

Stereo Frustums: A Siamese Pipeline for 3D Object Detection

Xi Mo · Usman Sajid · Guanghui Wang

Received: 7-19-2020 / Accepted: 10-27-2020

Abstract The paper proposes a light-weighted stereo frustums matching module for 3D objection detection. The proposed framework takes advantage of a high-performance 2D detector and a point cloud segmentation network to regress 3D bounding boxes for autonomous driving vehicles. Instead of performing traditional stereo matching to compute disparities, the module directly takes the 2D proposals from both the left and the right views as input. Based on the epipolar constraints recovered from the well-calibrated stereo cameras, we propose four matching algorithms to search for the best match for each proposal between the stereo image pairs. Each matching pair proposes a segmentation of the scene which is then fed into a 3D bounding box regression network. Results of extensive experiments on KITTI dataset demonstrate that the proposed Siamese pipeline outperforms the state-of-the-art stereo-based 3D bounding box regression methods.

Keywords stereopsis · LiDAR · stereo matching · epipolar constraint · segmentation · amodal regression

First Author: Xi Mo
Department of Electrical Engineering and Computer Science
School of Engineering, University of Kansas
Lawrence, Kansas, United States, 66045
E-mail: x618m566@ku.edu

Second Author: Usman Sajid
Department of Electrical Engineering and Computer Science
School of Engineering, University of Kansas
Lawrence, Kansas, United States, 66045
E-mail: usajid@ku.edu

Corresponding author: Dr. Guanghui Wang
Associate Professor
Department of Computer Science
Ryerson University
Toronto, ON, Canada M5B 2K3
E-mail: wangcs@ryerson.ca

1 Introduction

How to regress accurate 3D bounding boxes (bbox) for autonomous driving vehicles has become a pivotal topic recently. This technique can also benefit mobile robots and unmanned aerial vehicles with regard to scene understanding and reasoning. In this paper, we propose a Siamese pipeline method for 3D object detection.

Given a pair of stereo images and the point cloud data collected by velodyne [5], many approaches on a basis of deep-learning theories have been proposed to generate 3D bbox artifacts which can also be projected to a bird’s-eye view (BEV) of LiDAR data for localization evaluation. According to the number of image views these approaches utilized, they can be divided into three categories: monocular view [27, 16, 3, 23, 28, 11, 8, 21], binocular views [10, 7, 26, 2, 19], and non-view approaches [35, 31, 25, 13, 30, 32, 22] that only processes point cloud. Mono-view based approaches focus on sensor-fusion of the camera and LiDAR sensors in either a global or a local manner, while non-view approaches extract point cloud features from hand-crafted voxels or raw coordinates. Compared to the extensive development in both categories mentioned above, there are fewer stereo-based and stereopsis-LiDAR-fusion works for 3D object detection.

Considering the runtime of stereo matching, coarse disparity map generated by fast stereo matching and GPU acceleration achieves real-time frame-rate, yet less accurate 3D detection results [7] compared with that of coarse-to-fine disparity map [26]. However, it usually takes a few minutes to generate one panorama of coarse-to-fine disparity map before performing object detection tasks. Moreover, pixel-level stereo matching is sensitive to the error in the epipolar line calculated from camera calibration as stereo matching assumes all epipolar lines to be horizontal. We propose to reduce runtime by performing RoIs-level stereo matching instead of matching all pixels, and by a fast epipolar line searching strategy which calculates epipolar line from calibration data. We show in Section 4.4 that most of the runtime goes to point cloud processing.

Most stereo-based methods rely on stereo matching of stereopsis to generate depth maps for 3D object detection [7, 26, 2], and one Stereo R-CNN based method [10] directly regresses keypoints of 3D bbox from the left-right correspondence of regional proposals (RoIs). In stereo matching, very close objects are usually located on the border area in both views with very large disparities. Therefore, part of the same object can be missing in either view. In this case, stereo-based methods may be unable to locate matched keypoints by stereo matching. Furthermore, considering the perspective changes, the same object may appear distinctively in both views due to occlusions. These situations may cause intrinsic ambiguities in stereo matching. As shown in Fig. 1, the proposed method, by taking advantage of LiDAR-based 3D object detection methods, neither relies on stereo matching at pixel level nor predicts key points from corresponding RoIs.

To circumvent the matching ambiguities, several approaches introduce spatial constraints as 3D anchors [19] and regress point cloud proposals transformed from dense disparity maps [7, 26]. However, the design of a grid of anchors which has a proper density as well as high average precision (AP) for continuous space needs to be hand-crafted. Inspired by a novel single-frustum based method [27], we propose to solve the ambiguities by making full use of the spatial information for multi-modal regression. A novel module is proposed to directly map the point cloud onto the RoIs of the stereo image pairs and perform matching, by means of the normalized cross-correlation or the proposed 3D Intersection of Union (IoU) matching cost, and fast epipolar line search. This module correlates the 2D detection and synchronized 3D point cloud, and a novel network pipeline is proposed to accommodate this module.

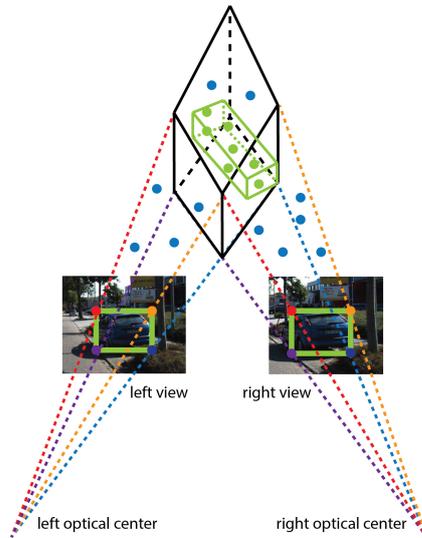


Fig. 1: Illustration of stereo frustums, valid inliers (green points) are segmented from inliers for predicting 3D bounding box (green). Given accurate 2D bboxes, the intersection of two frustums encloses the point cloud of the interested object with less ambiguity than single-frustum based and stereo-only methods, optimizing search space and being robust to perspective change.

The sparse nature of LiDAR points substantially reduces the number of points to be processed compared to the point cloud generated by the dense disparity map. Therefore, using LiDAR data can alleviate the computational burden of high-density matching at the pixel level. As to the perspective change, according to the setup of datum collecting vehicle, the baseline of stereo cameras is $0.54m$ such that both views trap the same set of 3D key-points from an object even though regions of the object appear distinctively. In addition, our method is robust to a slight disturbance on 2D bboxes (see Section 4.2).

The main contributions of this paper include:

- We propose an embedded light-weight matching module to generate 3D segmentation proposals by the RoIs from stereopsis.
- The proposed 3D IoU cost and epipolar line search algorithm are efficient in finding matches without Cython or GPU acceleration.
- The proposed Siamese architecture bridges the gap between stereopsis and real LiDAR points by integrating the point cloud segmentation network with 2D RoIs.

The proposed framework has been evaluated extensively on the KITTI dataset [4]. The experimental results outperform previous stereo-based 3D bbox regression approaches. When testing on KITTI validation set, our method maintains as high AP on car detection as F-PointNets (FPN) [16], and outperforms FPN on pedestrian detection. The proposed approach runs on average at 2-3 frames per second.

2 Related Works

2.1 Monocular Pipeline for 3D Detection

Du *et al.* [3] propose a general pipeline to detect cars from point cloud subsets constrained by monocular 2D detections. Three categories of normalized templates generalized from CAD models are fitted to 3D proposals in each subset. Each proposal is generated by RANSAC algorithm. The voxelized 3D proposals with the highest matching scores are fed to a two-stage refinement network.

