

UAMD-Net: A Unified Adaptive Multimodal Neural Network for Dense Depth Completion

Guancheng Chen, Junli Lin, Huabiao Qin*

South China University of Technology, Guangzhou, China
 eechengc@mail.scut.edu.cn, 1070403885@qq.com, eehbqin@scut.edu.cn

Abstract

Depth prediction is a critical problem in robotics applications especially autonomous driving. Generally, depth prediction based on binocular stereo matching and fusion of monocular image and laser point cloud are two mainstream methods. However, the former usually suffers from over-fitting while building cost volume, and the latter has a limited generalization due to the lack of geometric constraint. To solve these problems, we propose a novel multimodal neural network, namely **UAMD-Net**, for dense depth completion based on fusion of binocular stereo matching and the weak constrain from the sparse point clouds. Specifically, the sparse point clouds are converted to sparse depth map and sent to the *multimodal feature encoder (MFE)* with binocular image, constructing a cross-modal cost volume. Then, it will be further processed by the *multimodal feature aggregator (MFA)* and the depth regression layer. Furthermore, the existing multimodal methods ignore the problem of modal dependence, that is, the network will not work when a certain modal input has a problem. Therefore, we propose a new training strategy called **Modal-dropout** which enables the network to be adaptively trained with multiple modal inputs and inference with specific modal inputs. Benefiting from the flexible network structure and adaptive training method, our proposed network can realize unified training under various modal input conditions. Comprehensive experiments conducted on KITTI depth completion benchmark demonstrate that our method produces robust results and outperforms other state-of-the-art methods.

1 Introduction

Dense depth prediction is of great significance to the robotics applications such as autonomous driving. The ac-

quisition of depth information is the prerequisite for solving the tasks like obstacle avoidance, 3D object detection and 3D scene reconstruction [12]. Typically, there are two major application environments, indoors and outdoors. For the former, the mainstream method is using the depth camera to proactively acquire the depth information or utilizing the stereo vision in a passive way. But for the latter, it is better to apply the stereo vision or LiDAR sensors [19]. Besides, estimating the depth directly from monocular image [18, 6, 5, 21, 22, 7] is also an attempt although it is a morbid problem. Recently, stereo vision algorithms have achieved an impressive progress both in supervised [1, 10, 23, 4] and in self-supervised way [9, 26, 17, 14], but the problems of weak texture failure and over fitting are still unsolved, which will lead to limited accuracy. In contrast, LiDAR sensors can provide reliable and accurate depth sensing. But unfortunately, current LiDAR sensors only acquire sparse depth measurements which is not sufficient for real applications such as robotic navigation. Therefore, how to get both dense and accurate depth perception is still a challenging topic.

Many recent works on this topic turn on the trend of multimodal learning by fusing monocular image information and sparse depth measurements for depth completion. Cheng et al. [3, 2] proposed to utilize the convolutional spatial propagation network (CSPN, CSPN++) to assemble the features learning from monocular image and the corresponding LiDAR scans. Similarly, Park et al. [16] continued this idea and put forward the non-local spatial propagation network (NLSPN) that predicted non-local neighbors for each pixel, aiming to solve the mixed-depth problem. In contrast to the application of spatial propagation scheme, Tang et al. [19] proposed to fuse the LiDAR data and RGB image information by performing GuideNet which consists of learnable content-dependent and spatially-variant kernels. Zhao et al. [25] proposed to adopt the graph propagation to capture the observed spatial contexts. More recently, PENet [11] and FCFR-Net [13] were proposed to carried out the depth completion through a two-stage coarse-to-

*Corresponding author.

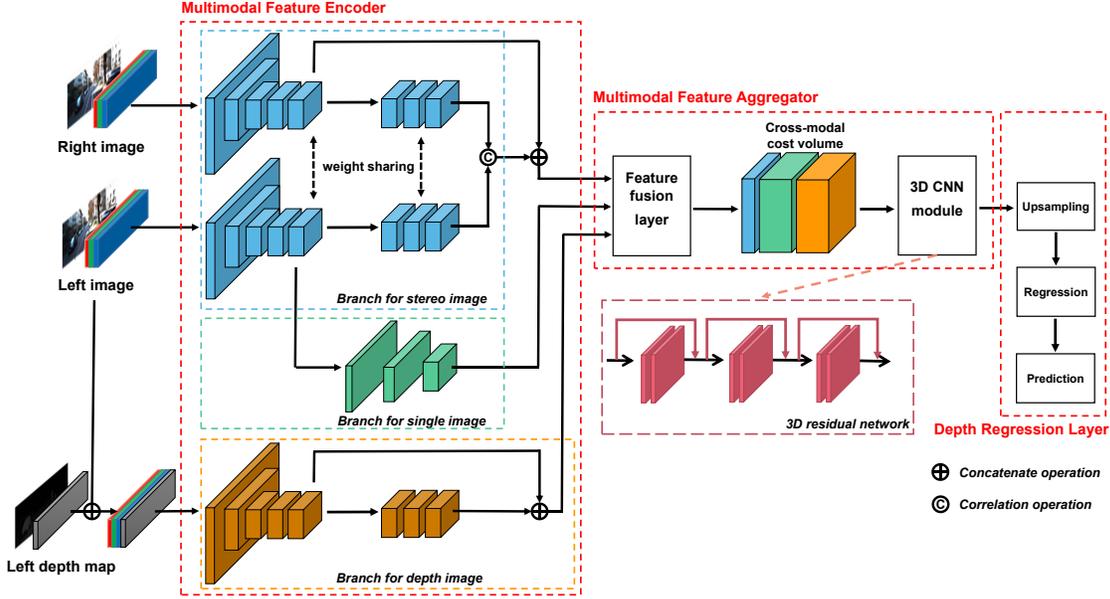


Figure 1. The network architecture of the proposed UAMD-Net, which consists of three main components: multimodal feature encoder (MFE), multimodal feature aggregator (MFA) and Depth Regression Layer (DRL).

