

UNIVERSITY OF TECHNOLOGY SYDNEY
Faculty of Engineering and Information Technology

Research on Point Cloud Segmentation and 3D Scene Understanding

by

Anan Du

Supervisor: Jian Zhang
Co-Supervisor: Qiang Wu

A THESIS SUBMITTED
IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE

Doctor of Philosophy

Sydney, Australia

April, 2022

Certificate of Authorship/Originality

I, Anan Du declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the School of Electrical and Data Engineering/Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

This research is supported by the Australian Government Research Training Program.

Signature: Production Note:
Signature removed prior to publication.

Date: 17/04/2022

Acknowledgements

During my doctoral study in UTS, I have received a lot of support from my supervisors, colleagues, and family. Here, I would like to express my sincere gratitude to them.

First and foremost, I would like to thank my Ph.D supervisor, Prof. Jian Zhang, for his guidance and support throughout my research and completion of this thesis. Prof. Zhang Jian has been an academic role model to me. He is rigorous and enthusiastic in work, has deep insights in academic, has rich skills on research writing and presentation and always keeps curiosity on new knowledge. All of these characters have guided me and will continue encouraging me to be a professional researcher. I also want to thank my co-supervisor, Associate Prof. Qiang Wu, for his constructive guidance regarding my research directions. Besides, his optimistic attitude to work and life has inspired me a lot.

I would like to thank my colleague, Dr. Xiaoshui Huang, who is also engaged in research on 3D point clouds. Dr. Xiaoshui has given me a lot of professional guidance during my entire Ph.D. study, from introducing me to the basic knowledge about 3D point cloud, to discussing my follow-up research works, to writing articles. He always works hard and efficiently. He also has a wide range of scientific research fields and often has novel ideas. It is an honor to work with him. I want to thank my co-author Dr. Shuchao Pang for supporting my study and life. In my Ph.D study, we often shared the articles we read, discussed our views, debugged the experiments, and polished each other's articles. In life, when I encounter setbacks, Dr. Shuchao provided me a lot of encouragement. I also want to thank my research collaborator Dr. Tianfei Zhou, for his help to my research. He has given me a lot of support in the method design, validation, and article writing. I also appreciate my colleague, Dr. Lingxiang Yao, for his help in polishing the articles and guidance in doing the

presentation with his excellent communication skills.

I want to thank all my colleagues supervised by Prof. Jian and Associate Prof. Qiang during my Ph.D study. Thanks to Dr. Jinsong Xu, Dr. Junjie Zhang, Dr. Muming Zhang, Dr. Lu Zhang, Dr. Yongshun Gong, Dr. Zhibin Li, Dr. Huaxi Huang, Dr. Litao Yu, Dr. Zongjian Zhang, Dr. Guofeng Mei, Dr. Wenbo Xu, Dr. Yan Huang, Dr. Peng Zhang, Dr. Qiang Li, Dr. Xunxiang Yao. I also thank all other colleagues I have studied with. They make my study and life in UTS beautiful.

At last, I want to thank my parents, brother and sister. They always show their patience, trust and love to me. I love them.

Anan Du
Sydney, Australia, 2022.

ABSTRACT

Recently, the amount of 3D data has been sharply increasing thanks to widely available 3D sensors like Lidar, Kinect and RealSense. Compared to 2D images, 3D data provides clear topology and geometric information, which is very important for many computer vision applications. Among many types of 3D data, the point cloud is widely used because of its easy availability and storage. This thesis summaries the works that have been conducted on understanding the high-level semantics and basic structure of 3D scenes based on point cloud data.

Firstly, a fully-supervised point cloud semantic segmentation framework is designed. Contextual information can help resolve ambiguity and improve the robustness of a recognition system. To capture long-range context, a long-short-term feature bank based framework is introduced to exploit the patch-wise relationship for point cloud semantic segmentation. This approach can capture context in an arbitrary range by costing a little extra computation than a standard segmentation model. Experiments demonstrate that the proposed approach outperforms other point cloud semantic segmentation approaches exploring long-range context.

Manually labeling point cloud datasets, especially point-wise annotations, for fully-supervised methods is expensive. To avoid manually annotating 3D keypoints, a weakly-supervised 3D keypoint extraction method for point cloud registration, called KPSNet, is proposed, which only uses the relative transformation matrices between input pairs of point clouds as weak labels to establish point-to-point correspondences on-the-fly for training. Moreover, KPSNet can simultaneously detect 3D keypoints and learn the representations. Experimental results reveal that our proposed method performs better in point cloud registration than other methods without keypoint annotations.

To cut the need for manual point-wise annotations, a weakly-supervised point

cloud semantic segmentation approach is proposed, which uses coarse labels only. The proposed approach reformulates this problem as a semi-supervised point cloud semantic segmentation problem with noisy pseudo labels. Then a three-branch network is proposed, which can reduce the impact of noises in pseudo labels and leverage the inner structure of point clouds to improve the segmentation performance. By evaluating on benchmark datasets, this method surpasses other weakly-supervised methods and fully-supervised method PointNet++, and narrows the gap with other outstanding fully-supervised approaches. Subsequently, another weakly-supervised point cloud semantic segmentation framework is introduced, which uses incomplete point-level labels. It explores contrastive learning with cross-sample contrast and low-level contrast and a pseudo label refinement module to mine more helpful supervision information from pseudo labels. This method gains state-of-the-art performance on benchmark datasets with several weakly-supervised annotation settings.