Another improvement of the monocular pipeline is FPN for the purpose of regressing amodal 3D bbox including bbox sizes, orientation, and 3D bbox center. To our best knowledge, this work firstly proposes the concept of frustum which assigns bins and scores to point cloud subsets constrained by 2D detections. Instead of voxelizing the point cloud [35] before feeding it into a segmentation network followed by a T-Net to regress offset of 3D bbox center, the authors of [16] design and apply the initial feature extraction networks PointNet(v1) [17] and PointNet++(v2) [18] that learn point coordinates directly. In following sections, we denote FPN with PointNet++ backbone as FPNv2.

RoarNet [23] points out that the performance of FPN degrades if the camera sensors and velodyne sensor are not synchronized. This work proposes a geometric agreement search by selecting the best projection from a 2D detection to its 3D bbox within single frustum with the help of spatial scattering to refine the location of 3D bbox. This improvement alleviates but can not solve the ambiguity of monocular back-projection (see Fig. 1). Recently, F-ConvNet [27] ranks leading position on KITTI benchmark. It proposes a sliding-window fashion along frustum to solve the localization ambiguity. As one of the state-of-the-art monocular detectors, F-ConvNet remedies improper hand-crafted divisions by concatenating point features from all windows, and learns valid objectness by a fully-convolutional network.

2.2 Stereo-Based Methods

Several works have discussed the possibilities of 3D bbox regression by 2D detection w/o auxiliary depth information. Li *et al.* [10] propose a new end-to-end approach based on stereo R-CNN to perform regional detection, and it is also integrated to a 2D-keypoint prediction network designed for vertex estimation of 3D bbox. Nevertheless, the predicted 2D-keypoints are lack of accuracy considering perspective changes. In addition, the orientation of the object is not included in the regression artifacts. 3DOP [2] trains a structured SVM to generate 3D bbox proposals, which learns weights for an energy function that incorporates the point cloud density, free space, height prior, and height contrast information. However, without taking advantage of dense epipolar constraints in raw stereopsis, 3DOP predicts less accurate 3D bboxes than our approach.

Recently, a triangulation learning network [19] aims at learning epipolar constraints. The network requires anchors grid and 3D bbox ground truth to train. Left and right RoIs are selected by frustum-like forward-projection of 3D bboxes before cosine similarity is imposed on their feature maps. By cosine coherence scores computed from the left-right RoI pairs, the reweighting process weakens the signals from noisy channels. This method does not utilize dense raw epipolar constraints, and heavily relies on sparse anchor grids for localization. Also, this method brings in ambiguity since it searches among all potential anchors captured by a single frustum.

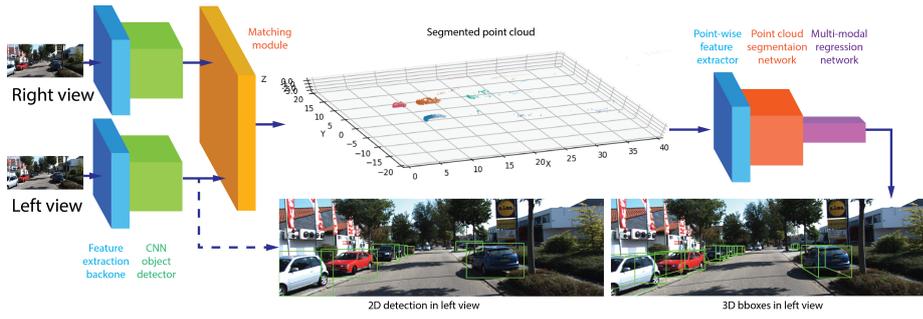


Fig. 2: Siamese pipeline for 3D object detection. Arrows signify flow directions, the dash line signifies an inspection of flow node. Feature extraction backbone and CNN object detector of sibling branches share the same weights and parameters.

RT3D [7] and Pseudo-LiDAR [26] estimate RGB-D image from stereopsis, then transform it into the point cloud, and regress 3D bbox with off-the-shelf clustering or LiDAR-based methods. RT3D presents a realtime detection scheme but comes with lower AP, while Pseudo-LiDAR predicts a less accurate depth map than real LiDAR data, which we will show in our experiments that, by the same FPN detector, our method achieves higher AP.

Another trend for object detection is through sensor fusion, i.e., to synchronize signals from multiple sensors. You *et al.* [33] proposed a scheme, namely PL++, to fuse sparse LiDAR points and the corresponding point cloud generated from the RGB-D image. As one of the top performers in 3D object detection, PL++ takes advantage of highly precise LiDAR points in localization. Zhang *et al.* [34] designed a deep network that fuses accurate depth maps, color images, and optical flow data. The model outperforms state-of-the-art works by a large margin. In order to further explore accurate RGB-D image in object detection, Tian *et al.* [24] proposed a novel representation for a 2D convolutional network that encodes depth map, multi-order depth template, and height difference map. This approach achieves realtime performance, as well as being robust to insufficient illumination and partial occlusion.

3 Stereo Frustums Pipeline

In this section, we propose a Siamese pipeline (SFPN) that takes advantage of over-constrained epipolar geometry. Section 3.1 describes the dense epipolar constraints that SFPN bases on. Section 3.2 introduces an overview of stereo frustum pipeline. Section 3.3 introduces four RoIs matching methods that our module has implemented. In this section, we assume the 2D object detections on both views are consistent with their ground truths.

3.1 Dense Epipolar Constraints

Notations. Let $F \in \mathbb{R}^{3 \times 3}$ be the fundamental matrix defined by the left-right camera coordinate systems (see Fig. 3) whose optical centers are O_2 and O_3 respectively. Let $\mathcal{S}_2 = \{l_i | i = 1, 2, \dots, m\}$ be the set of bboxes centered at l_i in the left view, $\mathcal{S}_3 = \{r_i | i = 1, 2, \dots, n\}$ be the set of bboxes centered at r_i in the right view, P_c denotes the point cloud set of scene, with $P_c(l_i)$ a subset of P_c by forward projecting P_c to the bbox region centered at l_i , N_{thres}

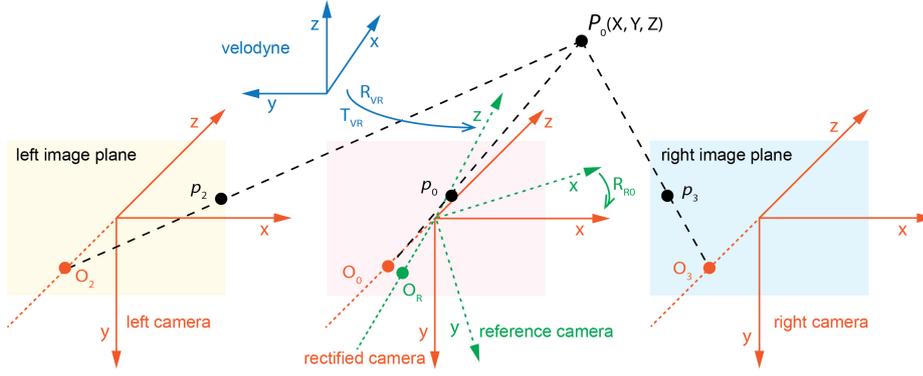


Fig. 3: Coordinate systems of the KITTI data-collection vehicle [4]. P_0 is a LiDAR point in velodyne coordinate system (blue), its forward-projections in left view and right view are denoted as p_2 and p_3 respectively. O_2 and O_3 are the optical centers (orange) of left color camera and right color camera, O_0 (orange) is the optical center of rectified camera coordinate system (gray-scale camera 0), and O_R (green) is the optical center of reference camera coordinate system. Rotation from reference camera system (green) to rectified camera system (orange) is denoted as R_{R0} .

denotes the minimum number of spacial points required for the multi-modal regression network, and $N(P_c(l_i))$ denotes the number of inliers of $P_c(l_i)$. SPFN is valid if the following two constraints are satisfied:

- (1) **One-to-one onto mapping.** $\exists \mathcal{S}'_2 \subseteq \mathcal{S}_2$ and $\exists \mathcal{S}'_3 \subseteq \mathcal{S}_3$ subject to, $\forall l_i \in \mathcal{S}'_2, \exists r_j \in \mathcal{S}'_3$ such that $r_j^T F l_i = 0$, and $\forall r_j \in \mathcal{S}'_3, \exists l_i \in \mathcal{S}'_2$ such that $r_j^T F l_i = 0$.
- (2) **Minimum intersection.** Given conditions (1), suppose N_{sec} is number of inliers in the stereo-frustums intersection as shown in Fig. 1, $\forall l_i \in \mathcal{S}'_2, \exists r_j \in \mathcal{S}'_3$ subject to $r_j^T F l_i = 0$, then $N_{sec} = N(P_c(l_i) \cap P_c(r_j)) \geq N_{thres}$.