fine mechanism. Although these methods have achieved remarkable results, its demand for a huge amount of label data and a long training time for convergence cannot be ignored. Besides, the hard fitting of monocular image and scene depth lacks geometric constraints, which will lead to scene dependence and limited generalization.

Instead of establishing the multimodal depth completion network based on monocular image and sparse point cloud, LiStereo [24] was the pioneer to make use of the multimodal learning of binocular image and sparse depth measurements. Except training on supervised mode, it also can be trained on semi-supervised mode benefited from the view synthesis scheme of binocular image and the weak constraint from the sparse point cloud. However, its feature fusion and aggregation modules are not sufficient for producing satisfactory results.

Besides, all these existing multimodal methods ignore the problem of modal dependence, which means the network will not work when a certain modal input has a problem.

To address the aforementioned issues, in this paper we propose a unified multimodal neural network, namely **UAMD-Net**, that is capable to fuse the feature learning of binocular image and sparse depth map. Specifically, it consists of the **MFE** and **MFA** module which can extract the cross-modal features to construct the 4D cost volume and then accomplish the feature aggregation based on 3D con-

volution. Besides, to solve the modal dependence problem, we propose a new training scheme called **Modal-dropout** which is capable to adaptively train the network with multiple modal inputs and inference with specific modal inputs. In particular, the flexible network structure and adaptive training method enable the network to realize unified training under various modal input conditions, including binocular stereo matching (*dual*), fusion of monocular image and sparse depth map (*mono Lidar*), and combination of binocular image and sparse depth map (*dual Lidar*). We conducted extensive experiments on KITTI depth completion benchmark and the results show that our technique achieves strong results and outperforms current state-of-the-art methods.

In short, the contributions of our research include:

- We propose a novel multimodal neural network for realizing depth completion, which we called **UAMD-Net**. It is capable to combine the advantages of binocular stereo matching and sparse point cloud constraint to get rid of the risk of over fitting and obtain better generalization performance.
- We propose a new training strategy called **Modal-dropout** to solve the modal dependence problem for multimodal learning. To the best of our knowledge, this is the first trial to provide a viable solution to the modal dependence problem.

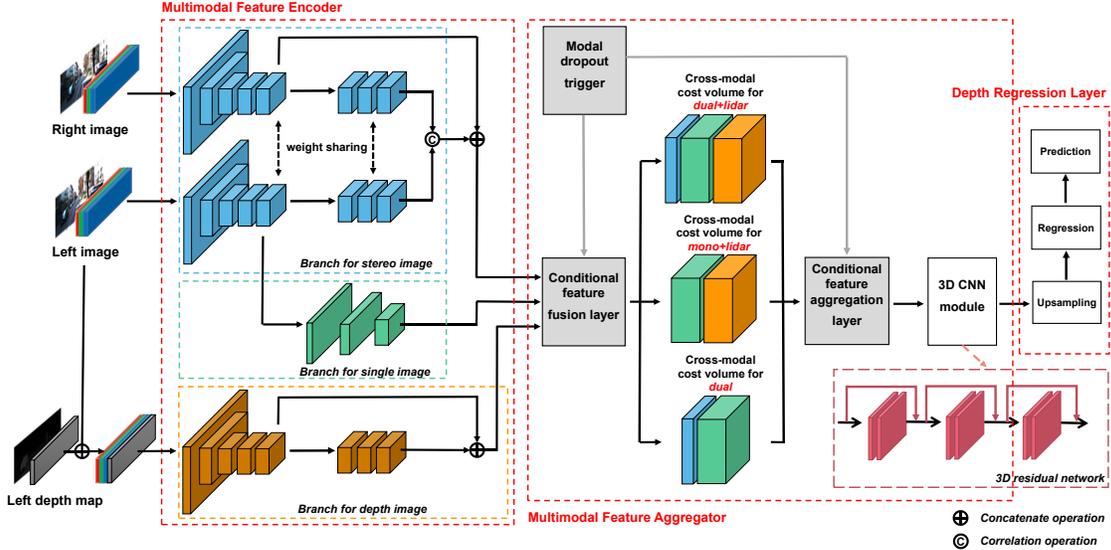


Figure 2. The extended network structure of UAMD-Net integrating three new modules including modal dropout trigger (MDT), conditional feature fusion layer (CFFL) and conditional feature aggregation layer (CFAL), which is designed for carrying out the proposed Modal-dropout training scheme.

- Our proposed network has great flexibility to realize unified training under various modal input conditions. Extensive experimental results on KITTI depth completion benchmark demonstrate the superiority of our proposed method quantitatively and qualitatively.

