

# 3dDepthNet: Point Cloud Guided Depth Completion Network for Sparse Depth and Single Color Image

Rui Xiang\*  
Department of Mathematics  
University of California, Irvine  
xiangr1@uci.edu

Feng Zheng, Huapeng Su, Zhe Zhang  
Trifo, Inc.  
{feng.zheng, huapeng.su, zhe.zhang}@trifo.com

## Abstract

In this paper, we propose an end-to-end deep learning network named 3dDepthNet, which produces an accurate dense depth image from a single pair of sparse LiDAR depth and color image for robotics and autonomous driving tasks. Based on the dimensional nature of depth images, our network offers a novel 3D-to-2D coarse-to-fine dual densification design that is both accurate and lightweight. Depth densification is first performed in 3D space via point cloud completion, followed by a specially designed encoder-decoder structure that utilizes the projected dense depth from 3D completion and the original RGB-D images to perform 2D image completion. Experiments on the KITTI dataset show our network achieves state-of-art accuracy while being more efficient. Ablation and generalization tests prove that each module in our network has positive influences on the final results, and furthermore, our network is resilient to even sparser depth.

## 1. Introduction

Depth sensing is a fundamental task for both indoor and outdoor applications, such as home robotics and autonomous driving. Indoor depth sensors, e.g. RGB-D camera, are usually more accurate because of the limited distance of indoor scenes and good illumination conditions [13, 38]. For outdoor scenes, LiDAR is usually the primary depth sensor due to its high accuracy and long sensing range. However, even high-resolution LiDAR albeit being extremely expensive<sup>1</sup> still produces sparse depth output. The absence of affordable dense depth sensors leads to a great research interest to develop methods to estimate smooth and dense depth from sparse samples.

There are two main challenges in this depth completion

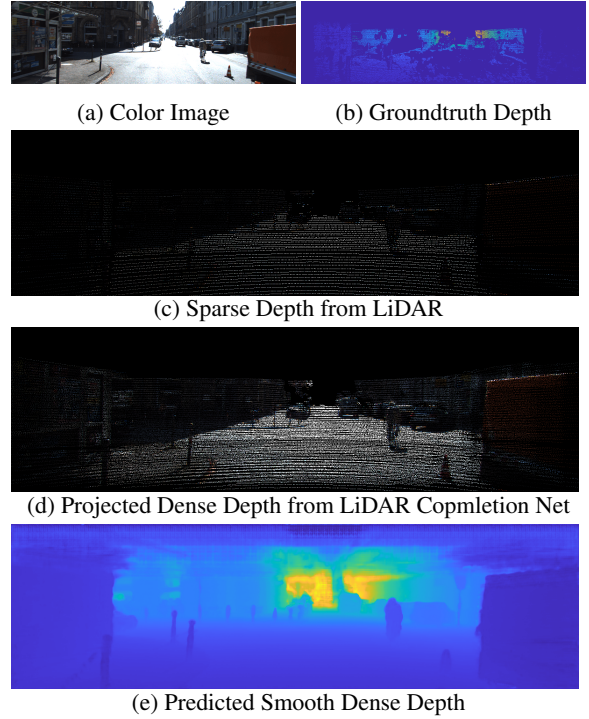


Figure 1: Our model takes sparse depth (c) as input to generate a dense point cloud by our LiDAR Completion Net (d), and then predict the smooth, dense depth (e) via our Depth Completion Net with the aid of the projected dense depth, and the original color image (a) and sparse depth (c).

task. The first challenge is feature extraction from highly sparse depth measurements and in particular, how to combine features from different input modalities (i.e. depth and color). Simple concatenation may not keep consistency well since depth and color values have different units. Also, the object boundary could be fuzzy due to the occlusion problem, which is the phenomenon that the depth image may not perfectly align with the color image. Secondly, for learning-based methods, dense groundtruth depth is diffi-

\*Paper is finished during author's internship at Trifo, Inc.

<sup>1</sup>Currently, the price of 16-line and 64-line Velodyne LiDARs is \$4k and \$75k, respectively [28].

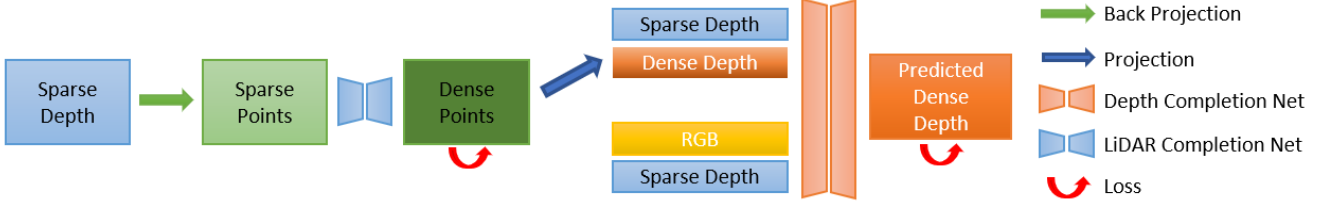


Figure 2: **Network architecture.** There are two important units in our network: LiDAR Completion Net and Depth Completion Net. LiDAR Completion Net performs point cloud completion in 3D space according to local geometry, while Depth Completion Net fuses projected dense depth from 3D completion and original sparse depth and color image to produce a smooth and dense depth image output.

cult to obtain. Compared to image classification or segmentation tasks, pixel-wise manual annotation of depth values is not reliable and extremely time-consuming.

The architecture of our proposed network is shown in Figure 2. We summarize our main contribution as follows:

1. We propose a novel 3D-to-2D coarse-to-fine depth completion network, which performs coarse densification in 3D space via point cloud completion, followed by 2D image space fine depth completion with the aid of the generated intermediate dense 3D point cloud.
2. Based on the property of sparse LiDAR points, we propose a LiDAR Completion Net, which performs patching, adjusting, and densification in 3D space. We illustrate the process in figure 3 and show that our LiDAR Completion Net performs a geometrically meaningful 3D completion.
3. We propose a modified encoder-decoder structure, called Depth Completion Net, to fuse the output from LiDAR Completion Net and the original sparse depth and color image to generate the final dense and smooth depth image.

