of 26.8°. Hence the depth maps obtained from the projection of the LiDAR 3D points are sparse. Due to rotary motion of

the Velodyne LiDAR sensor and movement of the vehicle while data recording was made, points that are far away had poor reflectivity. Therefore the extracted depth maps are sparser for scenes composed of far away objects.

#### B. Scale-Invariant Error

The sparse nature of the ground truth depth maps is considered in the design of the loss function. We have adopted the loss function as described by Eigen et al. [2] which is a  $l_2$ -loss with a scale-invariant term. There is a lot of uncertainty regarding the global scale associated with the image, since we consider only a single image for depth prediction. The scale-invariant loss considers this scaling effect and produces the same loss for two scenes that differ only by the scaling factor. Last linear layer in the third scale of the architecture predicts the depth, which is compared to the ground truth depth map. The loss function is defined by equation 1,

$$\text{Loss}(p, p^*) = \frac{1}{n} \sum_{i \in V} d_i^2 - \frac{1}{n^2} \left( \sum_{i \in V} d_i \right)^2 \quad (1)$$

where  $p$  is the pixel wise set of predictions from the neural network.  $p^*$  represents the ground truth depth map. Hence,  $d_i = p_i - p_i^*$  is the difference for pixel  $i$ . The ground truth depth map is sparse, i.e. not for all pixels exists an equivalent depth measurement. We define a set of valid pixels  $V \subset P^*$ , with  $V = \{p_1 \dots p_i \dots p_n\}$ , where  $n$  is the number of valid pixels within the ground truth depth map [2]:

$$\text{Loss}(p, p^*) = \frac{1}{n} \sum_{i \in V} (\log p_i - \log p_i^* + \alpha(p, p^*))^2. \quad (2)$$

For a given  $(p, p^*)$ , the error is minimized by  $\alpha$ . The value of  $\alpha$  is  $\alpha(p, p^*) = \frac{1}{n} \sum_{i \in V} (\log p_i^* - \log p_i)$ . The scale that best aligns to the ground truth is given by  $e^\alpha$  for any prediction  $p$ . The error is same across all the scalar multiples of  $p$ , hence the term scale invariance as mentioned in [2].

An equivalent form of metric was obtained by Eigen et al. [2] by setting  $d_i = \log p_i - \log p_i^*$  to be the difference between the prediction and ground truth at pixel  $i$ ,

$$\text{Loss}(p, p^*) = \frac{1}{n^2} \sum_{i, j \in V} ((\log p_i - \log p_j) - (\log p_i^* - \log p_j^*))^2 \quad (3)$$

The error is demonstrated in equation 3 by comparing the relationships between pairs of pixels  $i, j$  in the output: each pair of pixels in the prediction must differ in depth by an amount similar to that of the corresponding pair in the ground truth to have a low error. Our fisheye dataset is extremely sparse due to the nature of LiDAR sensors, the loss function is adapted to this sparsity. By masking out pixels that do not have a valid depth value, the loss is calculated only on pixels which have depth values. This facilitates efficient feature extraction by the neural network. In addition to the scale-invariant error, we evaluate our method using the error metrics used in [2], [3] as described in section IV.

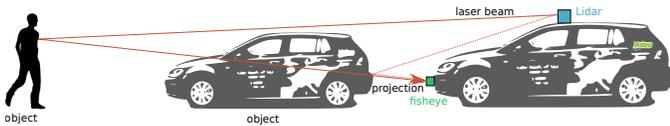


Fig. 2. Illustration of occlusion due to LiDAR’s viewpoint being higher than the fisheye camera’s viewpoint. 3D points from the object (person) will be mapped to image plane even though it is not visible from camera.