The one-to-one onto mapping constraint finds all spacial points that are forward-projected to matched the left and right RoIs, and the minimum intersection constraint is designed for the multi-modal regression network, ensuring that the regression will not fail due to too few spacial points as the input. It is widely known that the initial step after installing mounted devices is to get calibration information for all sensors including camera sensors, velodyne sensors, and IMU sensors. Then, fine-tune all sensors according to the calibration data until an acceptable error rate is observed. For the KITTI dataset, we assume the data collection system is well-calibrated, thus calibration files are reliable enough.

A necessary constraint for stereo matching based approaches is that, the stereo search can be performed horizontally. To relax this constraint, it is important to estimate the epipolar constraints from calibration files. In Section 3.3, we employ a searching algorithm based on the estimated epipolar constraints other than the horizontal search strategy.

Fundamental matrix. In KITTI dataset, the fundamental matrix between the left and the right image planes I_2, I_3 can be calculated as:

$$F = K_{c3}^{-T} [t_{c2}^{c3}]_{\times} K_{c2}^{-1}, \quad (1)$$

where K_{c_2}, K_{c_3} are non-singular 3×3 intrinsic matrices of the left and right cameras respectively, $\mathbf{t}_{c_3}^{c_2}$ is a 3D translation vector¹ from the left to the right camera coordinate system (see Fig. 3), $[\cdot]_{\times}$ indicates the antisymmetric matrix of a vector. Using homogeneous coordinates of image points, $\forall p_3 \in I_3$ and its correspondent image point $p_2 \in I_2$ subject to $p_3^T F p_2 = 0$, by which the epipolar line through p_3 or p_2 can be deduced directly. To compute $\mathbf{t}_{c_2}^{c_3}$ of Eq. (1), according to the KITTI calibration methodology [4], we know the projection from a 3D point \tilde{P}_0 to its image \tilde{p}_0 in each image plane can be denoted as $\tilde{p}_i = [K_{ci} | \mathbf{C}_i] \tilde{P}_0 = P_{rect}^{(i)} \tilde{P}_0, i \in \{0, 1, 2, 3\}$, where \mathbf{C}_i is the last column of projection matrix $P_{rect}^{(i)}$, and i denotes the i -th camera: value ‘2’ stands for the left camera, ‘3’ for the right camera, and ‘0’ for the gray-scale camera selected as the rectified camera. Then, we have the result

$$\mathbf{t}_{c_2}^{c_3} = K_{c_2}^{-1} \mathbf{C}_2 - K_{c_3}^{-1} \mathbf{C}_3. \quad (2)$$

Proof The coordinate systems of KITTI data-collection vehicle are depicted in Fig. 3. Assume P_0 is a inhomogeneous-coordinate point in rectified camera coordinate system, its images in left and right image planes are p_2 and p_3 (homogeneous-coordinates), respectively. According to the pinhole model, we have:

$$p_3 = K_{c_3} (R_{c_0}^{c_3} P_0 + \mathbf{t}_{c_0}^{c_3}) \Rightarrow p_3 = [K_{c_3} R_{c_0}^{c_3} | K_{c_3} \mathbf{t}_{c_0}^{c_3}] P_0, \quad (3)$$

where $R_{c_0}^{c_3}$ and $\mathbf{t}_{c_0}^{c_3}$ are 3×3 rotation matrix and 3×1 translation vector from the rectified camera coordinate system to right camera coordinate system. Likewise,

$$p_2 = K_{c_2} (R_{c_0}^{c_2} P_0 + \mathbf{t}_{c_0}^{c_2}) \Rightarrow p_2 = [K_{c_2} R_{c_0}^{c_2} | K_{c_2} \mathbf{t}_{c_0}^{c_2}] P_0. \quad (4)$$

According to the projection matrix $P_{rect}^{(i)}$ mentioned above, we know $p_i = [K_{ci} | \mathbf{C}_i] P_0$. Comparing the projection matrices $[K_{c_2} | \mathbf{C}_2]$ and $[K_{c_3} | \mathbf{C}_3]$ with the projection matrices in Eq. (3) and Eq. (4), respectively, we have the following observations:

$$R_{c_0}^{c_3} = R_{c_0}^{c_2} = I \Rightarrow R_{c_2}^{c_3} = I, \mathbf{t}_{c_0}^{c_i} = K_{ci}^{-1} \mathbf{C}_i, i = 2, 3, \quad (5)$$

which indicates the rotation between the left camera and right camera can be ignored. It can be derived from Eq. (5) that

$$\mathbf{t}_{c_2}^{c_3} = \mathbf{t}_{c_0}^{c_2} - \mathbf{t}_{c_0}^{c_3} = K_{c_2}^{-1} \mathbf{C}_2 - K_{c_3}^{-1} \mathbf{C}_3, \text{ QED.} \quad (6)$$

It is necessary to prove Eq. (1) as well. Consider Eq. (3), Eq. (4) again, we have

$$K_{c_3}^{-1} p_3 = R_{c_0}^{c_3} P_0 + \mathbf{t}_{c_0}^{c_3}, K_{c_2}^{-1} p_2 = R_{c_0}^{c_2} P_0 + \mathbf{t}_{c_0}^{c_2}, \quad (7)$$

and by using the observation $R_{c_0}^{c_3} = R_{c_0}^{c_2} = I$ from Eq. (5) and eliminating P_0 , we have

$$K_{c_3}^{-1} p_3 - \mathbf{t}_{c_0}^{c_3} = K_{c_2}^{-1} p_2 - \mathbf{t}_{c_0}^{c_2} \Rightarrow K_{c_2}^{-1} p_2 - K_{c_3}^{-1} p_3 = \mathbf{t}_{c_2}^{c_3}. \quad (8)$$

Premultiply $[\mathbf{t}_{c_2}^{c_3}]_{\times}$ to both sides of Eq. (8):

$$[\mathbf{t}_{c_2}^{c_3}]_{\times} K_{c_2}^{-1} p_2 = [\mathbf{t}_{c_2}^{c_3}]_{\times} K_{c_3}^{-1} p_3, \quad (9)$$

¹ Each pair of subscript and superscript indicates the from-to relation, the same notation applies to other translation vectors and rotation matrices in this paper.

by premultiplying $(K_{c3}^{-1} p_3)^T$ to both sides of Eq. (9), we have

$$(K_{c3}^{-1} p_3)^T [t_{c2}^{c3}]_{\times} K_{c2}^{-1} p_2 = 0. \quad (10)$$

which can be rearranged as

$$p_3^T (K_{c3}^{-T} [t_{c2}^{c3}]_{\times} K_{c2}^{-1}) p_2 = 0. \quad (11)$$

It is obvious that $F = K_{c3}^{-T} [t_{c2}^{c3}]_{\times} K_{c2}^{-1}$, and t_{c2}^{c3} can be calculated by Eq. (6), QED.

