

HMS-Net: Hierarchical Multi-scale Sparsity-invariant Network for Sparse Depth Completion

Zixuan Huang, Junming Fan, Shenggan Cheng, Shuai Yi, Xiaogang Wang, *Senior Member, IEEE*,
Hongsheng Li, *Member, IEEE*

Abstract—Dense depth cues are important and have wide applications in various computer vision tasks. In autonomous driving, LIDAR sensors are adopted to acquire depth measurements around the vehicle to perceive the surrounding environments. However, depth maps obtained by LIDAR are generally sparse because of its hardware limitation. The task of depth completion attracts increasing attention, which aims at generating a dense depth map from an input sparse depth map. To effectively utilize multi-scale features, we propose three novel sparsity-invariant operations, based on which, a sparsity-invariant multi-scale encoder-decoder network (HMS-Net) for handling sparse inputs and sparse feature maps is also proposed. Additional RGB features could be incorporated to further improve the depth completion performance. Our extensive experiments and component analysis on two public benchmarks, KITTI depth completion benchmark and NYU-depth-v2 dataset, demonstrate the effectiveness of the proposed approach. As of Aug. 12th, 2018, on KITTI depth completion leaderboard, our proposed model without RGB guidance ranks **1st** among all peer-reviewed methods without using RGB information, and our model with RGB guidance ranks **2nd** among all RGB-guided methods.

Index Terms—depth completion, convolutional neural network, sparsity-invariant operations.

1 INTRODUCTION

DEPTH completion, aiming at generating a dense depth map from the input sparse depth map, is an important task for computer vision and robotics. In Fig. 2 (a), (b), (e), we show one example input sparse depth map, its corresponding RGB image, and the depth completion result by our proposed method. Because of the limitation of current LIDAR sensors, the inputs of depth completion are generally sparse. For instance, the \$100,000 Velodyne HDL-64E has only a vertical resolution of $\sim 0.4^\circ$ and an azimuth angular resolution of 0.08° . It generates sparse depth maps, which might be insufficient for many real-world applications. Depth completion algorithms could estimate dense depth maps from sparse inputs and has great potential in practice. With an accurate depth completion algorithm, many high-level vision tasks, such as semantic segmentation, 3D object detection, visual odometry and SLAM with 3D point clouds, could be solved more effectively. Therefore, it becomes a hot research topic for self-driving cars and UAVs, and is listed as one of the ranked tasks in the KITTI benchmark [1].

Many different methods have been proposed for depth completion, which could be generally categorized into learning-based [1], [2], [3], [4] and non-learning-based methods [5], [6], [7]. Non-learning-based approaches generate

dense depth maps from sparse inputs based on hand-crafted rules. Therefore, the outputs of these algorithms are generated based on assumed prior by humans. As a result, they are not robust enough to sensor noises and are usually specifically designed for certain datasets. In addition, most non-learning-based methods ignore the correlations among sparse input depth points and might result in inaccurate object boundaries. An example of errors by a non-learning-based method [5] is shown in Fig. 2(e). The noises in the white box are not removed at all, and boundaries of the cars and trees in the yellow box are inaccurate.

For learning-based approaches, state-of-the-art methods are mainly based on deep neural networks. Previous methods mainly utilized deep convolutional neural networks (CNN) for generating dense depth maps from sparse inputs. Ma and Karaman [3] simply filled 0s to locations without depth inputs to create dense input maps, which might introduce ambiguity to very small depth values. Chodosh et al. [4] proposed to extract multi-level sparse codes from the inputs and used a 3-layer CNN for depth completion. However, those two methods used conventional convolution operations designed for dense inputs (see Fig. 2(c) for an example). Uhrig et al. [1] proposed sparsity-invariant convolution, which is specifically designed to process sparse maps and enables processing sparse inputs more effectively with CNN.

However, the sparsity-invariant convolution in [1] only mimics the behavior of convolution operations in conventional dense CNNs. Its feature maps of later stages lose much spatial information and therefore cannot effectively integrate both low-level and high-level features for accurate depth completion (see Fig. 1(a) for illustration). On the other hand, there exist effective multi-scale encoder-decoder net-

- Zixuan Huang is with Department of Computer Sciences, University of WisconsinMadison, United States
- Junming Fan, Shenggan Cheng and Shuai Yi are with SenseTime Research, Beijing, China.
- Xiaogang Wang and Hongsheng Li are with the Department of Electronic Engineering at Chinese University of Hong Kong, Hong Kong, China.
- The first two authors contribute equally to the paper. They finished the work while they were at SenseTime Research.
- E-mail: {hsli, xgwang}@ee.cuhk.edu.hk

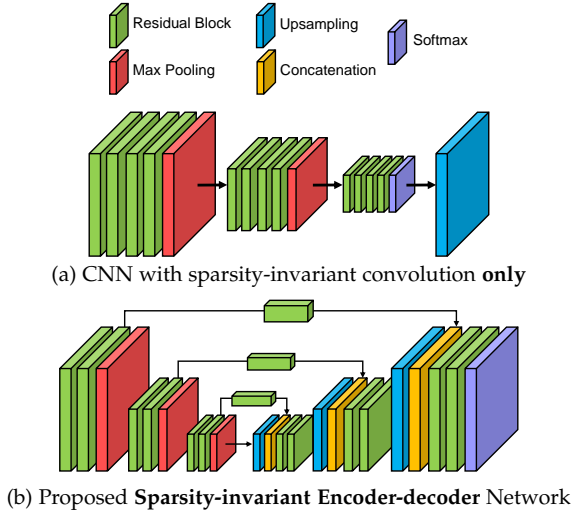


Fig. 1. (a) CNN with only sparsity-invariant convolution could only gradually downsample feature maps, which loses much resolution at later stages. (b) Our proposed encoder-decoder network with novel sparsity-invariant operations could effectively fuse multi-scale features from different layers for depth completion.

work structures for dense pixel-wise classification tasks (see Fig. 1(b)), such as U-Net [8], Feature Pyramid Network [9], Full Resolution Residual Network [10]. Direct integration of the sparsity-invariant convolution in [1] into the multi-scale structures is infeasible, as those structures also require other operations for multi-scale feature fusion, such as sparsity-invariant feature upsampling, average, and concatenation.

To overcome such limitation, we propose three novel sparsity-invariant operations to enable using encoder-decoder networks for depth completion. The three novel operations include sparsity-invariant upsampling, sparsity-invariant average, and joint sparsity-invariant concatenation and convolution. To effectively and efficiently handle sparse feature maps, sparsity masks are utilized at all locations of the feature maps. They record the locations of the sparse features at the output of each processing stage and guide the calculation of the forward and backward propagation. Each sparsity-invariant operation is designed to properly maintain and modify the sparsity masks across the network. The design of those operations are non-trivial and are the keys to using encoder-decoder structures with sparse features. Based on such operations, we propose a multi-scale encoder-decoder network, HMS-Net, which adopts a series of sparsity-invariant convolutions with downsampling and upsampling to generate multi-scale feature maps and shortcut paths for effectively fusing multi-scale features. Extensive experiments on KITTI [1] and NYU-depth-v2 [11] datasets show that our algorithm achieves state-of-the-art depth completion accuracy.

The main contributions of our work can be summarized to threefold. 1) We design three sparsity-invariant operations for handling sparse inputs and feature maps, which are important to handle sparse feature maps. 2) Based on the proposed sparsity-invariant operations, a novel hierarchical multi-scale network structure fusing information from different scales is designed to solve the depth completion task. 3) Our method outperforms state-of-the-art methods in depth completion. On KITTI depth completion benchmark,

our method without RGB information ranks *1st* among all peer-reviewed methods with RGB inputs, while our method with RGB guidance ranks *2nd* among all RGB-guided methods.

2 RELATED WORK

2.1 Depth completion

Depth completion is an active research area with a large number of applications. According to the sparsity of the inputs, current methods could be divided into two categories: sparse depth completion and depth enhancement. The former methods aim at recovering dense depth maps from spatially sparse inputs, while the later methods work on conventional RGB-D depth data and focus on filling irregular and relatively small holes in input dense depth maps. Besides, if the input depth maps are regularly sampled, the depth completion task could be regarded as a depth upsampling (also known as depth super-resolution) task. In other words, depth upsampling algorithms handle a special subset of the depth completion task. The inputs of depth upsampling are depth maps of lower resolution. According to whether RGB information is utilized, depth upsampling methods could be divided into two categories: guided depth upsampling and depth upsampling without guidance.

2.1.1 Sparse depth completion

This type of methods take sparse depth maps as inputs. To handle sparse inputs and sparse intermediate feature maps, Uhrig et al. [1] proposed sparsity-invariant convolution to replace the conventional convolution in convolution neural networks (CNN). The converted sparsity-invariant CNN keeps track of sparsity masks at each layer and is able to estimate dense depth maps from sparse inputs. Ku et al. [5] proposed to use a series of hand-crafted image processing algorithms to transform the sparse depth inputs into dense depth maps. The proposed framework first utilized conventional morphological operations, such as dilation and closure, to make the input depth maps denser. It then filled holes in the intermediate denser depth maps to obtain the final outputs. Eldesokey et al. [12] proposed an algebraically-constrained normalized convolution operation for handling sparse data and propagate depth confidence across layers. The regression loss and depth confidence are jointly optimized. Ren et al. [13] focused on another task, the efficiency of convolution, but also have the design of sparsity masks. However, their algorithm could not be used in the scenario where the sparsity of the input varies in a large range. The above mentioned methods did not consider image information captured by the calibrated RGB cameras.

