# Generative and Discriminative Learning for Visual Matching

**by Zhedong Zheng**

Thesis submitted in fulfilment of the requirements for the degree of

**Doctor of Philosophy**

under the supervision of Prof. Yi Yang and Dr. Liang Zheng

University of Technology Sydney
Faculty of Engineering and Information Technology

May 2021

# CERTIFICATE OF ORIGINAL AUTHORSHIP

I, Zhedong Zheng, declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the School of Computer Science at 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.

Production Note:
SIGNATURE: Signature removed prior to publication.

DATE: 24, May, 2021

# ABSTRACT

## Generative and Discriminative Learning
## for Visual Matching

by

Zhedong Zheng

Visual matching aims to establish image correspondences across viewpoints. Given a query image, the visual matching system seeks to retrieve images containing the object of interest from non-overlapping viewpoints according to the similarity score. The visual matching task remains challenging because objects captured by different viewpoints often contain significant intra-class variations caused by background, viewpoint, object pose, etc. In this thesis, I present my research on combining generative learning with discriminative learning to build one robust visual matching system. First, due to lack of sufficient data to enhance robustness against input variations, generative learning is aimed at letting the model potentially "see" these variations (particularly intra-class variations) during training. With recent progress in the generative adversarial networks (GANs), generative models have become appealing choices to introduce additional augmented data for free. Second, discriminative learning is designed to formulate visual matching as a metric learning problem and adopt the discriminative optimization objective to learn the distance. With these objectives in mind, it motivates us to enable Convolutional Neural Network (CNN) to learn the mapping function to discriminate between different objects. In this thesis, I investigate two scientific problems of combining two learning strategies: 1) How to obtain high-quality generated data for subsequential training? 2) How to leverage the generated data to promote discriminative learning? To study the two problems, I explore improving learned visual representations by better leveraging the data from the following three aspects.

First, we present a semi-supervised pipeline that integrates GAN-generated images into discriminative learning. It contains a generative adversarial model for unsupervised data generation and a discriminative convolutional neural network for semi-supervised learning. Second, we observe that the generative pipelines are typically presented as standalone models, which are relatively separate from the discriminative learning models. To make the best of the two worlds, we further propose a learning framework that couples discriminative and generative learning. This design leads to a unified framework that enables the interactions between generative and discriminative modules in an end-to-end manner. Third, we further investigate different discriminative learning approaches on various data sources. Specifically, we study the feasibility of borrowing the knowledge from real-world vehicle images collected on the web and propose a two-stage learning strategy to minimize the domain gap between the web data and real-world data. Furthermore, we also explore the possibility of learning from synthetic data simulated by 3D engines. We propose a new geo-localization benchmark and build a strong and flexible baseline to learn from multi-view multi-source data.

In summary, this thesis studies and solves the critical challenges of data limitation and robust representation learning in visual matching. We show the benefits of leveraging the generative and discriminative learning in deep learning, which achieves better performance than previous methods.

Dissertation directed by Professor Yi Yang
The Australian Artificial Intelligence Institute (AAII), School of Computer Science

# Acknowledgements

First, I would like to thank my supervisor, Prof. Yi Yang, for his encouragement, guidance, and patience. He has provided his perpetual supports to help me overcome difficulties and guided me to pursue challenging problems. He also taught me how to re-consider the same thing from different views, which is valuable to my future career. I have learned a lot under his caring supervision, and he has my deepest gratitude. Besides, I would like to thank my co-supervisor, Dr. Liang Zheng. He has provided me invaluable suggestions and consistent supports throughout my Ph.D. life. Thanks to Xiaodong Yang, Zhiding Yu, and Jan Kautz, my mentors when I interned at Nvidia Research, from whom I learned critical thinking and time management. Thanks to Mingyue Jiang and Xiao Tan, my mentors when I remotely interned at Baidu Research, from whom I learned the competition skills and knowledge. I also would like to thank all colleagues in the ReLER Lab at the University of Technology Sydney for their friendship and supports. I was fortunate to work with them.

Finally, I would like to thank my parents Yanmin Yang and Gaoxin Zheng, my grandma Wenlan Tang for loves and supports over the years.

<div align="right">

Zhedong Zheng

Sydney, Australia, 2021.

</div>

# List of Publications

**Journal Papers**

J-1. **Zhedong Zheng**, Liang Zheng, Yi Yang (2018). A discriminatively learned cnn embedding for person reidentification. ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM), 14(1), 1-20. [**ESI Highly-cited Paper**]

J-2. **Zhedong Zheng**, Liang Zheng, Yi Yang (2018). Pedestrian alignment network for large-scale person re-identification. IEEE Transactions on Circuits and Systems for Video Technology (TCSVT), 29(10), 3037-3045. [**ESI Highly-cited Paper**]

J-3. **Zhedong Zheng**, Liang Zheng, Michael Garrett, Yi Yang, Mingliang Xu, Yi-Dong Shen (2020). Dual-Path Convolutional Image-Text Embeddings with Instance Loss. ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM), 16(2), 1-23.