## 2. Related Work

**Point Cloud Completion** Raw LiDAR scan data is represented as point cloud in 3D space. Including LiDAR data, most real-world 3D data are often incomplete, resulting in a loss in geometric information. Incompleteness mainly comes from two major limitations of modern 3D sensors, the first one is the viewing angle limitation, and the second one is the resolution limitation. Current completion methods for 3D data can be categorized into traditional methods and learning-based methods.

For traditional methods, interpolation methods [2, 8, 33] are explored to fill holes or gaps in relatively small local surface. Symmetry detection methods [31, 32, 34] utilize symmetry property of nature objects to complete unobserved symmetrical parts. Such methods are limited by their strong

assumptions at the underlying geometry of the 3D objects. Alignment-based methods [16, 23, 20, 22] retrieve object or object parts from a library to match or assemble the target incomplete 3D object. Planes and quadrics are also taken into account for some alignment-based methods [3, 24]. These methods often have to solve a large-scale optimization problem and are limited by the content of the corresponding library.

For learning-based methods, most methods use voxels as the representation in order to take advantage of the voxel convolutional neural network. However, voxel-based methods are significantly limited by its computation cost both in the discretization step and training step. There are also learning methods [25] based on mesh deformation. [50] is a point cloud-based learning method dealing with specific 3D object completion. However, there are few methods designed for point cloud completion of large scale scenes, such as LiDAR scan data. The content complexity and information sparsity of LiDAR scan data make the completion task extremely difficult even for deep learning methods.

**Depth Completion** Depth completion aims to recover an accurate, smooth, and dense depth image, given a semi-dense or sparse depth image with or without the guidance of the corresponding color image. Structure light-based sensors and laser-based sensors have very different resolution and range. The sensor-dependent nature and different input modalities of depth completion make this task very challenging and application-specific. The sparser the input depth, the more difficult the problem becomes.

For indoor applications, with structured light sensor, typically less than 20% depth is missing. Traditional methods like depth inpainting [1, 4], depth super-resolution [29, 40, 49], and depth denoising [9] already provides good results. Deep learning methods also have been used for indoor depth completion. Zhang *et al.* [51] proposed adding surface normal information into learning network. Chen *et al.* [6] utilized an affinity matrix to guide the depth completion with a recurrent neural network.

The problem becomes more difficult when the outdoor scene is involved because the depth information is ex-

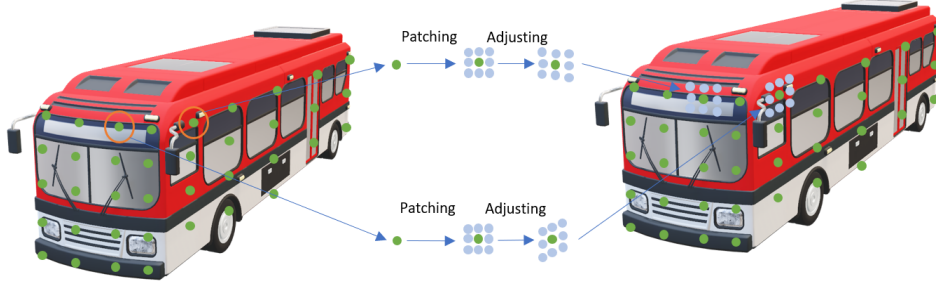


Figure 3: **Intuition of LiDAR Completion Net.** LiDAR Completion Net performs both patching of original points and then adjusting the patch to fit its local geometry. Green dots represent original LiDAR scan points, and blue dots represent new patch points.

tremely sparse. The typical sensor for outdoor application, e.g., autonomous driving, is LiDAR, which provides a very limited number of line scan depth information. Even high-resolution LiDAR can only provide 4% depth information in one image. Wavelet methods [17, 26] are used before the emerging of deep learning. Recently, Qiu *et al.* [37] proposed a two pathway deep learning network involving surface normal and attention mask to generate geometrically meaningful output. Ma *et al.* [28] designed an encoder-decoder network to fuse the information from RGB channel and depth channel with an extension on self-supervised learning. Uhrig *et al.* [43] proposed a sparse convolution filter to address the importance of data sparsity. Chodosh *et al.* [7] combined deep learning and dictionary learning into an integrated network to generate dense depth. Chen *et al.* [5] proposed a domain-specific network that applies 2D convolution on image pixels and continuous convolution on 3D points.

Most state-of-art methods tackle the problem in 2D image space via CNN. In fact, the essence of the depth image is 3D but not 2D. Our method takes advantage of this intrinsic property of depth image by performing completion in both 3D point cloud space and 2D image space sequentially, which makes the completion more geometrically meaningful. To our best knowledge, our 3dDepthNet is the first work that integrates 3D point cloud completion into 2D depth image completion problem.

**Depth Prediction** Depth prediction tasks predict depth from a single monocular color image. With the absence of depth measurements, depth prediction is even more challenging than depth completion. Handcrafted features [20, 21, 39] are used in the early study of this problem. Recent learning-based methods also integrated those features into novel loss functions to achieve a better result. For example, Zhou *et al.* [52] proposed to use the photometric loss as supervision during training. Mahjourian *et al.* [30] and Yin *et al.* [48] added 3D geometric constraints to predict depth and estimate optical flow. Qi *et al.* [36] combined

estimated surface normal and depth prediction into a two pathway neural network.

Even though most approaches achieve promising results on prediction, the scale ambiguity problem still exists because no accurate depth guidance is provided. When it comes to industrial applications, the performance of depth prediction based on a single monocular camera might not be stable.