There also exist works utilizing RGB images as additional information to achieve better depth completion. Schneider et al. [14] combined both intensity cues and object boundary cues to complete sparse depth maps. Liao et al. [2] utilized a residual neural network and combine the classification and regression losses for depth estimation with both RGB and sparse depth maps as inputs. Ma and Karaman [3] proposed to use a single deep regression network to learn directly from the RGB-D raw data, where the depth channel

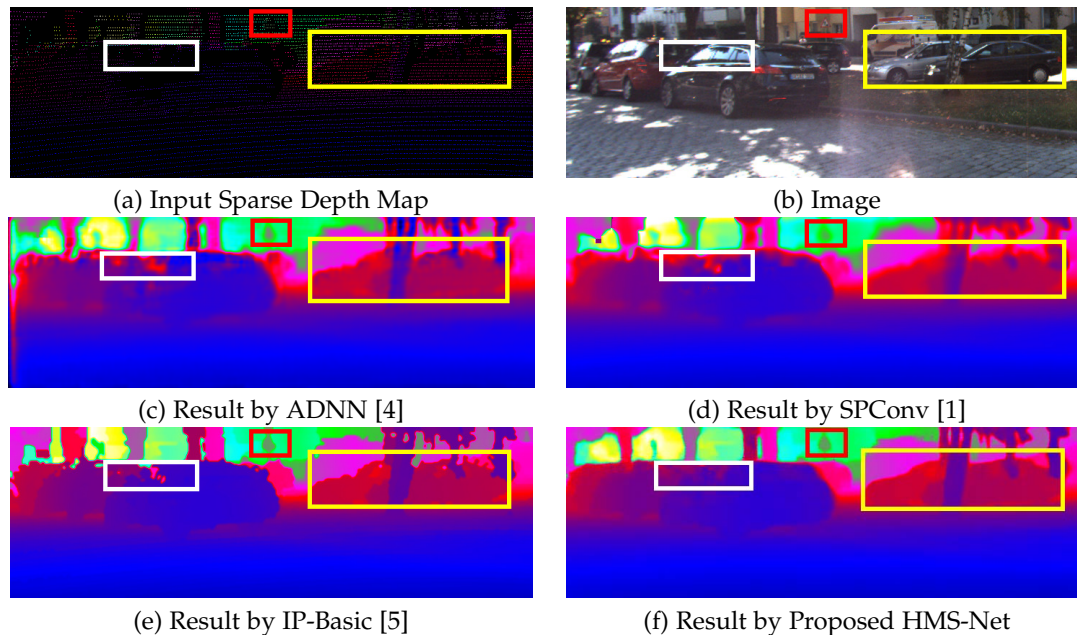


Fig. 2. Illustration of depth completion results by previous methods and our proposed HMS-Net. (a) An example input sparse depth map. (b) Corresponding RGB image. (c) Result by ADNN [4]. (d) Result by sparse convolution [1]. (e) Result by hand-crafted classical image processing method [5]. (f) Result by the proposed HMS-Net.

only has sparse depth values. However, the proposed algorithm mainly focused on the indoor scenes and was only tested on the indoor NYUv2 dataset. Instead, Van Gansbeke et al. [15] combine RGB and depth information through summation according to two predicted confidence maps. Chodosh et al. [4] utilized compressed sensing techniques and Alternating Direction Neural Networks to create a deep recurrent auto-encoder for depth completion. Sparse codes were extracted at the outputs of multiple levels of the CNN and were used to generate dense depth prediction. Zhang and Funkhouser [16] adopted a neural network to predict dense surface normals and occlusion boundaries from RGB images. These predictions were then used as auxiliary information to accomplish the depth completion task from sparse depth data. Qiu et al. [17] further extended this idea on outdoor datasets by generating surface normals as intermediate representation. Jaritz et al. [18] argued that by using networks with large receptive field, the networks were not required to have special treatment for sparse data. They instead trained networks with depth maps of different sparsities. Cheng et al. [19] proposed to learn robust affinities between pixels to spatially propagate depth values via convolution operations to fulfill the entire depth map, while Eldesokey et al. [20] used continuous confidence instead of validity maps and algebraically constrained filters to tackle the sparsity-invariance problem and also guidance from both RGB images and the output confidence produced by their unguided network. Yang et al. [21] proposed to yield the full posterior over depth maps with a Conditional Prior Network from their previous work [22]. And a self-supervised training framework was designed by Ma et al. [23], which explores temporal relations of video sequences to provide additional photometric supervisions for depth completion networks.

2.1.2 Depth enhancement

The inputs of depth enhancement or depth hole-filling methods are usually dense depth maps with irregular and rare small holes. The input depth maps are usually captured with RGB images. Matyunin et al. [24] used the depth from the neighborhoods of the hole regions to fill the holes, according to the similarity of RGB pixels. In [25], the missing depth values are obtained by iteratively applying a joint-bilateral filter to the hole regions' neighboring pixels. Yang et al. [26] propose an efficient depth image recovery algorithm based on auto-regressive correlations and recover high-quality multi-view depth frames with it. They also utilized color image, depth maps from neighboring views and temporal adjacency to help the recovery. Chen et al. [27] firstly created smooth regions around the pixels without depth, and then adopted a bilateral filter without smooth region constraints to fill the depth values. Yang et al. [28] proposed an adaptive color-guided autoregressive model for depth enhancement, which utilized local correlation in the initial depth map and non-local similarity in the corresponding RGB image.

2.1.3 Guided depth upsampling

Depth upsampling methods take low-resolution depth maps as inputs and output high-resolution ones. As the guidance signals, the provided RGB images bring valuable information (e.g., edges) for upsampling. Li et al. [29] proposed a CNN to extract features from the low-resolution depth map and the guidance image to merge their information for estimating the upsampled depth map. In [30], an anisotropic diffusion tensor is calculated from a high-resolution intensity image to serve as the guidance. An energy function is designed and solved to solve the depth upsampling problem. Hui et al. [31] proposed a convolution neural network, which fused the RGB guidance signals at

different stages. They trained the neural network in the high-frequency domain. Xie et al. [32] upsampled the low-resolution depth maps with the guidance of image edge maps and a Markov Random Field model. Jiang et al. [33] proposed a depth super-resolution algorithm utilizing transform domain regularization with an auto-regressive model, as well as a spatial domain regularization by injecting a multi-directional total variation prior.

2.1.4 Depth upsampling without guidance

Depth upsampling could also be achieved without the assistance of corresponding RGB images. As an MRF based method, the approach in [34] matched against the height field of each input low-resolution depth patch, and searched the database for a list of most appropriate high-resolution candidate patches. Selecting the correct candidates was then posed as a Markov Random Field labeling problem. Ferstl et al. [35] learned a dictionary of edge priors from an external database of high- and low-resolution depth samples, and utilized a novel variational sparse coding approach for upsampling. Xie et al. [36] used a coupled dictionary learning method with locality coordinate constraints to transform the original depth maps into high-resolution depth maps. Riegler et al. [37] integrated a variational method that modeled the piecewise affine structures in depth maps on top of a deep network for depth upsampling.

2.2 Multi-scale networks for pixelwise prediction

Neural networks that utilize multi-scale feature maps for pixelwise prediction (e.g., semantic segmentation) were widely investigated. Combining both low-level and high-level features was proven to be crucial for making accurate pixelwise prediction.

Ronneberger et al. [8] proposed a U-shaped network (U-Net). It consists of an iterative downsampling image encoder that gradually summarized image information into smaller but deeper feature maps, and an iterative upsampling decoder that gradually combined low-level and high-level features to output the pixelwise prediction maps. The low-level information from the encoder was passed to the high-level information of the decoder by direct concatenation of the feature maps of the same spatial sizes. Such a network has shown its great usefulness in many 2D and 3D segmentation tasks. Similar network structures, which include Hourglass [38] and Feature Pyramid Network [9], have also been investigated to tackle pixelwise prediction tasks. Recently, Pohlen et al. [10] proposed the full-resolution residual network (FRRN), which treated full-resolution information as a residual flow to pass valuable information across different scales for semantic segmentation. He et al. [39] designed a fully fused network to utilize both images and focal length to learn the depth map.

However, even with the sparsity-invariant convolution operation proposed in [1], the multi-scale encoder-decoder networks cannot be directly converted to handle sparse inputs. This is because there exist many operations that do not support sparse feature maps. Our proposed sparsity-invariant operations solve the problem and allow encoder-decoder networks to be used for sparse data.

3 METHOD

We introduce our proposed framework for depth completion in this section. In Section 3.1, we first review sparsity-invariant convolution proposed in [1]. In Section 3.2, we then introduce three novel sparsity-invariant operations, which are crucial for adopting multi-scale encoder-decoder networks to process sparse inputs. In Section 3.3, based on such sparsity-invariant operators, the hierarchical multi-scale encoder-decoder network (HMS-Net) for effectively combining multi-scale features is proposed to tackle the depth completion task.

