# Noise-Aware Unsupervised Deep Lidar-Stereo Fusion

*Xuelian Cheng[1,2], *Yiran Zhong[2,4,5], Yuchao Dai[1], Pan Ji[3], Hongdong Li[2,4]
[1]Northwestern Polytechnical University [2]Australian National University
[3]NEC Laboratories America, [4]ACRV, [5]Data61 CSIRO

## Abstract

*In this paper, we present LidarStereoNet, the first unsupervised Lidar-stereo fusion network, which can be trained in an end-to-end manner without the need of ground truth depth maps. By introducing a novel "Feedback Loop" to connect the network input with output, LidarStereoNet could tackle both noisy Lidar points and misalignment between sensors that have been ignored in existing Lidar-stereo fusion studies. Besides, we propose to incorporate a piecewise planar model into network learning to further constrain depths to conform to the underlying 3D geometry. Extensive quantitative and qualitative evaluations on both real and synthetic datasets demonstrate the superiority of our method, which outperforms state-of-the-art stereo matching, depth completion and Lidar-Stereo fusion approaches significantly.*

## 1. Introduction

Accurately perceiving surrounding 3D information from passive and active sensors is crucial for numerous applications such as localization and mapping [15], autonomous driving [18], obstacle detection and avoidance [25], and 3D reconstruction [10, 33]. However, each kind of sensors alone suffers from its inherent drawbacks. Stereo cameras are well-known for suffering from computational complexities and their incompetence in dealing with textureless/repetitive areas and occlusion regions [28], while Lidar sensors often provide accurate but relatively sparse depth measurements [5].

Therefore, it is highly desired to fuse measurements from Lidar and stereo cameras to achieve high-precision depth perception by exploiting their complementary properties. However, it is a non-trivial task as accurate Stereo-Lidar fusion requires a proper registration between Lidar and stereo images and noise-free Lidar points. Existing methods are not satisfactory due to the following drawbacks:

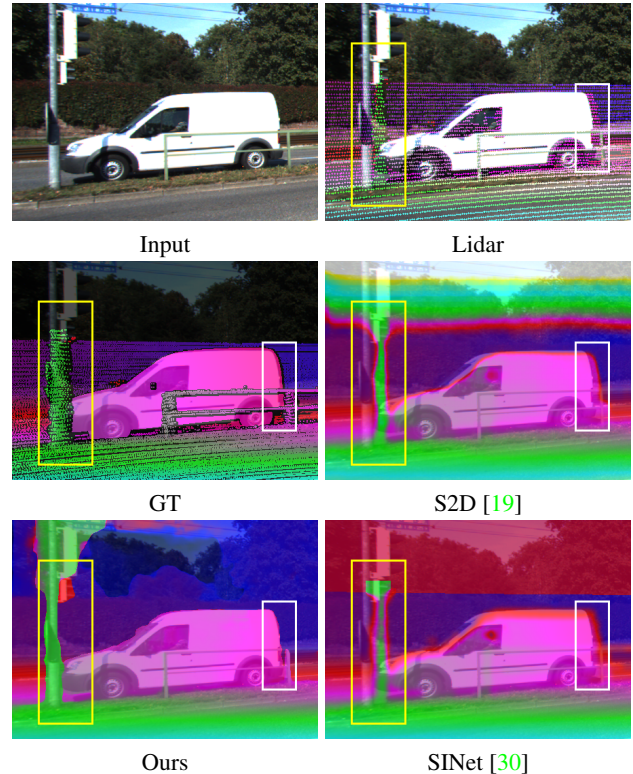- Existing deep neural network based Lidar-Stereo fu-



Figure 1. **Results on KITTI 2015.** We highlight the displacement error of Lidar points with bounding boxes. Lidar points are dilated for better visualization and we overlay our disparity maps to the colour images for illustration. Note the Lidar points for the foreground car and utility pole have been aligned to the background. Our method successfully recovers accurate disparities on tiny and moving objects while the other methods are misled by drifted and noisy Lidar points.

sion studies [2, 20, 26] strongly depend on the availability of large-scale ground truth depth maps, and thus their performance is fundamentally limited by their generalization ability to real-world applications.

- Due to rolling-shutter effects of Lidar and other calibration imperfections, a direct registration will introduce significant alignment errors between Lidar and stereo depth. Furthermore, existing methods tend to assume the Lidar measurements are noise-free [19,

---

*These authors contributed equally in this work.

30]. However, as illustrated in Fig. 1, the misalignment and noisy Lidar measurements cause significant defects in Stereo-Lidar fusion.

In this paper, we tackle the above challenges and propose a novel framework "LidarStereoNet" for accurate Stereo-Lidar fusion, which can be trained in an end-to-end unsupervised learning manner. Our framework is noise-aware in the sense that it explicitly handles misalignment and noise in Lidar measurements.

Firstly, we propose to exploit photometric consistency between stereo images, and depth consistency between stereo cameras and Lidar to build an unsupervised training loss, thus removing the need of ground truth depth/disparity maps. It enables a strong generalization ability of our framework to various real-world applications.

Secondly, to alleviate noisy Lidar measurements and slight misalignment between stereo cameras and Lidar, we present a novel training strategy that gradually removes these noisy points during the training process automatically. Furthermore, we have also presented a novel structural loss (named plane fitting loss) to handle the inaccurate Lidar measurements and stereo matching.

Under our problem setting, we make no assumption on the inputs such as the pattern/number of Lidar measurements, the probability distribution of Lidar points or stereo disparities. Our network allows the sparsity of input Lidar points to be varied, and can even handle an extreme case when the Lidar sensor is completely unavailable.

Experimental results on different datasets demonstrate that our method is able to recover highly accurate depth maps through Lidar-Stereo fusion. It outperforms existing stereo matching methods, depth completion methods and Lidar-Stereo fusion methods with a large margin (at least twice better than previous ones). To the best of our knowledge, there is no deep learning based method available that can achieve this goal under our problem setting.

## 2. Related Work

**Stereo Matching** Deep convolutional neural networks (CNNs) based stereo matching methods have recently achieved great success. Existing supervised deep methods either formulate the task as depth regression [21] or multi-label class classifications [31]. Recently, unsupervised deep stereo matching methods have also been introduced to relief from a large amount of labeled training data. Godard *et al.* [11] proposed to exploit the photometric consistency loss between left images and the warped version of right images, thus forming an unsupervised stereo matching framework. Zhong *et al.* [36] presented a stereo matching network for estimating depths from continuous video input. Very recently, Zhang *et al.* [34] extended the self-supervised stereo network [35] from passive stereo cameras to active stereo scenarios. Even though stereo matching has been greatly

advanced, it still suffers from challenging scenarios such as texture-less and low-lighting conditions.

**Depth Completion/Interpolation** Lidar scanners can provide accurate but sparse and incomplete 3D measurements. Therefore, there is a highly desired requirement in increasing the density of Lidar scans, which is crucial for applications such as self-driving cars. Uhrig *et al.* [30] proposed a masked sparse convolution layer to handle sparse and irregular Lidar inputs. Chodosh *et al.* [4] utilized compressed sensing to approach the sparsity problem for scene depth completion. With the guidance of corresponding color images, Ma *et al.* [19] extended the up-projection blocks proposed by [17] as decoding layers to achieve full depth reconstruction. Jaritz *et al.* [14] handled sparse inputs of various densities without any additional mask input.