### C. Training-Model

We train our model in a single pass in an end-to-end fashion compared to [1], [2] where the first two scales of the network were trained jointly. For each gradient step, the entire image area is considered for training. Pre-trained weights from AlexNet [13] are used. ConvNet is incorporated as a fixed feature extractor for our dataset and the last fully-connected layers are removed. The fully connected layers are initialized randomly with values from a normal truncated distribution. Scale 2 and Scale 3 are randomly initialized. The dataset contains 60 000 images from fisheye camera and sparse Velodyne LiDAR scans as ground truth with validation and test set of 5000 images each. We trained our model with a batch size of 20 using the Adam [27] optimization algorithm, with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$ . We adopt an exponential decay function to lower the learning rate as the training progresses, with an initial learning rate of  $\lambda = 10^{-4}$ . The function decays every 7500 steps with a base of 0.95. For the non-linearities in the network, we used swish [25] activation function instead of the commonly used rectified linear units (ReLU) [26] which tend to work better on deeper models. The swish function is defined as  $f(x) = x \cdot \sigma(x)$  [25], where  $\sigma(x) = (1 + e^{-x})^{-1}$  is the sigmoid function. The interesting aspect about the swish is that it does not monotonically increase compared to other activations functions like ReLU. The problem of *dead neurons* arises as the parameter will not be updated if the gradient is 0, since gradient descent being the parameter update algorithm. We initially experimented by adopting different proposed alternative activation functions such as scaled exponential linear units (SELU) [28], exponential linear units (ELU) [29] and leaky ReLU [30]. However, we found that swish performed best.

### D. Occlusion Correction

The sensor fusion of the data will be correct, if both camera and the Velodyne LiDAR scanner beholds the world from the same viewpoint. However, for technical reasons in our vehicle the fisheye camera is in the front and the LiDAR is placed at the top as seen in Fig. 2. LiDAR perceives the environment behind objects that occlude the view for the camera. This problem of occlusion results in wrong mapping of depth-points that are not visible to the camera. It is hard to solve, since occluded points are projected adjacently to unoccluded points [31].

To solve this problem, we adapted a distance based segmentation technique with morphological filters as shown in the Fig. 3. Instead of directly projecting points from

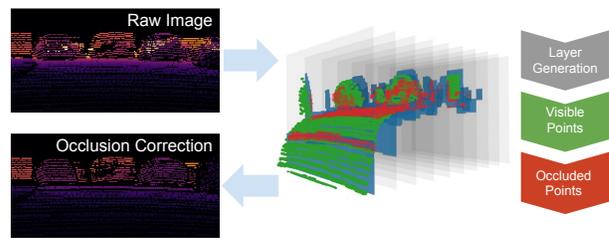


Fig. 3. Visualization of the distance based segmentation technique with morphological filters. LiDAR points are projected to corresponding layers and are removed if they are occluded by dilated parts of a neighboring layer.

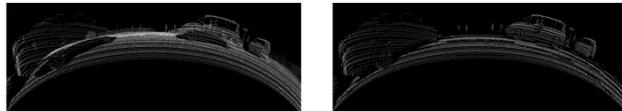


Fig. 4. Illustration of Occluded Velodyne Ground Truth (left) and Dis-Occluded Velodyne Ground Truth (right)

the LiDAR into the image plane of the fisheye camera, we introduce  $I$  layers within the camera view located at a distance  $d_i^{\text{layer}}$ ,  $i = 1, \dots, I$ . Each LiDAR point will be projected onto the layer next to it. We apply a morphological filter that dilates points within each layer to fill the sparse regions (in Fig. 3 dilated parts of the layers are colored blue). A point at a distance  $d^{\text{point}}$  is regarded as occluded, if a layer  $i$  exists with  $d_i^{\text{layer}} < d^{\text{point}}$ . Otherwise the valid point is projected onto the image plane of the fisheye camera.

## IV. RESULTS

The model is completely trained on our internal dataset. Our dataset contains 55 000 images obtained from raw fisheye camera and sparse Velodyne HDL-64E rotating 3D laser scanner as ground truth. Points without depth value are left unfilled without any post-processing. Eigen’s model [1] handles missing values by eliminating them in the loss function. The input images are down-sampled to  $576 \text{ px} \times 172 \text{ px}$  primarily to get faster inference and training times.

The ground truth depth for this dataset is captured at various intervals using a Velodyne HDL-64E rotating 3D laser scanner, and are sampled at irregularly spaced points. Conflicting values are found when constructing the ground truth depths for training, since sensor records data at a set maximum frequency of 10 Hz and the fisheye cameras record data at 30 Hz. Time synchronization is essential as the sensors capture data at different frequencies. Each spin of the LiDAR sensor is considered as a frame and carries a time-stamp associated with it. Similarly, each image frame recorded by the fisheye camera carries a time-stamp. For the purpose of synchronization, time-stamps provided with the recordings are used. We resolve conflicts by choosing the depth recorded closest to the RGB capture time in Intempora RTMaps (Real-Time Multisensor Applications) framework.