3.1 Sparsity-invariant Convolution

In this subsection, we first review sparsity-invariant convolution in [1], which modifies conventional convolution to handle sparse input feature maps. The sparsity-invariant convolution is formulated as

$$z(u, v) = \frac{\sum_{i,j=-k}^k \mathbf{m}_x(u+i, v+j) \mathbf{w}(i, j) \mathbf{x}(u+i, v+j)}{\sum_{i,j=-k}^k \mathbf{m}_x(u+i, v+j) + \epsilon} + b. \quad (1)$$

The sparsity-invariant convolution takes a sparse feature map \mathbf{x} and a binary single-channel sparsity mask \mathbf{m} as inputs, which both has the same spatial size $H \times W$. The convolution generates output features $z(u, v)$ for each location (u, v) . At each spatial location (u, v) , the binary sparsity mask $\mathbf{m}_x(u, v)$ records whether there are input features at this location, i.e., 1 for features existence and 0 otherwise. The convolution kernel \mathbf{w} is of size $(2k+1) \times (2k+1)$, and b represents a learnable bias vector. Note that the kernel weights \mathbf{w} and bias vector b are learned via back-propagation, while the sparsity mask \mathbf{m}_x is specified by the previous layer and is trained.

The key difference with conventional convolution is the use of binary sparsity mask \mathbf{m}_x for convolution calculation. The mask values on numerator of Eq. (1) denote that when conducting convolution, only input features at the valid or visible locations specified by the sparsity mask \mathbf{m}_x are considered. The mask values in denominator denote that, since only a subset of input features are involved, the output features should be normalized according to the number of valid input locations. ϵ represents a very small number and is used to avoid division by 0.

Note that the sparsity mask should always indicate the validity or sparsity of each location of the feature maps. Since the convolution layers in a neural network is generally stacked for multiple times, the output sparsity mask \mathbf{m}_z at each stage should be modified to match the output of each stage z . For each output feature location (u, v) , if there is at least one valid input location in its receptive field of the previous input, its sparsity mask $\mathbf{m}_z(u, v)$ should be updated to 1. In practice, the output sparsity mask is obtained by conducting max pooling on the input sparsity mask with the same kernel size of convolution $(2k+1) \times (2k+1)$.

3.2 Sparsity-invariant operations

The sparsity-invariant convolution successfully converts conventional convolution to handle sparse input features and is able to stack multiple stages for learning highly

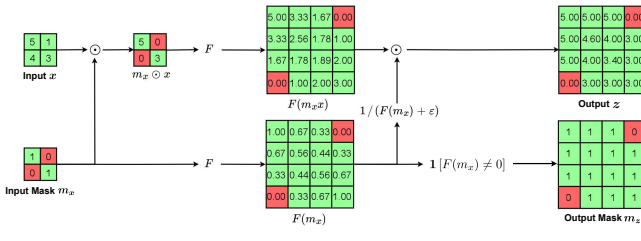


Fig. 3. Illustration of the proposed sparsity-invariant upsampling operation. F stands for Bilinear Upsampling.

non-linear functions. However, only modifying convolution operations is not enough if one tries to utilize state-of-the-art multi-scale encoder-decoder structure for pixelwise prediction. As shown in Fig. 1(b), average, upsampling, and concatenation are also common operations in a multi-scale encoder-decoder networks. Therefore, we propose the sparsity-invariant version of these operations: sparsity-invariant average, sparsity-invariant upsampling, and joint sparsity-invariant concatenation and convolution. The three sparsity-invariant operations allow effectively handling sparse feature maps across the entire encoder-decoder network. They are the foundations of complex building blocks in our overall framework.

Designing the sparsity-invariant operations is non-trivial. Our proposed sparsity-invariant operations share the similar spirit of sparsity-invariant convolution: using single-channel sparsity masks to track the validity of feature map locations. The sparsity masks could be used to guide and regularize the calculation of the operations.

3.2.1 Sparsity-invariant bilinear upsampling

One of the most important operations in encoder-decoder networks is the upsampling operation in the decoder part. We first propose the sparsity-invariant bilinear upsampling operation. Let x and m_x denote the input sparse feature map and the corresponding input sparsity mask of size $H \times W$. The operation generates the output feature map z and its corresponding sparsity mask m_z of size $2H \times 2W$. Let F represents the conventional bilinear upsampling operator, which bilinearly upsamples the input feature map or mask by two times. The proposed sparsity-invariant bilinear upsampling can be formulated as

$$z = \frac{F(m_x \odot x)}{F(m_x) + \epsilon}, \quad (2)$$

$$m_z = \mathbf{1}[F(m_x) \neq 0], \quad (3)$$

where \odot denotes the spatial elementwise multiplication, ϵ is a very small number to avoid division by zero, and $\mathbf{1}[\cdot]$ denotes the indicator function, i.e., $\mathbf{1}[\text{true}] = 1$ and $\mathbf{1}[\text{false}] = 0$. The proposed sparsity-invariant bilinear upsampling operation is illustrated in Fig. 3.

As shown by Eq. (2), the proposed operation first uses the input sparsity mask m_x to mask out the invalid features from the input feature maps x as $m_x \odot x$. The traditional bilinear upsampling F operator is then applied to upsample both the masked feature maps $m_x \odot x$ and the input sparsity mask m_x . The upsampled sparse features $F(m_x \odot x)$ are

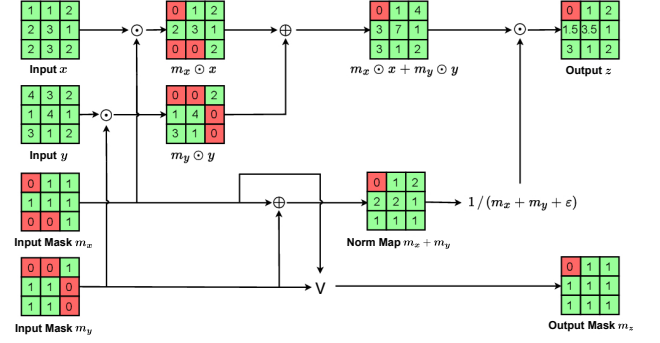


Fig. 4. Illustration of the proposed sparsity-invariant average.

then normalized at each location according to the upsampled sparsity mask values $F(m_x)$. The final sparsity mask m_z is obtained by identifying the non-zero locations of the upsampled sparsity mask $F(m_x)$.

Note that for sparsity-invariant max-pooling or down-sampling, it could be calculated the same as Eqs. (2) and (3) by replacing the upsampling function F with max-pooling or down-sampling operators.

3.2.2 Sparsity-invariant average

Pixelwise average of two feature maps of the same spatial sizes is needed for fusing features from different levels without increasing the output channels. For average of sparse input feature maps, however, specifically designed average operation is needed. We propose sparsity-invariant average, which takes two input sparse feature maps, x and y , with their corresponding sparsity masks, m_x and m_y , as inputs. It generates the fused sparse feature maps z with its corresponding sparsity mask m_z . Unlike the sparsity-invariant upsampling or downsampling, the key difference is that the average operation takes two sparse features as inputs.

We formulate the sparsity-invariant average as

$$z = \frac{m_x \odot x + m_y \odot y}{m_x + m_y + \epsilon}, \quad (4)$$

$$m_z = m_x \vee m_y, \quad (5)$$

where \vee denotes elementwise alternation (i.e., logical “or”) function, \odot represents elementwise multiplication, and ϵ is a very small number to avoid division by zero.

The sparsity-invariant average is illustrated in Fig. 4. The two input sparse features are first masked out by their corresponding masks to obtain $m_x \odot x$ and $m_y \odot y$. Both the masked features and the masks are then pixelwisely added. The added features are then normalized by added sparsity masks to calculate the output sparse features z . For the output sparsity mask, at each location, if the location is valid for either of the input feature maps, the mask is set to 1 for this location. At each location of the output feature map, the output feature vector is the mean of two input feature vectors if both input maps are valid at this location. If only one input feature is valid at a location, the valid single input feature vector is directly copied to the same location of the output feature maps. The two types of output feature vectors would have similar magnitude

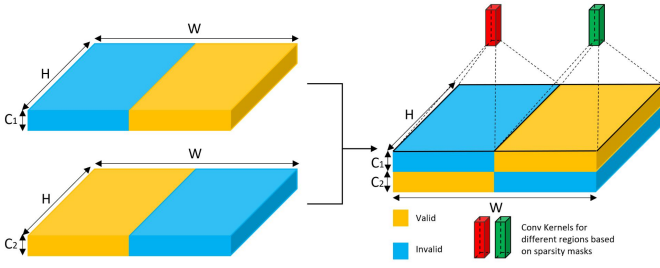


Fig. 5. Sparsity patterns vary from regions, thus we need several different kernels to deal with out feature maps after concatenation.