**Lidar-Stereo Fusion** Existing studies mainly focus on fusing stereo and time-of-flight (ToF) cameras for indoor scenarios [7, 23, 24, 12, 8], while Lidar-Stereo fusion for outdoor scenes has been seldom approached in the literature. Badino *et al.* [2] used Lidar measurements to reduce the searching space for stereo matching and provided predefined paths for dynamic programming. Later on, Maddern *et al.* [20] proposed a probabilistic model to fuse Lidar and disparities by combining prior from each sensor. However, their performance degrades significantly when the Lidar information is missing. To tackle this issue, instead of using a manually selected probabilistic model, Park *et al.* [26] utilized CNNs to learn such a model, which takes two disparities as input: one from the interpolated Lidar and the other from semi-global matching [13]. Compared with those supervised approaches, our unsupervised method can be end-to-end trained using stereo pairs and sparse Lidar points without using external stereo matching algorithms.

## 3. Lidar-Stereo Fusion

In this section, we formulate Lidar-Stereo fusion as an unsupervised learning problem and present our main ideas in dealing with the inherent challenges encountered by existing methods, *i.e.*, noise in Lidar measurements.

### 3.1. Problem Statement

Lidar-Stereo fusion aims at recovering a dense and accurate depth/disparity map from sparse Lidar measurements $S \in \mathbb{R}^{n \times 3}$ and a pair of stereo images $I_l$, $I_r$. We assume the Lidar and stereo camera have been calibrated with extrinsic matrix $T$ and the stereo camera itself is calibrated with intrinsic matrices $K_l$, $K_r$ and projection matrices $P_l$, $P_r$. We can then project sparse Lidar points $S$ onto the image plane of $I_l$ by $d_l^s = P_l T S$. Since disparity is used in the stereo matching task, we convert the projected depth $d_l^s$ to disparity $D_l^s$ using $D = Bf/d$, where $B$ is the baseline between the stereo camera pair and $f$ is the focal length. The
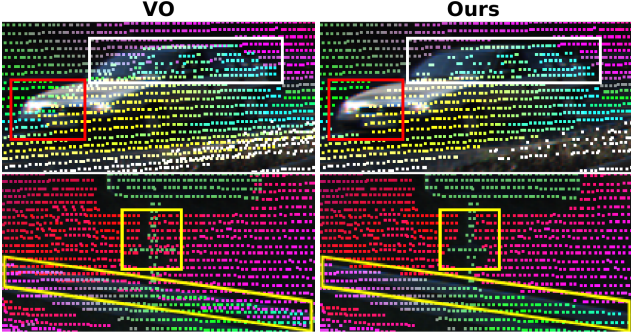
Figure 2. **KITTI VO Lidar points and our cleaned Lidar points.** Erroneous Lidar points on transparent/relective areas and tiny object surface have been successfully removed.

same process is applied to the right image as well. Mathematically this problem can be defined as:

$$(\widehat{D}_l, \widehat{D}_r) = \mathcal{F}(I_l, I_r, D_l^s, D_r^s; \Theta), \qquad (1)$$

where $\mathcal{F}$ is the learned Lidar-Stereo fusion model (a deep network in our paper) parameterized by $\Theta$, $\widehat{D}_l$, $\widehat{D}_r$ are the fusion outputs defined on the left and right coordinates.

Under our problem setting, we do not make any assumption on the Lidar points' configuration (*e.g.*, the number or the pattern) or error distribution of Lidar points and stereo disparities. Removing all these restrictions makes our method more generic and wider applicability.

### 3.2. Dealing with Noise

In Lidar-Stereo fusion, existing methods usually assume the Lidar points are noise free and the alignment between Lidar and stereo images is perfect. We argue that even for dedicated systems such as the KITTI dataset, the Lidar points are never perfect and the alignment cannot be consistently accurate, *c.f*. Fig. 1. The errors in Lidar scanning are inevitable for two reasons: (**1**) Even for well calibrated Lidar-stereo systems, *e.g.*, KITTI, and after eliminating the rolling shutter effect in Lidar scans by compensating the ego-motion, Lidar errors still persist even for stationary scenes, as shown in Fig. 2. According to the readme file in the KITTI VO dataset, the rolling shutter effect has already been removed in Lidar scans. However, we still find Lidar errors on transparent (white box) and reflective (red box) surfaces. Also, due to the displacement between Lidar and cameras, the Lidar can see through tiny objects as shown in the yellow box. (**2**) It is hard to perform motion-compensation on dynamic scenes, thus the rolling shutter effect will persist for moving objects.

It is possible to eliminate these errors by manually inserting 3D models and other post-processing steps [22]. However, lots of human efforts will be involved. Our method can automatically deal with these Lidar errors without the need of human power. Hence the problem we are tack-
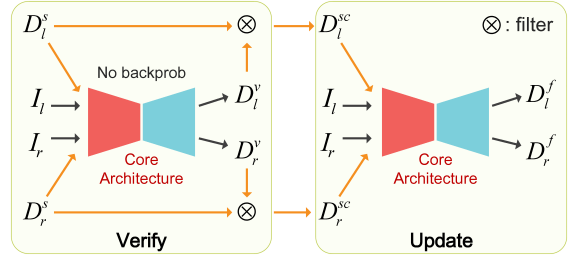


Figure 3. **The feedback loop.** For each iteration, the input stereo pair first computes initial disparities to filter errors in sparse Lidar points. At this stage, no *backprob* is taken place. So we call it the Verify phase. Then in the Update phase, the Core Architecture takes stereo pairs and cleaned sparse Lidar points as inputs to generate the final disparities. The parameters of the Core Architecture will be updated through *backprob* this time.

ling ("Noise-Aware Lidar-Stereo Fusion") is not a simple "system-level" problem that can be solved through "engineering registration".

It is known that the ability of deep CNNs to overfit or memorize the corrupted labels can lead to poor generalization performance. Therefore, we aim to deal with the noise in Lidar measurements properly to train deep CNNs.

Robust functions such the $\ell_1$ norm, Huber function or the truncated $\ell_2$ norm are natural choices in dealing with noisy measurements. However, these functions will not eliminate the effects caused by noises but only suppress them. Further, these errors also exist in the input. Automatically correct/ignore these erroneous points creates an extra difficulty for the network. To this end, we introduce a *feedback loop* in our network to allow the input also to depend on the output of the network. In this way, the input Lidar points can be cleaned before being fed into the network.

**The Feedback Loop**  We propose a novel framework to progressively detect and remove erroneous Lidar points during the training process and generate a highly accurate Lidar-Stereo fusion. Fig. 3 illustrates an unfolded structure of our network design, namely the feedback loop. It consists of two phases: "Verify" phase and "Update" phase. Each phase shares the same network structure of Core Architecture, and the details will be illustrated in Section 4.1.

