

3D Object Detection Method Based on YOLO and K-Means for Image and Point Clouds

Xuanyu Yin^{1,2} and Yoko Sasaki² and Weimin Wang² and Kentaro Shimizu¹

Abstract—Lidar based 3D object detection and classification tasks are essential for autonomous driving(AD). A lidar sensor can provide the 3D point cloud data reconstruction of the surrounding environment. However, real time detection in 3D point clouds still needs a strong algorithmic. This paper proposes a 3D object detection method based on point cloud and image which consists of three parts.(1)Lidar-camera calibration and undistorted image transformation. (2)YOLO-based detection and PointCloud extraction, (3)K-means based point cloud segmentation and detection experiment test and evaluation in depth image. In our research, camera can capture the image to make the Real-time 2D object detection by using YOLO, we transfer the bounding box to node whose function is making 3d object detection on point cloud data from Lidar. By comparing whether 2D coordinate transferred from the 3D point is in the object bounding box or not can achieve High-speed 3D object recognition function in GPU. The accuracy and precision get improved after k-means clustering in point cloud. The speed of our detection method is a advantage faster than PointNet.

Index Terms—3D Object Detection, Point Cloud Processing, Robot Vision, Machine Learning, Deep Learning

I. INTRODUCTION

Great progress has been made on 2D image understanding tasks, such as object detection and instance segmentation [1]. However, since the creation of 2D bounding boxes or pixel masks, real time detection on 3D point cloud data is becoming increasingly important in many applications areas, such as autonomous driving (AD) and augmented reality (AR). This paper presents our experiments on 3D object detection tasks, which are one of the most important tasks in 3D computer vision. Also presented is the analysis of the experimental results in precision, accuracy, recall, and time and possible future work to improve the 3D object detection average precision (AP).

In the AD field, the LIDAR sensor is the most common 3D sensor. It generates 3D point clouds and captures the 3D structure of scenes. The difficulty of point cloud-based 3D object detection mainly lies in the irregularity of the point clouds from LIDAR sensors [2]. Thus, state-of-art 3D object detection methods either leverage a mature 2D detection framework by projecting the point clouds into a bird's eye view or into a frontal view [2]. However, the information on the point cloud will suffer loss during the quantization process. Charles et al. at Stanford University published a paper on

CVPR in 2017, in which he proposed a deep learning network called PointNet that directly handles point clouds. This paper was a milestone, marking the point cloud processing entered a new stage. The reason is that before PointNet, we had no way to deal with point clouds directly. Because point clouds are three-dimensional, they are not smooth. Moreover, deep neural networks, which make many ordinary algorithms, do not work. Thus, researchers have come up with a variety of methods[1][2][3], such as flattening the point clouds into pictures (MVCNN), dividing the point clouds into voxels, and then dividing them into nodes and straightening them in order. Thus, the cloud domain has advanced from the “pre-PointNet era” to the “post-PointNet era” thanks to the development of this technology. After PointNet, PointCNN, SO-Net, etc., came out, the operation of these methods improved steadily.

PointNet[3] has achieved 83.7 percent mean accuracy. However, the speed is still a problem. Compared with two-dimensional data, point cloud data with an additional dimension are too large to achieve the requirements for real-time 3D object detection. This paper presents our extraction of every point that may be an object after transformation in a 2D bounding box, enabling high-speed 3D object detection to be achieved. First, we describe a device we constructed including six cameras and one LIDAR. Then, we present the experiments we conducted to show how 2D images are captured by cameras and how 3D point clouds using LIDAR store the data in a rosbag, which was reused in subsequent experiments. The image data needed to be distorted, but an undistorted transform process was also needed. After the undistorted transform process, five images are split and dropped into the you only look once (YOLO) detection process. The YOLO detection process returns the related bounding box and class label. We store the bounding box and the class label for later reading, reducing the coupling of project research. In the second step, a function is written to extract the different topic information of the rosbag. After extracting the point cloud file, we put it into the numpy matrix for future operations.

Data conversion based on external and internal parameters is performed for every point cloud by matching the 2D images corresponding to each point cloud by matching the fps of the cameras and LIDAR. For each different bounding box of each point cloud, we collect all the matching points and render different colors based on different class labels. Finally, unsupervised clustering of point clouds in different bounding boxes improves the detection performance by removing some of the noise. The results of 3D object recognition are presented at the end of the paper. The recognition results were saved in

¹Graduate School of Interdisciplinary Information Studies, The University of Tokyo, The Tokyo, Japan yinxuanyu@g.ecc.u-tokyo.ac.jp

²National Institute of Advanced Industrial Science and Technology, The Tokyo, Japan y-sasaki@aist.go.jp

a rosbag, and then 3D visualization was performed to check the experimental results.