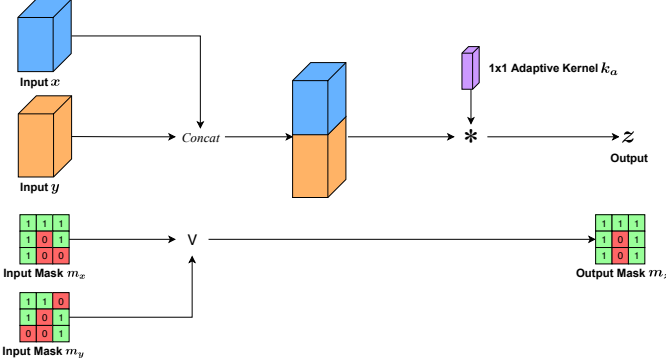


Fig. 6. Illustration of the proposed joint sparsity-invariant concatenation and convolution.

because the added features are properly normalized by the added sparsity mask.

Note feature addition was also explored in [1]. However, an output sparsity mask was not generated. Without such a mask, following convolution operations cannot handle sparse output features. Therefore, their operation could only be used as the last layer of a neural network.

3.2.3 Joint sparsity-invariant concatenation and convolution

Another commonly used approach of fusing two feature maps with the same spatial size is feature concatenation, i.e., concatenation along the channel dimension. However, different from aforementioned other operations, concatenation would introduce sparsity in both the spatial dimension ($H \times W$) and the feature dimension (C), and the latter actually prevent us from simply extending the idea of sparsity-invariant average in the previous subsection. Let's consider the scenario where two feature maps with shape $C_1 \times H \times W$ and $C_2 \times H \times W$ are now being concatenated into one with shape $(C_1 + C_2) \times H \times W$ as illustrated in Fig. 5. The concatenation is further followed by a convolution layer (1×1 for simplicity), which is a common operation in CNNs. However, we know that convolution performs filtering on a location by extracting one local feature vector with length C and summing all entries up into a number with learnable weights. Then, the convolution kernels iterate over the whole feature map, treating every location equally.

Recall the situation in sparsity-invariant convolution where the feature vector for a certain location only have two possible sparsity patterns—the whole vector of length

C is valid, or all the entries of this vector are zeros. Note that the latter situation would not influence the training as its contribution to the output as well as the gradient of kernels is zero. Therefore, it is enough for us to use one set of convolution kernels, equally for every valid location.

However, when we are convolving on the feature maps after concatenation, we have four different types of vectors or sparsity patterns for each location: the first C_1 feature channels of the vector is valid while the latter C_2 feature channels are not; or C_2 is valid while C_1 is not, or both of them are valid/invalid. Therefore, we need three different sets of kernels to tackle these four different sparsity patterns. In other words, to effectively handle different scenarios at different locations of the concatenated feature maps, we propose to use an adaptive-kernel version convolution to solve the difficulty and combine it with concatenation together.

Another advantage of combining them is that all convolution kernels would generate outputs with the same spatial sparsity patterns. Therefore the output mask is still of single channel, which is computationally efficient and reduces the model complexity significantly. Specifically, our joint sparsity-invariant concatenation and convolution is described and explained formally as following:

Given the two input sparse feature maps x and y with their sparsity masks m_x and m_y , the proposed joint concatenation and convolution operation is formulated as

$$z = [x; y] * k_a, \quad (6)$$

$$m_z = m_x \vee m_y, \quad (7)$$

where $[;]$ denotes the concatenation of two feature maps along the channel dimension, and $*$ denotes the conventional convolution operation. Note that the output sparsity mask is calculated exactly the same as that in sparsity-invariant average. The key of the proposed operation is a 1×1 convolution with an adaptive convolution kernel k_a that handles three different scenarios of concatenating sparse feature maps, which is formulated as

$$k_a(u, v) = \begin{cases} k_a^{(1)} & m_x(u, v) = 1, m_y(u, v) = 0; \\ k_a^{(2)} & m_x(u, v) = 0, m_y(u, v) = 1; \\ k_a^{(3)} & m_x(u, v) = 1, m_y(u, v) = 1, \end{cases} \quad (8)$$

where $k_a(u, v)$ are the 1×1 adaptive convolution kernel at location (u, v) of the concatenated feature maps $[x; y]$. $k_a^{(1)}$, $k_a^{(2)}$, $k_a^{(3)}$ are the three sets of learnable convolution weights for the three different feature concatenation scenarios: at each location (u, v) , either both input feature vectors are valid (i.e., $m_x(u, v) = 1$ and $m_y(u, v) = 1$), or only one of the input feature vectors is valid (i.e., either $m_x(u, v) = 1$ or $m_y(u, v) = 1$). The key reason for using different sets of kernel weights instead of the same set of convolution weights, as we briefly introduced before, is to avoid involving invalid input features in the concatenated feature maps into feature learning process. Illustrating with our notation above, if the current 1×1 convolution kernel is on the location (u, v) and find that the first mask here $m_x(u, v)$ is one and the second mask $m_y(u, v)$ is zero, it would choose the first set of kernel weights $k_a^{(1)}$ to use during forward pass. So is the back propagation. In this case, we know the second chunk of the feature vector, of which the length is fixed, is always zero. And because it's consistently processed by the first

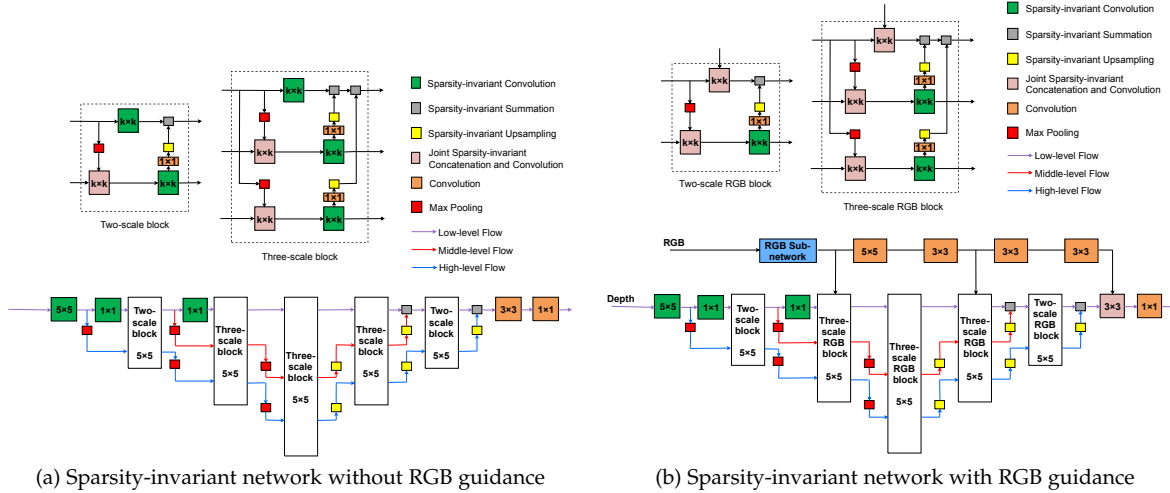


Fig. 7. Illustration of our multi-scale encoder-decoder network structure for depth completion based on the proposed sparsity-invariant operations. (a) Proposed network with RGB guidance. (b) Proposed network without RGB guidance.

kernel, this kernel would naturally learn how to adapt to this pattern.

In other words, by adopting the proposed adaptive convolution kernel k_a , the three sets of kernel weights $k_a^{(1)}$, $k_a^{(2)}$, $k_a^{(3)}$ are able to handle different sparse feature concatenation scenarios. With joint training, the different kernels are learned to best adapt each other to generate appropriate feature representations for further processing. In this way, the sparse feature maps could be effectively fused with proposed concatenation.

3.3 Hierarchical Multi-scale Network (HMS-Net) for depth completion

Multi-scale encoder-decoder neural networks for dense inputs are widely investigated for pixelwise prediction tasks. Those networks have the advantages of fusing both low-level and high-level features for accurate pixel prediction. However, with only the sparsity-invariant convolution in [1], encoder-decoder networks cannot be converted to handle sparse inputs. On the other hand, for frequently studied pixelwise prediction tasks, such as semantic segmentation, global high-level information usually shows greater importance to the final performance, while the full-resolution low-level features are less informative and generally go through fewer non-linearity layers compared with high-level features. However, we argue that depth completion is a low-level vision task. The low-level features in this task should be non-linearly transformed and fused with mid-level and high-level features for more times to achieve satisfactory depth completion accuracy.

Based on this motivation, we propose the Hierarchical Multi-scale encoder-decoder Network (HMS-Net) with our proposed sparsity-invariant operations for sparse depth completion. The network structure without RGB information is illustrated in Fig. 7(a). We propose two basic building blocks, a two-scale block and a three-scale block, consisting proposed sparsity-invariant operations. The two-scale block has an upper path that non-linearly transforms the full-resolution low-level features by a $k \times k$ sparsity-invariant convolution. The lower path takes downsampled low-level

features as inputs for learning higher-level features with another $k \times k$ convolution. We empirically set $k=5$ in all our experiments according to our hyperparameter study. The resulting higher-level features are then upsampled and added back to the full-resolution low-level features. Compared with the two-scale block, the three-scale block fuses features from two higher levels into the upper low-level feature path to utilize more auxiliary global information. In this way, the full-resolution low-level features are effectively fused with higher-level information and are non-linearly transformed multiple times to learn more complex prediction functions. All feature maps in our network are of 16 channels regardless of scales.