In the Verify phase, the network takes stereo image pairs $(I_l, I_r)$ and noisy Lidar disparities $(D_l^s, D_r^s)$ as input, and generates two disparity maps $(D_l^v, D_r^v)$. No back-propagation takes place in this phase. We then compare $(D_l^v, D_r^v)$ and $(D_l^s, D_r^s)$ and retain the sparse Lidar points $(D_l^{sc}, D_r^{sc})$ that are consistent in both stereo matching and Lidar measurements. In the Update phase, the network takes both stereo pairs $(I_l, I_r)$ and cleaned sparse Lidar points $(D_l^{sc}, D_r^{sc})$ as the inputs to recover dense disparity maps $(D_l^f, D_r^f)$. All loss functions are evaluated on the final disparity outputs $(D_l^f, D_r^f)$ only. Once the network is
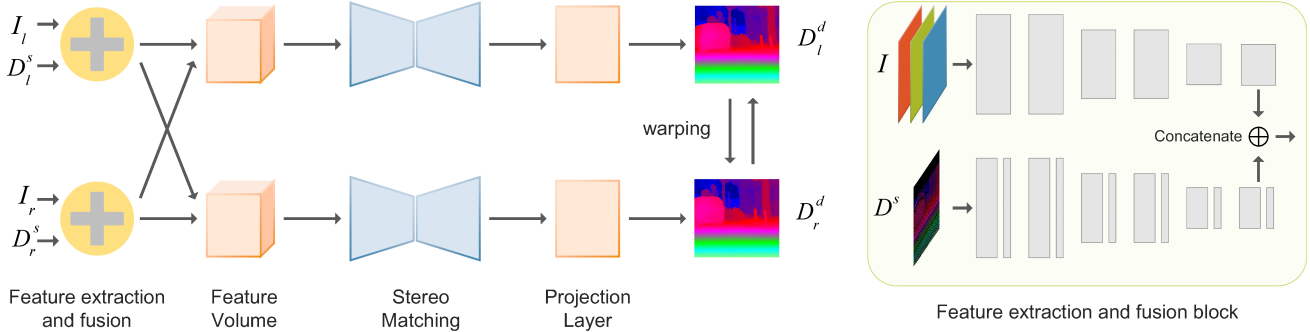
Figure 4. **Core Architecture of our LidarStereoNet.** It consists of a feature extraction and fusion block, a stack-hourglass type feature matching block and a disparity computing layer. Given a stereo pair $I_l, I_r$ and corresponding projected Lidar points $D_l^s, D_r^s$, the feature extraction block produces feature maps separately for images and Lidar points. The feature maps are then concatenated to form final input features which are aggregated to form a feature volume. The feature matching block learns the cost of feature-volume. Then we use the disparity computing layer to obtain disparity estimation. Details of the feature extraction and fusion block is illustrated on the right.

trained, we empirically find that there is no performance drop if we directly feed the Core Architecture with noisy Lidar points. Therefore, we remove the feedback loop module and only use the Core Architecture in testing.

Our feedback loop detects erroneous Lidar points by measuring the consistency between Lidar points and stereo matching. Lidar and stereo matching are active and passive depth acquisition techniques. Hence, it is less likely that they would make the same errors. It may also filter out some correct Lidar points at the first place but we have image warping loss and other regularization losses to keep the network training on the right track.

## 4. Our Network Design

In this section, we present our "LidarStereoNet" for Lidar-Stereo fusion, which can be learned in an unsupervised end-to-end manner. To remove the need of large-scale training data with ground truth, we propose to exploit the photometric consistency between stereo images, and the depth consistency between stereo cameras and Lidar. This novel network design enables the following benefits: **1**) A wide generalization ability of our framework in various real-world scenarios; **2**) Our network design allows the sparsity of input Lidar points to be varied, and can even handle the extreme case when the Lidar sensor is completely unavailable. Furthermore, to alleviate noisy Lidar measurements and the misalignment between Lidar and stereo cameras, we incorporate the "Feedback Loop" into the network design to connect the output with the input, which enables the Lidar points to be cleaned before fed into the network.

### 4.1. Core Architecture

We illustrate the detailed structure of the Core architecture of our LidarStereoNet in Fig. 4. LidarStereoNet consists of the following blocks: **1**) Feature extraction and fusion; **2**) Feature matching and **3**) Disparity computing. The

general information flow of our network is similar to [35] but has some crucial modifications in the feature extraction and fusion block. In view of different characteristics between dense colour images and sparse disparity maps, we leverage different convolution layers to extract features from each of them. For colour images, we use the same feature extraction block from [3] while for sparse Lidar inputs, the sparse invariant convolution layer [30] is used. The final feature maps are produced by concatenating stereo image features and Lidar features. Feature maps from left and right branches are concatenated to form a 4D feature volume with a maximum disparity range of 192. Then feature matching is processed through an hourglass structure of 3D convolutions to compute matching cost at each disparity level. Similar to [16], we use the soft-argmin operation to produce a 2D disparity map from the cost volume.

**Dealing with dense and sparse inputs** To extract features from sparse Lidar points, Uhrig *et al*. [30] proposed a sparsity invariant CNN after observing the failure of conventional convolutions. However, Ma *et al*. [19] and Jaritz *et al*. [14] argued that using a standard CNN with special training strategies can achieve better performance and also handle varying input densities. We compared both approaches and realized that standard CNNs can handle sparse inputs and even get better performance but they request much deeper network (ResNet38 encoded VS 5 Convolutional layers) with 500 times more trainable parameters ($13675.25K$ VS $25.87K$). Using such a "deep" network as a feature extractor will make our network not feasible for end-to-end training and hard to converge.

In our network, we choose sparsity invariant convolutional layers [30] to assemble our Lidar feature extractor which can handle varying Lidar points distribution elegantly. It consists of 5 sparse convolutional layers with a stride of 1. Each convolution has an output channel of 16 and is followed by a ReLU activation function. We attached

a plain convolution with a stride of 4 to generate the final 16 channels Lidar features in order to make sure the Lidar features compatible with the image features.

## 4.2. Loss Function

Our loss function consists of two data terms and two regularization terms. For data terms, we directly choose the image warping error $\mathcal{L}_w$ as a dense supervision for every pixel and discrepancy on filtered sparse Lidar points $\mathcal{L}_l$. For regularization terms, we use colour weighted smoothness term $\mathcal{L}_p$ and our novel slanted plane fitting loss $\mathcal{L}_p$. Our overall loss function is a weighted sum of the above loss terms:

$$\mathcal{L} = \mathcal{L}_l + \mu_1 \mathcal{L}_w + \mu_2 \mathcal{L}_s + \mu_3 \mathcal{L}_p, \quad (2)$$

we empirically set $\mu_1 = 1, \mu_2 = 0.001, \mu_3 = 0.01$.

### 4.2.1 Image Warping Loss

We assume photometric consistency between stereo pairs such that corresponding points between each pair should have similar appearance. However, in some cases, this assumption does not hold. Hence, we also compare the difference between small patches' Census transform as it is robust for photometric changes. Our image warping loss is defined as follow:

$$\mathcal{L}_w = \mathcal{L}_i + \lambda_1 \mathcal{L}_c + \lambda_2 \mathcal{L}_g, \quad (3)$$

where $\mathcal{L}_i$ stands for photometric loss, $\mathcal{L}_c$ represents Census loss and $\mathcal{L}_g$ is the image gradient loss. We set $\lambda_1 = 0.1, \lambda_2 = 1$ to balance different terms.

The photometric loss is defined as the difference between the observed left (right) image and the warped left (right) image, where we have weighted each term with the observed pixels to account for the occlusion:

$$\mathcal{L}_i = \left[ \sum_{i,j} \varphi \left( \hat{I}(i,j) - I(i,j) \right) \cdot O(i,j) \right] / \sum_{i,j} O(i,j), \quad (4)$$

where $\varphi(s) = \sqrt{s^2 + 0.001^2}$ and the occlusion mask $O$ is computed through left-right consistency check.

To further improve the robustness in evaluating the image warping error, we used the Census transformation to measure the difference:

$$\mathcal{L}_c = \left[ \sum_{i,j} \varphi \left( \hat{C}(i,j) - C(i,j) \right) \cdot O(i,j) \right] / \sum_{i,j} O(i,j). \quad (5)$$

Lastly, we have also used the difference between image gradients as an error metric:

$$\mathcal{L}_g = \left[ \sum_{i,j} \varphi \left( \nabla \hat{I}(i,j) - \nabla I(i,j) \right) \cdot O(i,j) \right] / \sum_{i,j} O(i,j). \quad (6)$$

### 4.2.2 Lidar Loss

The cleaned sparse Lidar points after our feedback verification can also be used as a sparse supervision for generating disparities. We leverage the truncated $\ell_2$ function to handle noises and errors in these sparse Lidar measurements,

$$\mathcal{L}_l = ||M(\hat{D} - D^{sc})||_\tau, \quad (7)$$

where $M$ is the mask computed in the Verify phase. The truncated $\ell_2$ fuction is defined as:

$$|| \cdot ||_\tau = \begin{cases} 0.5x^2, & |x| < \epsilon \\ 0.5\epsilon^2, & \text{otherwise.} \end{cases} \quad (8)$$

### 4.2.3 Smoothness Loss

The smoothness term in the loss function is defined as:

$$\mathcal{L}_s = \sum \left( e^{-\alpha_1 |\nabla I|} |\nabla d| + e^{-\alpha_2 |\nabla^2 I|} |\nabla^2 d| \right) / N, \quad (9)$$

where $\alpha_1 = 0.5$ and $\alpha_2 = 0.5$. Note that previous studies [11, 35] often neglect the weights $\alpha_1, \alpha_2$, which actually play a crucial role in colour weighted smoothness term.

