# Learning Based Active Rejection of Environmental Disturbances for Underwater Robots

by

Tianming Wang

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

at the
Centre for Autonomous Systems
Faculty of Engineering and Information Technology
**University of Technology Sydney**

June 2020

# Certificate of Original Authorship

I, Tianming Wang declare that this thesis, is submitted in fulfillment of the requirements for the award of Doctor of Philosophy, in the Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise reference 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.

Signed:
Production Note:
Signature removed prior to publication.

Date:          02 Jun 2020

# Learning Based Active Rejection of Environmental Disturbances for Underwater Robots

by

Tianming Wang

A thesis submitted in partial fulfilment of the requirements for the

degree of Doctor of Philosophy

# *Abstract*

Underwater robots in shallow waters usually suffer from turbulent flows and strong waves. Such disturbances may frequently exceed the robot's control constraints, thus severely destabilize robot during task operation. Conventional disturbance observer and model predictive control are not particularly effective since they heavily rely on a sufficiently accurate dynamics model. Learning-based controllers are able to alleviate model dependency and achieve high computational efficiency. Learned control policies normally specialize on one dynamics model, and may not directly generalize to other models. Transfer learning offers a pathway to bridge the mismatch between different dynamics models. In this thesis, reinforcement learning algorithms are applied that enables optimal control of underwater robots under unobservable excessive time-correlated disturbances, and transfer learning algorithms are implemented for control policy adaptation under dynamics model mismatch.

History Window Reinforcement Learning (HWRL) and Disturbance Observer Network (DOB-Net) are developed for disturbance rejection control. Both algorithms jointly optimize a disturbance observer and a motion controller, and implicitly learn embedding of disturbance waveforms from motion history of robot. A modular design of learning disturbance rejection controller is also developed. A Generalized Control Policy (GCP) is trained over a wide range of disturbance waveforms, an Online Disturbance Identification

Model (ODI) exploits motion history of robot to predict the disturbance waveforms, which served as input to GCP. Together, GCP-ODI provides robust control across a wide variety of disturbances.

Transfer learning algorithms are applied to address the mismatch between a mathematical model of system dynamics developed from the fundamental principles of dynamics and an empirical model of system dynamics derived from real-world experimental data. Hybrid Policy Adaptation (HPA) is first proposed where learning a model-free policy under the empirical model is accelerated by pre-training a model-based policy with the mathematical model. Transition Mismatch Learning (TML) is then proposed that learns a compensatory policy based on the modular architecture of GCP-ODI through minimizing transition mismatch between the mathematical model and the empirical model.

Numerical simulations on a pose regulation task have demonstrated that HWRL, DOB-Net and GCP-ODI can successfully stabilize the underwater robot across a wide range of disturbance waveforms, and outperform conventional controllers and classical RL policies. Both HPA and TML achieve satisfactory control performance when deployed under the empirical model, with high sample efficiency and avoidance of initial exploratory actions.

# *Acknowledgements*

# Contents

# List of Figures

# List of Tables

# Acronyms & Abbreviations

**UTS**  University of Technology Sydney

**CAS**  Centre for Autonomous Systems

**AUV**  Autonomous Underwater Vehicle

**ROV**  Remotely Operated Vehicle

**SPIR**  Submerged Pile Inspection Robot

**IMU**  Inertial Measurement Unit

**DOF**  Degree of Freedom

**TCM**  Thruster Control Matrix

**PWM**  Pulse Width Modulation

**DOB**  Disturbance Observer

**DOBC**  Disturbance Observer Based Control

**RISE**  Robust Integral of Sign Error

**MPC**  Model Predictive Control

**LQR**  Linear Quadratic Regulator

**iLQR**  iterative Linear Quadratic Regulator

**EKF**  Extended Kalman Filter

**MDP**  Markov Decision Process

**POMDP**      Partially Observable Markov Decision Process

**TD**             Temporal Difference

**DQN**          Deep Q-Network

**DDPG**        Deep Deterministic Policy Gradient

**A2C**           Advantage Actor Critic

**TRPO**         Trust Region Policy Optimization

**PPO**           Proximal Policy Optimization

**SGD**          Stochastic Gradient Descent

**RNN**          Recurrent Neural Network

**GRU**          Gated Recurrent Unit