### 3. The Proposed Method

Unlike state-of-art depth completion works which function in 2D image space via CNN [37, 28, 44, 46], we lift the problem to 3D since depth image, in essence, is three dimensional. Specifically, we propose an end-to-end deep learning network to complete the depth in 3D point cloud space and image space sequentially. The network takes sparse depth image, RGB image, and corresponding calibration information<sup>2</sup> as input to generate a dense depth image.

Completion is first performed in 3D space by patching sparse 3D LiDAR points and adjusting the position of corresponding patch points (LiDAR Completion Net). After projecting the dense point cloud onto image space, we pass it through a specially designed encoder-decoder structure (Depth Completion Net) to fuse the information from the original sparse depth image, projected dense depth image, and RGB image to produce a smooth dense depth image. The whole structure is illustrated in Figure 2.

#### 3.1. LiDAR Completion Net

Completion in 3D space is a long-studied subject. Learning-based methods [47, 50, 41] mainly focus on single object point cloud completion instead of the entire scene in our scenario. In our work, the sparse depth information serves as landmark points in 3D space, and we aim to complete the missing neighboring depth information around

<sup>2</sup>Since the calibration information is already used when projecting raw LiDAR data to the image space, we do not consider it as extra information.

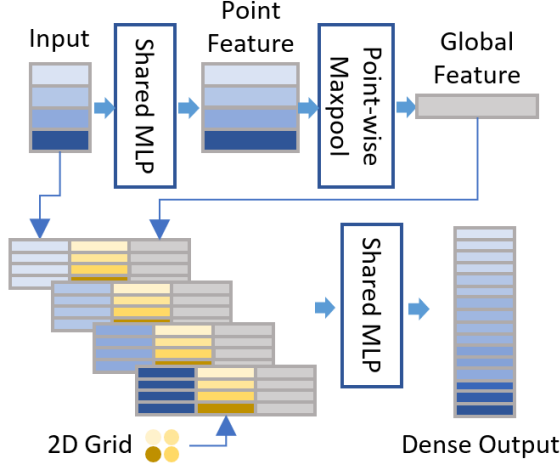


Figure 4: **LiDAR Completion Net.** We first extract global features by PointNet, and then concatenate the 2D grid coordinates with repeated local point coordinates and global features. Next, combined local and global features are fed into a shared-MLP to produce dense point cloud output. Each colored row in the above figure stands for one point in the corresponding feature space; the same color represents the same point.

those anchor points. Inspired by the Point Completion Net[50] (PCN) which is designed for single object point cloud completion, we propose a LiDAR Completion Net by adjusting the structure of PCN, making it more suitable for sparse LiDAR points completion task. The shared-MLP forming coarse level points in PCN is removed, where coarse level points are directly replaced by sparse input points. Moreover, the second half of the encoder in PCN is removed. These changes make our LiDAR Completion Net more efficient and effective in completing sparse LiDAR points, as shown in Figure 5.

As illustrated in Figure 4 and 3, sparse LiDAR points are first fed into a PointNet [35] to extract global feature vector. Meanwhile, raw LiDAR points are also served as coarse level landmarks. Global feature vector, the local landmark, and patch points around the corresponding landmark are then combined into a multi-level shared-MLP to generate dense point clouds.

By directly using raw LiDAR points as coarse level landmarks, the number of parameters used in our LiDAR Completion Net is much less than PCN since most parameters are used for generating coarse level landmarks in original PCN, PointNet encoder and last shared-MLP only have very fewer parameters. In our experiments, LiDAR Completion Net only has around 0.2 million parameters while original PCN has 6.85 million parameters.

### 3.2. Depth Completion Net

The Depth Completion Net is always the core part of methods solving depth completion problem in image space. We propose a network architecture that fuses the information of projected dense depth image from our LiDAR Completion Net, original sparse depth image, and RGB image, as shown in Figure 6.

Inspired by [37], our Depth Completion Net consists of two pathways; one is a dual-depth pathway that integrates completed depth from LiDAR Completion Net and sparse depth, the other one is a RGB-D pathway which takes original RGB image and sparse depth as input. The left pathway is the main pathway, while the right pathway provides guidance during the up-projection stage to generate smoother results. Encoders in both pathways consist of ResNet blocks to obtain a 1/16 downsized feature. The decoder then gradually integrates features from both pathways and increases feature resolution back to the original resolution. As suggested in [37], we concatenate features from the right pathway but sum the features from the left path. The decoder is encouraged to keep consistency and learn more from the completed depth and the original sparse depth when features from the left pathway are added to it [4]. Different from [37] which has three encoder-decoder structures, we only need one encoder decoder, which results in a significant decrease in the total number of parameters. Comparing with one encoder-decoder network such as [28] with ResNet34 [18] as building blocks, we only use ResNet18 which is half of the parameter numbers used in ResNet34. Therefore, our Depth Completion Net is more lightweight and computationally efficient.

### 3.3. Loss Function

Our loss function is defined as:

$$L = \{L_{pt}, L_{im}\} \quad (1)$$

where  $L_{im}$  is the standard masked MSE loss between generated dense depth image and ground truth depth image, and  $L_{pt}$  is the Chamfer Distance [12] between the dense point cloud generated from LiDAR Completion Net and ground truth point cloud.