### 4.2.4 Plane Fitting Loss

We also introduce a slanted plane model into deep learning frameworks to enforce structural constraint. This model has been commonly used in conventional Conditional Random Field (CRF) based stereo matching/optical flow algorithms. It assumes that all pixels within a superpixel lie on a 3D plane. By leveraging this piecewise plane fitting loss, we could enforce strong regularization on 3D structure. Although our slanted plane model is defined on disparity space, it has been proved that a plane in disparity space is still a plane in 3D space [29]. Mathematically, the disparity $d_p$ of each pixel $p$ is parameterized by a local plane,

$$d_p = a_p u + b_p v + c_p, \quad (10)$$

where $(u, v)$ is the image coordinate, the triplet $(a_p, b_p, c_p)$ denotes the parameters of a local disparity plane.

Define $P$ as the matrix representation of pixel's homogeneous coordinates within a SLIC superpixel [1] with a dimension of $N \times 3$ where $N$ is number of pixels within a segment, and denote $\mathbf{a}$ as the planar parameters. Given the current disparity predictions $\hat{\mathbf{d}}$, we can estimate the plane parameter in closed-form via $\mathbf{a}^* = (P^T P)^{-1} P^T \hat{\mathbf{d}}$. With the estimated plane parameter, the fitted planar disparities $\widetilde{\mathbf{d}} \in \mathbb{R}^N$ can be computed as $\widetilde{\mathbf{d}} = P \mathbf{a}^* = P(P^T P)^{-1} P^T \hat{\mathbf{d}}$.

Our plane fitting loss then can be defined as

$$\mathcal{L}_p = \|\hat{\mathbf{d}} - \widetilde{\mathbf{d}}\| = \|[I - P(P^T P)^{-1} P^T]\hat{\mathbf{d}}\|. \quad (11)$$

Table 1. **Quantitative results on the selected KITTI 141 subset.** We compare our LidarStereoNet with various state-of-the-art Lidar-Stereo fusion methods, where our proposed method outperforms all the competing methods with a wide margin.

| Methods | Input | Supervised | Abs Rel | $> 2$ px | $> 3$ px | $> 5$ px | $\delta < 1.25$ | Density |
|---|---|---|---|---|---|---|---|---|
| Input Lidar | Lidar | - | - | 0.0572 | 0.0457 | 0.0375 | - | 7.27% |
| S2D [19] | Lidar | Yes | 0.0665 | 0.0849 | 0.0659 | 0.0430 | 0.9626 | 100.00% |
| SINet [30] | Lidar | Yes | 0.0659 | 0.0908 | 0.0660 | 0.0456 | 0.9576 | 100.00% |
| Probabilistic fusion [20] | Stereo + Lidar | No | - | - | 0.0591 | - | - | 99.6% |
| CNN Fusion [26] | Stereo + Lidar | Yes | - | - | 0.0484 | - | - | 99.8% |
| Our method | Stereo | No | **0.0572** | **0.0540** | **0.0345** | **0.0220** | **0.9731** | 100.00% |
| Our method | Stereo + Lidar | No | **0.0350** | **0.0287** | **0.0198** | **0.0126** | **0.9872** | 100.00% |



(a) Input image    (b) Input lidar disparity    (c) Ground truth    (d) Ours

(e)S2D [19]    (f) SINet [30]    (g) Probabilistic fusion [20]    (h) CNN fusion [26]
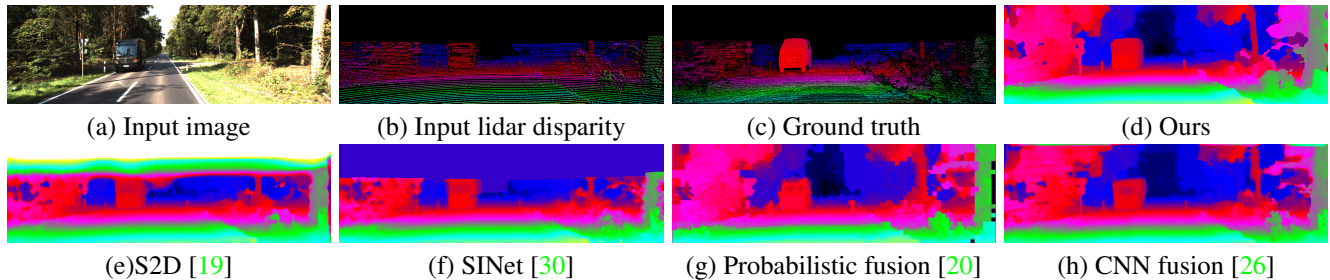
Figure 5. **Qualitative results of the methods from Tab. 1.** Our method is trained on KITTI VO dataset and tested on the selected unseen KITTI 141 subset without any finetuning.

## 5. Experiments

We implemented our LidarStereoNet in Pytorch. All input images were randomly cropped to $256 \times 512$ during training phases while we used their original size in inference. The typical processing time of our net was about 0.5 fps on Titan XP. We used the Adam optimizer with a constant learning rate of 0.001 and a batch size of 1. We performed a series of experiments to evaluate our LidarStereoNet on both real-world and synthetic datasets. In addition to analyzing the accuracy of depth prediction in comparison to previous work, we also conducted a series of ablation studies on different sensor fusing architectures and investigate how each component of the proposed losses contributes to the performance.

### 5.1. KITTI Dataset

The KITTI dataset [9] is created to set a benchmark for autonomous driving visual systems. It captures depth information from a Velodyne HDL-64E Lidar and corresponding stereo images from a moving platform. They use a highly accurate inertial measurement unit to accumulate 20 frames of raw Lidar depth data in a reference frame and serves as ground truth for the stereo matching benchmark. In KITTI 2015 [22], they also take moving objects into consideration. The dynamic objects are first removed and then re-inserted by fitting CAD models to the point clouds, resulting in a clean and dense ground truth for depth evaluation.

**Dataset Preparation** After these processes, the raw Lidar points and the ground truth differ significantly in terms of outliers and density as shown in Fig. 1. In raw data, due to the large displacement between the Lidar and the stereo cameras [29], boundaries of objects may not perfectly align when projecting Lidar points onto image planes. Also, since Lidar system scans depth in a line by line order, it will create a rolling shutter effect on the reference image, especially for a moving platform. Instead of heuristically removing measurements, our method is able to omit these outliers automatically which is evidently shown in Fig. 1 and Fig. 2.