The evaluation results of the recognition experiment were completed using the method of depth images; the point cloud detection results were transferred to 32*1024 depth images. The final evaluation experiment was done for a comparison with ground truth in every pixel. The aforementioned is the rough research process of this article.

II. 3D OBJECT DETECTION METHOD

A. Overview

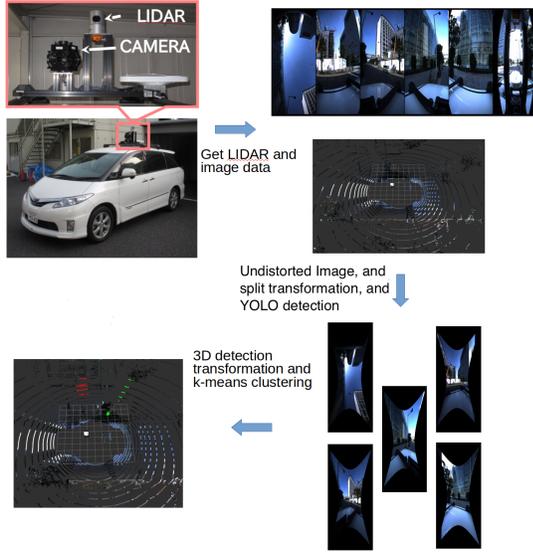


Fig. 1. Overview of the proposed object detection system

Fig.1 shows the overview of the proposed system. This research was basically divided into six parts. The first part mainly focused on the calibration of the cameras and the structural design of the testing equipment. The second part was to convert the distorted images into undistorted ones. The third part was YOLOv3 object recognition with 2D images. We mainly applied YOLOv1 tiny and YOLOv3 methods in doing the experiments, using keras to reproduce YOLO. The fourth part was the extraction of point clouds. We used rosbag to store the data and RVIZ for point cloud visualization. The fifth part was the unsupervised clustering of k-means, which were used to optimize the detection results of the basic experiments and to improve the detection accuracy of 3D object recognition.

B. Lidar-camera calibration

Here is the main information on the equipment used in this experiment and the external reference of the cameras. This experiment used Velodyne lidar(HDL-32e) with omni-directional cameras(PointGrey Ladybag5) to achieve 360 no dead angle monitoring, which is shown in Fig. 2.

A geometric model of camera imaging must be established during the image measurement process and machine vision application to determine the relationship between the three-dimensional geometric position of a point on the surface of a



Fig. 2. Structure of detection device: LIDAR and omni-directional camera

space object and its corresponding point in the image. These geometric model parameters are camera parameters. Under most conditions, these parameters must be obtained through experiments and calculations[4]. This process of solving parameters is called camera calibration. The internal parameters of the five cameras obtained at the end of this study are shown in Fig. 3.

```
Internal_parameters = [[(262.437031921347, 0., 521.935575263155, 0., 519.453341046090, 1252.10255470112, 0., 0., 1.),
(2.421699227481561e+02, 0., 5.171454569398794e+02, 0., 4.803203353454073e+02, 1.23807874e+03, 0., 0., 1.),
(236.722566162119, 0., 527.119344859111, 0., 470.499233566595, 1276.67270316542, 0., 0., 1.),
(236.356345549995, 0., 543.141363901400, 0., 470.51655348912, 1289.45734138472, 0., 0., 1.),
(244.22777457936, 0., 512.499661250746, 0., 487.087713904067, 1258.1831389491, 0., 0., 1.)]]
External_parameters = [[(7.023846171215485157e+00, -1.420039003532723054e+00, 8.5181478702986689e-01,
6.480931488336687682e-02, -1.672540390406322686e-01, 5.129173066285740434e-01,
(-1.564770985936392522e+00, -3.287539592246310744e-01, -3.126095115904812425e+00,
-1.300167248877406867e-01, -2.2982052755709237e-01, 3.039150183240042422e-02),
(-1.559532112765356526e+00, 9.447475452007216834e-01, -3.126403042750293793e+00,
-8.010125548818973484e-02, -2.363659554470468327e-01, -1.409445185079227780e-01),
(1.563879228745117755e+00, 9.33727244088895033e-01, 3.0833002121740866780e-03,
1.13918020955745548e-01, -2.318417592047426523e-01, -1.58808294958492050e-01),
(1.576105513946523586e+00, -3.074588238495890025e-01, 1.51922864785923845e-02,
1.744782300756334281e-01, -2.373308756671665409e-01, 8.360234406313583888e-02)]]
```