2 Related Work

2.1 Monocular Depth Estimation

The first work of monocular depth estimation can be traced back to 2005, when Saxena et al. constructed a Gaussian MRF probabilistic model by handcraft features to directly perform RGB-to-depth regression. With the popularity of convolutional neural network (CNN), a series of CNN-based monocular depth estimation networks have been proposed [6, 5, 21, 22, 7, 15], and achieved a constant improvement in accuracy. However, since monocular depth estimation is an ill-conditioned problem that relies heavily on the learning of scene texture and structure information, it is difficult to solve the problem of scene generalization.

2.2 Depth Estimation Based on Stereo Match

Unlike monocular depth estimation, binocular depth estimation can use the geometric constraints of stereo matching for depth prediction. Generally, it can be divided into two learning paradigms. The first one is supervised learning. This kind of methods usually apply a two-stream neural

network to extract binocular image features, then construct the matching cost volume for further depth regression, like [1, 10, 23, 4]. The second one is self-supervised learning. This kind of methods perform self-supervised learning through the mechanism of view synthesis without the need for labeled data, such as [9, 26, 17, 14]. Although binocular depth estimation can utilize the geometric prior information of the scene, the constraint based on the similarity of the matching pixels in the left and right views is not strong, which easily leads to the problem of overfitting.

2.3 Depth Completion Based on Multimodal Learning

Recently, depth completion based on multimodal learning receives raising spotlight. CSPN [3] started a fashion of fusing of features learning from monocular image and the corresponding depth measurement through the convolutional spatial propagation network. Soon after, CSPN++ [2] further improves the effectiveness and efficiency of CSPN by learning adaptive convolutional kernel sizes and the number of iterations for the propagation. Park et al. [16] put forward the non-local spatial propagation network (NL-SPN) that predicted non-local neighbors for each pixel, aiming to solve the mixed-depth problem. Except the application of spatial propagation scheme, Tang et al. [19] proposed GuideNet to fuse the LiDAR data and RGB image information by performing learnable content-dependent and spatially-variant kernels. Zhao et al. [25] proposed to adopt

the graph propagation to model the observed spatial contexts with depth values, so as to better guide the recovery of the unobserved pixels’ depth. More recently, PENet [11] and FCFR-Net [13] were proposed to carry out the depth completion through a two-stage coarse-to-fine mechanism and achieved impressive results. However, these methods need be trained in supervised learning mode, which require a large amount of annotated label. In order to get rid of the limit of supervised learning, Zhang et al. proposed LiStereos [24] to realize depth completion by accomplishing the multimodal learning of binocular image and sparse depth measurements. It can be trained on semi-supervised mode based on the view synthesis scheme of binocular image and the weak constraint from the sparse point cloud. However, since the sparse point cloud is not able to offer enough constraint like the label data, and its feature fusion and aggregation modules are not sufficient enough, it still has a large performance gap compared with supervised learning methods. In this paper, we propose a novel multimodal neural network which tries to combine the advantages of binocular stereo matching and sparse point cloud constraint, aiming at improving the performance of both supervised and semi-supervised learning modes.

3 Proposed Method

In this section, we firstly describe the network architecture of the proposed **UAMD-Net**, shown in Fig. 1, which is mainly divided into three components: **MFE**, **MFA** and **DRL**. Next, we detail the proposed adaptive multimodal training strategy **Modal-dropout** and the extended network structure of **UAMD-Net** which is designed for unified training under various modal input conditions. Finally, we introduce the objective function design.

3.1 Multimodal Neural Network for Depth Completion

MFE: For multimodal inputs, we design three branches to extract the specific modal features. The branch for stereo image and the branch for depth map have the same configuration, consisting of multiple convolutional layers with ReLUs as the activation functions. We use the branch for stereo image to extract the features from both left and right image, and then obtain the cross-modal features by accomplishing the correlation operation. Besides, we concatenate the image features from the middle layer to enhance the feature representation ability. We use the branch for depth map to extract the cross-modal features from the concatenation of monocular image and the corresponding sparse depth map. Moreover, we design a branch for single image in order to enhance the features from the image domain. More specific settings are specified in the supplementary material.

MFA: After acquiring the multimodal features from different branches, we design a feature fusion layer to construct the cross-modal 4D cost volume by fusing the multimodal features. Then, inspired by [1], we establish a simple yet effective 3D CNN module by stacking six $3 \times 3 \times 3$ 3D convolutional layers and three residual blocks, which can aggregate the features from both spatial and channel dimensions.

DRL: After feature aggregation, we apply the trilinear interpolation on the disparity feature map to recover the resolution to $H \times W \times D$ (H, W represent the height and width of the image, D denotes the disparity range which we set as 192). Then, we adopt the softmax operation to carry out the disparity regression. In this case, the features in D dimension are considered to be the probability of the corresponding disparity. Finally, the disparity map will be transformed to depth map according to the stereo constraint: $d = fl/disp$, where d denotes the depth map, b denotes the length of baseline, fl denotes the focal length of the camera, and $disp$ denotes the predicted disparity map.