We used the KITTI VO dataset [9] as our training set. We sorted all 22 KITTI VO sequences and found 7 frames from sequence 17 and 20 having corresponding frames in the KITTI 2015 training set. Therefore we excluded these two sequences and used the remaining 20 stereo sequences as our training dataset. Our training dataset contains 42104 images with a typical image resolution of $1241 \times 376$. To obtain sparse disparities inputs, we projected raw Lidar points onto left and right images using provided extrinsic and intrinsic parameters and converted the raw Lidar depths to disparities. Maddern *et al*. [20] also traced 141 frames from KITTI raw dataset that have corresponding frames in the KITTI 2015 dataset and reported their results on this subset. For consistency, we used the same subset to evaluate our performance and utilize the 6 frames from KITTI VO dataset as our validation set (we excluded 1 frame that overlaps the KITTI 141 subset from our validation).
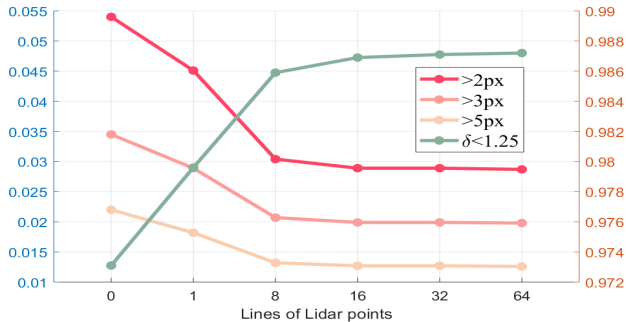
Figure 6. **Test results of our network on the selected KITTI 141 subset with varying levels of input Lidar points sparsity.** Left column: lower is better; right column: higher is better.



Figure 7. **Gradually cleaned input Lidar points.** From top to bottom, left column: left image, cleaned Lidar points at the $2^{nd}$ epoch, cleaned Lidar points at the $5^{th}$ epoch; Right column: raw Lidar points, error points find at the $2^{nd}$ epoch, error points find at the $5^{nd}$ epoch. Note that the error measurements on the right car have been gradually removed.

**Comparisons with State-of-the-Art** We compared our results with depth completion methods and Lidar-stereo fusion methods using depth metrics from [6] and bad pixel ratio disparity error from KITTI [22]. We also provide a comparison of our method and stereo matching methods in the supplemental material.

For depth completion, we compared with S2D [19] and SINet [30]. In our implementation of S2D and SINet, we trained them on KITTI depth completion dataset [30]. From 151 training sequences, we excluded 28 sequences that overlaps with KITTI 141 dataset and used the remaining 123 sequences to train these networks from scratch in a supervised manner. As a reference, we computed the error rate of the input Lidar. It is worth noting that our method increases the disparity density from less than 7.3% to 100% while reducing the error rate by a half.

We also compared our method with two existing Lidar-Stereo fusion methods: Probabilistic fusion [20] and CNN fusion [26] and outperforms them with a large margin. Quantitative comparison between our method and the competing state-of-the-art methods is reported in Tab. 1. We can clearly see that our self-supervised LidarStereoNet achieves the best performance throughout all the metrics evaluated. Note that, our method even outperforms recent supervised CNN based fusion method [26] with a large margin. More qualitative evaluations of our method in challenging scenes are provided in Fig. 11. These results demonstrate the superiority of our method that can effectively leverage the complementary information between Lidar and stereo images.

**On Input Sparsity** Thanks to the nature of deep network and sparsity invariant convolution, our LidarStereoNet can handle Lidar input of varying density, ranging from no Lidar input to 64 lines input. To see this trend, we downsampled the vertical and horizontal resolution of the Lidar points. As shown in Fig. 6, our method performs equally well when using 8 or more lines of Lidar points. Note that even when there are no Lidar points as input (in this case, the problem becomes a pure stereo matching problem), our method still
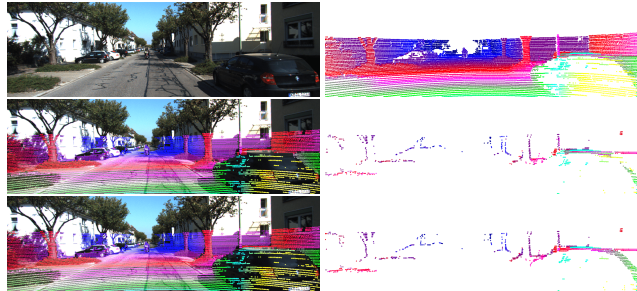
outperforms SOTA stereo matching methods.

Table 2. **Ablation study on the feedback loop** Type 1 and Type 2 show the performance only use the Core Architecture without and with removing error Lidar points from the input, while Full model means our proposed feedback loop.

| Methods | Abs Rel | > 2 px | > 3 px | > 5 px |
|---|---|---|---|---|
| Type 1 | 0.0539 | 0.0411 | 0.0310 | 0.0229 |
| Type 2 | 0.0468 | 0.0401 | 0.0302 | 0.0226 |
| Full model | **0.0350** | **0.0287** | **0.0198** | **0.0126** |

### 5.2. Ablation Study

In this section, we perform ablation studies to evaluate the importance of our feedback loop and proposed losses. Notably, all ablation studies on losses and fusion strategies are evaluated on Core Architecture only in order to reduce the randomness introduced by our feedback loop module.

**Importance of the feedback loop** We evaluate the importance of the feedback loop in two aspects. One is to remove the error points from the back-end, *i.e.* the loss computation part. The other is to remove them from the input. In our baseline model (*Type 1*), we use raw Lidar as our input and compute the Lidar loss on them. For *Type 2* model, we also use the raw Lidar as input but compute the Lidar loss only on cleaned Lidar points. Our full model uses clean Lidar points in both parts. As shown in Tab. 2, removing errors in the back-end can improve the performance by $2.58\%$. However, using cleaned Lidar points as input can boost the performance in $34.44\%$ in $> 3px$ metric, which demonstrates the importance of our feedback loop module.

**Comparing different loss functions** Tab. 3 shows the performance gain with different losses. As we can see, when only using Lidar points as supervision, its performance is affected by the outliers in Lidar measurements.

Adding a warping loss can reduce the error rate from $4.71\%$ to $3.02\%$. Adding our proposed plane fitting loss can further reduce the metric from $3.02\%$ to $2.73\%$. In the supplementary material, we further compare our soft slanted plane model and a hard plane fitting model. The soft one achieves better performance.

Table 3. **Evaluation of different loss functions.** $\mathcal{L}_w$, $\mathcal{L}_s$, $\mathcal{L}_l$ and $\mathcal{L}_p$ represent warping loss, smoothness loss, Lidar loss and plane fitting loss separately.

| Loss | Abs Rel | $> 2$ px | $> 3$ px | $> 5$ px |
|---|---|---|---|---|
| $\mathcal{L}_l$ | 0.0555 | 0.0733 | 0.0471 | 0.0296 |
| $\mathcal{L}_w + \mathcal{L}_s$ | 0.0628 | 0.0940 | 0.0637 | 0.0405 |
| $\mathcal{L}_w + \mathcal{L}_s + \mathcal{L}_l$ | 0.0565 | 0.0401 | 0.0302 | 0.0226 |
| $\mathcal{L}_w + \mathcal{L}_s + \mathcal{L}_l + \mathcal{L}_p$ | **0.0468** | **0.0393** | **0.0276** | **0.0201** |

**Comparing different fusion strategies** Considering the problem of utilizing sparse depth information, one no-fusion approach will be directly using Lidar measurements for supervisions. As shown in Tab. 4, its performance is affected by the misaligned Lidar points and it has a relatively high error rate of $4.10\%$. The second method is to leverage the depth as a fourth channel additionally to the RGB images. We term it an early fusion strategy. As shown in Tab. 4, it has the worst performance among the baselines. This may be due to the incompatible characteristics between RGB images and depth maps thus the network is unable to handle well within a common convolution layer. Our late fusion strategy achieves the best performance among them.

### 5.3. Generalizing to Other Datasets

To illustrate that our method can generalize to other datasets, we compare our method to several methods on the Synthia dataset [27]. Synthia contains 5 sequences under different scenarios. And for each scenario, they capture images under different lighting and weather conditions such as Spring, Winter, Soft-rain, Fog and Night. We show quantitative results of experiments in Tab. 5 and qualitative results are provided in the supplementary material.