Keywords: Point cloud, semantic segmentation, 3D keypoint, weakly-supervised, 3D scene understanding

Contents

| | |
|---|----------|
| Certificate | i |
| Acknowledgments | ii |
| Abstract | iv |
| List of Figures | x |
| List of Tables | xii |
| List of Publications | xiv |
| 1 Introduction | 1 |
| 1.1 Background | 1 |
| 1.2 Research Challenges and Objectives | 3 |
| 1.3 Research Contributions | 4 |
| 1.4 Structure of Thesis | 6 |
| 2 Literature Review | 8 |
| 2.1 Deep Learning on Point Clouds | 8 |
| 2.1.1 Multi-view based Methods | 8 |
| 2.1.2 Voxel based Methods | 10 |
| 2.1.3 Raw Point based Methods | 11 |
| 2.1.4 Other Methods | 13 |
| 2.2 Point Cloud Semantic Segmentation | 14 |
| 2.2.1 Multi-view based Deep Architectures | 14 |

| | | |
|-------|--|----|
| 2.2.2 | Voxel based Deep Architecture | 16 |
| 2.2.3 | Point based Deep Architecture | 21 |
| 2.2.4 | Other Architecture | 23 |
| 2.3 | 3D Keypoint Detection | 26 |
| 2.3.1 | 3D Keypoint Detector | 26 |
| 2.3.2 | 3D Keypoint Descriptor | 27 |
| 2.4 | Weakly Supervised Learning for Point Cloud Understanding | 28 |
| 2.4.1 | Incomplete Supervision | 29 |
| 2.4.2 | Inexact Supervision | 31 |
| 2.5 | Contrastive Learning for Dense Prediction | 31 |

3 Exploring Long-Short-Term Context for Point Cloud

Semantic Segmentation 33

| | | |
|-------|--|----|
| 3.1 | Introduction | 33 |
| 3.2 | Methods | 36 |
| 3.2.1 | The Encoder Module | 36 |
| 3.2.2 | The Long-Short-Term Feature Bank Establishment | 37 |
| 3.2.3 | The Feature Bank Fusion Module | 38 |
| 3.2.4 | Decoder and Predictor Modules | 39 |
| 3.2.5 | Implementation Details | 40 |
| 3.3 | Experiments | 40 |
| 3.3.1 | Datasets and Evaluation Metrics | 40 |
| 3.3.2 | Experimental Settings | 41 |
| 3.3.3 | Performance on S3DIS | 42 |
| 3.3.4 | Qualitative Results and Discussions | 44 |

| | | |
|----------|---|-----------|
| 3.4 | Conclusions | 46 |
| 4 | Weakly-supervised 3D Keypoint Detection for Point Cloud Registration | 48 |
| 4.1 | Introduction | 48 |
| 4.2 | Methods | 50 |
| 4.2.1 | Architecture of KPSNet | 50 |
| 4.2.2 | Model Training | 52 |
| 4.3 | Experimental Setup and Results | 55 |
| 4.3.1 | Datasets and Experimental Settings | 55 |
| 4.3.2 | Performance on 3DMatch Benchmark | 57 |
| 4.3.3 | Qualitative Results and Discussions | 59 |
| 4.4 | Conclusions | 64 |
| 5 | Weakly-supervised Point Cloud Semantic Segmentation with Coarse-grained Labels | 65 |
| 5.1 | Introduction | 65 |
| 5.2 | Pipeline for Weakly-supervised Point Cloud Semantic Segmentation . . | 68 |
| 5.2.1 | Phase-1: Pseudo Label Generation With Cloud-level Weak Annotations | 70 |
| 5.2.2 | Phase-2: Semantic Segmentation Network Learning With Selected Pseudo Labels | 71 |
| 5.3 | Experimental Setup and Results | 75 |
| 5.3.1 | Datasets and Experimental Settings | 75 |
| 5.3.2 | Performance on ScanNet | 78 |
| 5.3.3 | Performance on S3DIS | 82 |

| | | |
|----------|--|------------|
| 5.3.4 | Ablation Study | 85 |
| 5.4 | Conclusions | 89 |
| 6 | Point Contrast and Labeling for Weakly-supervised Point Cloud Semantic Segmentation | 90 |
| 6.1 | Introduction | 90 |
| 6.2 | Methods | 93 |
| 6.2.1 | Overview of PCL | 93 |
| 6.2.2 | Cross-sample Contrast | 96 |
| 6.2.3 | Low-level Contrast | 99 |
| 6.2.4 | Pseudo Label Refinery | 99 |
| 6.3 | Experimental Settings and Results | 100 |
| 6.3.1 | Datasets and Experimental Setup | 101 |
| 6.3.2 | Performance on ScanNet | 103 |
| 6.3.3 | Performance on S3DIS | 106 |
| 6.3.4 | Ablation Study | 108 |
| 6.3.5 | Qualitative Results on ScanNet and S3DIS | 110 |
| 6.4 | Conclusions | 113 |
| 7 | Conclusions and Future Work | 114 |
| 7.1 | Conclusions | 114 |
| 7.2 | Future Work | 116 |
| | Bibliography | 119 |