The Chamfer Distance between two point clouds  $P_1$  and  $P_2$  is defined as:

$$CD(P_1, P_2) = \frac{1}{|P_1|} \sum_{x \in P_1} \min_{y \in P_2} \|x - y\|_2 + \frac{1}{|P_2|} \sum_{y \in P_2} \min_{x \in P_1} \|x - y\|_2 \quad (2)$$

It computes the average closest point distance between the  $P_1$  and  $P_2$ . Since the parameter number and loss magnitude are very different in  $L_{pt}$  and  $L_{im}$ , we train our network in



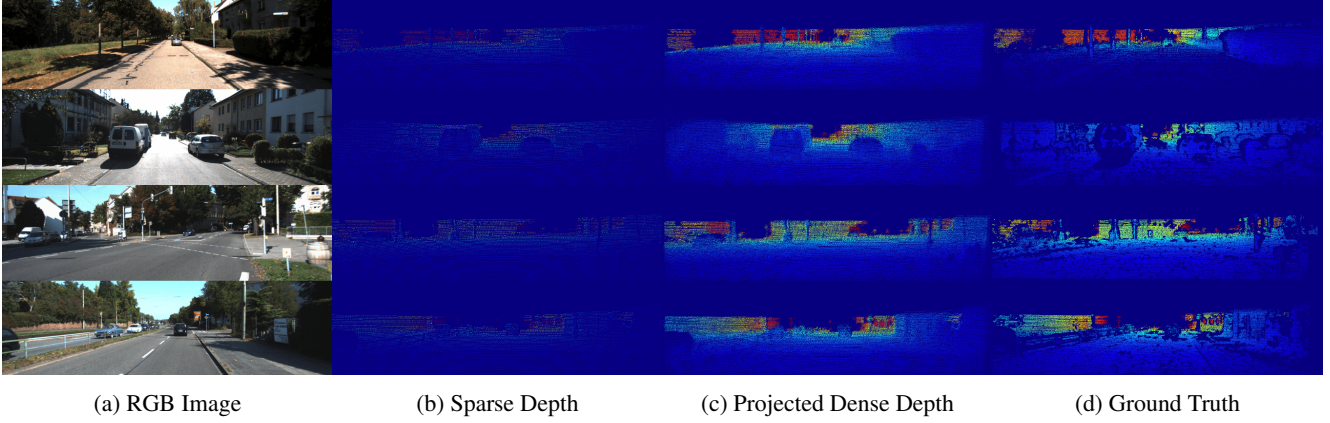


Figure 5: **Visualization of the LiDAR Completion Net.** To visualize the result of LiDAR Completion Net, we project the dense point cloud to the image space and compare it with groundtruth depth image. The third column is the projection result of our dense point cloud generated by the LiDAR Completion Net.

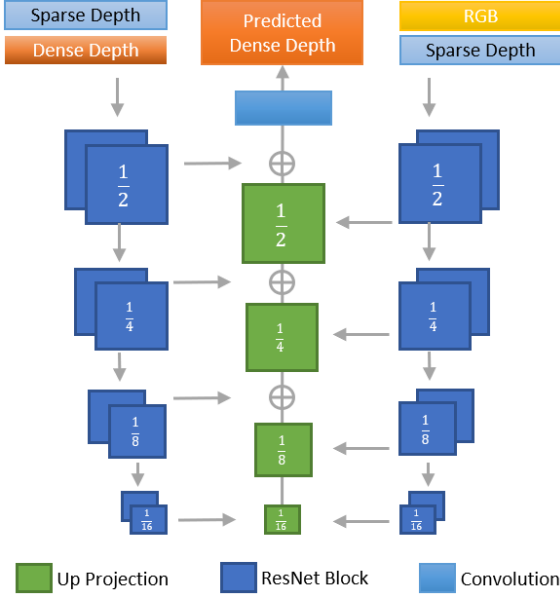


Figure 6: **Depth Completion Net.** Similar to [37], we use late fusion strategy, where different encoder features are only combined in the decoder stage. To emphasize the importance of depth information, we sum the dual-depth channel and concatenate RGB-D channel in the decoder.

two steps. We first freeze parameters in the Depth Completion Net and train parameters in the LiDAR Completion Net for 1 epoch, and then we freeze parameters in the LiDAR completion Net and train parameters in Depth Completion Net for 11 epochs. We use Adam as our optimizer with an initial learning rate 0.001,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and drop learning rate to half for every 2000 steps during training LiDAR Completion Net, for every 5 epoch during

training Depth Completion Net.

## 4. Experiments

In this section, we first describe the training dataset and device. Next, we compare our method against existing state-of-art methods and conduct an ablation study of our method. Finally, we show the potential generalization ability of our model.

### 4.1. Training Data

Our model is trained on KITTI depth completion benchmark [15]. Original KITTI data set provides RGB images, sparse depth images, raw LiDAR scans, and calibration files. We generate point cloud files from raw sparse depth images with corresponding calibration files. We use 1 TITAN RTX GPU for training and 1 RTX 2080Ti GPU for testing. Since we have limited computation resources, the training is on one camera’s data from scratch and fine-tune on both cameras’ data.

### 4.2. Comparison with State-of-art Methods

**Metrics.** Following the KITTI benchmark, we use 4 standard metrics to measure the quality of our method: root mean square error in  $mm$  (RMSE), mean absolute error in  $mm$  (MAE), root mean squared error of the inverse depth in  $1/km$  (iRMSE) and mean absolute error of the inverse depth in  $1/km$  (iMAE). iRMSE and iMAE compute the mean error of inverse depth and concentrate more on close objects while RMSE and MAE directly measure the accuracy of all estimated depth and focus more on further objects. RMSE is used as the dominant metric in the KITTI benchmark leader board.