The training set was collected by driving around Paris, France and various parts of Bavaria, Germany. The training set includes scenes from the *city*, *residential* and *sub-urban*

TABLE I

QUANTITATIVE RESULTS OF LEADERBOARD ALGORITHMS ON KITTI 2015 [32] DATASET AND OUR APPROACH ON VALEO'S FISHEYE DATASET

Approach	Supervised	cap	RMSE	RMSE (log)	ARD	SRD	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
			lower is better				higher is better		
Mancini et al. [33]	Yes	0 – 100 m	7.508	0.524	0.196	-	0.318	0.617	0.813
Eigen et al. [2] coarse $28 \times 144$	Yes	0 – 80 m	7.216	0.273	0.194	1.531	0.679	0.897	0.967
Eigen et al. [2] fine $27 \times 142$	Yes	0 – 80 m	7.156	0.270	0.190	1.515	0.692	0.899	0.967
Liu et al. [34] DCNF-FCSP FT	Yes	0 – 80 m	6.986	0.289	0.217	1.841	0.647	0.882	0.961
Ma et al. [35]	Yes	0 – 100 m	6.266	-	0.208	-	0.591	0.900	0.962
Kuznetsov et al. [6]	Yes	0 – 50 m	3.531	0.183	0.117	0.597	0.861	0.964	0.989
Zhou et al. [36] (w/o explainability)	No	0 – 50 m	5.452	0.273	0.208	1.551	0.695	0.900	0.964
Zhou et al. [36]	No	0 – 50 m	5.181	0.264	0.201	1.391	0.696	0.900	0.966
Godard et al. [5]	No	0 – 50 m	4.471	0.232	0.140	0.976	0.818	0.931	0.969
Ours fine $80 \times 284$	Yes	0 – 50 m	1.717	0.236	0.160	0.397	0.816	0.934	0.969

categories of our raw dataset. These are randomly shuffled and fed to the network. We train the entire model for 80 epochs and test prediction takes 3.45s/batch with a batch size of 20 images (0.17s/image).

The evaluation of accuracy in our method in depth prediction is using the 3D laser ground truth on the test images. We use the depth evaluation metrics used by Eigen et al. [2]. Exemplary predictions are shown in figure 5. The qualitative results show that image regions without sufficient large ground truth data points (e.g. sky), the model fails to predict reasonable values.

A protocol evaluation is applied and results are shown by discarding ground-truth depth below 0 m and above 50 m while capping the predicted depths into 0 m – 50 m depth interval. This implies, we set predicted depths to 0 m and 50 m if they are below 0 m or above 50 m, respectively.

In Table I, we show how our approach performs on Valeo's fisheye dataset. Furthermore the results of leaderboard algorithms on KITTI 2015 [32] are reproduced. For lack of a better comparison, we use this as a proxy to illustrate that we obtained comparable RMSE on our sparse fisheye dataset. It should be noted that although we predict a dense depth map, the sparse dataset only allows us to take a fraction of the predicted values into consideration for error calculation. To tackle this problem we plan to refine our model on a synthetic dataset, close to our Valeo's fisheye dataset, that allows a full verification of the predicted depth. First tests show promising results with excluded sky.

## V. CONCLUSION

Even though the camera/LiDAR setups are different, the results provide a reasonable comparison to KITTI on performance of monocular depth regression using sparse LiDAR input. In future work, we aim to improve the results by using more consecutive frames which can exploit the motion parallax and better CNN encoders. We also plan to augment the supervised training with synthetic data and unsupervised training techniques.

## REFERENCES

[1] D. Eigen and R. Fergus, "Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture,"

in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2650–2658.

[2] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," in *Advances in neural information processing systems*, 2014, pp. 2366–2374.

[3] F. Liu, C. Shen, G. Lin, and I. Reid, "Learning depth from single monocular images using deep convolutional neural fields," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 10, pp. 2024–2039, Oct. 2016. [Online]. Available: <https://doi.org/10.1109/TPAMI.2015.2505283>

[4] I. P. Howard, *Perceiving in depth, volume 1: basic mechanisms*. Oxford University Press, 2012.

[5] C. Godard, O. Mac Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *CVPR*, vol. 2, no. 6, 2017, p. 7.

[6] Y. Kuznetsov, J. Stückler, and B. Leibe, "Semi-supervised deep learning for monocular depth map prediction," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6647–6655.

[7] V. R. Kumar, S. Milz, C. Witt, M. Simon, K. Amende, J. Petzold, S. Yogamani, and T. Pech, "Near-field depth estimation using monocular fisheye camera: A semi-supervised learning approach using sparse lidar data," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, Deep Vision: Beyond Supervised learning*, 2018.