For sparse disparity inputs, we randomly selected 10% of full image resolution. As discussed before, projected Lidar points have misalignment with stereo images in KITTI dataset. To simulate the similar interference, we add various density levels of Gaussian noise to sparse disparity maps. As shown in Fig. 8, our proposed LidarStereoNet adapts

Table 4. **Comparison of different fusion strategies.**

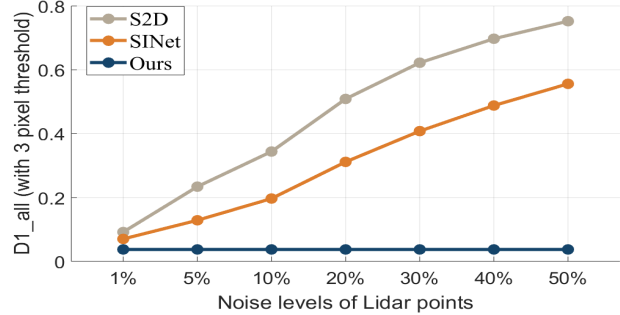| Methods | Abs Rel | $> 2$ px | $> 3$ px | $> 5$ px |
|---|---|---|---|---|
| No Fusion | 0.0555 | 0.0733 | 0.0471 | 0.0296 |
| Early fusion | 0.0644 | 0.0667 | 0.0526 | 0.0398 |
| Our method | **0.0468** | **0.0393** | **0.0276** | **0.0201** |



Figure 8. **Ablation study on noise resistance on Synthia dataset.** Our method has a consistent performance while the others have a notable performance drop.

well to the noisy input disparity maps, while S2D [19] fails to recover disparity information.

Table 5. **Quantitative results on the Synthia dataset.**

| Methods | Abs Rel | $> 2$ px | $> 3$ px | $> 5$ px |
|---|---|---|---|---|
| SPS-ST [32] | 0.0475 | 0.0980 | 0.0879 | 0.0713 |
| S2D [19] | 0.0864 | 0.5287 | 0.4414 | 0.270 |
| SINet [30] | **0.0290** | 0.0642 | 0.0472 | **0.0283** |
| Our method | 0.0334 | **0.0446** | **0.0373** | 0.0299 |

## 6. Conclusion

In this paper, we have proposed an unsupervised end-to-end learning based Lidar-Stereo fusion network "LidarStereoNet" for accurate 3D perception in real world scenarios. To effectively handle noisy Lidar points and misalignment between sensors, we presented a novel "Feedback Loop" to sort out clean measurements by comparing output stereo disparities and input Lidar points. We have also introduced a piecewise slanted plane fitting loss to enforce strong 3D structural regularization on generated disparity maps. Our LidarStereoNet does not need ground truth disparity maps for training and has good generalization capabilities. Extensive experiments demonstrate the superiority of our approach, which outperforms state-of-the-art stereo matching and depth completion methods with a large margin. Our approach can reliably work even when Lidar points are completely missing. In the future, we plan to extend our method to other depth perception and sensor fusion scenarios.

# References

[1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Susstrunk. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE Trans. Pattern Anal. Mach. Intell.*, 34(11):2274–2282, Nov. 2012. 5

[2] H. Badino, D. Huber, and T. Kanade. Integrating lidar into stereo for fast and improved disparity computation. In *International Conference on 3D Imaging, Modeling, Processing, Visualization and Transmission*, pages 405–412, May 2011. 1, 2

[3] Jia-Ren Chang and Yong-Sheng Chen. Pyramid stereo matching network. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 5410–5418, 2018. 4, 11, 13

[4] Nathaniel Chodosh, Chaoyang Wang, and Simon Lucey. Deep convolutional compressed sensing for lidar depth completion. *arXiv preprint arXiv:1803.08949*, 2018. 2

[5] Jeffrey S Deems, Thomas H Painter, and David C Finnegan. Lidar measurement of snow depth: a review. *Journal of Glaciology*, 59(215):467–479, 2013. 1

[6] David Eigen, Christian Puhrsch, and Rob Fergus. Depth map prediction from a single image using a multi-scale deep network. In *Proc. Adv. Neural Inf. Process. Syst.*, NIPS'14, pages 2366–2374, Cambridge, MA, USA, 2014. MIT Press. 7

[7] Samir El-Omari and Osama Moselhi. Integrating 3d laser scanning and photogrammetry for progress measurement of construction work. *Automation in Construction*, 18(1):1 – 9, 2008. 2

[8] V. Gandhi, J. ech, and R. Horaud. High-resolution depth maps based on tof-stereo fusion. In *IEEE International Conference on Robotics and Automation*, pages 4742–4749, May 2012. 2

[9] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2012. 6

[10] Andreas Geiger, Julius Ziegler, and Christoph Stiller. Stereoscan: Dense 3d reconstruction in real-time. In *IEEE Intelligent Vehicles Symposium*, pages 963–968. Ieee, 2011. 1

[11] Clément Godard, Oisin Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, volume 2, page 7, 2017. 2, 5

[12] Alastair Harrison and Paul Newman. Image and sparse laser fusion for dense scene reconstruction. In Andrew Howard, Karl Iagnemma, and Alonzo Kelly, editors, *Field and Service Robotics*, pages 219–228, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg. 2

[13] Heiko Hirschmuller. Stereo processing by semiglobal matching and mutual information. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30(2):328–341, Feb. 2008. 2

[14] Maximilian Jaritz, Raoul De Charette, Emilie Wirbel, Xavier Perrotton, and Fawzi Nashashibi. Sparse and dense data with cnns: Depth completion and semantic segmentation. In *International Conference on 3D Vision*, 2018. 2, 4

[15] Jonathan Kelly and Gaurav S Sukhatme. Visual-inertial sensor fusion: Localization, mapping and sensor-to-sensor self-calibration. *International Journal of Robotics Research*, 30(1):56–79, 2011. 1

[16] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. In *Proc. IEEE Int. Conf. Comp. Vis.*, Oct 2017. 4

[17] Iro Laina, Christian Rupprecht, Vasileios Belagiannis, Federico Tombari, and Nassir Navab. Deeper depth prediction with fully convolutional residual networks. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 239–248. IEEE, 2016. 2

[18] Jesse Levinson, Jake Askeland, Jan Becker, Jennifer Dolson, David Held, Soeren Kammel, J Zico Kolter, Dirk Langer, Oliver Pink, Vaughan Pratt, et al. Towards fully autonomous driving: Systems and algorithms. In *IEEE Intelligent Vehicles Symposium (IV)*, pages 163–168. IEEE, 2011. 1

[19] Fangchang Ma and Sertac Karaman. Sparse-to-dense: Depth prediction from sparse depth samples and a single image. In *IEEE International Conference on Robotics and Automation*, pages 1–8. IEEE, 2018. 1, 2, 4, 6, 7, 8, 12

[20] W. Maddern and P. Newman. Real-time probabilistic fusion of sparse 3d lidar and dense stereo. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 2181–2188, Oct 2016. 1, 2, 6, 7

[21] Nikolaus Mayer, Eddy Ilg, Philip Hausser, Philipp Fischer, Daniel Cremers, Alexey Dosovitskiy, and Thomas Brox. A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pages 4040–4048, 2016. 2, 13

[22] Moritz Menze, Christian Heipke, and Andreas Geiger. Joint 3d estimation of vehicles and scene flow. In *ISPRS Workshop on Image Sequence Analysis (ISA)*, 2015. 3, 6, 7

[23] Peyman Moghadam, Wijerupage Sardha Wijesoma, and Dong Jun Feng. Improving path planning and mapping based on stereo vision and lidar. In *International Conference on Control, Automation, Robotics and Vision*, pages 384–389. IEEE, 2008. 2

[24] Kevin Nickels, Andres Castano, and Christopher M. Cianci. Fusion of lidar and stereo range for mobile robots. *International Conference on Advanced Robotics (ICAR)*, 1:65–70, 2003. 2