**Evaluation on KITTI Test set.** KITTI test set contains 1,000 RGB images and corresponding sparse depth. The

	RMSE	MAE	iRMSE	iMAE
HMS-Net [19]	841.78	253.47	2.73	1.13
MSFF-Net [45]	836.69	241.54	2.63	1.07
NConv-CNN [11]	829.98	233.26	2.60	1.03
Sparse-to-Dense [28]	814.73	249.95	2.80	1.21
PwP [46]	777.05	235.17	2.42	1.13
DeepLiDAR [37]	<b>758.38</b>	226.50	2.56	1.15
Ours	798.44	<b>226.27</b>	<b>2.36</b>	<b>1.02</b>

Table 1: **Comparison on KITTI Test Set.** We selected several published state-of-art methods listed on the KITTI leaderboard. The evaluation is done on KITTI testing server. Our method outperforms these methods on iRMSE, MAE, iMAE and achieves comparable RMSE

	Device	Runtime	Converted Run Time
DeepLiDAR [37]	GTX 1080Ti	0.07s	0.07s
Sparse-to-Dense [28]	Tesla V100	0.08s	0.16s
PwP [46]	Tesla V100	0.1s	0.19s
Ours	RTX 2080Ti	0.03s	0.05s

Table 2: **Runtime Comparison on the KITTI validation set.** We list inference speeds and corresponding devices of several state-of-art methods and ours. According to [27], RTX 2080Ti has 1.62 times average speedup over GTX 1080Ti on FP16 computation, and Tesla V100 has 1.98 times average speedup. Hence, our method is the most efficient one with the consideration of different devices.

groundtruth is not provided; all test results will be evaluated on the KITTI test server to prevent overfitting.

Quantitative results of our method and other state-of-art methods are listed in Table 1. Qualitative results are shown in Figure 7. Our method produces more stable and accurate details on close objects than most of the other methods as indicated by lower iRMSE and iMAE.

Moreover, since our LiDAR Completion Net and Depth Completion Net are more lightweight compared to other state-of-art methods, our network is more efficient. We list the device and runtime comparison with several state-of-art methods in Table 2.

**Evaluation on KITTI Validation set.** We further compare our method with some other methods which are not on the KITTI benchmark on the KITTI validation set. KITTI validation set also contains 1,000 RGB images and corresponding sparse depth image with the ground truth dense depth image provided. Compared methods include bilateral filter using color (Bilateral) [42], fast bilateral (Fast) [10], optimization using total variance (TGV) [14], and deep depth completion for indoor scene [51]. The first three methods are non-learning based methods, which do not perform very well because of the huge complexity of the depth completion model. Indoor learning-based methods also do

	RMSE	MAE	iRMSE	iMAE
Bilateral[42]	2989.02	1200.56	9.67	5.08
Fast[10]	3548.87	1767.80	26.48	9.13
TGV[14]	2761.29	1068.69	15.02	6.28
Zhang <i>et al.</i> [51]	1312.10	356.60	4.29	1.41
Ours	<b>693.23</b>	<b>208.96</b>	<b>2.37</b>	<b>0.98</b>

Table 3: **Comparison on KITTI Validation Set.** Our method outperforms three non-learning based methods and one indoor learning based method in all metrics.

not perform very well because of the scene change. The quantitative result is in Table 3.

### 4.3. Ablation Study

To better understand the impact of our network modules, we conduct a systematical ablation study by numerically presenting the influence of disabling specific components in our network. The quantitative result is listed in Table 4, and the qualitative result of LiDAR Completion Net is shown in Figure 8.

**Effectiveness of LiDAR Completion Net.** The dense depth image generated from the LiDAR Completion Net is an essential part of our network. To verify the effectiveness of it, we replace the dense depth in the dual depth channel by the original sparse depth. The Depth Completion Net will be the only component left, so it is equivalent to study the ability of Depth Completion Net. We train the modified network with the same training strategy described in Section 3.3. All the metrics drop significantly. This indicates that the dense depth image generated from LiDAR Completion Net provides additional important information which sparse depth image and RGB image can not provide. This extra information helps improve the final dense depth in all metrics.

Furthermore, in order to show our modification to PCN as described in Section 3.1 is effective, we compare the performance of the PCN [50] and our LiDAR Completion Net on scene point cloud completion. For simplicity, we only train both of them to perform 3D point cloud completion with randomly selected 5,000 pairs of data from KITTI training set for 10 epochs. The Depth Completion Net is removed and only the Chamfer Distance loss is back-propagated. Completed point clouds are projected to image space for visualization. As shown in figure 8, our LiDAR Completion Net is much more suitable for completing sparse point cloud of an entire scene.

**Effectiveness of Depth Completion Net.** We verify the positive effects of our Depth Completion Net by replacing it with the image encoder-decoder structure describe in [28]. The significant difference of [28] with our Depth Completion Net is that it only has one encoder pathway, and all the skip connection is implemented by concatenation. We