List of Figures

| | | |
|-----|--|----|
| 1.1 | Thesis contributions: supervised and weakly-supervised methods for 3D scene understanding. | 5 |
| 2.1 | Illustrations of sparse convolution and submanifold sparse convolution | 18 |
| 2.2 | Difference between submanifold sparse convolution and generalized sparse convolution | 19 |
| 2.3 | Architecture of a 3D-UNet | 20 |
| 2.4 | Illustrations of popular point based deep architectures. | 22 |
| 2.5 | Illustrations of training with different types of annotations. | 29 |
| 3.1 | Architecture of our Long-Short-Term Context framework (LSTC) . . | 37 |
| 3.2 | Feature Bank Fusion Module. | 39 |
| 3.3 | Visualization of segmentation results on S3DIS. | 45 |
| 4.1 | Architecture of our KPSNet | 51 |
| 4.2 | Per-scene performance on 3DMatch benchmark | 58 |
| 4.3 | Visualization of our KPSNet applied for point cloud registration . . . | 60 |
| 4.4 | Visualization of keypoint features. | 61 |
| 4.5 | Cases with poor point cloud registration results applying our proposed KPSNet | 63 |

| | | |
|-----|--|-----|
| 5.1 | Overview of our weakly-supervised point cloud semantic segmentation method | 69 |
| 5.2 | Architecture of the three-branch framework in Phase-2 | 72 |
| 5.3 | Visualization of our segmentation results on ScanNet validation set . | 80 |
| 5.4 | Visualization of our segmentation results on S3DIS testing set. . . . | 84 |
| 6.1 | Illustration of our PCL framework | 94 |
| 6.2 | Two types of point-level relationships for our cross-sample contrastive loss and similarity loss | 97 |
| 6.3 | Performance on ScanNet using different losses. | 108 |
| 6.4 | Qualitative results on ScanNet validation set under different weak annotation settings | 110 |
| 6.5 | Qualitative results on S3DIS under different weak annotation settings | 112 |

List of Tables

| | | |
|-----|--|-----|
| 3.1 | Comparison results on the S3DIS | 43 |
| 4.1 | Overall performance on the 3D-match benchmark | 57 |
| 5.1 | Segmentation results on ScanNet validation set. We provide the mIoU (%) and per-class IoUs (%). Bold represents the best results across all listed approaches. | 79 |
| 5.2 | Comparisons on ScanNet testing set | 81 |
| 5.3 | Comparison results on S3DIS testing set (Area-5) | 83 |
| 5.4 | Segmentation performance of pseudo labels and selected version on ScanNet training set. We report the mIoU (%) and per-class IoUs (%). | 85 |
| 5.5 | Comprehensive evaluation of the pseudo labels on ScanNet training set. | 86 |
| 5.6 | Performance of our method with different thresholds θ on ScanNet validation set | 87 |
| 5.7 | Impact of individual losses on ScanNet validation set. | 88 |
| 6.1 | Comparisons with existing weakly-supervised methods for point cloud semantic segmentation on ScanNet testing set | 104 |
| 6.2 | Performance of our fully-supervised baseline on ScanNet testing dataset. | 105 |
| 6.3 | Comparisons with existing methods on S3DIS. | 106 |

| | | |
|-----|--------------------------------------|-----|
| 6.4 | Evaluation on Pseudo labels. | 109 |
|-----|--------------------------------------|-----|

List of Publications

Conference Papers

1. **Anan Du**, Shuchao Pang, Xiaoshui Huang, Jian Zhang, Qiang Wu, "Exploring Long-Short-Term Context For Point Cloud Semantic Segmentation", In *2020 IEEE International Conference on Image Processing (ICIP)*, pp. 2755-2759. IEEE, 2020.
2. **Anan Du**, Xiaoshui Huang, Jian Zhang, Lingxiang Yao, Qiang Wu, "KPSNet: Keypoint Detection and Feature Extraction for Point Cloud Registration", In *2019 IEEE International Conference on Image Processing (ICIP)*, pp. 2755-2759. IEEE, 2019.

Journal Papers

1. **Anan Du**, Shuchao Pang, Xiaoshui Huang, Jian Zhang, Qiang Wu, "Weakly-supervised Point Cloud Semantic Segmentation with Coarse-grained Labels." In *IEEE Transactions on Multimedia*. (Submitted)
2. **Anan Du**, Tianfei Zhou, Shuchao Pang, Jian Zhang, Qiang Wu, "Point Contrast and Labelling for Weakly-supervised Point Cloud Semantic Segmentation." In *IEEE Transactions on Image Processing*. (Under submission)