[25] Florin Oniga and Sergiu Nedevschi. Processing dense stereo data using elevation maps: Road surface, traffic isle, and obstacle detection. *IEEE Transactions on Vehicular Technology*, 59(3):1172–1182, 2010. 1

[26] Kihong Park, Seungryong Kim, and Kwanghoon Sohn. High-precision depth estimation with the 3d lidar and stereo fusion. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 2156–2163. IEEE, 2018. 1, 2, 6, 7

[27] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio Lopez. The SYNTHIA Dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016. 8

[28] Daniel Scharstein and Richard Szeliski. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *Int. J. Comp. Vis.*, 47(1-3):7–42, 2002. 1

[29] Nick Schneider, Lukas Schneider, Peter Pinggera, Uwe Franke, Marc Pollefeys, and Christoph Stiller. Semantically guided depth upsampling. In *German Conference on Pattern Recognition*, pages 37–48. Springer, 2016. 5, 6

[30] Jonas Uhrig, Nick Schneider, Lukas Schneider, Uwe Franke, Thomas Brox, and Andreas Geiger. Sparsity invariant cnns. In *International Conference on 3D Vision*, 2017. 1, 2, 4, 6, 7, 8, 12

[31] Jure Žbontar and Yann LeCun. Stereo matching by training a convolutional neural network to compare image patches. *J. Mach. Learn. Res.*, 17(1):2287–2318, Jan. 2016. 2, 13

[32] Koichiro Yamaguchi, David McAllester, and Raquel Urtasun. Efficient joint segmentation, occlusion labeling, stereo and flow estimation. In *Proc. Eur. Conf. Comp. Vis.*, pages 756–771. Springer, 2014. 8, 12, 13

[33] Chi Zhang, Zhiwei Li, Yanhua Cheng, Rui Cai, Hongyang Chao, and Yong Rui. Meshstereo: A global stereo model with mesh alignment regularization for view interpolation. In *Proc. IEEE Int. Conf. Comp. Vis.*, pages 2057–2065, 2015. 1

[34] Yinda Zhang, Sameh Khamis, Christoph Rhemann, Julien Valentin, Adarsh Kowdle, Vladimir Tankovich, Michael Schoenberg, Shahram Izadi, Thomas Funkhouser, and Sean Fanello. Activestereonet: End-to-end self-supervised learning for active stereo systems. In *Proc. Eur. Conf. Comp. Vis.*, September 2018. 2

[35] Yiran Zhong, Yuchao Dai, and Hongdong Li. Self-supervised learning for stereo matching with self-improving ability. *arXiv preprint arXiv:1709.00930*, 2017. 2, 4, 5, 13

[36] Yiran Zhong, Hongdong Li, and Yuchao Dai. Open-world stereo video matching with deep rnn. In *Proc. Eur. Conf. Comp. Vis.*, September 2018. 2

# Noise-Aware Unsupervised Deep Lidar-Stereo Fusion –Supplementary Material –

## Abstract

*In this supplementary material, we provide our detailed network structure, qualitative comparison of hard and soft slanted plane constraint, qualitative and quantitative comparison to stereo matching algorithms and qualitative results on the Synthia dataset.*

## 1. Detailed Network Structure

The core architecture of our LidarStereoNet contains three blocks: 1) Feature extraction and fusion; 2) Feature matching, and 3) Disparity computing. We provide the detailed structure of the feature extraction and fusion block in Table 6. The feature matching block and disparity computing block share the same structures with PSMnet [3].

Table 6. Feature extraction and fusion block architecture, where **k**, **s**, **chns** represent the kernel size, stride and the number of the input and the output channels. We use "+" to represent feature concatenation.

| layer | k , s | chns | input |
|---|---|---|---|
| **Lidar feature extraction** | | | |
| conv_s1 | 11×11, 1 | 1/16 | disparity |
| conv_s2 | 7×7, 2 | 16/16 | conv_s1 |
| conv_s3 | 5×5, 1 | 16/16 | conv_s2 |
| conv_s4 | 3×3, 2 | 16/16 | conv_s3 |
| conv_s5 | 3×3, 1 | 16/16 | conv_s4 |
| conv_mask | 1×1, 1 | 17/16 | conv_s5+mask |
| **Stereo feature extraction** | | | |
| conv0_1 | 3×3, 2 | 3/32 | image |
| conv0_2 | 3×3, 1 | 32/32 | conv0_1 |
| conv0_3 | 3×3, 1 | 32/32 | conv0_2 |
| conv1_n | [3×3, 1 ; 3×3, 1] ×3 | 32/32 | conv0_3 |
| conv2_1 | [3×3, 2 ; 3×3, 1] | [32/64 ; 64/64] | conv1_3 |
| conv2_n | [3×3, 1 ; 3×3, 1] ×15 | 64/64 | conv2_1 |
| conv3_1 | [3×3, 1 ; 3×3, 1] | [64/128 ; 128/128] | conv2_16 |
| conv3_n | [3×3, 1 ; 3×3, 1] ×2 | 128/128 | conv3_1 |
| conv4_n | [3×3, 1 ; 3×3, 1] ×3 | 128/128 | conv3_3 |
| branch1 | 64×64, 64 | 128/32 | conv4_3 |
| branch2 | 32×32, 32 | 128/32 | conv4_3 |
| branch3 | 64×16, 16 | 128/32 | conv4_3 |
| branch4 | 8×8, 8 | 128/32 | conv4_3 |
| lastconv | [3×3, 1 ; 1×1, 1] | [320/128 ; 128/32] | conv2_16+conv4_3 +branch1+branch2 +branch3+branch4 |
| **Feature fusion** | | | |
| lastconv + conv_mask | | | |

## 2. Hard versus Soft Plane Fitting

There are two kinds of plane fitting constraints. Conventional CRF based methods use one slanted plane model to describe all disparities in one segment, *i.e.*, disparities insides one segment exactly obeys one slanted plane model. We term it as "Hard" plane fitting constraint. Our method, on the other hand, only applies this term as part of the whole optimization target. In other words, we only require the recovered disparities to fit a plane in a segment if possible but it can still be balanced by other loss terms.