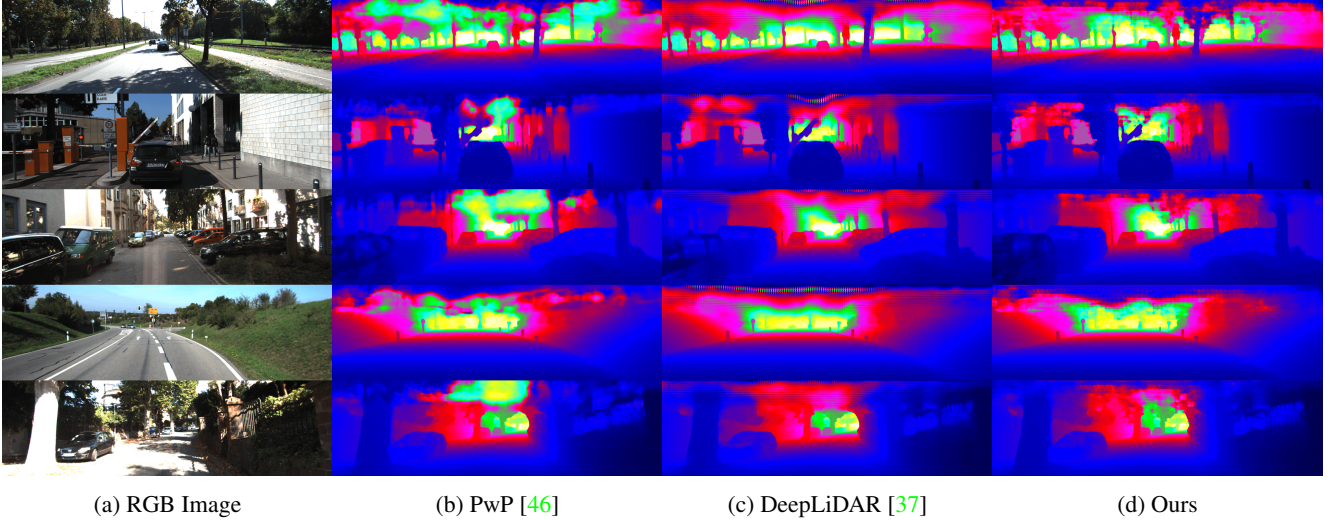


Figure 7: **Visualization Comparison on KITTI Test Set** We list the visualization results of PwP [46], DeepLiDAR [37] and our method of several test images. Our method is more accurate on close objects as indicated by lower iRMSE and iMAE. For areas without groundtruth depth guidance, such as upper part in color image, our network produces more meaningful results than PwP and has less aliasing than DeepLiDAR.

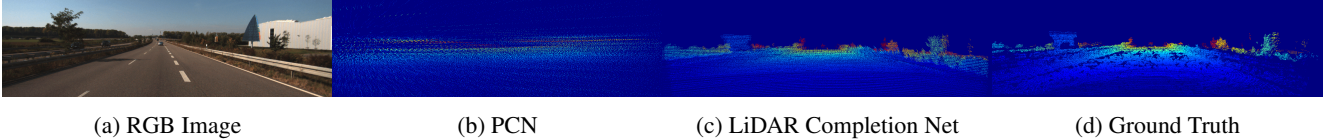


Figure 8: **Comparison of PCN and LiDAR Completion Net.** We visualize the results of the retrained PCN and LiDAR Completion Net for 3D scene point cloud completion. LiDAR Completion Net produces more meaningful projected dense depth.

change the sparse depth input into a two-channel depth by concatenating the dense depth generated from LiDAR Completion Net. The training strategy is the same as described in Section 3.3.

#### 4.4. Generalization

We further test the generalization ability of our method by providing sparser input. In the KITTI data set, the original LiDAR depth provides around 4% depth information which is close to 17,000 pixels with depth in a 1216 by 352 image. To test the generalization ability of our network, We sub-sample the depth information in sparse depth data from LiDAR to 1/4, 1/16, 1/64, and 1/256. The quantitative result of RMSE and iRMSE comparing with several other methods is illustrated in Figure 9. The performance of our network drops as the sparsity increases as expected. Our network outperforms other methods when the sampling ratio becomes lower. The patching process in LiDAR Completion Net helps densify the sparse LiDAR data, and this phenomenon amplifies when sampling density drops. This

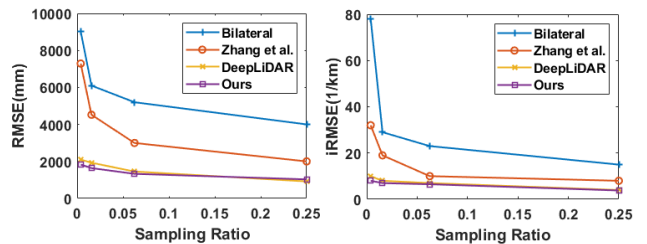


Figure 9: **Generalization test on different sparsity input.** Due to our novel design of the 3D-to-2D coarse-to-fine densification pipeline, in particular our LiDAR Completion Network, our model generalizes well on sub-sampled sparse depth.

explains why our network is more resilient to even sparser depth.



	LiDAR Completion Net	Depth Completion Net	RMSE	MAE	iRMSE	iMAE
Model 1	No	Yes	767.68	249.07	3.67	1.14
Model 2	Yes	No	742.28	221.29	2.52	1.10
Full	Yes	Yes	693.23	208.96	2.37	0.98

Table 4: Ablation study on KITTI validation set.

## 5. Discussion

**3D point cloud completion.** The LiDAR Completion Net is an important part of our network architecture. How to perform a better 3D completion task for large scale point cloud data is essential to generate a dense and smooth depth image. Our LiDAR Completion Net performs *patching* and *adjusting* functions according to global feature and the coordinate of each sparse point. To further improve the quality of LiDAR Completion Net, local geometry information and semantic information can be taken into consideration. Object segmentation information in 3D space can also help perform an object-specific point cloud completion. Specifically, points describing people or cars should be treated differently in LiDAR Completion Net.

**Penetration Problem.** Once the completion is finished in 3D space, the dense point cloud will be projected into image space to generate a dense depth image. However, since objects described by point cloud are not solid, points describing faraway objects may penetrate through close object point cloud onto the projection plane, as illustrated in Figure 10. The original sparse depth image does not have this problem because the LiDAR scan can not penetrate through solid objects. In our network, the Depth Completion Net learns to solve the penetration problem. However, we feel like there is no geometric explanation and theoretical guarantee on how good this problem is solved by learning. We aim to solve it more elegantly in future work.

**Combination of 2D and 3D completion.** Most previous methods to depth completion employ 2D image encoder-decoder structures that take color and sparse depth image as input. In our case, with the generated intermediate dense point cloud from Lidar Completion Net, how to design an encoder-decoder structure that efficiently utilizes RGB image, sparse depth, and intermediate dense point cloud, becomes essential. Our Depth Completion Net with dual depth channel plus RGB-D channel design is modified from 2D image encoder-decoder structure, which only takes image input. We will explore possibilities of designing a new encoder-decoder structure that directly uses information from 3D point cloud and 2D images.