[8] A. Saxena, M. Sun, and A. Y. Ng, "Make3d: Learning 3d scene structure from a single still image," *IEEE transactions on pattern analysis and machine intelligence*, vol. 31, no. 5, pp. 824–840, 2009.

[9] D. Hoiem, A. A. Efros, and M. Hebert, "Automatic photo pop-up," in *ACM transactions on graphics (TOG)*, vol. 24, no. 3. ACM, 2005, pp. 577–584.

[10] L. Ladický, J. Shi, and M. Pollefeys, "Pulling things out of perspective," in *Proceedings of the 2014 IEEE Conference on Computer Vision and Pattern Recognition*, ser. CVPR '14. Washington, DC, USA: IEEE Computer Society, 2014, pp. 89–96. [Online]. Available: <http://dx.doi.org/10.1109/CVPR.2014.19>

[11] K. Karsch, C. Liu, and S. B. Kang, "Depth extraction from video using non-parametric sampling," in *European Conference on Computer Vision*. Springer, 2012, pp. 775–788.

[12] C. Liu, J. Yuen, A. Torralba, J. Sivic, and W. T. Freeman, "Sift flow: Dense correspondence across different scenes," in *European conference on computer vision*. Springer, 2008, pp. 28–42.

[13] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[14] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014. [Online]. Available: <http://arxiv.org/abs/1409.1556>

[15] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.

[16] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler, "Joint training of a convolutional network and a graphical model for human pose estimation," in *Advances in neural information processing systems*, 2014, pp. 1799–1807.