3.2 Adaptive Multimodal Training Strategy and the Extended Unified Network Structure

We get inspiration from the universal training method **Dropout**, which randomly discards nodes to prevent network overfitting. Similarly, we propose to randomly drop the specific modal inputs during training while inference with fixed modal inputs, so as to prevent the network from being limited to specific modal inputs, which solves the modal dependence problem. Naturally, we name this training strategy **Modal-dropout**. To carry out this training scheme, we need to further extend the network structure. As shown in Fig. 2, we design the **MDT** component to guide the **CFFL** and the **CFAL** to adaptively accomplish the feature fusion and aggregation, respectively.

MDT: The key of **MDT** is to realize the random sampling of three modal input combinations: *dual Lidar*, *mono Lidar*, and *dual*, which can be formulated as follow:

$$X \sim P \{X = k\} = 1/3, k = 1, 2, 3 \quad (1)$$

where X denotes the sample variable of three cases.

Therefore, the formulation of the network trained on supervised mode can be described as:

$$\begin{cases} d = f(I_l, I_r, D_l) & X = 1 \\ d = f(I_l, D_l) & X = 2 \\ d = f(I_l, I_r) & X = 3 \end{cases} \quad (2)$$

and the formulation of the network trained on semi-supervised mode can be described as:

$$\begin{cases} d = f(I_l, I_r, D_l, D_r) & X = 1 \\ d = f(I_l, I_r) & X = 2 \end{cases} \quad (3)$$

Table 1. Ablation study on weights of loss for semi-supervised learning. Our UAMD-Net is trained in semi-supervised mode with the modal input *dual Lidar*.

Loss weights	$w_l = 1$ $w_p = 0$	$w_l = 1$ $w_p = 0.2$	$w_l = 1$ $w_p = 0.4$	$w_l = 1$ $w_p = 0.6$	$w_l = 1$ $w_p = 0.8$	$w_l = 1$ $w_p = 1.0$	$w_l = 1$ $w_p = 1.3$	$w_l = 1$ $w_p = 1.5$	$w_l = 0$ $w_p = 1.0$
RMSE (mm)	1725.587	1513.612	1373.387	1337.512	1270.628	1305.006	1267.047	1381.939	2587.277

Table 2. Ablation study on different learning mode for various modal input combinations.

	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)
<i>supervised</i>				
dual_lidar	1.747	1.007	669.166	252.580
mono_lidar	1.938	1.098	918.067	346.224
dual	6.400	4.545	1163.147	603.036
<i>semi-supervised</i>				
dual_lidar	4.954	1.864	1267.047	460.880
dual	4.639	2.031	2587.277	731.396

where d denotes the predicted depth map, f denotes the network, I_l, I_r, D_l, D_r denote the left image, right image, left sparse depth map, right sparse depth map, respectively.

CFFL: To realize the adaptive multimodal training, we construct a *conditional feature fusion layer*. It receives the command from the **MDT** to adaptively produce three forms of 4D cost volume, corresponding to three different modal input combinations: *dual lidar*, *mono lidar*, and *dual*.

CFAL: Since the feature dimension of 4D cost volume is changeable, we establish a *conditional feature aggregation layer* to adaptively accomplish the feature aggregation. It receives the same command as the **CFFL** from the **MDT**. We realize it by constructing two 3D convolutional layers with ReLUs as the activation functions. The input feature dimension is variant while the output feature dimension is fixed.

3.3 Objective Function Design

Our network can be trained in both supervised and semi-supervised mode benefited from the view synthesis scheme of stereo vision. For supervised learning, we optimize **UAMD-Net** by minimizing the following L_2 loss,

$$Loss_{sup} = \frac{1}{N} \sum_{p \in P_v} \|d_p^{gt} - d_p\|^2 \quad (4)$$

where P_v represents the set of the valid pixels. d_p^{gt} and d_p denote the ground truth and predicted depth at the pixel p , respectively. N is the number of valid pixels.

For semi-supervised learning, we only provide sparse ground truth depth map for supervision. Since the density

of sparse depth map is low and L_2 is more sensitive to the outliers, we adopt the following L_1 loss.

$$Loss_{lidar} = \frac{1}{N} \sum_{p \in P_v} \|d_p^{sp} - d_p\| \quad (5)$$

where d_p^{sp} denotes the ground truth sparse depth map.

Besides, we follow [1] to use a combination of an L_1 and single scale SSIM term as our photometric image reconstruction loss, which compares the input image I_p^r and its reconstruction \tilde{I}_p^r .

$$Loss_{photometric} = \frac{1}{N} \sum_{p \in P_v} \alpha \cdot SSIM(I_p^r, \tilde{I}_p^r) + (1 - \alpha) \cdot \|I_p^r - \tilde{I}_p^r\| \quad (6)$$

Here, we use a simplified SSIM with a 3×3 block filter and set $\alpha = 0.85$. We noted that the correlation of stereo feature map can construct the cost volume for both left and right disparity, so in semi-supervised training, we construct the cost volume for right disparity with the reconstruction constraint of right image, while construct the cost volume for left disparity in the inference. Moreover, we apply the following objective for encouraging the depth to be locally smooth with an L_1 penalty on the depth gradients ∂_d .

$$Loss_{gradient} = \frac{1}{N} \sum_{p \in P_v} |\partial_x d_p^r| \cdot e^{-\|\partial_x I_p^r\|} + |\partial_y d_p^r| \cdot e^{-\|\partial_y I_p^r\|} \quad (7)$$

Furthermore, inspired by [27, 20], we employ the noise label learning strategy. We generate the noise depth map by the traditional depth estimation method Semi-Global Matching (SGM) for convenience. Since the accuracy of the noise depth map is limited and L_2 is more sensitive to the outliers, we adopt the following L_1 loss.

$$Loss_{noise} = \frac{1}{N} \sum_{p \in P_v} \|d_p^n - d_p\| \quad (8)$$

where d_p^n denotes the generated noise depth map.