## 6. Conclusion

We propose a new approach to the depth completion problem by introducing 3D completion into the learning

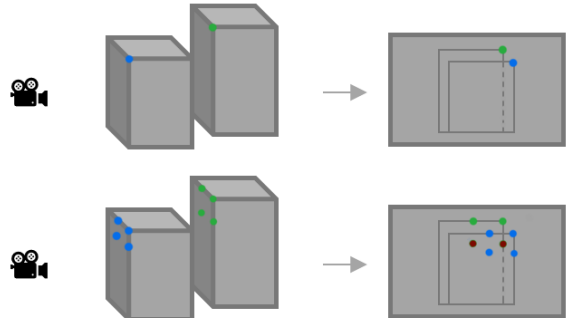


Figure 10: **Penetration Problem.** The first row demonstrates the process from LiDAR scan points (first column) to depth image (second column), which has no penetration problem since the LiDAR sensor can not penetrate through solids. The second row is the dense point cloud generated by LiDAR Completion Net. After projection, the two red dots will exist in the depth image since blue points can not block them. However, for LiDAR scan depth image, red dots are not allowed to show up in the smaller rectangular region according to the perspective principle.

network. Our novel deep learning network offers a 3D-to-2D coarse-to-fine dual densification design. By taking advantage of the dimensional nature of depth image, we propose a novel LiDAR Completion Net to do completion in 3D space and pass it to a specially designed Depth Completion Net which integrates projected dense depth, sparse depth and RGB image to produce smooth dense depth image. Experiments show our network is efficient and achieves state-of-art accuracy. Ablation and generalization tests prove that each module in our network has positive influences on the final results, and furthermore, our network is resilient to even sparser depth.

## References

- [1] Jonathan T Barron and Ben Poole. The fast bilateral solver. In *European Conference on Computer Vision*, pages 617–632. Springer, 2016. 2
- [2] Matthew Berger, Andrea Tagliasacchi, Lee Seversky, Pierre Alliez, Joshua Levine, Andrei Sharf, and Claudio Silva. State of the art in surface reconstruction from point clouds. 2014. 2
- [3] Anne-Laure Chauve, Patrick Labatut, and Jean-Philippe Pons. Robust piecewise-planar 3d reconstruction and com-



- pletion from large-scale unstructured point data. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 1261–1268. IEEE, 2010. 2
- [4] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1907–1915, 2017. 2, 4
- [5] Yun Chen, Bin Yang, Ming Liang, and Raquel Urtasun. Learning joint 2d-3d representations for depth completion. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 10023–10032, 2019. 3
- [6] Xinjing Cheng, Peng Wang, and Ruigang Yang. Depth estimation via affinity learned with convolutional spatial propagation network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 103–119, 2018. 2
- [7] Nathaniel Chodosh, Chaoyang Wang, and Simon Lucey. Deep convolutional compressed sensing for lidar depth completion. In *Asian Conference on Computer Vision*, pages 499–513. Springer, 2018. 3
- [8] James Davis, Stephen R Marschner, Matt Garr, and Marc Levoy. Filling holes in complex surfaces using volumetric diffusion. In *Proceedings. First International Symposium on 3D Data Processing Visualization and Transmission*, pages 428–441. IEEE, 2002. 2
- [9] James Diebel and Sebastian Thrun. An application of markov random fields to range sensing. In *Advances in neural information processing systems*, pages 291–298, 2006. 2
- [10] Frédo Durand and Julie Dorsey. Fast bilateral filtering for the display of high-dynamic-range images. In *ACM transactions on graphics (TOG)*, volume 21, pages 257–266. ACM, 2002. 6
- [11] Abdelrahman Eldesokey, Michael Felsberg, and Fahad Shahbaz Khan. Propagating confidences through cnns for sparse data regression. *arXiv preprint arXiv:1805.11913*, 2018. 6
- [12] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 605–613, 2017. 4
- [13] Sean Ryan Fanello, Julien Valentin, Christoph Rhemann, Adarsh Kowdle, Vladimir Tankovich, Philip Davidson, and Shahram Izadi. Ultrastereo: Efficient learning-based matching for active stereo systems. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6535–6544. IEEE, 2017. 1
- [14] David Ferstl, Christian Reinbacher, Rene Ranftl, Matthias Rüther, and Horst Bischof. Image guided depth upsampling using anisotropic total generalized variation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 993–1000, 2013. 6
- [15] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013. 5
- [16] Feng Han and Song-Chun Zhu. Bottom-up/top-down image parsing with attribute grammar. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(1):59–73, 2008. 2
- [17] Simon Hawe, Martin Kleinsteuber, and Klaus Diepold. Dense disparity maps from sparse disparity measurements. In *2011 International Conference on Computer Vision*, pages 2126–2133. IEEE, 2011. 3
- [18] 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*, pages 770–778, 2016. 4
- [19] Zixuan Huang, Junming Fan, Shuai Yi, Xiaogang Wang, and Hongsheng Li. Hms-net: Hierarchical multi-scale sparsity-invariant network for sparse depth completion. *arXiv preprint arXiv:1808.08685*, 2018. 6
- [20] Kevin Karsch, Ce Liu, and S Kang. Depth extraction from video using non-parametric sampling-supplemental material. In *European conference on Computer Vision*. Citeseer, 2012. 2, 3
- [21] Kevin Karsch, Ce Liu, and Sing Bing Kang. Depth transfer: Depth extraction from video using non-parametric sampling. *IEEE transactions on pattern analysis and machine intelligence*, 36(11):2144–2158, 2014. 3
- [22] Vladimir G Kim, Wilmot Li, Niloy J Mitra, Siddhartha Chaudhuri, Stephen DiVerdi, and Thomas Funkhouser. Learning part-based templates from large collections of 3d shapes. *ACM Transactions on Graphics (TOG)*, 32(4):70, 2013. 2
- [23] Yangyan Li, Angela Dai, Leonidas Guibas, and Matthias Nießner. Database-assisted object retrieval for real-time 3d reconstruction. In *Computer Graphics Forum*, volume 34, pages 435–446. Wiley Online Library, 2015. 2
- [24] Yangyan Li, Xiaokun Wu, Yiorgos Chrysathou, Andrei Sharf, Daniel Cohen-Or, and Niloy J Mitra. Globfit: Consistently fitting primitives by discovering global relations. *ACM transactions on graphics (TOG)*, 30(4):52, 2011. 2
- [25] Or Litany, Alex Bronstein, Michael Bronstein, and Ameesh Makadia. Deformable shape completion with graph convolutional autoencoders. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1886–1895, 2018. 2
- [26] Lee-Kang Liu, Stanley H Chan, and Truong Q Nguyen. Depth reconstruction from sparse samples: Representation, algorithm, and sampling. *IEEE Transactions on Image Processing*, 24(6):1983–1996, 2015. 3
- [27] Lambda LLC. Deep learning gpu benchmarks - tesla v100 vs rtx 2080 ti vs gtx 1080 ti vs titan v, 2018. 6
- [28] Fangchang Ma, Guilherme Venturéli Cavalheiro, and Sertac Karaman. Self-supervised sparse-to-dense: self-supervised depth completion from lidar and monocular camera. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 3288–3295. IEEE, 2019. 1, 3, 4, 6
- [29] Oisín Mac Aodha, Neill DF Campbell, Arun Nair, and Gabriel J Brostow. Patch based synthesis for single depth image super-resolution. In *European conference on computer vision*, pages 71–84. Springer, 2012. 2
- [30] Reza Mahjourian, Martin Wicke, and Anelia Angelova. Un-supervised 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*, pages 5667–5675, 2018. 3