3.2 Pipeline For Stereo Frustums

Each matched RoIs pair from left-right views encodes dense epipolar constraints between pixels and their corresponding 3D spacial points. Though fewer LiDAR points are forward-projected to front-view, we observe that $\sim 10,000$ points are captured by most bboxes. Among these points, there are many outliers that do not belong to the interested objects. As pointed out by Zhou *et al.* [35] that it enforces high computational cost to FPN, our matching module effectively filters out 15% to 49% of points by dense epipolar constraints and thus can greatly speed up the networks that process point cloud.

The proposed Siamese pipeline is shown in Fig. 2. The stereopsis is fed to a 2D object detector. The detector can be any of traditional ones, such as the CNN-based detectors as Faster R-CNN [20], RetinaNet [12], the leading PC-CNN-V2 [3] on KITTI object detection benchmark, and others [14, 9]. An accurate detector is crucial before populating the RoIs pairs into the matching module. Theoretically, each RoI in the left view should have its counterpart in the right view since one-to-one onto mapping constraint will be inconsistent otherwise. Also, If one salient object fails to be detected simultaneously in both views, it would lower both detection precision and recall. Note that only the objects in the left view have their ground truths, thus we focus on matching RoIs of the right view to the left RoIs. The matching module follows the one-to-one mapping constraint or its relaxed form and takes all right RoIs as candidates. The goal is to find all matches to objects in the left view by the matching costs, which is illustrated in Section 3.3.

In the light of valid pairs proposed by the matching module and forward-projection matrices, the scene is segmented into 3D RoIs subject to minimum intersection constraint, whose semantic information and objectness scores are inherited from the 2D detections. A segmentation network refines the 3D RoIs with the extracted point-wise features. Furthermore, multi-modal regressor estimates the orientation, sizes of 3D bbox, and localization. To verify the validity of SFPN, we employ Faster R-CNN and Mask R-CNN [6] as the 2D detectors, FPNv2 as the point cloud segmentation and multi-modal regression networks.

3.3 Matching Algorithms

Notations. Consider the centers of bboxes as potential matches, let D_{thres} be the threshold of distance to the epipolar line, P_{thres} the threshold of regional similarity, P_{l_i} is the RoI image patch centered at l_i , $Prob_{ij} = NCC(P_{l_i}, P_{r_j})$ denotes the normalized cross-correlation operator, $dist(r_j, e_i)$ is the operator to compute the distance from r_j to the epipolar line e_i of l_i , $IoU_{ij} = N(P_c(l_i) \cap P_c(r_j)) / N(P_c(l_i) \cup P_c(r_j))$ as 3D IoU matching cost, $P_{3d.thres}$ the probability threshold of the 3D IoU cost. We then formulate four matching algorithms as:

- (1) RoIs matching by 3D IoU cost and epipolar line search (3DCES, Alg. 1), which relaxes the one-to-one onto mapping constraint to left view only for alleviating computational burden.
- (2) RoIs matching by 3D IoU cost and method of exhaustion (3DCME, Alg. 2), which relaxes the one-to-one onto mapping constraint.
- (3) RoIs matching by regional similarity cost (RSC, Alg. 3), which relaxes the one-to-one onto mapping constraint.
- (4) RoIs matching by regional similarity cost and left-right consistency check (RSCCC, Alg. 4), which strictly follows the one-to-one onto mapping constraint.

The epipolar line constraint facilitates searching for a fast pixel-level matching. As for RoI-level matching, we adapt this approach for faster search to RSC, RSCCC, and 3DCES. Experiments on both regional similarity and 3D IoU costs show the effectiveness of epipolar line search. The matching module is designed using Python without any Cython and GPU acceleration, or paralleled programming. On account of the compatibility purpose, 3DCME and 3DCES merely return matches as RSC and RSCCC do. We will show in the experiments that SFPN directly processes raw point cloud with high efficiency, and RSC achieves the shortest runtime among the four modules. As previous works done by Li *et al.* [10] and Qin *et al.* [19] have designed learning-based methods to find RoIs matches as RSC and RSCCC, we are seeking the possibilities to design learning-based 3DCME and 3DCES.

In order to deploy the most appropriate module for different tasks, it is recommended to inspect the ground-truth of the labeled dataset. In terms of a stereo camera configuration, if the ground-truths of both views are provided, then RSCCC may be implemented for the reason that 2D detection in both views can be described as ‘reliable’, which reduces unreliable matches by left-right consistency check while maintains high efficiency in RoIs matching. If the ground-truth of one view is unavailable, the view without ground-truth is thereby unreliable. In this case, although implementation of RSC is the fastest in runtime, 3DCES may be the best choice as a trade-off of runtime and detection performance. Moreover, 3DCME is better implemented in case of a very sparse LiDAR scene for its ability to find accurate matches of frustums. Reader may refer to Section 4.2 for comparisons on the performance of the proposed matching modules and Section 4.4 for details on the runtime analysis.

Algorithm 1 3DCES

Input:

$\mathcal{S}_2, \mathcal{S}_3, P_{thres}, D_{thres}, P_c$;

Steps:

- 1: Compute set of epipolar lines $\mathcal{E} = \{Fl_i | \forall l_i \in \mathcal{S}_2\}$;
 - 2: $\forall l_i \in \mathcal{S}_2$, search leftwards in right view along $e_i \in \mathcal{E}$, calculate $\mathcal{R}_i = \{r_j | \forall r_j \in \mathcal{S}_3, dist(r_j, e_i) < D_{thres}\}$;
 - 3: Compute costs $\mathcal{C}_i = \{IoU_{ij} | \forall r_j \in \mathcal{R}_i, \exists IoU_{ij} \neq 0\}$;
 - 4: Find match. If $\mathcal{C}_i \neq \emptyset$ and $\exists r_j \in \mathcal{R}_i$ s.t. $IoU_{ij} \geq \max(P_{3d_thres}, \max(\mathcal{C}_i))$, let $m_i = (l_i, r_i)$;
 - 5: **Output:** \mathcal{M} ;
-

4 Experiments

In this section, we present both validation and testing results on publicly available KITTI dataset. The dataset consists of 7,481 annotated (2D/3D bbox labeled in left view only)

Algorithm 2 3DCME**Input:** $\mathcal{S}_2, \mathcal{S}_3, P_{3d_thres}, P_c;$ **Steps:**

- 1: $\forall l_i \in \mathcal{S}_2$, let $\mathcal{C}_i = \{IoU_{ij} | \forall r_j \in \mathcal{S}_3, \exists IoU_{ij} \neq 0\}$;
- 2: Find match. If $\mathcal{C}_i \neq \emptyset$ and $\exists r_j \in \mathcal{S}_3$ s.t. $IoU_{ij} \geq \max(P_{3d_thres}, \max(\mathcal{C}_i))$, let $m_i = (l_i, r_i)$;
- 3: **Output:** \mathcal{M} ;

Algorithm 3 RSC**Input:** $\mathcal{S}_2, \mathcal{S}_3, P_{thres}, D_{thres};$ **Steps:**

- 1: Compute set of epipolar lines $\mathcal{E} = \{Fl_i | \forall l_i \in \mathcal{S}_2\}$;
- 2: $\forall l_i \in \mathcal{S}_2$, search leftwards in right view along $e_i \in \mathcal{E}$, calculate $\mathcal{R}_i = \{r_j | \forall r_j \in \mathcal{S}_3, \text{dist}(r_j, e_i) < D_{thres}\}$;
- 3: Perform RoI alignment to P_{l_i} and P_{r_j} , compute set of costs $\mathcal{C}_i = \{Prob_{ij} | \forall r_j \in \mathcal{R}_i, \mathcal{R}_i \neq \emptyset\}$;
- 4: Find match. If $\mathcal{C}_i \neq \emptyset$ and $\exists r_j \in \mathcal{R}_i$ s.t. $Prob_{ij} \geq \max(P_{thres}, \max(\mathcal{C}_i))$, let $m_i = (l_i, r_i)$;
- 5: **Output:** \mathcal{M} ;