Finally, the above objective will be combined to train in a multi-task learning fashion. w_l, w_p, w_g, w_n represent the corresponding weight parameters, which will be fine-tuned according to the training feedback.

$$Loss_{semi} = w_l \cdot Loss_{lidar} + w_p \cdot Loss_{photometric} + w_g \cdot Loss_{gradient} + w_n \cdot Loss_{noise} \quad (9)$$

Table 3. Ablation study on Modal-dropout training scheme: both training and validating with various modal input combinations.

	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)	Performance attenuation
<i>supervised</i>					
dual_lidar	1.873	1.095	730.327	284.489	↓9.14%
mono_lidar	9.423	5.986	1306.826	610.729	↓12.35%
dual	2.086	1.234	973.540	385.330	↓6.04%
avg.	-	-	-	-	↓9.18%
<i>semi-supervised</i>					
dual_lidar	4.269	1.862	1419.445	515.308	↓13.02%
dual	4.343	2.014	3087.818	744.682	↓19.35%
avg.	-	-	-	-	↓16.19%

Table 4. Ablation study on Modal-dropout training scheme: training with various modal input combinations while validating with specific modal input combination. For example, the first section means training with {dual_lidar, mono_lidar, dual} while validating with dual_lidar.

	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)	Performance attenuation
<i>supervised</i>					
<i>dual_lidar</i>					
dual_lidar	1.888	1.065	747.816	285.691	↓11.75%
mono_lidar	7.135	4.908	1319.096	631.323	↓13.48%
dual	1.165	2.035	991.919	373.009	↓8.04%
avg.	-	-	-	-	↓11.09%
<i>mono_lidar</i>					
dual_lidar	1.810	1.036	709.975	269.336	↓6.10%
mono_lidar	4.758	3.011	1319.360	526.489	↓13.43%
dual	2.054	1.139	967.731	364.451	↓5.41%
avg.	-	-	-	-	↓ 8.31%
<i>dual</i>					
dual_lidar	1.888	1.056	731.881	278.886	↓9.37%
mono_lidar	6.138	3.760	1369.923	527.990	↓17.78%
dual	2.130	1.154	988.307	370.867	↓7.65%
avg.	-	-	-	-	↓11.60%
<i>semi-supervised</i>					
<i>dual_lidar</i>					
dual_lidar	5.129	1.887	1320.503	492.800	↓4.22%
dual	5.692	2.156	7733.671	955.438	-
avg.	-	-	-	-	-
<i>dual</i>					
dual_lidar	4.663	1.811	1351.973	488.343	↓6.70%
dual	4.140	2.020	2902.286	742.696	↓12.18%
avg.	-	-	-	-	↓9.44%

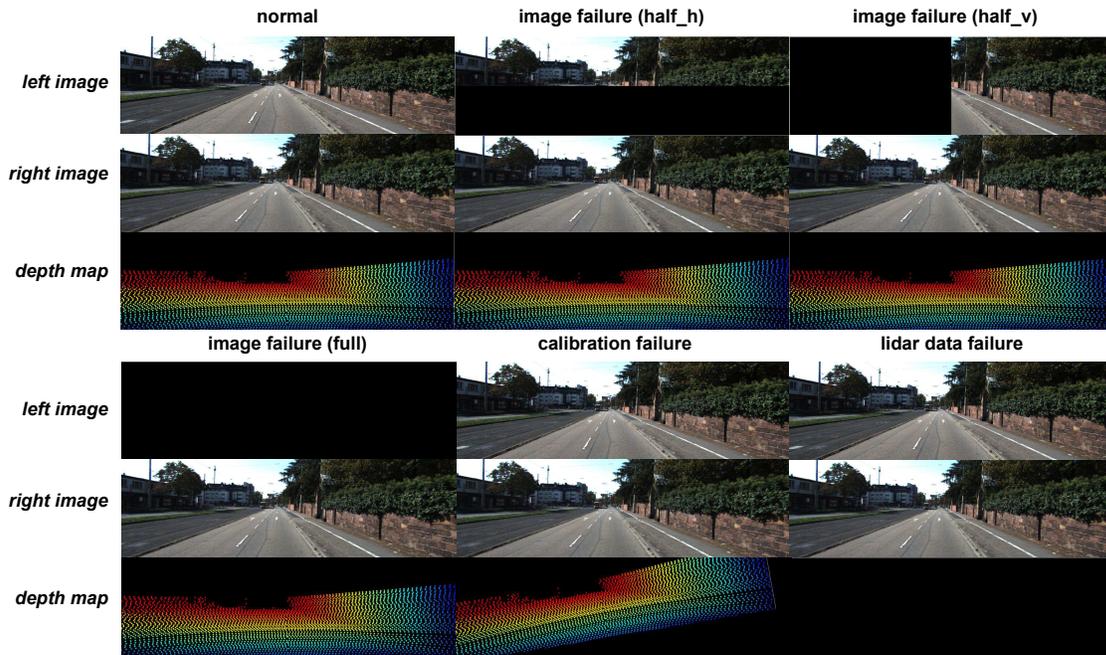


Figure 3. Modal failure situations: image failure (three forms including half of horizontal pixels, half of the vertical pixels and full pixels), camera and LiDAR calibration failure, and LiDAR data failure.