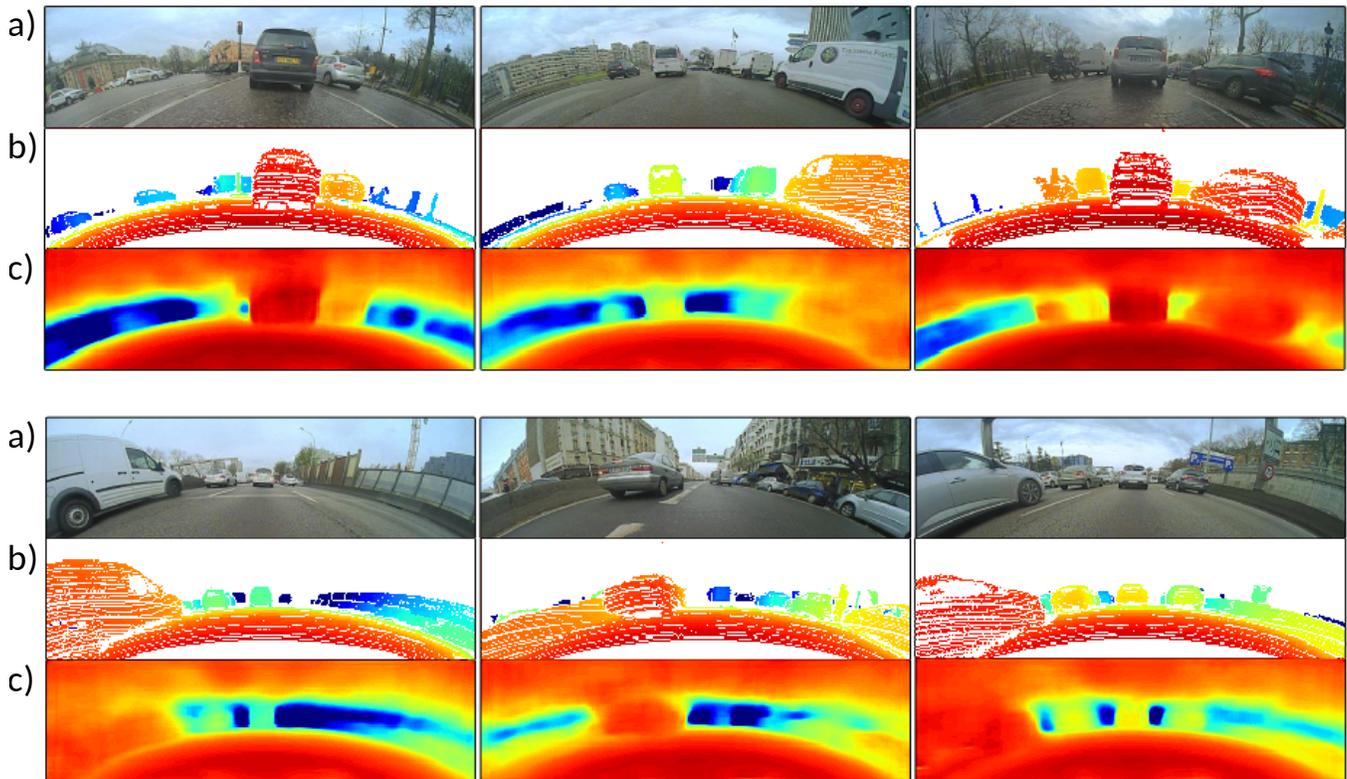


Fig. 5. Qualitative results: Exemplary predictions by the proposed CNN network. For each image, we show (a) RGB Input (b) LiDAR Ground Truth (c) Predicted Depth Map [The sky is considered to be invalid pixel i.e masked as zero while training. We have not considered disparity depth for ground truth generation as compared to KITTI [32]. The depth values are in 8-bit intensity range (0 - 255)]

- [17] J. Zbontar and Y. LeCun, "Computing the stereo matching cost with a convolutional neural network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1592–1599.
- [18] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik, "Learning rich features from rgb-d images for object detection and segmentation," in *European Conference on Computer Vision*. Springer, 2014, pp. 345–360.
- [19] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab, "Deeper depth prediction with fully convolutional residual networks," in *3D Vision (3DV), 2016 Fourth International Conference on*. IEEE, 2016, pp. 239–248.
- [20] K. Karsch, C. Liu, and S. B. Kang, "Depth extraction from video using non-parametric sampling," in *European Conference on Computer Vision*. Springer, 2012, pp. 775–788.
- [21] M. Liu, M. Salzmann, and X. He, "Discrete-continuous depth estimation from a single image," in *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*. IEEE, 2014, pp. 716–723.
- [22] B. Liu, S. Gould, and D. Koller, "Single image depth estimation from predicted semantic labels," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. IEEE, 2010, pp. 1253–1260.
- [23] A. Levin, R. Fergus, F. Durand, and W. T. Freeman, "Image and depth from a conventional camera with a coded aperture," *ACM transactions on graphics (TOG)*, vol. 26, no. 3, p. 70, 2007.
- [24] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [25] P. Ramachandran, B. Zoph, and Q. V. Le, "Searching for activation functions," *CoRR*, vol. abs/1710.05941, 2017. [Online]. Available: <http://arxiv.org/abs/1710.05941>
- [26] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ser. ICML'10. USA: Omnipress, 2010, pp. 807–814. [Online]. Available: <http://dl.acm.org/citation.cfm?id=3104322.3104425>
- [27] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *CoRR*, vol. abs/1412.6980, 2014. [Online]. Available: <http://arxiv.org/abs/1412.6980>
- [28] G. Klambauer, T. Unterthiner, A. Mayr, and S. Hochreiter, "Self-normalizing neural networks," in *Advances in Neural Information Processing Systems*, 2017, pp. 972–981.
- [29] D. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," *CoRR*, vol. abs/1511.07289, 2015. [Online]. Available: <http://arxiv.org/abs/1511.07289>
- [30] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *Proc. icml*, vol. 30, no. 1, 2013, p. 3.
- [31] P. Biasutti, J.-F. Aujol, M. Brédif, and A. Bugeau, "Disocclusion of 3D LiDAR point clouds using range images," in *ISPRS International Society for Photogrammetry and Remote Sensing (CMRT)*, Hannover, Germany, Jun. 2017. [Online]. Available: <https://hal.archives-ouvertes.fr/hal-01522366>
- [32] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 3354–3361.
- [33] M. Mancini, G. Costante, P. Valigi, and T. A. Ciarfuglia, "Fast robust monocular depth estimation for obstacle detection with fully convolutional networks," in *Intelligent Robots and Systems (IROS), 2016 IEEE/RSJ International Conference on*. IEEE, 2016, pp. 4296–4303.
- [34] F. Liu, C. Shen, and G. Lin, "Deep convolutional neural fields for depth estimation from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5162–5170.
- [35] F. Ma and S. Karaman, "Sparse-to-dense: Depth prediction from sparse depth samples and a single image," *CoRR*, vol. abs/1709.07492, 2017. [Online]. Available: <http://arxiv.org/abs/1709.07492>
- [36] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, "Unsupervised learning of depth and ego-motion from video," in *CVPR*, vol. 2, no. 6, 2017, p. 7.