Fig. 3. Internal and external parameters

C. Image undistorted transform

In photography, wide-angle lenses are generally believed to be prone to barrel distortion, while telephoto lenses are prone to pincushion distortion. If a camera uses a short focal length wide-angle lens, the resulting image will be more susceptible to barrel distortion because the magnification of the lens gradually decreases as the distance increases, causing the image pixels to surround the center point radially. Fig. 4 shows raw images, and Fig. 5 shows images after the undistorted transform using OpenCV for image correction and camera calibration.

D. YOLO-based detection

YOLO is a fast target detection algorithm that is very useful for tasks with very high real-time requirements. The YOLO authors launched YOLOv3 version in 2018. After training on



Fig. 4. Raw images of five cameras

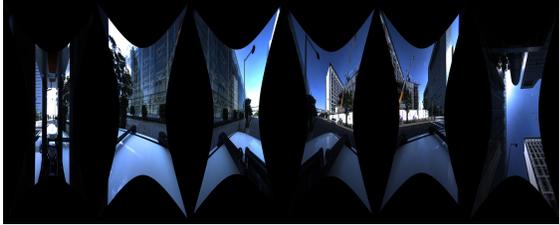


Fig. 5. Images of five cameras after undistorted transform

Titan X, v3 is 3.8 times faster than RetinaNet regarding mean average precision (mAP), and it can create a 320320 picture in 22 ms. The objective score is 51.5, which is comparable to the accuracy of the single shot detector (SSD), but it is three times faster. Thus, YOLOv3 is very fast and accurate. In the case of IoU=0.5, it is equivalent to the mAP value of Focal Loss, but it is four times faster. We utilized YOLOv3 as a 2D object detection algorithm. Fig. 6 shows an example from a camera image, and Fig. 7 shows an example of 3D object detection from a point cloud.

We present a total of all the classes that can be identified in the coco dataset, including people, bicycles, cars, motorbikes, airplanes, buses, trains, trucks, boats, traffic, lights, fires, hydrants, and stop lights. In total, 80 classes are presented.

However, the most frequent classes are trucks, people, and cars. Thus, training a new YOLOv3 neural network to detect these three classes may be useful in saving detection time and enabling real-time functionality.

Because the images are largely black after undistorted conversion, we must remove the noise that exceeds its maximum bounding box when doing a conversion. Because the boundingbox of more than a quarter of the image size contains black parts, this part is invariably noisy. Fig 6 shows the YOLO detection examples from the camera image.

E. Point cloud extraction

We mainly used rosbag to read and output the data. The read data contained undistorted images and point cloud images. The output was mainly the result of point cloud detection.

F. K-means based point cloud segmentation

Because 2D is converted into 3D, all points that can be mapped into the bounding box in a certain direction are marked with different label colors. To improve the experiment, we utilized the unsupervised clustering method for k-means machine learning. The detection was faster, and the accuracy

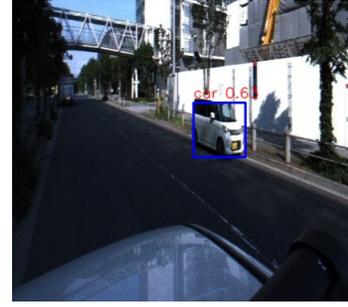


Fig. 6. YOLO detection example from camera image

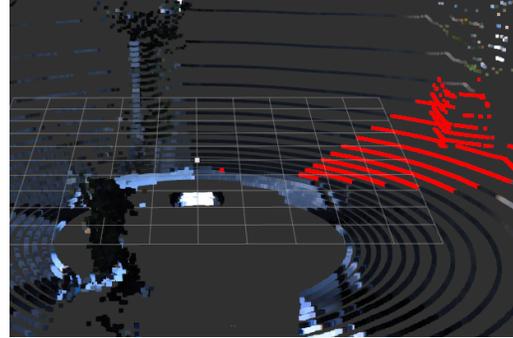


Fig. 7. 3D object detection example of point cloud

of the points substantially improved, removing most of the noise points. The k-means clustering graph is shown in Fig. 8, and the 3D object detection examples by this method is shown in Fig. 7.

