# Certificate of Original Authorship Template

**Graduate research students are required to make a declaration of original authorship when they submit the thesis for examination and in the final bound copies. Please note, the Research Training Program (RTP) statement is for all students.** The Certificate of Original Authorship must be placed within the thesis, immediately after the thesis title page.

**Required wording for the certificate of original authorship**

CERTIFICATE OF ORIGINAL AUTHORSHIP

I, Zishuo Cheng, declare that this thesis is submitted in fulfilment 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 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: 6 March 2023

# Differential Privacy in Multi-agent Reinforcement Learning

by

Zishuo Cheng

A thesis submitted in fulfilment of the requirements for the

degree of Doctor of Philosophy

# *Abstract*

Due to the increasing demand for autonomous systems, e.g., multi-robot systems and autonomous driving systems, multi-agent reinforcement learning, technology attracts considerable attention. In the field of multi-agent reinforcement learning, agent advising by allowing agents to ask for or give advice to others has become an increasingly important topic as it can significantly promote agents' learning speed with negligible computation overheads. However, there are some critical challenging issues in agent advising learning, particularly performance and privacy issues. Differential privacy is a promising privacy-preserving model with several valuable properties. Its privacy-preserving property provides a provable guarantee to privacy protection while its randomisation property helps to resist the inference in machine learning schemes. Therefore, this thesis aims to explore the feasibility of adopting differential privacy mechanisms to resolve the two research challenges in multi-agent reinforcement learning. In summary, this thesis consists of the following four contributions:

- A differential advising method is proposed, which allows agents to use a piece of advice in various states. This method is implemented by using differential privacy technology to mask the difference in states. In this way, agents' learning performance can be remarkably improved while their communication overheads can be significantly reduced.

- A differential knowledge transfer method is proposed, which stimulates the learning performance in a homogeneous multi-agent reinforcement learning system. This method jointly utilises the randomisation property of differential privacy and relevance weight to mitigate

the interference of negative transfer in a multi-agent reinforcement learning environment. By acquiring a greater number of relevant sample sets, agents' learning rates can be largely improved.

- A novel time-drive and privacy-preserving navigation learning scheme is proposed for multi-agent vehicular communication. This learning scheme, which is particularly suitable for multi-agent deep reinforcement learning systems, consists of three parts: 1) a function for estimating the average traffic flow of a segment; 2) a function for calculating a valid route; 3) confidence weights to estimate the accuracy of the estimated traffic flow. Second, we adopt a customised $\epsilon$-differentially private mechanism for the *TDPP* model. To the best of our knowledge, this is the first navigation system for vehicle-to-vehicle systems to provably guarantee the location privacy of vehicles. Third, we theoretically prove that *TDPP* satisfies the definition of $\epsilon$-differential privacy, accompanied by extensive experiments examining its performance.

- A novel multi-agent reinforcement learning model that jointly adopts deep reinforcement learning and differential privacy is proposed for evolutionary game theory, which promotes cultivating more cooperators while protecting agents' sensitive information. First, this is the first evolutionary cooperation method which takes the privacy of agents into account. By adopting differential privacy, *NNDP* protects agents' private information while still encouraging cooperation with negligible performance reduction. Additionally, compared to encryption methods, differentially private mechanisms use much less time and computational resources. Second, *NNDP* pioneers the use of deep learning to promote cooperation. By adopting a deep reinforcement learning algorithm, the *NNDP* offers a more adaptive, generalised framework and a framework with greater stability in dynamic situations. Third, we theoretically prove that *NNDP* satisfies the definition of differential privacy, accompanied by extensive experiments to examine its performance.

# *Acknowledgements*

I would like to show my deepest gratitude to my supervisor, Prof. Tianqing Zhu, Prof. Wanlei Zhou and Dr. Dayong Ye. They give me a lot of help during the course of this course. Especially Prof. Tianqing Zhu, She gave me a lot of support in my study. I am also very grateful to my parents for their support over the past three years. Finally, I would like to express my heartfelt gratitude once again to those who have helped me.

# Contents

# List of Figures

# List of Tables

# List of Publications

1. **Z. Cheng**, D. Ye, T. Zhu, W. Zhou, P. Yu, 'Multi-Agent Reinforcement Learning via Knowledge Transfer with Differentially Private Noise', International Journal of Intelligent Systems, published

2. **Z. Cheng**, D. Ye, T. Zhu, W. Zhou, T. Zhang, 'Evolutionary Cooperation in Neural Network Games with Differentially Private Mechanisms', IEEE Trans on System, Man and Cybernetics, in the second round review after major revision.

3. **Z. Cheng**, D. Ye, T. Zhu, C. Zhu, W. Zhou, 'Time-Driven and Privacy-Preserving Navigation Model for V2V Communications', IEEE Trans on Vehicular and Technology, under review.

4. D. Ye, T. Zhu, **Z. Cheng**, W. Zhou, P. Yu, 'Differentially-Private Knowledge Transfer for Reinforcement Learning', IEEE Trans on Cybernetics, published.

5. T. Zhu, D. Ye, **Z. Cheng**, W. Zhou, "Learning Games for Defending Advanced Persistent Threats in Cyber Systems", IEEE Transactions on Systems, Man and Cybernetics: Systems, in the second round review after major revision.

6. T. Zhu, W. Zhou, D. Ye, **Z. Cheng**, J. Li, 'Resource Allocation in IoT Edge Computing via Concurrent Federated Reinforcement Learning', IEEE Internet of Things, published.

7. S. Abahussein, **Z. Cheng**, T. Zhu, D. Ye, W. Zhou, 'Privacy-preserving in Double Deep Q-Network with Differential Privacy in Continuous Spaces', 2021 International Conference on Web-based Learning (ICWL), published.

8. S. Abahussein, D. Ye, **Z. Cheng**, T. Zhu, W. Zhou, 'Differential privacy to protect privacy in Dueling Deep-Q-Network', Concurrency and Computation Practice and Experience, under review.