Fig. 9 illustrates the difference between our soft constraint and the CRF-style hard constraint in a recovered disparity map. As can be seen in Fig. 9, strictly applying the slanted plane model in recovered disparity map decreases its performance from 3.27% to 3.97% and it is very sensitive to segments as well. By switching segments from Stereo SLIC to SLIC, its performance further decreases from 3.97% to 4.52%.
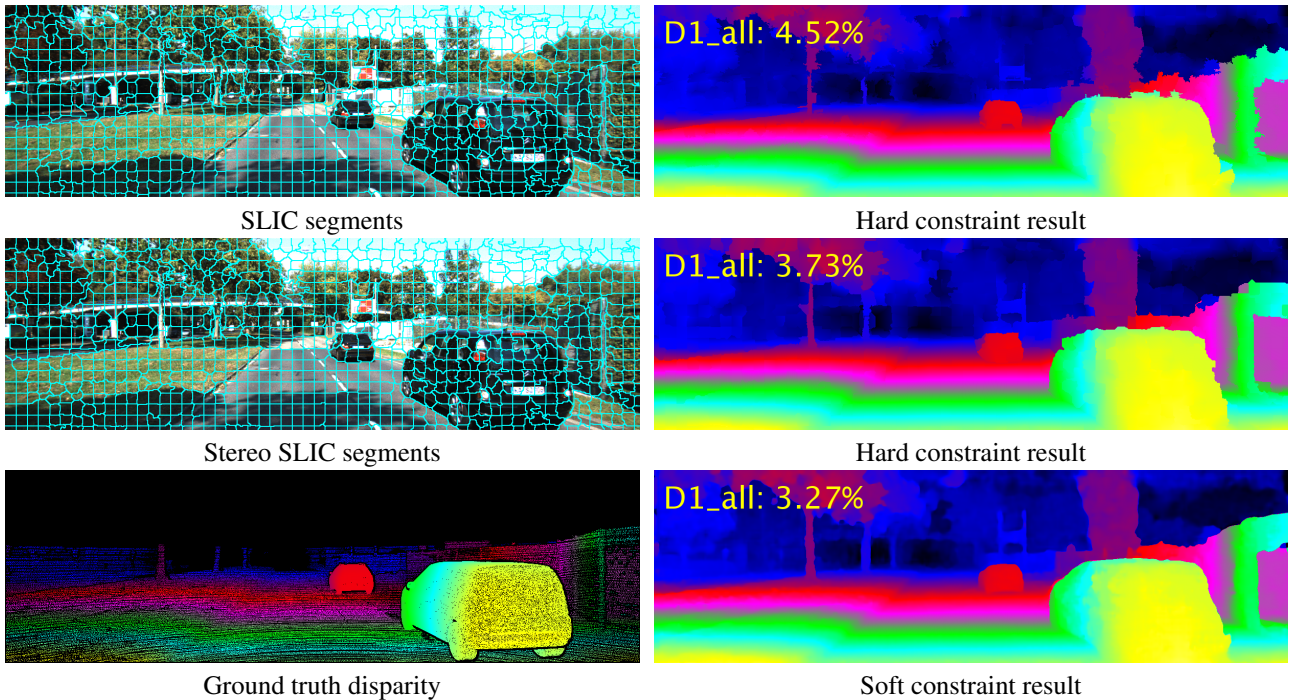


| SLIC segments | Hard constraint result |
| Stereo SLIC segments | Hard constraint result |
| Ground truth disparity | Soft constraint result |

Figure 9. **Comparison of soft and hard constraints on slanted plane model with different superpixel segmentation methods.** Note that our recovered disparity map has more aligned boundaries with the color image.
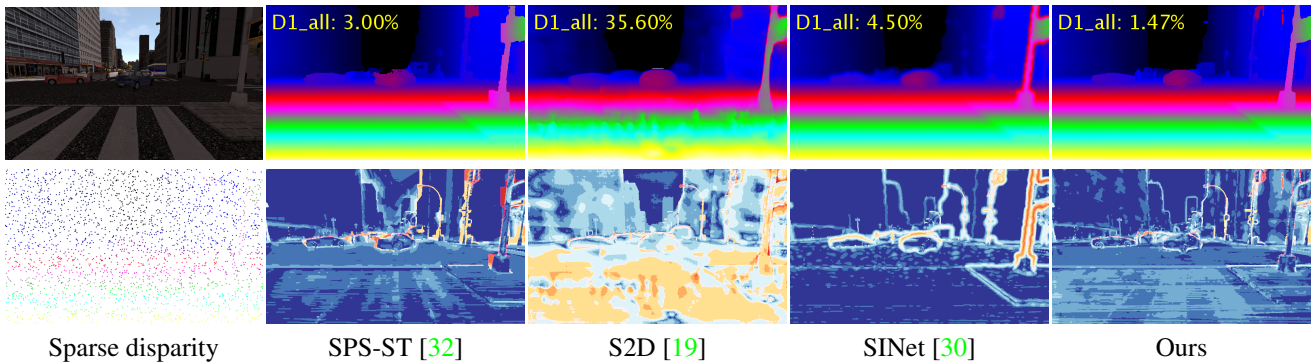


| Sparse disparity | SPS-ST [32] | S2D [19] | SINet [30] | Ours |

Figure 10. **Qualitative results on the Synthia dataset.** The first raw is the colorized disparity results, and the second row is the corresponding error maps.

# 3. Comparisons with STOA stereo matching methods

For the sake of completeness, we provide qualitative and quantitative comparisons with state-of-the-art stereo matching methods. We choose SPS-ST [32], MC-CNN[31], PSMnet [3] and SsSMnet[35] for reference. Note that the SPS-ST method is a traditional (non-deep) method, and its meta-parameters were tuned on KITTI dataset. For deep MC-CNN we used a model which was firstly trained on Middlebury dataset and for PSMnet we used the model that was trained on SceneFlow [21] dataset and the model ("-ft") that we fine-tuned on KITTI VO dataset. We also compared our method with state-of-the-art self-supervised stereo matching network SsSMnet [35].

Table 7. **Quantitative comparison on the selected KITTI 141 subset.** We compare our LidarStereoNet with various state-of-the-art stereo matching methods, where our proposed method outperforms all the competing methods with a wide margin.

| Methods | Input | Supervised | Abs Rel | $> 2$ px | $> 3$ px | $> 5$ px | $\delta < 1.25$ | Density |
|---|---|---|---|---|---|---|---|---|
| MC-CNN [31] | Stereo | Yes | 0.0798 | 0.1070 | 0.0809 | 0.0555 | 0.9472 | 100.00% |
| PSMnet [3] | Stereo | Yes | 0.0807 | 0.2480 | 0.1460 | 0.0639 | 0.9399 | 100.00% |
| PSMnet-ft [3] | Stereo | Yes | 0.0609 | 0.0635 | 0.0410 | 0.0277 | 0.9689 | 100.00% |
| SPS-ST [32] | Stereo | No | 0.0633 | 0.0702 | 0.0413 | 0.0265 | 0.9660 | 100.00% |
| SsSMnet [35] | Stereo | No | 0.0619 | 0.0743 | 0.0498 | 0.0334 | 0.9633 | 100.00% |
| Our method | Stereo | No | **0.0572** | **0.0540** | **0.0345** | **0.0220** | **0.9731** | 100.00% |
| Our method | Stereo + Lidar | No | **0.0350** | **0.0287** | **0.0198** | **0.0126** | **0.9872** | 100.00% |



(a) Input image    (b) Input lidar disparity    (c) Ground truth    (d) Ours

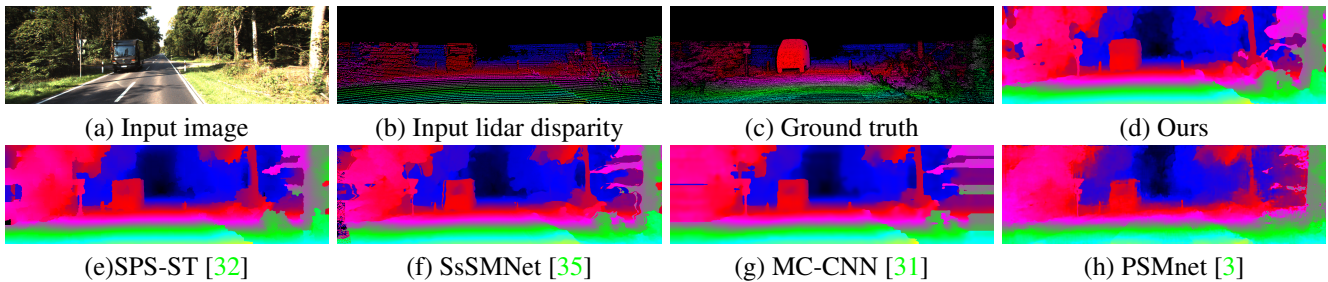(e)SPS-ST [32]    (f) SsSMNet [35]    (g) MC-CNN [31]    (h) PSMnet [3]

Figure 11. **Qualitative results of the methods from Table 7.** Our method is trained on KITTI VO dataset and tested on the selected unseen KITTI 141 subset without any finetuning.

# 4. Qualitative results on Synthia dataset

In Fig. 10, we show qualitative comparison results on Synthia dataset. Our method achieves the lowest bad pixel ratio.