# Object Detection and Localization in 2D & 3D Environment

**by Zhihao Cui**

Thesis submitted in fulfilment of the requirements for the degree of

**Doctor of Philosophy**

under the supervision of Forest Zhu

University of Technology Sydney
Faculty of Engineering and IT

26th August 2021

# Certificate of Authorship/Originality

CERTIFICATE OF ORIGINAL AUTHORSHIP

I, *zhihao cui* declare that this thesis, is submitted in fulfilment of the requirements for the award of *Doctor of Philosophy*, in the *FEIT* 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.

Signature:
Production Note:
Signature removed prior to publication.

Date: 23/8/2021

# ABSTRACT

# OBJECT DETECTION AND LOCALIZATION IN 2D & 3D ENVIRONMENT

Computer vision is a science that studies how to make machines "see." It refers to utilizing vision sensors and computers to identify, locate, and track objects. Under this topic, this thesis proposed three frameworks to improve 2D and 3D object detection and localization performance. In the first 3D object detection framework, we investigated the bilateral convolution layers' feasibility to alternate the widely used point cloud voxelization process. The second framework explored the voxel-wise and point-wise proposal fusions method to improve 3D object detection performance. For the 2D instance segmentation, the framework formed an NMS-free and anchor-free detector designed explicitly for the eye-to-hand robotic system.

In existing works, most of the state-of-the-art 3D object detection approaches are based on the point clouds' voxelization method to sample the point cloud into a subdivide voxel space. Although it provides an efficient way to process point cloud data, its lack of feature relationship on voxel-level limits the model's detection accuracy. Furthermore, the voxel sizes hyperparameters tuning increased the model complexity, resulting in a fluctuated model performance. To this end, we aim to simplify the process by re-projecting the point cloud data onto a lattice hyper-plane that saves point cloud processing time while maintaining the model accuracy. The proposed framework Bilateral Lattice Point Network (BLPNet) is provided in chapter three.

In the second framework, Point and Voxel Fusion Net (PVF-Net) is proposed to further push the 3D object detection performance forward. In two-stage approaches, increasing the first stage proposals recall rate positively influences the model final prediction performance. Therefore, in the PVF-Net, we proposed a twofold proposal fusion architecture to extract and fuse the voxel-level and point-level features of the point clouds. The model details are in chapter four, mainly consisting of two novel modules: the Twofold Proposal Fusion (TPF) module and the ROI Deep Fusion (RDF) module.

Lastly, it is well-known that 3D and 2D sensors jointly depict the real world. In chapter five, 2D object detection will become the next goal for improvement. So far, the existing 2D instance segmentation algorithms developed significantly and reached a saturated performance. However, there is no solid solution for heavy occluded or diagonally arranged objects, especially in the vision-guided robot picking system. To solve the problem above, we proposed a real-time occlusion and oblique friendly instance segmentation framework, terms as Keypoint-Mask, assisting the robotic system to handle the complicated detection scenario.

# Contents

# List of Figures

# List of Tables