Our final network utilizes a 5×5 sparsity-invariant convolution at the first layer. The resulting features then go through three of the proposed multi-scale blocks followed by sparsity-invariant max-pooling, and are then upsampled three times to generate the full-resolution feature maps. The final feature maps are then transformed by one 1×1 convolution layers to generate the final per-pixel prediction. The output depth predictions are supervised by the Mean-square Error (MSE) loss function with ground-truth annotations.

Another way to understand the network structure we proposed is that our structure could be considered as a backbone encoder-decoder CNN shown in Fig. 1(b) but with special design for the depth completion task. By looking at the details in Fig. 7(a) and compare it with commonly used encoder-decoder networks [8], [9], [38], it's easy to find that there are several key shortcuts (importance demonstrated later in the ablation study) between different flows uniquely in our network for a better fusion as we described before. Except for this aspect, the low-level features in the mentioned encoder-decoder networks go through very few non-linear transformations, while our proposed network emphasizes much more on the low-level features. Furthermore, the total depth of our network is also much shallower. Compared with Full-Resolution Residual Network [10] which also has multiple shortcuts, the latter's full-resolution low-level features only serve as residual signals. In addition, it does not consider the fusion of multiple-scale features at

the same time as our three-scale block does. We further compared the proposed network with the commonly used encoder-decoder network structures in our experimental studies.

3.4 RGB-guided multi-scale depth completion

The LIDAR sensors are usually paired with RGB cameras to obtain the aligned sparse depth maps and RGB images. RGB images could therefore act as auxiliary guidance for depth completion.

To integrate RGB features into our multi-scale encoder-decoder network, we added an RGB feature path to our proposed network. The network structure is illustrated in Fig. 7(b). The input image is first processed by an RGB sub-network to obtain mid-level RGB features. The structure of the sub-network follows the first six blocks of the ERFNet [40]. It consists of two downsampling blocks and four residual blocks. The downsampling block has a 2×2 convolution layer with stride 2 and a 2×2 max-pooling layer. The input features are input into the two layers simultaneously, and their results are concatenated along the channel dimension to obtain the $1/2$ size feature maps. The main path of the residual block has two sets of 1×3 conv \rightarrow BN \rightarrow ReLU $\rightarrow 3 \times 1$ conv \rightarrow BN \rightarrow ReLU. Because the obtained mid-level RGB features are downsampled to $1/4$ of its original size, they are upsampled back to the input image’s original size. The upsampled RGB features are then transformed by a series of convolutions. They act as additional guidance signals and are concatenated to the low-level sparse depth feature maps of different multi-scale blocks. Our experimental results show including the additional RGB mid-level features as guidance further improves the depth completion accuracy.

3.5 Training scheme

We adopt the mean squared error (MSE) loss function to train our proposed encoder-decoder networks. Since some datasets could only provide sparse ground-truth depth maps, the loss function is only evaluated at locations with ground-truth annotations, which could be formulated as

$$L(x, y) = \frac{1}{|\mathbf{V}|} \sum_{u,v \in \mathbf{V}} |o(u, v) - t(u, v)|^2 \quad (9)$$

where \mathbf{V} is the set containing coordinates with ground-truth depth values, $|\mathbf{V}|$ calculates the total number of valid points in \mathbf{V} , and o and t are the predicted and ground-truth depth maps.

For network training, all network parameters except those of the RGB sub-network are randomly initialized. We adopt the ADAM optimizer [41] with an initial learning rate of 0.01. The network is trained for 50 epochs. To gradually decrease the learning rate, it is decayed according to the following equation,

$$\text{learning rate} = 0.01 \times \left(1 - \frac{\text{iter_epoch}}{50}\right)^{0.9}, \quad (10)$$

where iter_epoch denotes the current epoch iteration.

Also note that sparsity masks are generated for every input directly depending on the network structure, without

TABLE 1
Depth completion errors by different methods on the test set of KITTI depth completion benchmark.

Methods	RGB info?	RMSE	MAE	iRMSE	iMAE
SparseConvs [1]	×	1601.33	481.27	4.94	1.78
IP-Basic [5]	×	1288.46	302.60	3.78	1.29
NConv-CNN [12]	×	1268.22	360.28	4.67	1.52
Spade-sD [18]	×	1035.29	248.32	2.60	0.98
Sparse-to-Dense(d) [23]	×	954.36	288.64	3.21	1.35
Ours w/o RGB	×	937.48	258.48	2.93	1.14
<hr/>					
Bilateral NN [7]	✓	1750.00	520.00	-	-
ADNN [4]	✓	1325.37	439.48	59.39	3.19
CSPN [19]	✓	1019.64	279.46	2.93	1.15
Spade-RGBsD [18]	✓	917.64	234.81	2.17	0.95
Sparse-to-Dense(gd) [23]	✓	814.73	249.95	2.80	1.21
Ours w/ RGB	✓	841.78	253.47	2.73	1.13

any learnable parameters. The values of input masks are set to 1 for all valid spatial locations and 0 for invalid locations. The mask in one layer is propagated to the following layer. In other words, they purely depend on the network structure and the current input. During training, they filter out both invalid feature points and the gradient for invalid spatial locations.

For our RGB sub-network, we use the first six blocks of the ERFNet [40]. Its initial network parameters are copied from the network pretrained on the CityScapes dataset [42]. Both paths are then end-to-end trained until convergence.

4 EXPERIMENTS

We conduct experiments on the KITTI depth completion dataset [1] and NYU-depth-v2 dataset [11] for evaluating the performance of our proposed approach.

4.1 KITTI depth completion benchmark

4.1.1 Data and evaluation metrics

We first evaluate our proposed approach on the KITTI depth completion benchmark [1]. Following the experimental setup in [1], 85,898 depth maps are used for training, 1,000 for validation and 1,000 for test. The LIDAR depth maps are aligned with RGB images by projecting the depth map into the image coordinates according to the two sensors’ essential matrix. The input depth maps generally contains $< 10\%$ sparse points with depth values and the top $1/3$ of the input maps do not contain any depth measurements. One example is shown in Fig. 2(a).

According to the benchmark, all algorithms are evaluated according to the following metrics, root mean square error (RMSE in mm), mean absolute error (MAE in mm), root mean squared error of the inverse depth (iRMSE in $1/\text{km}$), and mean absolute error of the inverse depth (iMAE

in 1/km), i.e.,

$$\text{RMSE} = \left(\frac{1}{|\mathcal{V}|} \sum_{u,v \in \mathcal{V}} |\mathbf{o}(u,v) - \mathbf{t}(u,v)|^2 \right)^{0.5}, \quad (11)$$

$$\text{MAE} = \frac{1}{|\mathcal{V}|} \sum_{u,v \in \mathcal{V}} |\mathbf{o}(u,v) - \mathbf{t}(u,v)|, \quad (12)$$

$$\text{iRMSE} = \left(\frac{1}{|\mathcal{V}|} \sum_{u,v \in \mathcal{V}} \left| \frac{1}{\mathbf{o}(u,v)} - \frac{1}{\mathbf{t}(u,v)} \right|^2 \right)^{0.5}, \quad (13)$$

$$\text{iMAE} = \frac{1}{|\mathcal{V}|} \sum_{u,v \in \mathcal{V}} \left| \frac{1}{\mathbf{o}(u,v)} - \frac{1}{\mathbf{t}(u,v)} \right|, \quad (14)$$

where \mathbf{o} and \mathbf{t} represent the output of our approach and ground-truth depth values.

For RMSE and MAE, RMSE is more sensitive to large errors compared. This is because even a small number of large errors would be magnified by the square operation and dominate the overall loss value. RMSE is therefore chosen as the main metric for ranking different algorithms in the KITTI leaderboard. Since large depth values usually have greater errors and might dominate the calculation of RMSE and MAE, iRMSE and iMAE are also evaluated, which calculate the mean of inverse of depth errors. In this way, large depth values' errors would have much lower weights on the two metrics. The two metrics focus more on depth points near the LIDAR sensors.

4.1.2 Comparison with state-of-the-arts

The performance of our proposed approaches and state-of-the-art depth completion methods are recorded in Table 1.

SparseConvs represents the 6-layer convolution neural network with only sparsity-invariant convolution proposed in [1]. It only supports convolution and max-pooling operations and therefore loses much high-resolution information. IP-Basic represents the method in [5], a well-designed algorithm hand-crafted rules based on several traditional image processing algorithms. NConv-CNN [12] proposed a constrained convolution layer and propagating confidence across layers for depth completion. Bilateral NN [7] uses RGB images as guidance and integrate bilateral filters into deep neural networks. It was modified to handle sparse depth completion following [1]. Spade-sD and Spade-RGBsD [18] do not have special treatment for sparse data. They utilize conventional dense CNN but adopt different loss function and training strategy. CSPN [19] iteratively learns inter-pixel affinities with RGB guidance via recurrent convolution operations. The affinities could then be used to spatially propagate depth values between different locations. Sparse-to-dense(d) and Sparse-to-dense(gd) [23] explore additional temporal information from sequential data to apply additional supervisions based on the photometric loss between neighboring frames.