- [31] Niloy J Mitra, Leonidas J Guibas, and Mark Pauly. Partial and approximate symmetry detection for 3d geometry. In *ACM Transactions on Graphics (TOG)*, volume 25, pages 560–568. ACM, 2006. 2
- [32] Niloy J Mitra, Mark Pauly, Michael Wand, and Duygu Ceylan. Symmetry in 3d geometry: Extraction and applications. In *Computer Graphics Forum*, volume 32, pages 1–23. Wiley Online Library, 2013. 2
- [33] Andrew Nealen, Takeo Igarashi, Olga Sorkine, and Marc Alexa. Laplacian mesh optimization. In *Proceedings of the 4th international conference on Computer graphics and interactive techniques in Australasia and Southeast Asia*, pages 381–389. ACM, 2006. 2
- [34] Mark Pauly, Niloy J Mitra, Johannes Wallner, Helmut Pottmann, and Leonidas J Guibas. Discovering structural regularity in 3d geometry. In *ACM transactions on graphics (TOG)*, volume 27, page 43. ACM, 2008. 2
- [35] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 652–660, 2017. 4
- [36] Xiaojuan Qi, Renjie Liao, Zhengzhe Liu, Raquel Urtasun, and Jiaya Jia. Geonet: Geometric neural network for joint depth and surface normal estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 283–291, 2018. 3
- [37] Jiaxiong Qiu, Zhaopeng Cui, Yinda Zhang, Xingdi Zhang, Shuaicheng Liu, Bing Zeng, and Marc Pollefeys. DeepLidar: Deep surface normal guided depth prediction for outdoor scene from sparse lidar data and single color image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3313–3322, 2019. 3, 4, 5, 6, 7
- [38] Sean Ryan Fanello, Christoph Rhemann, Vladimir Tankovich, Adarsh Kowdle, Sergio Orts Escolano, David Kim, and Shahram Izadi. Hyperdepth: Learning depth from structured light without matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5441–5450, 2016. 1
- [39] Ashutosh Saxena, Sung H Chung, and Andrew Y Ng. Learning depth from single monocular images. In *Advances in neural information processing systems*, pages 1161–1168, 2006. 3
- [40] Elham Shabaninia, Ahmad Reza Naghsh-Nilchi, and Shohreh Kasaei. High-order markov random field for single depth image super-resolution. *IET Computer Vision*, 11(8):683–690, 2017. 2
- [41] Minhyuk Sung, Vladimir G Kim, Roland Angst, and Leonidas Guibas. Data-driven structural priors for shape completion. *ACM Transactions on Graphics (TOG)*, 34(6):175, 2015. 3
- [42] Carlo Tomasi and Roberto Manduchi. Bilateral filtering for gray and color images. In *Iccv*, volume 98, page 2, 1998. 6
- [43] Jonas Uhrig, Nick Schneider, Lukas Schneider, Uwe Franke, Thomas Brox, and Andreas Geiger. Sparsity invariant cnns. In *2017 International Conference on 3D Vision (3DV)*, pages 11–20. IEEE, 2017. 3
- [44] Wouter Van Gansbeke, Davy Neven, Bert De Brabandere, and Luc Van Gool. Sparse and noisy lidar completion with rgb guidance and uncertainty. *arXiv preprint arXiv:1902.05356*, 2019. 3
- [45] Benzhang Wang, Yiliu Feng, and Hengzhu Liu. Multi-scale features fusion from sparse lidar data and single image for depth completion. *Electronics Letters*, 54(24):1375–1377, 2018. 6
- [46] Yan Xu, Xinge Zhu, Jianping Shi, Guofeng Zhang, Hujun Bao, and Hongsheng Li. Depth completion from sparse lidar data with depth-normal constraints. *arXiv preprint arXiv:1910.06727*, 2019. 3, 6, 7
- [47] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 206–215, 2018. 3
- [48] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1983–1992, 2018. 3
- [49] Lap-Fai Yu, Sai-Kit Yeung, Yu-Wing Tai, and Stephen Lin. Shading-based shape refinement of rgb-d images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1415–1422, 2013. 2
- [50] Wentao Yuan, Tejas Khot, David Held, Christoph Mertz, and Martial Hebert. Pcn: Point completion network. In *2018 International Conference on 3D Vision (3DV)*, pages 728–737. IEEE, 2018. 2, 3, 4, 6
- [51] Yinda Zhang and Thomas Funkhouser. Deep depth completion of a single rgb-d image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 175–185, 2018. 2, 6
- [52] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1851–1858, 2017. 3