Algorithm 4 RSCCC**Input:** $\mathcal{S}_2, \mathcal{S}_3, P_{thres}, D_{thres};$ Matches of left-to-right RoIs pairs, $\mathcal{M}_l = \emptyset$;Matches of right-to-left RoIs pairs, $\mathcal{M}_r = \emptyset$;**Steps:**

- 1: Compute sets of epipolar lines
 $\mathcal{E}_r = \{Fl_i | \forall l_i \in \mathcal{S}_2\}$ and $\mathcal{E}_l = \{r_l^T F | \forall r_l \in \mathcal{S}_3\}$;
- 2: $\forall l_i \in \mathcal{S}_2$, search leftwards in right view along $e_i \in \mathcal{E}_r$, calculate $\mathcal{R}_i = \{r_j | \forall r_j \in \mathcal{S}_3, \text{dist}(r_j, e_i) < D_{thres}\}$,
 $\forall r_l \in \mathcal{S}_3$, search rightwards in left view along $e_l \in \mathcal{E}_l$, calculate $\mathcal{L}_i = \{l_j | \forall l_j \in \mathcal{S}_2, \text{dist}(l_j, e_l) < D_{thres}\}$;
- 3: Perform RoI alignment to P_{l_i} and P_{r_j} , compute two sets of costs $\mathcal{C}_{li} = \{Prob_{ij} | \forall r_j \in \mathcal{R}_i, \mathcal{R}_i \neq \emptyset\}$ and
 $\mathcal{C}_{ri} = \{Prob_{ji} | \forall l_j \in \mathcal{L}_i, \mathcal{L}_i \neq \emptyset\}$;
- 4: Find matches from either view. If $\mathcal{C}_{li} \neq \emptyset$ and $\exists r_j \in \mathcal{R}_i$ s.t. $Prob_{ij} \geq \max(P_{thres}, \max(\mathcal{C}_{li}))$, let $\mathcal{M}_l = \mathcal{M}_l \cup \{(l_i, r_i)\}$. If $\mathcal{C}_{ri} \neq \emptyset$ and $\exists l_j \in \mathcal{L}_i$ s.t. $Prob_{ji} \geq \max(P_{thres}, \max(\mathcal{C}_{ri}))$, let $\mathcal{M}_r = \mathcal{M}_r \cup \{(l_j, r_i)\}$;
- 5: **Output:** $\mathcal{M} = \mathcal{M}_l \cap \mathcal{M}_r$;

image pairs with corresponding point clouds and calibrations, 7518 unlabeled testing image pairs with corresponding point clouds and calibrations. In our experiments, 20% of 7,481 randomly shuffled training samples are divided as the validation set, 2 categories of objects - ‘Car’ and ‘Pedestrian’ are examined for 3D/BEV detections. Also, we compare proposed SPFN with the state-of-the-arts including Pseudo-LiDAR [26], Stereo R-CNN [10], TLNet [19], et. al.

4.1 Experiments Setup

The 2D detector Faster R-CNN [20] (VGG-16 backbone) is trained on KITTI training set to provide shared weights and parameters for the sibling branches as shown in Fig. 2, and Mask R-CNN [6] (ResNet backbone + Feature Pyramid Network) is pre-trained on COCO dataset. We train FPNv2 with raw LiDAR data other than segmented point cloud from the matching module. In this way, each module pairs with FPNv2 as 3DCES-FPNv2, 3DCME-FPNv2, RSC-FPNv2, RSCCC-FPNv2. Before testing, several parameters need to be manually set: the percentage of bbox size $S_{enlarge}$ to be enlarged, distance threshold D_{thres} from the bbox

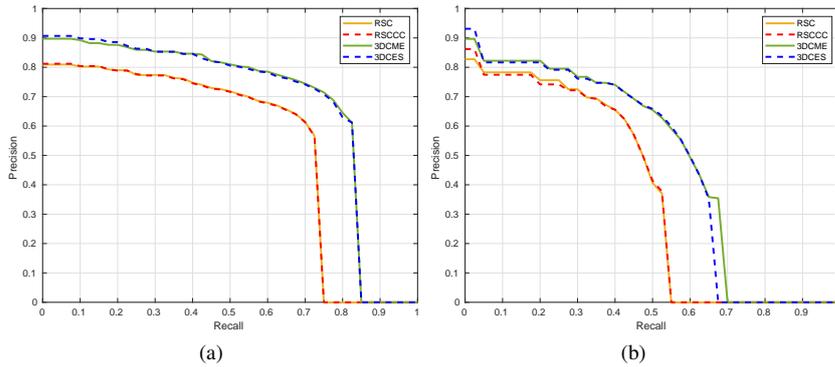


Fig. 4: P-R curves of 3D object detection results. **Moderate** level is presented. (a) Car category at IoU threshold 0.7. (b) Pedestrian category at IoU threshold 0.5.

center to the epipolar line, similarity threshold P_{thres} for the template matching, the minimum number of inliers of stereo-frustums intersection N_{thres} and P_{3d_thres} the 3D IoU threshold. We set $S_{enlarge}$ to 8%, N_{thres} to 5, $D_{thres}=30px$, $P_{thres}=0.4$, and $P_{3d_thres}=0.5$. In this section, we enforce enlargement to all bboxes unless specified by ‘NE’ (w/o enlargement).

2D Detector	AP _{2D}	AP _{BEV}	AP _{3D}
Faster R-CNN(NE)	65.9	56.6	55.5
Mask R-CNN(NE)	85.8	74.2	73.5
Faster R-CNN	65.9	57.7	56.7
Mask R-CNN	85.8	83.0	82.2

Table 1: 3D object detection on car category using two 2D detectors, APs(%) at **Moderate** level with IoU=0.5. Enlarged bboxes of both views by 8% on height and width. Only 3DCES-FPNv2 is tested.

It should be noted that the reported 2D detection results come from 2D detection stage other than the projection of the predicted 3D bounding box to the front view. Tab. 1 shows that when the 2D detection AP is low, enlarging bboxes slightly increases AP_{BEV} and AP_{3D} since the original bboxes are not precise enough to capture all keypoints. On the other hand, in terms of SFPN’s serial pipeline structure, higher 2D detection AP may result in better 3D detection. We also observe that, if setting P_{3d_thres} to a larger threshold, the number of total matches drops greatly and vice versa, which behaves similar to 2D IoU during the detection phase.

4.2 Quantitative Evaluation on Validation Set

The results of BEV and 3D object detection are presented in Tab. 2, Tab. 3, Tab. 4, and Tab. 5. Both RSC and RSCCC require RoI alignment during matching. RSC is designed to relax the one-to-one onto mapping constraint by merely finding matches to RoIs in the left view while

RSCCC strictly follows this constraint. As shown in the detection results, RSCCC-FPNv2 achieves better performance than RSC-FPNv2 with IoU=0.7, which provides evidence for the effectiveness of this constraint. 3DCME depends on the proposed 3D IoU without RoI alignment or epipolar line search, which achieves better AP than 3DCES with epipolar line search. We will show in Section 4.4 that RSC runs faster than the other three matching methods, and 3DCME is the slowest but with the highest AP.

Overview on comparisons 3DCES-FPNv2 and 3DCME-FPNv2 achieve competitive performance among all state-of-the-art stereo-based methods because of highly precise LiDAR data. Although almost all listed stereo-based methods yield higher 2D detection AP (IoU=0.7) on either the validation set or the test set, 3DCES-FPNv2 and 3DCME-FPNv2 regress disproportional accurate 3D bboxes when projected to BEV. It is interesting that RSCCC-FPNv2 and RSC-FPNv2 outperform their learning-based versions (without LiDAR data) Stereo-RCNN (~10%) and TLNet(~35%) by a large margin.