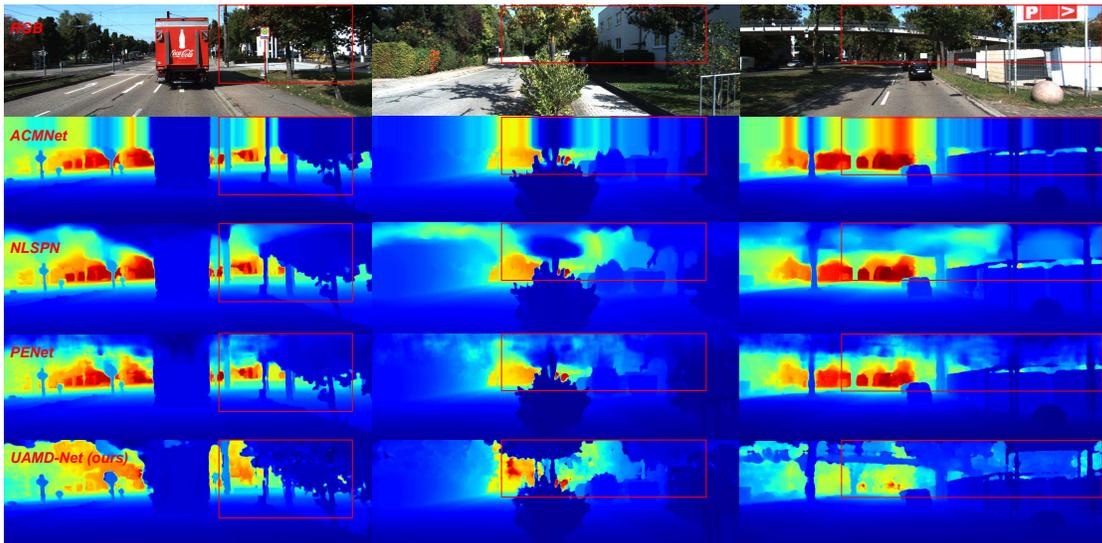


Figure 4. Qualitative results of different methods. From top to down are the input images, results of ACMNet [25], NLSPN [16], PENet [11] and our UAMD-Net respectively.

Table 5. Ablation study on noise label learning. Our UAMD-Net is trained in semi-supervised mode with the modal input *dual_lidar*.

w_n	0	0.1	0.2	0.3	0.4
RMSE (mm)	1267.05	1225.85	1270.06	1264.00	1269.72

4 Experiments

4.1 Experiments Settings

Benchmark Dataset: The KITTI depth completion dataset [8] contains 42949 pair depth maps for training, 3426 pair depth maps for validation, and 1000 frames for testing. Specially a selection set of 1000 frames is also provided. These ground truth depth maps are generated by registering 11 LiDAR scans temporally and further refined with the corresponding stereo image pairs. Since the test set does not contain stereo images, we split the validation set into two sub-sets, 2426 pairs of stereo images for testing (split-test) and 1000 for validation (split-val) according to the selection set, which guarantees the fairness of the experimental comparison with other methods.

Evaluation Metrics: We follow the KITTI benchmark and exiting methods [19, 25, 11] to use four standard metrics for evaluation: root mean squared error (RMSE), mean absolute error (MAE), root mean squared error of the inverse depth (iRMSE) and mean absolute error of the inverse depth (iMAE). Among them, RMSE and MAE measure the depth accuracy directly, while iRMSE and iMAE compute the mean error of the inverse depth, giving less weight to the far-away points. Since RMSE is more sensitive to the outliers, it is chosen as the dominant metric to rank the submissions on the KITTI leaderboard.

Implementation Details: Our network is implemented using PyTorch framework. For supervised learning, the learning rate begins at $1e-4$ and is decayed by 0.5 at 10 epochs, 0.1 at 14 epochs and 0.01 at 17 epochs. For semi-supervised learning, the learning rate begins at $1e-4$ and is decayed by 0.1 at $10e3$ iterations, 0.01 at $14e3$ iterations and 0.01 at $16e3$ iterations. The batch size is set to 4 for training on 2 NVIDIA GTX 1080Ti GPUs for all models. We report the experimental results based on the validation set (split-val) for ablation studies of our proposed method. While compared with other start-of-the-art methods, we report the experimental results conducted on the test set (split-test). In all tables, bold values indicate the best performance, underlined values indicate the suboptimal performance.

