UNIVERSITY OF TECHNOLOGY SYDNEY

Faculty of Engineering and Information Technology

# Advanced Techniques of Cross Domain Translation Learning

by

**Wanming Huang**

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

**Doctor of Philosophy**

Sydney, Australia

2020

# Certificate of Original Authorship

I, Wanming Huang, declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the Faculty of Engineering and IT 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: 17/09/2020

# ABSTRACT

**Advanced Techniques of Cross Domain Translation Learning**

by

Wanming Huang

Cross domain translation, such as image captioning, fashion synthesis from text descriptions, music composition in a particular style, has attracted considerable interest in the deep learning community lately. Despite significant progress in this field, certain drawbacks in previous methods have been identified. First, although the attention mechanism has been widely applied to domain transfer and achieved remarkable outcomes, cross-domain translation remains an open research question on cross domain transfer learning because of the different data structures. Second, most domain translation algorithms address only a pair of domains, and there is a need for $2 \times \binom{N}{2}$ transfer functions given $N$ image domains. This makes training prohibitively unmanageable. We have proposed a set of solutions to solve these two problems, as described in detail in Chapter 3. Third, most generative model based domain-transfer algorithms uses single-mode distribution to model the latent space. This does not work well on datasets that contain diversified samples that form multiple clusters. Our study applies mixture models to cross-domain generation, of which the effects and properties are illustrated in Chapter 4. Finally, cross-domain translation models usually suffer from long training time and are difficult to converge. Indeed, this applies to most deep neural network training that involves complex network designs and large datasets. Our work in Chapter 5 accelerates deep neural network training with a specially designed mini-batch sampling strategy.

Dissertation directed by Associate Professor Richard Yi Da Xu

School of Electrical and Data Engineering

# Acknowledgements

First and foremost, I would like to express deep feelings of gratitude to my supervisor, Associate Professor Richard Yi Da Xu, who has guided me through my PhD journey. He has taught me how to conduct research, and how interesting and amazing machine learning can be. He has spent time with me and my colleagues from day one, taught us calculus, statistics and all other mathematical theories required to conduct machine learning research. He has also shared his insights and the latest breakthrough in many areas of machine learning, which enlightens our own research. Further, he has provided guidance to each one of my paper and my thesis, not only on research ideas, but also on the structuring and writing. He helped me overcome each obstacle I met in my research. I appreciate all his time, effort, and funding that supported my PhD research. I am also thankful for his dedication and hardworking towards research, which exemplified moral and ethical research.

I would also like to thank my external supervisor Dr. Ian Oppermann. Ian and I have arranged meetings regularly to discuss my research progress and exchange ideas on the application of machine-learning algorithms to industry projects. I am deeply grateful for these meetings because they provided insights into my research. Ian also shared his own experience in research and discussed with me possible applications of my work in the industry, which shed light upon my future career path. Discussions with Ian further allowed me to examine problems in a more systematic way and realised what could be improved and what was left to be addressed. He also provided me with valuable advice on PhD study and thesis structuring.

The members of the research group have also contributed immensely to my research and personal life. Discussions with my colleagues, Shuai Jiang, Xuan Liang, Chen Deng, Yi Huang, Ying Li, Wei Huang and all other group members have allowed me to conquer numerous challenges. We have collaborated in multiple

projects, and I greatly appreciate the inspiration and suggestions they provided from multiple fields. Shuai Jiang helped me a great deal with optimisation and statistics; Xuan Liang has enlightened me with his experience in neural translation models. I would also like to thank Haodong Chang and Wei Huang for their help in finding a proofreader.

Regarding my research in latent space modelling of generative adversarial networks with applications in anomaly detection, I would like to thank Dr. Hongbo Xie from the University of Queensland. He not only provided me with his knowledge and experience from a statistical perspective, he has also provided me with valuable advice on paper writing.

Further, I would like to thank Elite Editing for their effort in proofreading this thesis, and editorial intervention was restricted to Standards D and E of the *Australian Standards for Editing Practice*.

Finally, I would like to thank my parents for always being supportive of my research and my life.

<div align="right">

Wanming Huang
Sydney, Australia, 2020.

</div>

# List of Publications

**Conference Papers**

C-1. W. Huang, R. Y. D. Xu, and I. Oppermann, "Realistic image generation using region-phrase attention," in Proceedings of The Eleventh Asian Conference on Machine Learning, ser. Proceedings of Machine Learning Research, W. S. Lee and T. Suzuki, Eds., vol. 101. Nagoya, Japan: PMLR, 17–19 Nov 2019, pp. 284–299.

C-2. W. Huang, R. Y. D. Xu, and I. Oppermann, "Efficient diversified mini-batch selection using variable high-layer features," in Proceedings of The Eleventh Asian Conference on Machine Learning, ser. Proceedings of Machine Learning Research, W. S. Lee and T. Suzuki, Eds., vol. 101. Nagoya, Japan: PMLR, 17–19 Nov 2019, pp. 300–315.

C-3. "GAN-based Gaussian Mixture Model Responsibility Learning" accepted by ICPR 2020

# Contents

# 4   Latent Space Modelling in Style Transfer     75

# 5   Training Acceleration in Style Transfer     117

# List of Figures

# Abbreviation

3D - three-dimensional

AUC - area under the receiver operating characteristic curve

bi-LSTM - bidirectional long short-term memory

CUB - Caltech-USCD Bird

CNN - convolutional neural network

DCGAN - deep convolutional generative adversarial network

DM-SGD - Diversified Mini-Batch SGD

DNN - deep neural network

DP - Dirichlet process

DPGMM - Dirichlet process Gaussian mixture model

DPP - determinantal point process

FC - fully connected

FID - Fréchet Inception Distance

GAN - generative adversarial network

GAWWN - generative adversarial what-where network

GLU: gated linear unit

GM-GAN: Gaussian mixture generative adversarial network

GMM: Gaussian Mixture Model

GRU: gated recurrent unit

IS: inception score

LSTM - long short-term memory

MSCOCO - Microsoft COCO dataset

NF - normalizing flow

NMT - neural machine translation

NLP - natural language processing

PCM - posterior consistency module

RCM - responsibility consistency module

RCNN - region-based convolutional neural network

ReLU: rectified linear unit

RNN - recurrent neural network

ROC - receiver operating characteristic

RoI - region of interest

SGD - stochastic gradient descent

SMC - sequential Monte Carlo

SVRG - stochastic variance reduced gradient

VAE - variational autoencoder

# Nomenclature and Notation

Capital letters denote matrices.

Lower-case alphabets denote column vectors.

$(.)^T$ denotes the transpose operation.

$I_n$ is the identity matrix of dimension $n \times n$.

$0_n$ is the zero matrix of dimension $n \times n$.

$\mathbb{R}$, $\mathbb{R}^+$ denote the field of real numbers, and the set of positive reals, respectively.