For methods without RGB guidance, our proposed network without RGB guidance outperforms all other peer-reviewed methods in terms of RMSE (the main ranking metric in KITTI leaderboard). Spade-sD has better MAE, iRMSE and iMAE, which mean that this method performs better on nearby objects but is more likely to generate large errors than our proposed method. Note that we utilize the

TABLE 2
Component analysis of our proposed method on the validation set of KITTI depth completion benchmark.

Method	RMSE	MAE
Baseline w/o sparseconv	1819.81	426.84
Baseline w/ sparseconv	1683.22	447.93
Baseline + MS (Up only)	1185.02	323.41
Baseline + MS (Down only)	1192.43	322.95
Baseline + MS (Mid-level flow removed)	1166.87	317.74
Baseline + MS (Full)	1137.42	315.32
Baseline + MS + SO	994.14	262.41
Baseline + MS + SO + RGB	883.74	257.11

$L2$ loss function to deliberately minimize RMSE. If other metrics are considered to be more important, different loss functions could be adopted for training.

For methods with RGB guidance, our method ranks 2nd behind Sparse-to-dense(gd) [23] in terms of RMSE. However, Sparse-to-dense(gd) utilized additional supervisions from temporal information, while our proposed method only uses supervisions from individual frames.

4.1.3 Ablation study

We investigate individual components in our framework to see whether they contribute to the final performance. The investigated components include, multi-scale structure, sparsity-invariant operations, and RGB guidance. We choose our full-resolution low-level feature path without the mid-level or high-level flow in our network (i.e., the upper path in Fig. 7(a)) as the baseline model for this section. The baseline model does not include RGB guidance. The analysis results on KITTI validation set are shown in Table 2. And the baseline model without multi-scale feature fusion generates large depth estimation errors.

Multi-scale structure. Using our multi-scale encoder-decoder structure (denoted as *Baseline+MS (Full)*) in addition to the baseline model enables fusing low-level and high-level information from different scales. The multi-scale structure provides much larger receptive fields for neurons in the last convolution layer. Therefore, even if some regions in the input depth map are very sparse, the model could still predict depth values for every grid points. Using multi-scale features generally results in clearer boundaries and shows higher robustness to noises. *Baseline+MS (Full)* significantly decreases RMSE from 1819.81 to 1137.42. Also, our fusing skip connections connect different scales also enable a better fusion for the information in our network. By removing either the up fusing skip connection (shown by up arrows within the blocks in Fig. 6(a)) or the down fusing skip connection, our network performs worse as shown in Table 2 (Up only and Down only entries). Also, the mid-level flow is making the performance better (Mid-level flow removed entry). An example showing the differences between *Baseline+MS (Full)* and *Baseline* is in Fig. 8.

Sparsity-invariant operations. The *Baseline+MS* utilizes dense convolution. It has difficulty in handling sparse input data and sparse feature maps, especially for regions where there are very sparse points. The *Baseline+MS+SO* uses our proposed sparsity-invariant operations to maintain a correct mask flow and then converts conventional operations in

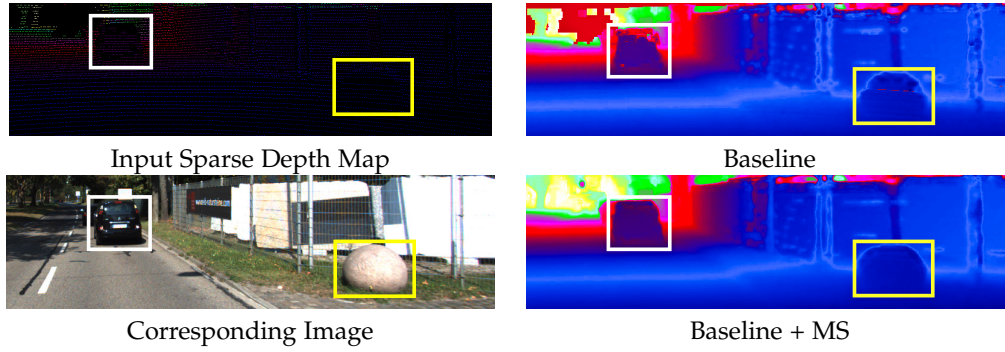


Fig. 8. An example in the KITTI validation set to show the results by the baseline model (*Baseline*) vs. baseline model with multi-scale encoder-decoder structure (*Baseline + MS*). See rectangles for better boundary regions by *Baseline + MS*.

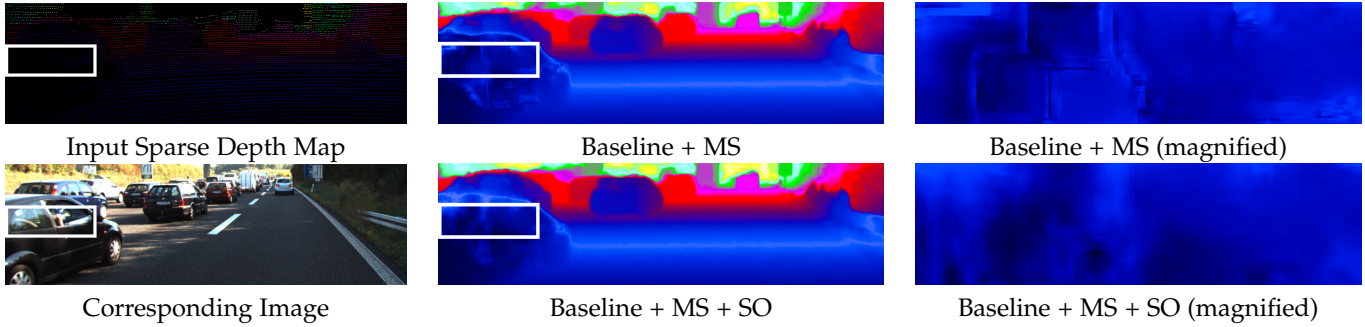


Fig. 9. An example in the KITTI validation set to show the effectiveness with (*Baseline + MS*) and without proposed sparsity-invariant operations (*Baseline + MS + SO*). Rectangles show *Baseline + MS + SO* better handles input depth maps with very sparse depth points.

TABLE 3
Comparison with commonly used encoder-decoder networks without RGB guidance on KITTI depth completion validation set.

Method	RMSE	MAE
U-Net [8]	1387.35	445.73
FRRN [10]	1148.27	338.56
PSPNet [43]	1185.39	354.21
FPN [44]	1441.82	473.65
He et al. [39]	1056.39	293.86
Ours w/o RGB	994.14	262.41

Baseline+MS into sparsity-invariant ones to handle very sparse inputs. RMSE by *Baseline+MS+SO* improves from 1137.42 to 994.14. An example is shown in Fig. 9, which shows *Baseline+MS+SO* better handles regions with very sparse inputs.

Deep fusion of RGB features. By incorporating RGB features into *Baseline+MS+SO+RGB*, the network utilizes useful additional guidance from RGB images for further improving the depth completion accuracy. An example comparing *Baseline+MS+SO+RGB* and *Baseline+MS+SO* is shown in Fig. 10, where the resulting depth maps with RGB guidance are much sharper at image boundaries.

4.1.4 Comparison with other encoder-decoder structures

To evaluate the effectiveness of our proposed HMS-Net structure, we conduct experiments to test other commonly used encoder-decoder network structures. We modify those structures with the sparsity-invariant convolution and our proposed sparsity-invariant operations to handle sparse inputs. The experimental results on the KITTI validation set

are reported in Table 3. We compared our network structure without RGB guidance with modified U-Net [8], FRRN [10], PSPNet [43], FPN [44] and He et al. [39] without focal length. Our proposed network structure achieves the lowest errors in terms of RMSE and MAE.

4.2 Robustness testing on KITTI benchmark

4.2.1 Robustness to depth noises

Since the sparse depth maps are obtained by LIDAR scans, inevitably, there would be noises in acquired depth values. As a result, the robustness of depth completion algorithms with regard to different noise levels is important in practice. We conduct experiments to test the robustness of our proposed model without RGB guidance and compare with SparseConvs [1] and IP-Basic [5]. Note that all models in this section are trained on original data and directly tested on noisy data.

Scene-level Gaussian noises. For this experiment, we add Gaussian noises on randomly selected 10% depth points among all points. Once the 10% points are selected, for a specific noise standard deviation level from 5-50 meters, we sample additive noise values from the zero-mean Gaussian distribution. The negative additive noises could simulate occlusion from raindrops, snowflakes or fog, while the positive additive noises mimic the laser going through glasses that mistakenly measures objects behind glasses. Noisy points whose depth values are smaller than 1 meter are set to 1 meter to simulate the minimum range of the LIDAR sensor. The RMSEs by three methods on the KITTI validation set with Gaussian noises are shown in Fig. 11(a).