4.2 Ablation Study on Weight of Loss for Semi-supervised Learning

The photometric loss is essential for the model to be trained in a self-supervised manner, so we investigate the influence of the weight of the photometric loss w_p . According to Eq. 9, we keep $w_l = 1$, $w_g = 0.01$, $w_n = 0$, and set w_p between 0 to 1.5, and the results are shown in Table 1. It is clear that the photometric loss does help complete sparse depth map, reducing RMSE from 1725mm to 1267mm with $w_p = 0$ and $w_p = 1.3$ respectively. However, too much weight on photometric loss worsens the results. So we set $w_p = 1.3$ throughout the experiments.

4.3 Ablation Study on Different Learning Modes for Various Modal Input Combinations

In this section, we report the performance of our network trained with various modal input combinations in supervised and semi-supervised mode. As shown in Table 2, *dual_lidar* achieves the lowest RMSE while *dual* acquires the highest in both supervised and semi-supervised mode, proving that our network can solve the problem of overfitting well. Besides, *mono_lidar* has a better performance than *dual* benefited from the constraint of sparse point cloud.

4.4 Ablation Study on Two Different Modal-out Training Schemes

In this section, we study the performance of our two proposed **Modal-dropout** training strategies. The first one is both training and validating with various modal input combinations, the results are shown in Table 3. The second one is training with various modal input combinations while validating with specific modal input combination, the results are shown in Table 4. It is obvious that the model performance drops while trained with the **Modal-dropout** scheme, with the lowest average 8.31%. However, the model then will obtain the ability to inference with different modal input combinations, solving the modal dependence problem, which greatly improves its robustness against modal input failure situations.

4.5 Ablation Study on Noise Label Learning

We introduce the noise label learning scheme to further improve the performance of semi-supervised learning. As shown in Table 5, an appropriate proportion of noise labels is conducive to improve the performance, reducing RMSE from 1267 to 1225 with $w_n = 0.1$.

Table 6. Robustness against different modal failure situations. In order to deal with the situations of image failure, our UAMD-Net switches to work with the modal input *mono Lidar*.

	image failure (half_h)				image failure (half_v)				image failure (full)			
	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)
PENet [11]	4.448	1.626	1630.624	463.079	3.545	1.320	1779.584	421.798	5.173	1.854	2531.038	656.476
ACMNet [25]	3.255	1.086	1131.010	274.131	2.592	0.971	1151.714	261.263	3.307	1.111	<u>1439.257</u>	319.110
UAMD-Net (mono_lidar)	4.821	3.009	<u>1348.657</u>	525.945	4.821	3.009	<u>1348.657</u>	525.945	4.821	3.009	1348.657	525.945

Table 7. Robustness against different modal failure situations. In order to deal with the situations of rotation failure and LiDAR data failure, our UAMD-Net switches to work with the modal input *dual*.

	rotation failure				LiDAR data failure			
	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)
PENet [11]	524.298	430.731	17517.829	12812.575	-	-	-	-
ACMNet [25]	276.951	247.391	<u>16856.885</u>	12277.822	-	-	-	-
UAMD-Net (dual)	2.115	1.151	956.168	361.121	2.115	1.151	956.168	361.121

Table 8. Comparison with other supervised methods. Our UAMD-Net is trained in supervised mode with the modal input *dual Lidar*.

	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)
PENet [11]	2.159	0.903	756.667	209.369
ACMNet [25]	2.099	0.868	765.210	205.503
GuideNet [19]	-	-	763.3	-
NLSPN [16]	2.0	0.8	776.3	198.5
UAMD-Net (ours)	1.729	0.999	677.132	254.056

Table 9. Comparison with other semi-supervised methods. Our UAMD-Net is trained in semi-supervised mode with the modal input *dual Lidar*.

	iRMSE (1/km)	iMAE (1/km)	RMSE (mm)	MAE (mm)
Sparse2dense [15]	4.08	1.61	1301.05	352.22
LiStereo [24]	3.84	1.32	1278.87	326.10
UAMD-Net (ours)	4.71	1.82	1241.10	464.38

4.6 Robustness against Different Modal Failure Situations

In order to prove the effectiveness of our proposed **Modal-dropout** training strategy, we simulate the situations when the input modalities are problematic: image failure (three forms including half of horizontal pixels, half of the vertical pixels and full pixels), camera and LiDAR calibration failure, and LiDAR data failure, as shown in Fig. 3. The experimental results reported in Table 6 and 7 shows that current multimodal depth completion models like PENet [11] and ACMNet [25] will have a large performance penalty for the case of image failure. Besides, they will complete failure for the case of calibration failure and LiDAR data failure. On the contrary, our model can switch to proper inference mode to maintain stable performance, demonstrating the great robustness of our method against modal input failure situations. More discussions are presented in the supplementary material.

4.7 Comparison with State-of-the-Arts

As shown in Table 8 and 9, our method surpasses other state-of-the-art methods no matter in supervised or semi-supervised learning fashion. Especially, our method achieves an 80mm RMSE gap with the closest competitor, PENet [11].

4.8 Qualitative Results

The qualitative results are reported in Fig. 4. From top to down are the input images, results of ACMNet [25], NL-

SPN [16], PENet [11] and our **UAMD-Net** respectively. It shows that the depth map predicted by our method has great advantages in preserving edge details.