Method	Type	IoU=0.5			
		AP _{2D}	Easy	Mode	Hard
RSC-FPNv2(ours)	Stereo+LiDAR	85.8	71.9	68.4	61.5
RSCCC-FPNv2(ours)	Stereo+LiDAR	85.8	71.9	68.2	61.4
3DCME-FPNv2(ours)	Stereo+LiDAR	85.8	82.9	82.9	75.7
3DCES-FPNv2(ours)	Stereo+LiDAR	85.8	83.1	83.0	75.7
PSMNET-AVOD [26]	Stereo	-	89.0	77.5	68.7
Stereo R-CNN [10]	Stereo	-	87.1	74.1	58.9
3DOP [2]	Stereo	-	55.0	41.3	34.6
TLNet [19]	Stereo	-	62.5	46.0	41.9

Table 2: BEV results on validation set. Car category is evaluated with IoU=0.5 and all APs in ‘%’. **Moderate** AP_{2D} on validation set is reported.

Method	Type	IoU=0.7			
		AP _{2D}	Easy	Mode	Hard
RSC-FPNv2(ours)	Stereo+LiDAR	54.2	67.1	58.0	51.3
RSCCC-FPNv2(ours)	Stereo+LiDAR	54.2	68.4	58.6	51.8
3DCME-FPNv2(ours)	Stereo+LiDAR	54.2	79.8	72.9	65.7
3DCES-FPNv2(ours)	Stereo+LiDAR	54.2	78.4	72.2	65.1
PSMNET-AVOD [26]	Stereo	-	74.9	56.8	49.0
Stereo R-CNN [10]	Stereo	88.3	68.5	48.3	41.5
3DOP [2]	Stereo	-	55.0	9.5	7.6
TLNet [19]	Stereo	-	29.22	21.9	18.8

Table 3: BEV results on validation set. Car category is evaluated with IoU=0.7 and all APs in ‘%’. **Moderate** AP_{2D} on validation set is reported.

Comparison with Pseudo-LiDAR Tab. 6 shows the comparison with a recently published Pseudo-LiDAR [26] method. Compared with FPNv2 on validation set and the car category, 3DCES-FPNv2 and 3DCME-FPNv2 achieve competitive results although they suffer from lower recall (see Fig. 4a) considering mismatches in both views. In order to explore potentials of proposed method, we specifically select PSMNET-AVOD that achieves the best

Method	Type	IoU=0.5			
		AP _{2D}	Easy	Mode	Hard
RSC-FPNv2(ours)	Stereo+LiDAR	85.8	70.7	67.1	54.0
RSCCC-FPNv2(ours)	Stereo+LiDAR	85.8	70.4	66.9	53.9
3DCME-FPNv2(ours)	Stereo+LiDAR	85.8	82.5	82.3	75.0
3DCES-FPNv2(ours)	Stereo+LiDAR	85.8	82.4	82.2	74.9
PSMNET-AVOD [26]	Stereo	-	88.5	76.4	61.2
Stereo R-CNN [10]	Stereo	-	85.8	66.3	57.2
3DOP [2]	Stereo	-	46.0	34.6	30.1
TLNet [19]	Stereo	-	59.5	43.7	38.0

Table 4: 3D detection results on validation set. Car category is evaluated with IoU=0.5. 3DCME-FPNv2 achieves the best performance on validation set out of stereo-based methods.

Method	Type	IoU=0.7			
		AP _{2D}	Easy	Mode	Hard
RSC-FPNv2(ours)	Stereo+LiDAR	54.2	58.5	53.7	47.4
RSCCC-FPNv2(ours)	Stereo+LiDAR	54.2	58.7	53.9	47.5
3DCME-FPNv2(ours)	Stereo+LiDAR	54.2	73.0	67.3	61.2
3DCES-FPNv2(ours)	Stereo+LiDAR	54.2	72.5	66.7	60.9
PSMNET-AVOD [26]	Stereo	-	61.9	45.3	39.0
Stereo R-CNN [10]	Stereo	88.3	54.1	36.7	31.1
3DOP [2]	Stereo	-	6.6	5.1	4.1
TLNet [19]	Stereo	-	18.2	14.3	13.7

Table 5: 3D detection results on validation set. Car category is evaluated with IoU=0.7. 3DCME-FPNv2 achieves the best performance on validation set out of stereo-based methods.

performance among the methods proposed in [26]. As mentioned in Section 1, dense disparity map generates a coarse estimation of the scene which is sensitive to the focal length and baseline of the stereo cameras. According to Pseudo-LiDAR [26], its point cloud detector AVOD and FPN are trained by fine-grained, sparse LiDAR points, therefore its performance may not compete with the state-of-the-art methods based on real-LiDAR data.

Method	AP _{2D}	Easy	Mode	Hard
3DCME-FPNv2(ours)	57.4	57.6	47.1	40.7
3DCES-FPNv2(ours)	57.4	57.8	47.3	40.8
FPNv2(our results)	57.4	55.6	49.1	42.6
PSMNET-FPN [26]	-	33.8	27.4	24.0
FPN [26]	-	64.7	56.5	49.9

Table 6: 3D detection results on validation set. Pedestrian category is evaluated.

Tab. 7 shows the BEV results compared to FPNv2 on validation set, along with another similar comparison with [26]. 3DCME-FPNv2 achieves better performance than FPNv2 by 4.1% at Easy level, 2.2% at Moderate level, and 1.8% at Hard level. The improvements are mainly derived from enlarged bboxes since they enrich sparse segmentation captured by stereo frustums. These enlarged bboxes are especially efficient for pedestrians whose poses vary. For 3D object detection, 3DCES-FPNv2 achieves higher AP at Easy level, but is less

Method	AP _{2D}	Easy	Mode	Hard
3DCME-FPNv2(ours)	57.4	62.6	53.7	46.6
3DCES-FPNv2(ours)	57.4	60.8	52.8	45.9
FPNv2(our results)	57.4	58.5	51.5	44.8
PSMNET-FPN [26]	-	41.3	34.9	30.1
FPN [26]	-	69.7	60.6	53.4

Table 7: BEV detection results on validation set. Pedestrian category is evaluated. Results of FPN in the table are presented in [27].

precise than FPNv2 at Moderate and Hard levels. To conclude, when comparing margins to corresponding FPN detection results in Tab. 7 and Tab. 6, our methods outperform pseudo-LiDAR based PSMNET-FPN by significant margins.