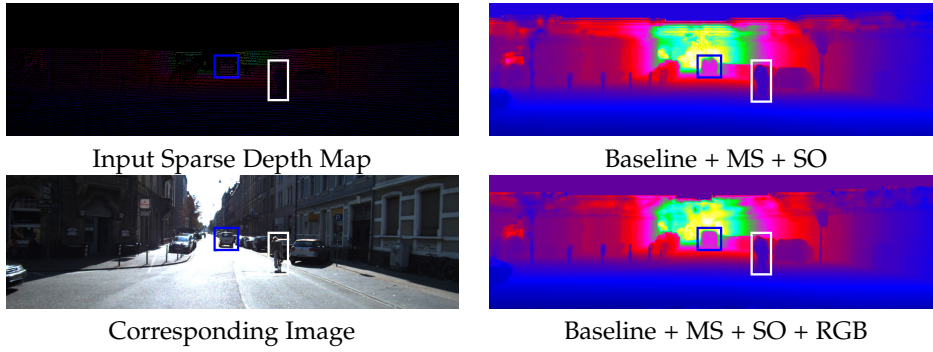


Fig. 10. An example in the KITTI validation set to show the effectiveness without RGB guidance (*Baseline + MS + SO*) and with RGB guidance (*Baseline + MS + SO + RGB*). Rectangles show *Baseline + MS + SO + RGB* generates clearer boundaries of the cyclist and the far-away vehicle.

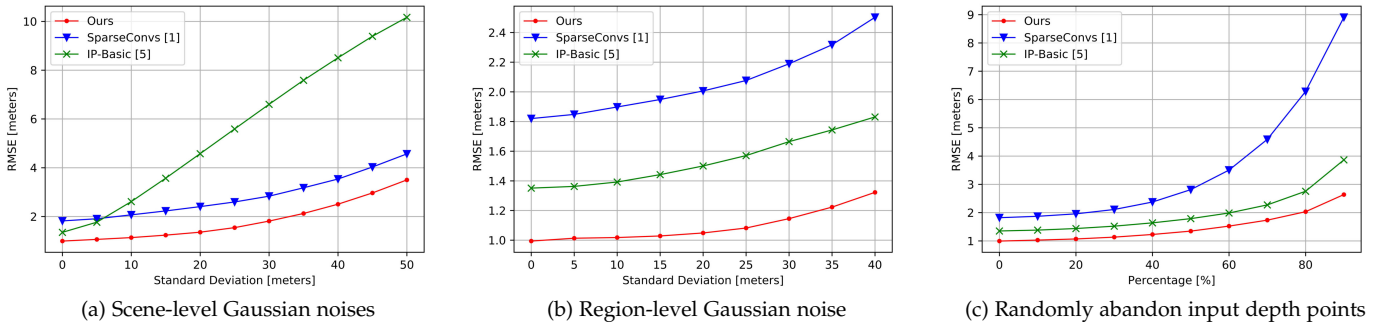


Fig. 11. Testing robustness on the validation set of KITTI depth completion benchmark with manually corrupted data. (a) Scene-level Gaussian noises on 10% depth points. (b) Region-level Gaussian noise on eight 25×25 regions. (c) Randomly abandon input depth points with varying probabilities.

TABLE 4
Comparison with other methods on NYU-depth-v2 depth dataset with varying N values.

Method	20 samples		50 samples		200 samples	
	RMSE	REL	RMSE	REL	RMSE	REL
Ma et al. [3] w/o RGB	0.461	0.110	0.347	0.076	0.259	0.054
Jaritz et al. [18] w/o RGB	0.476	0.114	0.358	0.074	0.246	0.051
Ma et al. [23] w/o RGB	0.481	0.113	0.352	0.073	0.245	0.049
He et al. [39] w/o RGB	0.478	0.113	0.355	0.072	0.238	0.045
Ours w/o RGB	0.449	0.110	0.344	0.073	0.233	0.044
Ma et al. [3] w/ RGB	0.351	0.078	0.281	0.059	0.230	0.044
Ours w/ RGB	0.350	<u>0.079</u>	0.274	<u>0.059</u>	0.212	0.041

Our method outperforms both SparseConvs [1] and IP-Basic [5] on different noisy depth values.

Region-level Gaussian noises. We randomly select eight regions of size 25×25 pixels in every input depth map. In each region, 50% of depth points are randomly selected to add Gaussian noises of zero mean and different standard deviation values. Noisy points whose depth values are smaller than 1 meter are set to 1 meter to simulate the minimum range of the LIDAR sensor. The region-level noises are used to simulate the cases where large glasses or mirrors exist. Those regions would reflect most laser and leave large holes in the obtained depth map. The RMSE by different methods on the KITTI validation set are shown in Fig. 11(b). Our method again outperforms the two compared methods, because of its capability of fusing low-level and high-level information with the proposed multi-scale encoder-decoder structure.

4.2.2 Robustness to sparsity

Robustness to sparsity is also essential to depth completion algorithms. We conduct experiments on testing different levels of sparsity of the input depth maps. For each input depth map, we randomly abandon 10%-90% of valid depth points. Note again that all methods are trained on original training data and are not finetuned to adapt the sparser inputs. The results on the KITTI validation set by the our model without RGB guidance and two compared methods are shown in Fig. 11(c). Our proposed method shows the highest tolerance against different sparsity levels.

4.3 NYU-depth-v2 dataset

4.3.1 Data, experimental setup, and evaluation metrics

We also evaluate our proposed method on the NYU-depth-v2 dataset [11] with its official train/test split. Each RGB image in the dataset is paired with a spatially aligned dense

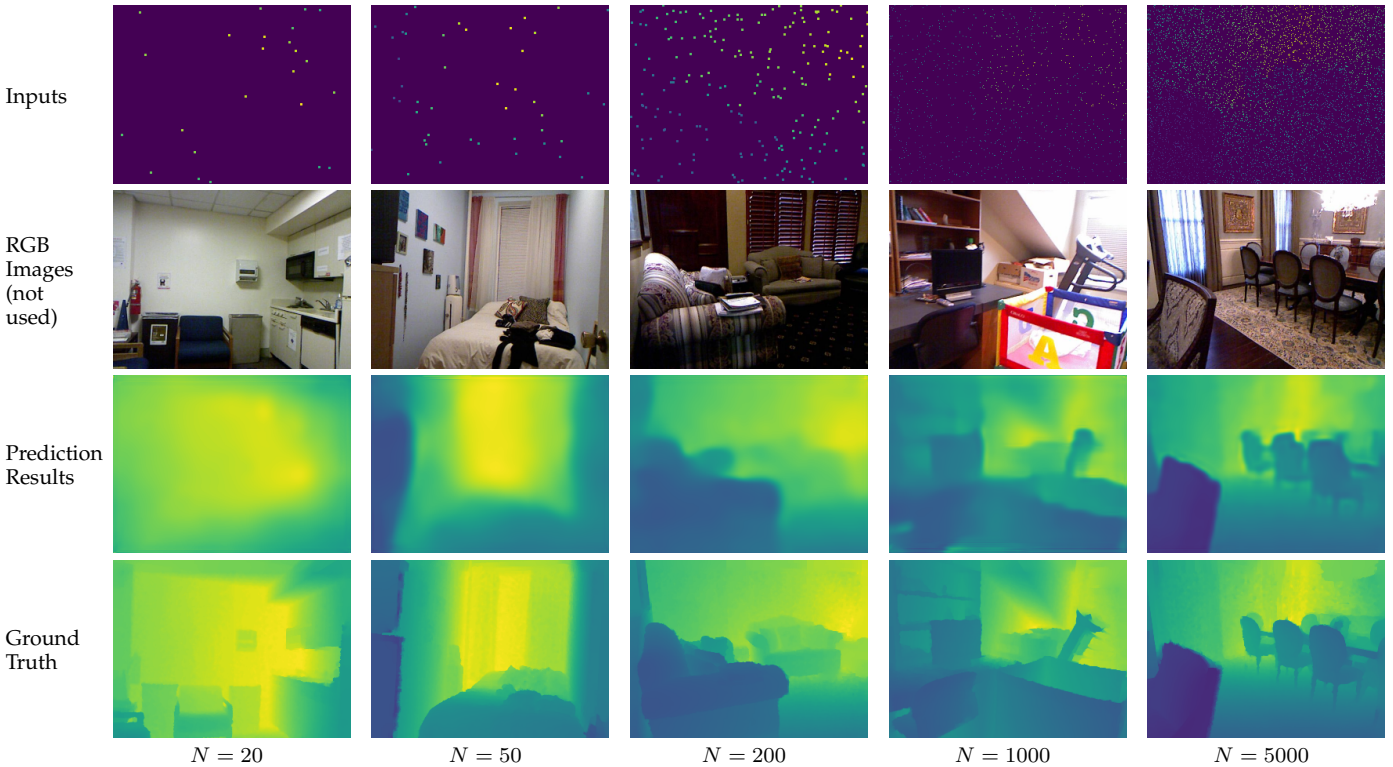


Fig. 12. Depth completion examples in NYU-depth-v2 dataset [11] by our proposed method with vary N values. (First row) Input sparse depth maps. (Second row) Corresponding RGB images (not used as algorithm inputs). (Third row) Predicted dense depth maps by our proposed method. (Fourth row) Ground truth dense depth maps.

depth map. The original depth maps are dense and the dataset is not originally proposed for sparse depth completion. Following the experimental setup in [3], synthetic sparse depth maps could be created to test the performance of depth completion algorithms. Only N depth points in each depth map are randomly kept as input depth maps for depth completion. The training set consists of depth and image pairs from 249 scenes, and 654 images are selected for evaluating the final performance according to the setup in [3], [45], [46]. RMSE (Eq. (11)) and mean absolute relative error (REL in meters) are adopted as the evaluation metrics. REL is calculated as

$$REL = \frac{1}{|\mathbf{V}|} \sum_{u,v \in \mathbf{V}} \left| \frac{o(u,v) - t(u,v)}{t(u,v)} \right|, \quad (15)$$

where o and t are the outputs of our network and the ground truth dense depth maps.