G. Evaluation of prediction results in depth images

To create a comparison with ground truth and to make the results easier to observe, we converted the point clouds into a 3*1024 panoramic depth image that is shown in Fig. 9 with different colors corresponding to different categories. Depth image generation properties are shown in Table I. In this step, handmade 0.1K ground truth images were created using the LabelMe annotation tool. The final evaluation results also included these 0.1K pictures.

TABLE I
DEPTH IMAGE GENERATION PROPERTIES

Description	Specifications
Number of detector pairs	32
Limitation in horizontal scanning	1024
Vertical scanning range	+10.67 to -30.67
Angular resolution in vertically	1.33
Angular resolution in horizontally	0.2
Model located area	0 to 25,000[mm] from origin
Depth image angle	0 to 360 degrees
Output depth image size	32 1024 1

III. EXPERIMENTS AND RESULTS

An in-vehicle sensor was used, and a large number of data were collected. The identification experiment consisted of two

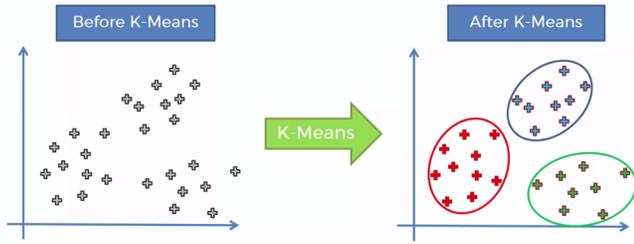


Fig. 8. K-means clustering



Fig. 9. 32*1024 depth image

parts. The first included visualization and quantity statistics of point cloud identification with and without k-means clustering. The experimental results were then evaluated. The second part included evaluation criteria the accuracy, precision, and recall by comparing the results after conversion to depth images with ground truth.

A. Experiments with different classes

In this study, we first made 3D prediction results that were directly converted. Basically, all the points on the image that could be mapped to the bounding box were recognized. This led to a lot of noise, making it impossible to identify the 3D bounding box, but enabled recognizing the 3D radar data. A specific category exists in a certain direction, thereby completing the rough 3D recognition function. A 2D YOLO example is shown in Fig. 10, and related experiment results with and without k-means are shown in Fig. 11.

The proofreading here is mainly for the naked eye. In accordance with the specific bounding box data and the data under the camera, judgments were made as to whether or not the recognition results were correct.



Fig. 10. 2D YOLO experimental results

B. Experiments with k-means

K-means clustering is a method of vector quantization, originally from signal processing, which is popular for cluster analysis in data mining. K-means clustering aims to partition n

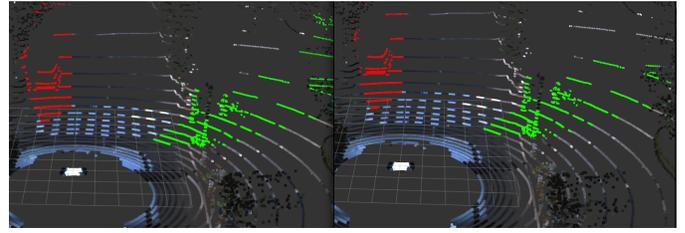


Fig. 11. 3D experiment results without(left) and with(right) k-means

observations into k clusters in which each observation belongs to the cluster with the nearest mean, serving as a prototype of the cluster. This results in a partitioning of the data space into Voronoi cells [4].

Because k-means pre-determines the number of categories and the maximum number of cluster iterations, the center point is randomly selected at first, so the results of each cluster become biased. However, the expected value of each category did not change too much at the end of the experiment.

After using YOLO to complete the 2D object recognition and to convert the data into 3D point cloud data, we mainly used k-means to cluster further the points of the corresponding bounding box that had been acquired so as to remove some noise and to make the recognition results more accurate.

C. Results with and without unsupervised learning

The results of the experiment with and without k-means are shown in Fig. 11. The total number of point clouds in each group was 46,464. The maximum number of points was 4989. This occurred when a car or truck got very close. The lowest number was 0, which means that no target object could be recognized around that time. In other words, it was an empty space with no pre-trained class label objects.