5 Conclusion

In this paper, we propose a unified multimodal neural network called **UAMD-Net** for depth completion task, aiming to combine the advantages of binocular stereo matching and sparse point cloud constraint to get rid of the risk of over fitting and obtain better generalization performance. Besides, to address the modal dependence problem, we further propose a new training strategy named **Modal-dropout**. The flexible network structure and adaptive training strategy enable the network realize unified training under various modal input conditions, which greatly improves the robustness of the multimodal neural network in the case of a modal input failure. Our experimental results demonstrate that the proposed method not only overcomes the modal dependence problem but also achieve better quantitative and qualitative performance compared with state-of-the-art methods.

References

- [1] J.-R. Chang and Y.-S. Chen. Pyramid stereo matching network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5410–5418, 2018.
- [2] X. Cheng, P. Wang, C. Guan, and R. Yang. Cspn++: Learning context and resource aware convolutional spatial propagation networks for depth completion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10615–10622, 2020.
- [3] X. Cheng, P. Wang, and R. Yang. Learning depth with convolutional spatial propagation network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(10):2361–2379, 2019.
- [4] X. Cheng, Y. Zhong, M. Harandi, Y. Dai, X. Chang, H. Li, T. Drummond, and Z. Ge. Hierarchical neural architecture search for deep stereo matching. *Proceedings of the Advances in Neural Information Processing Systems*, 33:22158–22169, 2020.
- [5] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2650–2658, 2015.
- [6] D. Eigen, C. Puhrsch, and R. Fergus. Depth map prediction from a single image using a multi-scale deep network. *Proceedings of the Advances in Neural Information Processing Systems*, 27, 2014.
- [7] H. Fu, M. Gong, C. Wang, K. Batmanghelich, and D. Tao. Deep ordinal regression network for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2002–2011, 2018.
- [8] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE, 2012.
- [9] C. Godard, O. Mac Aodha, and G. J. Brostow. Un-supervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 270–279, 2017.
- [10] X. Guo, K. Yang, W. Yang, X. Wang, and H. Li. Group-wise correlation stereo network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3273–3282, 2019.
- [11] M. Hu, S. Wang, B. Li, S. Ning, L. Fan, and X. Gong. Penet: Towards precise and efficient image guided depth completion. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 13656–13662. IEEE, 2021.
- [12] H. Laga, L. V. Jospin, F. Boussaid, and M. Benamoun. A survey on deep learning techniques for stereo-based depth estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [13] L. Liu, X. Song, X. Lyu, J. Diao, M. Wang, Y. Liu, and L. Zhang. Fcfr-net: Feature fusion based coarse-to-fine residual learning for monocular depth completion. *arXiv preprint arXiv:2012.08270*, pages arXiv–2012, 2020.
- [14] Y. Luo, J. Ren, M. Lin, J. Pang, W. Sun, H. Li, and L. Lin. Single view stereo matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 155–163, 2018.
- [15] F. Ma, G. V. Cavalheiro, and S. 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.

- [16] J. Park, K. Joo, Z. Hu, C.-K. Liu, and I. So Kweon. Non-local spatial propagation network for depth completion. In *European Conference on Computer Vision*, pages 120–136. Springer, 2020.
- [17] A. Pilzer, D. Xu, M. Puscas, E. Ricci, and N. Sebe. Unsupervised adversarial depth estimation using cycled generative networks. In *2018 International Conference on 3D Vision (3DV)*, pages 587–595. IEEE, 2018.
- [18] A. Saxena, S. Chung, and A. Ng. Learning depth from single monocular images. *Proceedings of the Advances in Neural Information Processing Systems*, 18, 2005.
- [19] J. Tang, F.-P. Tian, W. Feng, J. Li, and P. Tan. Learning guided convolutional network for depth completion. *IEEE Transactions on Image Processing*, 30:1116–1129, 2020.
- [20] A. Tonioni, M. Poggi, S. Mattoccia, and L. Di Stefano. Unsupervised domain adaptation for depth prediction from images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(10):2396–2409, 2019.
- [21] P. Wang, X. Shen, Z. Lin, S. Cohen, B. Price, and A. L. Yuille. Towards unified depth and semantic prediction from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2800–2809, 2015.
- [22] D. Xu, E. Ricci, W. Ouyang, X. Wang, and N. Sebe. Multi-scale continuous crfs as sequential deep networks for monocular depth estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5354–5362, 2017.
- [23] F. Zhang, X. Qi, R. Yang, V. Prisacariu, B. Wah, and P. Torr. Domain-invariant stereo matching networks. In *European Conference on Computer Vision*, pages 420–439. Springer, 2020.
- [24] J. Zhang, M. S. Ramanagopal, R. Vasudevan, and M. Johnson-Roberson. Listereo: Generate dense depth maps from lidar and stereo imagery. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 7829–7836. IEEE, 2020.
- [25] S. Zhao, M. Gong, H. Fu, and D. Tao. Adaptive context-aware multi-modal network for depth completion. *IEEE Transactions on Image Processing*, 30:5264–5276, 2021.
- [26] Y. Zhong, Y. Dai, and H. Li. Self-supervised learning for stereo matching with self-improving ability. *arXiv preprint arXiv:1709.00930*, 2017.
- [27] S. Zhu, G. Brazil, and X. Liu. The edge of depth: Explicit constraints between segmentation and depth. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13116–13125, 2020.