4.3.2 Comparison with state-of-the-art

We compare our method to methods proposed in Ma et al. [3], [23]¹, Jaritz et al. [18] and He et al. [39]. Since the input depth maps are much sparser than the depth maps in KITTI dataset [1], we added a 2×2 max-pooling layer following the first 5×5 convolution layer, and added a Batch Normalization [47] layer after each convolution layer. Each input depth map has $N = 20, 50, 200$ randomly kept depth points. RMSE and REL of different N values are reported in Table 4, which demonstrates that our proposed method outperforms all other methods without using RGB

information. We also show the result with RGB guidance in TABLE 4 as well as examples of different N values and the depth completion results in Fig. 12.

5 CONCLUSIONS

In this paper, we proposed several novel sparsity-invariant operations for handling sparse feature maps. The novel operations enable us to design a novel sparsity-invariant encoder-decoder network, which effectively fuses multi-scale features from different CNN layers for accurate depth completion. RGB features for better guiding depth completion is also integrated into the proposed framework. Extensive experiment results and component analysis show advantages of the proposed sparsity-invariant operations and the encoder-decoder network structure. The proposed method outperforms state-of-the-arts and demonstrate great robustness against different levels of data corruption.

REFERENCES

- [1] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox, and A. Geiger, "Sparsity invariant CNNs," *International Conference on 3D Vision*, 2017.
- [2] Y. Liao, L. Huang, Y. Wang, S. Kodagoda, Y. Yu, and Y. Liu, "Parse geometry from a line: Monocular depth estimation with partial laser observation," *IEEE International Conference on Robotics and Automation*, pp. 5059–5066, 2017.
- [3] F. Ma and S. Karaman, "Sparse-to-dense: Depth prediction from sparse depth samples and a single image," *IEEE International Conference on Robotics and Automation*, 2018.
- [4] N. Chodosh, C. Wang, and S. Lucey, "Deep convolutional compressed sensing for lidar depth completion," *European Conference on Computer Vision*, 2018.

1. The code released by the authors were utilized.

- [5] J. Ku, A. Harakeh, and S. L. Waslander, "In defense of classical image processing: Fast depth completion on the cpu," *arXiv preprint arXiv:1802.00036*, 2018.
- [6] E. A. Nadaraya, "On estimating regression," *Theory of Probability & Its Applications*, vol. 9, no. 1, pp. 141–142, 1964.
- [7] J. T. Barron and B. Poole, "The fast bilateral solver," *European Conference on Computer Vision*, pp. 617–632, 2016.
- [8] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," *International Conference on Medical image computing and computer-assisted intervention*, pp. 234–241, 2015.
- [9] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [10] T. Pohlen, A. Hermans, M. Mathias, and B. Leibe, "Full-resolution residual networks for semantic segmentation in street scenes," *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4151–4160, 2017.
- [11] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from rgb-d images," *European Conference on Computer Vision*, pp. 746–760, 2012.
- [12] A. Eldesokey, M. Felsberg, and F. S. Khan, "Propagating confidences through cnns for sparse data regression," *British Machine Vision Conference (BMVC)*, 2018.
- [13] M. Ren, A. Pokrovsky, B. Yang, and R. Urtasun, "Sbnet: Sparse blocks network for fast inference," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8711–8720.
- [14] N. Schneider, L. Schneider, P. Pinggera, U. Franke, M. Pollefeys, and C. Stiller, "Semantically guided depth upsampling," *German Conference on Pattern Recognition*, pp. 37–48, 2016.
- [15] W. Van Gansbeke, D. Neven, B. De Brabandere, and L. Van Gool, "Sparse and noisy lidar completion with rgb guidance and uncertainty," *arXiv preprint arXiv:1902.05356*, 2019.
- [16] Y. Zhang and T. Funkhouser, "Deep depth completion of a single rgb-d image," *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 175–185, 2018.
- [17] J. Qiu, Z. Cui, Y. Zhang, X. Zhang, S. Liu, B. Zeng, and M. 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*, 2019, pp. 3313–3322.
- [18] M. Jaritz, R. de Charette, E. Wirbel, X. Perrotton, and F. Nashashibi, "Sparse and dense data with cnns: Depth completion and semantic segmentation," *International Conference on 3D Vision*, 2018.
- [19] X. Cheng, P. Wang, and R. Yang, "Depth estimation via affinity learned with convolutional spatial propagation network," *European Conference on Computer Vision*, 2018.
- [20] A. Eldesokey, M. Felsberg, and F. S. Khan, "Confidence propagation through cnns for guided sparse depth regression," *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [21] Y. Yang, A. Wong, and S. Soatto, "Dense depth posterior (ddp) from single image and sparse range," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3353–3362.
- [22] Y. Yang and S. Soatto, "Conditional prior networks for optical flow," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 271–287.
- [23] F. Ma, G. V. Cavalheiro, and S. Karaman, "Self-supervised sparse-to-dense: Self-supervised depth completion from lidar and monocular camera," *arXiv:1807.00275*, 2018.
- [24] S. Matyunin, D. Vatolin, Y. Berdnikov, and M. Smirnov, "Temporal filtering for depth maps generated by kinect depth camera," *3DTV Conference: The True Vision-Capture, Transmission and Display of 3D Video*, 2011.
- [25] M. Camplani and L. Salgado, "Efficient spatio-temporal hole filling strategy for kinect depth maps," *Three-dimensional image processing and applications*, 2012.
- [26] J. Yang, X. Ye, and P. Frossard, "Global auto-regressive depth recovery via iterative non-local filtering," *IEEE Transactions on Broadcasting*, 2018. [Online]. Available: <http://infoscience.epfl.ch/record/253660>
- [27] L. Chen, H. Lin, and S. Li, "Depth image enhancement for kinect using region growing and bilateral filter," *International Conference on Pattern Recognition*, pp. 3070–3073, 2012.
- [28] J. Yang, X. Ye, K. Li, C. Hou, and Y. Wang, "Color-guided depth recovery from rgb-d data using an adaptive autoregressive model," *IEEE Transactions on Image Processing*, vol. 23, no. 8, pp. 3443–3458, 2014.
- [29] Y. Li, J.-B. Huang, N. Ahuja, and M.-H. Yang, "Deep joint image filtering," *European Conference on Computer Vision*, pp. 154–169, 2016.
- [30] D. Ferstl, C. Reinbacher, R. Ranftl, M. Rütther, and H. Bischof, "Image guided depth upsampling using anisotropic total generalized variation," *International Conference on Computer Vision*, pp. 993–1000, 2013.
- [31] T.-W. Hui, C. C. Loy, and X. Tang, "Depth map super-resolution by deep multi-scale guidance," *European Conference on Computer Vision*, pp. 353–369, 2016.
- [32] J. Xie, R. S. Feris, and M.-T. Sun, "Edge-guided single depth image super resolution," *IEEE Transactions on Image Processing*, vol. 25, no. 1, pp. 428–438, 2016.
- [33] C. Guo, C. Li, J. Guo, R. Cong, H. Fu, and P. Han, "Hierarchical features driven residual learning for depth map super-resolution," *IEEE Transactions on Image Processing*, vol. 28, no. 5, pp. 2545–2557, 2018.
- [34] O. Mac Aodha, N. D. Campbell, A. Nair, and G. J. Brostow, "Patch based synthesis for single depth image super-resolution," *European Conference on Computer Vision*, pp. 71–84, 2012.
- [35] D. Ferstl, M. Ruther, and H. Bischof, "Variational depth super-resolution using example-based edge representations," *IEEE International Conference on Computer Vision*, pp. 513–521, 2015.
- [36] J. Xie, R. S. Feris, S.-S. Yu, and M.-T. Sun, "Joint super resolution and denoising from a single depth image," *IEEE Transactions on Multimedia*, vol. 17, no. 9, pp. 1525–1537, 2015.
- [37] G. Riegler, M. Rütther, and H. Bischof, "Atgv-net: Accurate depth super-resolution," 2016, pp. 268–284.
- [38] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," *European Conference on Computer Vision*, 2016.
- [39] L. He, G. Wang, and Z. Hu, "Learning depth from single images with deep neural network embedding focal length," *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4676–4689, 2018.
- [40] E. Romera, J. M. Alvarez, L. M. Bergasa, and R. Arroyo, "Efficient convnet for real-time semantic segmentation," *IEEE Intelligent Vehicles Symposium*, 2016.
- [41] D. Kinga and J. B. Adam, "A method for stochastic optimization," *International Conference on Learning Representations*, 2015.
- [42] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [43] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2881–2890, 2017.
- [44] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," *IEEE Conference on Computer Vision and Pattern Recognition*, vol. 1, no. 2, p. 4, 2017.
- [45] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab, "Deeper depth prediction with fully convolutional residual networks," *International Conference on 3D Vision*, pp. 239–248, 2016.
- [46] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," *Advances in neural information processing systems*, pp. 2366–2374, 2014.
- [47] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *International Conference on Machine Learning*, 2015.