Fig. 12 and Fig.13 represents the ratio and number of dropped data after clustering using the k-means method on the point cloud dataset. The highest was 49.15 percent in 46 frames. In the case of not recognizing any objects, the lowest was only 0.85 percent in 21 frames.

When clustering was not performed, the experimental results stained all the points mapped to the bounding box in that direction. After clustering, if there was an object, the number of point clouds in that part significantly increased, while other places that had no object obviously had a noticeable decrease in the number of clouds. Therefore, unsupervised clustering could significantly change the results, generally removing 31.6 percent of the point clouds.

D. Evaluation of prediction results in depth images

We made a 0.1K ground truth by labeling, reading, and outputting operations on LabelMe software¹. In the previous stage, we got the prediction results of the point clouds. After converting the point clouds to the same 3*1024 depth image, the 3*1024*1 image file containing the object category labels were saved too. In the final evaluation test, accuracy, precision,

¹labelme:<http://labelme.csail.mit.edu/Release3.0/>

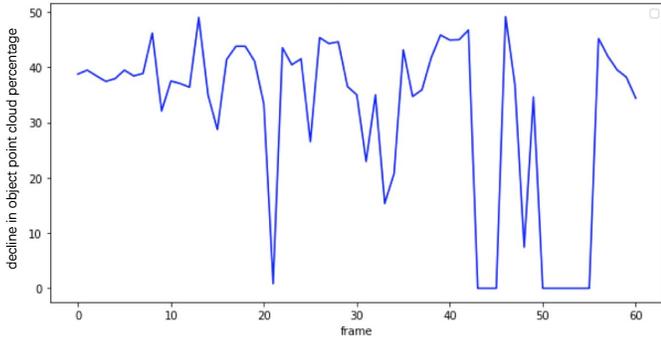


Fig. 12. Rate of point cloud decline with k-means function

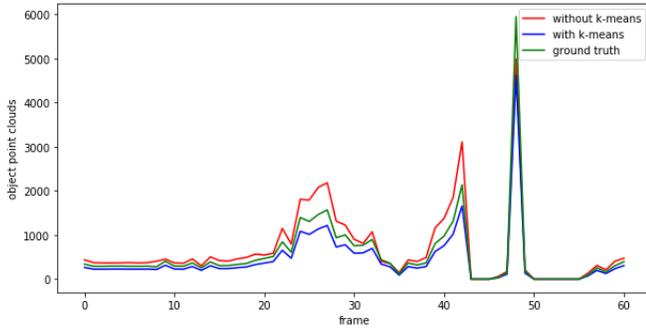


Fig. 13. Experiment results of number of detection point clouds with and without k-means function

and recall were calculated by comparing the ground truth and the numpy file holding the experimental prediction results. Fig. 14 is the Depth image which is converted from point cloud data. The ground truth is shown in Fig. 17, the detection results are in Fig. 16, and the final results after k-means clustering are in Fig. 17. The accuracy and consumed time results for the total prediction process are shown in Table II. Accuracy, precision, and recall with and without the k-means clustering were recorded. The results shown in Table III reveal YOLO 2Ds detection was successful.

TABLE II
PREDICTION RESULTS WITH AND WITHOUT K-MEANS CLUSTERING

precision method	accuracy	time
PointNet	0.837	1.00
image(YOLO)+pointcloud(without clustering)	0.7288	0.848
image(YOLO)+pointcloud(k-means clustering)	0.7301	0.856

The "time" stand for the consumed time to handle one million point clouds (around 1K objects)

TABLE III
PREDICTION RESULTS AFTER YOLO

precision method	accuracy	precision	recall
without clustering	0.900	0.7753	0.6011
K-means clustering	0.9023	0.7867	0.6103



Fig. 14. Depth image



Fig. 15. Ground truth



Fig. 16. Prediction result without k-means



Fig. 17. Prediction results with k-means function

IV. CONCLUSIONS

The conclusion of our study are as follows:

1.The method adopted by this paper is to directly convert the 3D point cloud to 2D image data, from the recognition of the 2D boundingbox to the dyeing of the 3D point cloud. Since the YOLO algorithm is adopted, the real-time performance is very strong, and the unsupervised clustering is used too. A lot of noise will be removed. It makes the recognition better.

2.This paper mainly wants to find a way to quickly and accurately determine whether there are objects and objects in a certain direction. This will contribute to the success of the unmanned field, allowing the car to obtain more information to make more judgments.