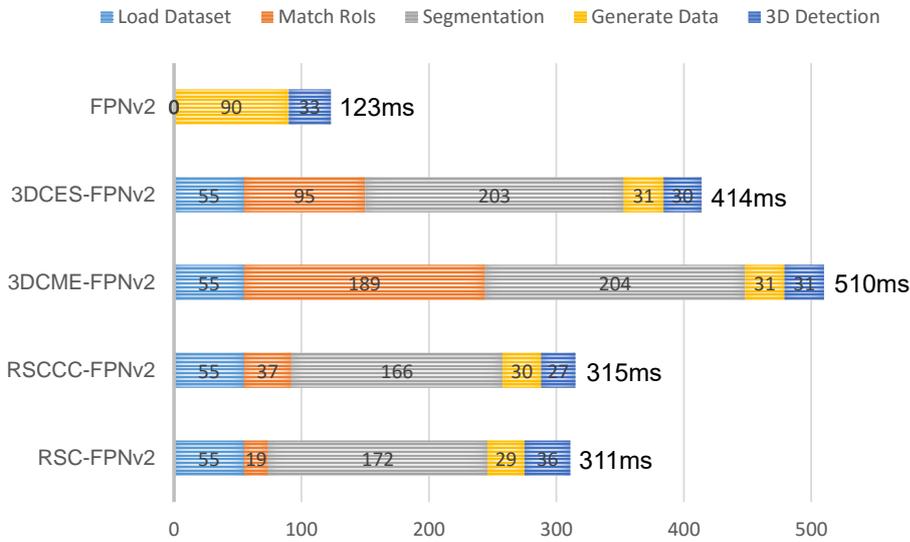


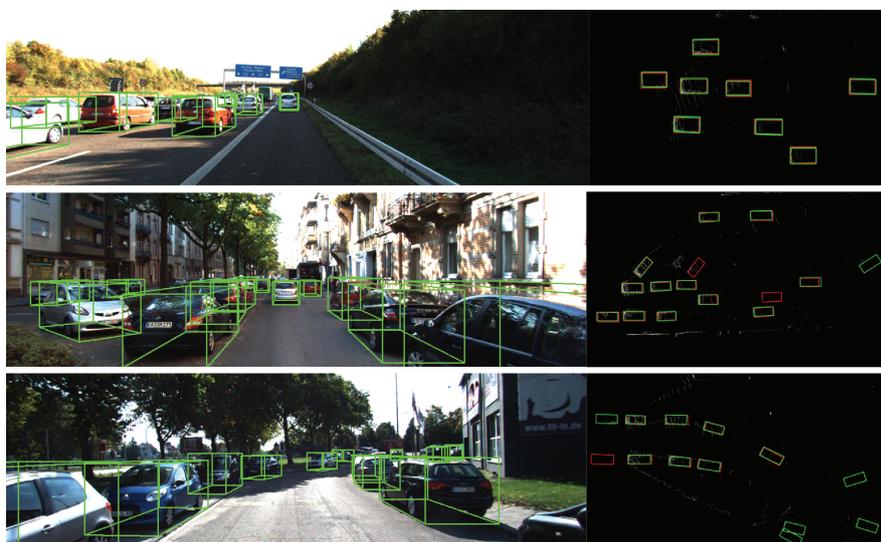
Fig. 5: Runtime comparison between matching modules and original FPNv2, best viewed in colour mode. Processing time (millionsecond) is calculated by averaging runtime of validating 1,496 samples.

4.3 Quantitative Evaluation on Test Set

In this subsection, we present the performance of 3DCES-FPNv2 on KITTI leaderboard [5] among all stereo-based submissions. We take Mask R-CNN [6] as the 2D detector as it has a higher detection AP over Faster R-CNN (see Section 4.1). FPNv2 is trained on our training set other than all 7,481 samples. Tab. 8 and Tab. 9 present 3D object detection results, Tab. 10 and Tab. 11 present BEV detection results.



(a) 3D(left) and BEV(right) detections. More persons than cars.



(b) 3D(left) and BEV(right) detections. Car category only.

Fig. 6: Qualitative evaluation on validation set using Mask R-CNN detector. Red boxes in BEV indicate ground truths, and green boxes indicate detections. Cyclist category is not evaluated in our experiments. However, Mask R-CNN [6] misclassifies cyclists as pedestrians, therefore significantly increases false positive rate in detecting pedestrians.

It can be observed that from Tab. 9 and Tab. 11, proposed method outperforms stereo-based method in detecting pedestrians. However, as there are more samples in the cyclist category than that of the validation set, while Mask R-CNN mistakes the cyclist category as the pedestrian category, we believe this is one of the main factors that lower the detection APs. In terms of the performance in 3D object detection, 3DCES-FPNv2 achieves competitive performance among stereo-based methods. It should be noted that untrained Mask R-CNN provides much lower 2D detection AP at 46.68% which is 38.47% less than 85.15% of PL++(SDN+GDC) [33]. Nonetheless, the proposed method maintains disproportionate performance in 3D object detection, which proves the effectiveness of the proposed pipeline.

Method	Type	AP _{2D}	Easy	Mode	Hard
3DCES-FPNv2(ours)	Stereo+LiDAR	46.68	58.88	51.92	44.59
Pseudo-LiDAR [26]	Stereo	67.79	54.53	34.05	28.25
Pseudo-LiDAR++ [33]	Stereo	82.90	61.11	42.43	36.99
PL++(SDN+GDC) [33]	Stereo+LiDAR	85.15	68.38	54.88	49.16
ZoomNet [29]	Stereo	83.92	55.98	38.64	30.97
Stereo R-CNN [10]	Stereo	85.98	47.58	30.23	23.72
StereoFENet [1]	Stereo	85.70	29.14	18.41	14.20
OC Stereo [15]	Stereo	74.60	55.15	37.60	30.25
RT3DStereo [7]	Stereo	45.81	29.90	23.28	18.96
TLNet [19]	Stereo	63.53	7.64	4.37	3.74

Table 8: 3D detection APs(%) on test set[5]. Car category is evaluated. **Moderate** 2D detection APs are reported. With much lower 2D detection AP, proposed method has the potential to outperform PL++(SDN+GDC).

Method	Type	AP _{2D}	Easy	Mode	Hard
3DCES-FPNv2(ours)	Stereo+LiDAR	51.83	37.16	29.77	26.61
OC Stereo [15]	Stereo	30.79	29.79	20.80	18.62
RT3Dstereo [7]	Stereo	29.30	4.72	3.65	3.00

Table 9: 3D detection APs(%) on test set [5]. Pedestrian category is evaluated.

Method	Type	AP _{2D}	Easy	Mode	Hard
3DCES-FPNv2(ours)	Stereo+LiDAR	46.68	74.20	65.74	58.35
Pseudo-LiDAR [26]	Stereo	67.79	67.30	45.00	38.40
Pseudo-LiDAR++ [33]	Stereo	82.90	78.31	58.01	51.25
PL++(SDN+GDC) [33]	Stereo+LiDAR	85.15	84.61	73.80	65.59
ZoomNet [29]	Stereo	83.92	72.94	54.91	44.14
Stereo R-CNN [10]	Stereo	85.98	61.92	41.31	33.42
StereoFENet [1]	Stereo	85.70	49.29	32.96	25.90
OC Stereo [15]	Stereo	74.60	68.89	51.47	42.97
RT3Dstereo [7]	Stereo	45.81	58.81	46.82	38.38
TLNet [19]	Stereo	63.53	13.71	7.69	6.73

Table 10: BEV detection APs(%) on test set [5]. Car category is evaluated.

Method	Type	AP _{2D}	Easy	Mode	Hard
3DCES-FPNv2(ours)	Stereo+LiDAR	51.83	31.61	24.84	21.96
OC Stereo [15]	Stereo	30.79	24.48	17.58	15.60
RT3Dstereo [7]	Stereo	29.30	3.28	2.45	2.35

Table 11: BEV detection APs(%) on test set [5]. Pedestrian category is evaluated.

4.4 Runtime

With a fixed number of LiDAR points (e.g., 1024) fed into FPNv2, 3D detection phases have almost the same computation costs (see Fig. 5). Therefore, it is necessary to show runtime efficiency of the stereo frustums pipeline. As depicted in Fig. 5, all shown phases are tested with a 2.5Ghz CPU except for the 3D detection phase which is tested with a P100 GPU. RSC-FPNv2 is the fastest among the four methods, RSCCC-FPNv2 doubles the matching time of RSC-FPNv2 but only slightly increases the overall preparation time, 3DCME significantly increases RoIs matching time while 3DCES reduces the time almost by half. Thus, 3DCES is efficient in matching as well as maintaining detection AP as 3DCME does. Most runtime of data-preparation phase goes to point cloud processing, which remains to be optimized by high-performance computing techniques.

4.5 Qualitative Results

The results of 3D detection and BEV detection using SPFN on the validation set are visualized in Fig. 6. The BEV detections are depicted on a sparse point cloud generated by the 3DCES matching module. Fig. 6a shows that two pedestrians are missed due to unsatisfying illumination condition that invalids the 2D detector, another two pedestrians are missed due to large occlusion and too small in dimensions. Despite the missed detections, SFPN is capable of regressing highly precise bboxes that enclose irregular objects as pedestrians. There is a very close object in Fig. 6b which is not detected due to too few clues of objectness in that region. Also, largely occluded objects in Fig. 6b can hardly be detected. Nevertheless, RoIs locate both very near and near (around 40m, which almost reaches the valid-range limit of LiDAR sensor) objects precisely. Some faraway objects located at 50-70m can be detected by SFPN but with less precision. To further study segmented point cloud of the faraway objects, the number of points in each segmented point cloud is usually less than 100, we believe this number is around the borderline threshold N_{thres} .

