# Motion Segmentation Based Robust RGB-D SLAM

*Author:*
Youbing Wang

*Supervisor:*
A/Prof. Shoudong Huang

*A thesis submitted in fulfilment of the requirements
for the degree of Doctor of Philosophy*

*in the*

Centre for Autonomous Systems
School of Elec, Mech and Mechatronic Systems
Faculty of Engineering and Information Technology

May 2015

# Declaration of Authorship

I, Youbing WANG, declare that this thesis titled, 'Motion Segmentation Based Robust RGB-D SLAM' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.

- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.

- Where I have consulted the published work of others, this is always clearly attributed.

- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.

- I have acknowledged all main sources of help.

- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

# *Abstract*

## Motion Segmentation Based Robust RGB-D SLAM

by Youbing Wang

While research on simultaneous localisation and mapping (SLAM) in static environments can be regarded as a significant success due to intensive work during the last several decades, conducting SLAM, especially vision-based SLAM, in dynamic scenarios is still at its early stage. Although it seems like just one step further, the dynamic elements have brought in many unanticipated challenges, including motion detection, segmentation, tracking and $3D$ reconstruction of both the static environments and the moving objects, in addition to the handling of motion blur.

Solely based on RGB-D data with no prior knowledge available, this work centres upon proposing new practical solution frameworks for conducting SLAM in dynamic environments with efficient and robust motion segmentation methods serving as the basis. After a detailed review of the related achievements for SLAM in static environments as well as dynamic ones, and an analysis of the unaddressed challenges, four different motion segmentation methods, which include two 2-view sparse feature based motion segmentation algorithms, a 2-view semi-dense motion segmentation algorithm and an extended n-view dense moving object segmentation algorithm, are firstly proposed and their advantages, disadvantages and feasibility for different practical SLAM application scenarios are evaluated.

Based on the proposed motion segmentation methods, two kinds of solution frameworks for performing SLAM in dynamic scenarios are then put forward: the first one is formulated by integrating our sparse feature based motion segmentation techniques with the available pose-graph SLAM framework; and the other one is built upon dense moving object segmentation and tailored for dense SLAM. Related simulation and experimental results have demonstrated the effectiveness of our approaches.

# Acknowledgements

First of all, I would like to offer my sincere gratitude to my supervisor, A/Prof. Shoudong Huang, who has supported me during my PhD study continuously with his patience, motivation, enthusiasm and immense knowledge while allowing me the room to work in my own way. Without his encouragement, guidance, and efforts, this work would not have been possible.

Besides my supervisor, I would like to thank Prof. Dikai Liu and Prof. Gamini Dissanayake for their invaluable support for my study at the Centre for Autonomous System.

Sincere thanks also go to the other staff members of CAS, Dr. Jack Jianguo Wang, A/Prof. Sarath Kodagoda, A/Prof. Guang Hong, A/Prof. Jaime Valls Miro, Dr. Gabriel Aguirre-Ollinger, Dr. Liang Zhao and Dr. Lei Shi. They have given me support for my study and life in various ways.

In my daily work, I have been blessed with a friendly and cheerful group of fellow students: Gibson Hu, Kasra Khosoussi, Lakshitha Dantanarayana, Yue Wang, a visiting student from Zhejiang University, and many other people, for the stimulating discussions, for the sleepless nights we were working together before deadlines, and for all the fun we have had in the last several years.

Last but not the least, I would like to thank my family and all my friends, for their unconditional support to me spiritually throughout my life.

# Contents

# List of Figures

# List of Tables

# Abbreviations

| | |
|---|---|
| **SLAM** | **S**imultaneous **L**ocalization **A**nd Mapping |
| **KF** | **K**alman **F**ilter |
| **EKF** | **E**xtended **K**alman **F**ilter |
| **EIF** | **E**xtended **I**nformation **F**ilter |
| **SfM** | **S**tructure **f**rom Motion |
| **MSaM** | **M**ultiple **S**tructure **a**nd Motion |
| **RANSAC** | **RAN**dom **SA**mple **C**onsensus |
| **GPS** | **G**lobal **P**ositioning **S**ystem |
| **IMU** | **I**nertial **M**easurement **U**nit |
| **RMSE** | **R**oot-**M**ean-**S**quare **E**rror |
| **RPE** | **R**elative **P**ose **E**rror |
| **ATE** | **A**bsolute **T**ranslational **E**rror |
| **SMSaM** | **S**imultaneous **M**ultibody **S**tructure **a**nd Motion |
| **MDL** | **M**inimum **D**escription **L**anguage |
| **KLT** | **K**anade-**L**ucas-**T**omasi |

*Dedicated to my family...*