3.The final experimental results, in the case of using two 1080Ti GPUs, basically ensure that the experiment without clustering consumes 0.19 seconds per frame and 0.192 seconds after k-means clustering in 5 threads. The fast identification process ensures the real-time detection of the surrounding conditions in unmanned driving. If parallel, distributed computing and other technologies are used, the recognition speed will be faster.

4.The speed is very fast. However, the accuracy is not very high due to a front yolo recognition accuracy needs to be considered. The recall for detection is not high too.

V. FUTURE WORK

In the future, robots will be added, and semantic mapping from running mobile robots will form the core of the next step. Then, we will consider not only the k-means function but also handling methods for directed point clouds like PointNet and FCN or more clustering methods like point cloud based depth clustering to figure out a faster method to complete 3d object detection using images and lidar. We will develop automatic labeling functions using our method for training data generation of LIDAR-based 3D objects.

ACKNOWLEDGMENT

We would like to thank Wonjik Kim and Ryusei Hasegawa for providing us with their conversion tools, which transferred the point cloud data to 32*1024 depth images.

REFERENCES

- [1] Qi, C. Ruizhongtai, W. Liu, C. Wu, H. Su and L. J. Guibas. "Frustrum PointNets for 3D Object Detection from RGB-D Data, 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (2018): 918-927.
- [2] S. Shaoshuai, W. Xiaogang and Li. Hongsheng "PointRCNN: 3D Object Proposal Generation and Detection from Point Cloud, CoRR abs/1812.04244 (2018): n. pag.
- [3] Qi, C. Ruizhongtai, H. Su, K. Mo and L. J. Guibas. "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation, 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017): 77-85.
- [4] K. Jason, M. Mozifian, J. Lee, A. Harakeh and S. Lake Waslander. "Joint 3D Proposal Generation and Object Detection from View Aggregation, 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (2018): 1-8.
- [5] W. Kim, M. Tanaka, M. Okutomi, Y. Sasaki, "Automatic Labeled LiDAR Data Generation based on Precise Human Model," in IEEE International Conference on Robotics and Automation (ICRA).IEEE, 2019.
- [6] Mousaviarz, Arsalan, D. Anguelov, J. Flynn and J. Kosecka. "3D Bounding Box Estimation Using Deep Learning and Geometry, 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017): 5632-5640.
- [7] K. Shin, Y. Paul Kwon, M. Tomizuka, "RoarNet: A Robust 3D Object Detection based on RegiOn Approximation Refinement," arXiv:1811.03818 [cs.CV].
- [8] Y. Zhou, O. Tuzel, "VoxelNet: End-to-End Learning for Point Cloud Based 3D Object Detection," 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)2017.
- [9] Wikipedia, "K-means clustering," https://en.wikipedia.org/wiki/K-means_clustering
- [10] C. Xiaozhi, M. Huimin, W. Ji, Li. Bo and X. Tian. "Multi-view 3D Object Detection Network for Autonomous Driving, 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017): 6526-6534.
- [11] Wikipedia, "Tutorial Camera Calibration," https://boofcv.org/index.php?title=Tutorial_Camera_Calibration.
- [12] C. Xiaozhi, K. Kundu, Z. Zhang, H. Ma, S. Fidler and R. Urtasun, "Monocular 3D Object Detection for Autonomous Driving, 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016): 2147-2156.
- [13] J. Sun, M. Ovsjanikov, and L. Guibas. "A concise and provably informative multi-scale signature based on heat diffusion," In Computer graphics forum, volume 28, pages 13831392. Wiley Online Library, 2009.
- [14] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. "3d shapenets: A deep representation for volumetric shapes," In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 19121920, 2015.
- [15] L. Yi, V. G. Kim, D. Ceylan, I-C. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, and L. Guibas. "A scalable active framework for region annotation in 3d shape collections," SIGGRAPH Asia, 2016.
- [16] H. Ling and D. W. Jacobs. "Shape classification using the inner-distance," IEEE transactions on pattern analysis and machine intelligence, 29(2):286299, 2007
- [17] M. Aubry, U. Schlickewei, and D. Cremers. "The wave kernel signature: A quantum mechanical approach to shape analysis," In Computer Vision Workshops (ICCV Workshops), 2011 IEEE International Conference on, pages 16261633. IEEE, 2011.
- [18] Y. Fang, J. Xie, G. Dai, M. Wang, F. Zhu, T. Xu, and E. Wong. "3d deep shape descriptor," In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 23192328, 2015.