5 Conclusion

In this paper, we have proposed four matching modules to bridge the gap between 2D object detection on stereopsis and real LiDAR data via dense epipolar geometry constraints: the one-to-one onto mapping and minimum intersection. To accommodate the proposed matching modules, we have proposed a stereo frustum pipeline for 3D object detection where the 2D detection results are fed to the matching module to generate matches to segment the point cloud of the scene, and the 3D segmentation proposals are then fed to a refinement network for more precise objectness segmentation, followed by a multi-modal regression network.

By integrating with the F-PointNets, our stereo frustum pipeline can achieve 2-3 frames per second without coding optimization. Although this frame rate is not yet applicable for real-time applications, its speed can be further increased if we adopt techniques like GPU acceleration. More efficient matching algorithms and 2D detection models are also expected. The proposed pipeline outperforms the state-of-the-art stereo-based approaches with a lower 2D detection average precision, it has the potential to outperform the state-of-the-art LiDAR and stereo fusion approaches if better 2D detection models are adopted.

We are currently working to design end-to-end matching modules for 3DCES and 3DCME to achieve more effective representation to encode the sparse point cloud. Some future work includes leveraging the reliability of both views for better performance in the detection accuracy and recall, optimizing runtime not only at the coding level, but also seeking into the possibilities of implementing distributed parallel computing techniques.

References

1. Bao, W., Xu, B., Chen, Z.: Monofenet: Monocular 3d object detection with feature enhancement networks. *Transactions on Image Processing* (2019)
2. Chen, X., Kundu, K., Zhu, Y., Ma, H., Fidler, S., Urtasun, R.: 3d object proposals using stereo imagery for accurate object class detection. *Transactions on Pattern Analysis and Machine Intelligence* **40**(5), 1259–1272 (2017)
3. Du, X., Ang, M.H., Karaman, S., Rus, D.: A general pipeline for 3d detection of vehicles. In: *IEEE International Conference on Robotics and Automation*, pp. 3194–3200 (2018)
4. Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research* **32**(11), 1231–1237 (2013)
5. Geiger, A., Lenz, P., Urtasun, R.: Official KITTI benchmark. URL <http://www.cvlibs.net/datasets/kitti/>. Accessed: 2019-11-19
6. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: *IEEE International Conference on Computer Vision*, pp. 2961–2969 (2017)
7. Königshof, H., Salscheider, N.O., Stiller, C.: Realtime 3 d object detection for automated driving using stereo vision and semantic information. In: *IEEE International Conference on Intelligent Transportation Systems* (2019)
8. Ku, J., Mozifian, M., Lee, J., Harakeh, A., Waslander, S.L.: Joint 3d proposal generation and object detection from view aggregation. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 1–8 (2018)
9. Li, K., Ma, W., Sajid, U., Wu, Y., Wang, G.: 2 object detection with convolutional neural networks. *Deep Learning in Computer Vision: Principles and Applications* **30**(31), 41 (2020)
10. Li, P., Chen, X., Shen, S.: Stereo r-cnn based 3d object detection for autonomous driving. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7644–7652 (2019)
11. Liang, M., Yang, B., Wang, S., Urtasun, R.: Deep continuous fusion for multi-sensor 3d object detection. In: *European Conference on Computer Vision*, pp. 641–656 (2018)
12. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: *IEEE International Conference on Computer Vision*, pp. 2980–2988 (2017)
13. Luo, W., Yang, B., Urtasun, R.: Fast and furious: Real time end-to-end 3d detection, tracking and motion forecasting with a single convolutional net. In: *IEEE conference on Computer Vision and Pattern Recognition*, pp. 3569–3577 (2018)
14. Ma, W., Wu, Y., Cen, F., Wang, G.: Mdfn: Multi-scale deep feature learning network for object detection. *Pattern Recognition* **100**, 107149 (2020)
15. Pon, A.D., Ku, J., Li, C., Waslander, S.L.: Object-centric stereo matching for 3d object detection. In: *IEEE International Conference on Robotics and Automation*, pp. 8383–8389 (2020)
16. Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3d object detection from rgb-d data. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 918–927 (2018)
17. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 652–660 (2017)
18. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: *Advances in Neural Information Processing Systems*, pp. 5099–5108 (2017)
19. Qin, Z., Wang, J., Lu, Y.: Triangulation learning network: from monocular to stereo 3d object detection. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7607–7615 (2019)

20. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: *Advances in Neural Information Processing Systems*, pp. 91–99 (2015)
21. Shi, S., Wang, X., Li, H.: Pointcnn: 3d object proposal generation and detection from point cloud. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–779 (2019)
22. Shi, S., Wang, Z., Shi, J., Wang, X., Li, H.: From points to parts: 3d object detection from point cloud with part-aware and part-aggregation network. *Transactions on Pattern Analysis and Machine Intelligence* (2020)
23. Shin, K., Kwon, Y.P., Tomizuka, M.: Roarnet: A robust 3d object detection based on region approximation refinement. In: *IEEE Intelligent Vehicles Symposium (IV)*, pp. 2510–2515 (2019)
24. Tian, L., Li, M., Hao, Y., Liu, J., Zhang, G., Chen, Y.Q.: Robust 3-d human detection in complex environments with a depth camera. *Transactions on Multimedia* **20**(9), 2249–2261 (2018)
25. Wang, B., An, J., Cao, J.: Voxel-fpn: multi-scale voxel feature aggregation in 3d object detection from point clouds. *Sensors* **20**(3), 704 (2020)
26. Wang, Y., Chao, W.L., Garg, D., Hariharan, B., Campbell, M., Weinberger, K.Q.: Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8445–8453 (2019)
27. Wang, Z., Jia, K.: Frustum convnet: Sliding frustums to aggregate local point-wise features for amodal. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 1742–1749 (2019)
28. Xu, D., Anguelov, D., Jain, A.: Pointfusion: Deep sensor fusion for 3d bounding box estimation. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 244–253 (2018)
29. Xu, Z., Zhang, W., Ye, X., Tan, X., Yang, W., Wen, S., Ding, E., Meng, A., Huang, L.: Zoomnet: Part-aware adaptive zooming neural network for 3d object detection. In: *AAAI*, pp. 12557–12564 (2020)
30. Yang, B., Luo, W., Urtaun, R.: Pixor: Real-time 3d object detection from point clouds. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7652–7660 (2018)
31. Yang, Z., Sun, Y., Liu, S., Shen, X., Jia, J.: Ipod: Intensive point-based object detector for point cloud. *arXiv preprint arXiv:1812.05276* (2018)
32. Yang, Z., Sun, Y., Liu, S., Shen, X., Jia, J.: Std: Sparse-to-dense 3d object detector for point cloud. In: *IEEE International Conference on Computer Vision*, pp. 1951–1960 (2019)
33. You, Y., Wang, Y., Chao, W.L., Garg, D., Pleiss, G., Hariharan, B., Campbell, M., Weinberger, K.Q.: Pseudo-lidar++: Accurate depth for 3d object detection in autonomous driving. In: *International Conference on Learning Representations* (2019)
34. Zhang, G., Liu, J., Li, H., Chen, Y.Q., Davis, L.S.: Joint human detection and head pose estimation via multistream networks for rgb-d videos. *Signal Processing Letters* **24**(11), 1666–1670 (2017)
35. Zhou, Y., Tuzel, O.: Voxelnet: End-to-end learning for point cloud based 3d object detection. In: *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4490–4499 (2018)