J-4. **Zhedong Zheng**, Tao Ruan, Yunchao Wei, Yi Yang, Mei Tao (2020). VehicleNet: Learning Robust Visual Representation for Vehicle Re-identification. IEEE Transactions on Multimedia (TMM).

J-5. **Zhedong Zheng**, Yi Yang (2021). Rectifying Pseudo Label Learning via Uncertainty Estimation for Domain Adaptive Semantic Segmentation. International Journal of Computer Vision (IJCV), 1-15.

J-6. Zhun Zhong, Liang Zheng, **Zhedong Zheng**, Shaozi Li, Yi Yang (2018). Camstyle: A novel data augmentation method for person re-identification. IEEE Transactions on Image Processing (TIP), 28(3), 1176-1190. [**ESI Highly-cited Paper**]

J-7. Yan Huang, Jinsong Xu, Qiang Wu, **Zhedong Zheng**, Zhaoxiang Zhang, Jian

Zhang (2018). Multi-pseudo regularized label for generated data in person re-identification. IEEE Transactions on Image Processing (TIP), 28(3), 1391-1403.

J-8. Yutian Lin, Liang Zheng, **Zhedong Zheng**, Yu Wu, Zhilan Hu, Chenggang Yan, Yi Yang (2019). Improving person re-identification by attribute and identity learning. Pattern Recognition, 95, 151-161. [**ESI Highly-cited Paper**]

J-9. Qingji Guan, Yaping Huang, Zhun Zhong, **Zhedong Zheng**, Liang Zheng, Yi Yang (2020). Thorax disease classification with attention guided convolutional neural network. Pattern Recognition Letters, 131, 38-45.

J-10. Yutian Lin, **Zhedong Zheng**, Hong Zhang, Chenqiang Gao, Yi Yang (2020). Bayesian query expansion for multi-camera person re-identification. Pattern Recognition Letters, 130, 284-292.

J-11. Bingwen Hu, **Zhedong Zheng**, Ping Liu, Wankou Yang, Mingwu Ren (2020). Unsupervised eyeglasses removal in the wild. IEEE Transactions on Cybernetics, 1 - 13.

J-12. Tingyu Wang, **Zhedong Zheng**, Chenggang Yan, Jiyong Zhang, Yaoqi Sun, Bolun Zheng, Yi Yang (2021). Each Part Matters: Local Patterns Facilitate Cross-view Geo-localization. IEEE Transactions on Circuits and Systems for Video Technology (TCSVT).

**Conference Papers**

C-1. **Zhedong Zheng**, Liang Zheng, Yi Yang (2017). Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In Proceedings of the IEEE International Conference on Computer Vision (ICCV) (pp. 3754-3762). [Spotlight]

C-2. **Zhedong Zheng**, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, Jan Kautz (2019). Joint discriminative and generative learning for person re-identification. In Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR) (pp. 2138-2147). [Oral]

C-3. **Zhedong Zheng**, Yunchao Wei, Yi Yang (2020). University-1652: A Multi-view Multi-source Benchmark for Drone-based Geo-localization. In Proceedings of the ACM Multmieda (ACM MM) (pp. 1395–1403).

C-4. **Zhedong Zheng**, Yi Yang (2020). Unsupervised Scene Adaptation with Memory Regularization in vivo. In Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI) (pp. 1076-1082).

C-5. Yawei Luo, **Zhedong Zheng**, Liang Zheng, Tao Guan, Junqing Yu, Yi Yang (2018). Macro-micro adversarial network for human parsing. In Proceedings of the European Conference on Computer Vision (ECCV) (pp. 418-434).

C-6. Zhun Zhong, Liang Zheng, **Zhedong Zheng**, Shaozi Li, Yi Yang (2018). Camera style adaptation for person re-identification. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (pp. 5157-5166).

C-7. Zhikun Huang, **Zhedong Zheng**, Chenggang Yan, Hongtao Xie, Yaoqi Sun, Jianzhong Wang, Jiyong Zhang (2020). Real-World Automatic Makeup via Identity Preservation Makeup Net. In Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI) (pp. 652-658).

C-8. Chuchu Han, **Zhedong Zheng**, Changxin Gao, Nong Sang, Yi Yang (2021). Decoupled and Memory-Reinforced Networks: Towards Effective Feature Learning for One-Step Person Search. In Proceedings of the AAAI Conference on Artificial Intelligence (AAAI).

The source codes of most works can be found at `https://github.com/layumi` .

# Contents

# 5   Two-stage Progressive Learning                                   61

# 6   Joint Discriminative and Generative Learning                    86

# List of Figures