UNIVERSITY OF TECHNOLOGY SYDNEY

Faculty of Engineering and Information Technology

# HUMAN GAIT RECOGNITION UNDER CHANGES OF WALKING CONDITIONS

by

**Lingxiang Yao**

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

**Doctor of Philosophy**

Sydney, Australia

2021

# Certificate of Authorship/Originality

I, Lingxiang Yao, declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the School of Electrical and Data Engineering, 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.

This research is supported by the Australian Government Research Training Program.

Signature: _Production Note: Signature removed prior to publication._

Date: _Dec/14/2021_

# Acknowledgements

It is a cherished memory to pursue PhD in UTS. I would like to express my sincere gratitude to all those who help me complete my doctoral study.

Foremost, I would like to express my sincere gratitude to my principal supervisor, A/Prof. Jian Zhang, for his professional guidance and persistent support throughout my PhD study. His enthusiasm, attitude, and devotion towards academy have deeply influenced me, which provides instructions to my future career.

I would also like to appreciate my co-supervisor A/Prof. Qiang Wu for his advice and guidance in both research and lives in Sydney. He always encourages me to focus on the advanced techniques that motivate the research going deeper.

I also deeply appreciate my friend Worapan Kusakunniran in Mahidol University for his significant support. He always spends a lot of time to help me polish my paper, and gives me valuable advice to improve my academic skills.

And, I appreciate the help, support and friendship from my dear colleagues and friends during my doctoral study. Thanks to Jingsong Xu, Yucheng Wang, Xiaoshui Huang, Yazhou Yao, Peng Zhang, Junjie Zhang, Muming Zhao, Zhibin Li, Yongshun Gong, Lu Zhang, Huaxi Huang, Anan Du, Xunxiang Yao, Yan Huang, Guofeng Mei, and all other labmates for their collaboration and discussion. I am also grateful to all my friends in Sydney, Razia Osman, Yumeng Yao, Xin Ba, Jiancheng Li, and Xiao Li for their encouragement and companion.

Finally, I would like to express my sincere thanks to my parents and girlfriend Xia Zhang for their support, trust, encouragement and love throughout my PhD studies..

Lingxiang Yao

August 2021 @ UTS.

# List of Publications

**Published Papers**

C-1. **Lingxiang Yao**, Worapan Kusakunniran, Qiang Wu, Jian Zhang and Jingsong Xu. "Part-based Collaborative Spatio-temporal Feature Learning for Cloth-changing Gait Recognition." 2020 25th International Conference on Pattern Recognition (ICPR) (2021): 2057-2064.

J-1. **Lingxiang Yao**, Worapan Kusakunniran, Qiang Wu and Jian Zhang. "Gait recognition using a few gait frames." PeerJ Computer Science 7 (2021).

J-2. **Lingxiang Yao**, Worapan Kusakunniran, Qiang Wu, Jian Zhang, Zhen-min Tang and Wankou Yang. "Robust gait recognition using hybrid descriptors based on Skeleton Gait Energy Image." Pattern Recognition Letters 150 (2021): 289-296.

C-2. Junyi Wu, **Lingxiang Yao**, Y. Huang, J. Xu, Qiang Wu and Liqin Huang. "Improving Person Re-Identification Performance Using Body Mask Via Cross-Learning Strategy." 2019 IEEE Visual Communications and Image Processing (VCIP) (2019): 1-4.

C-3. Anan Du, Xiaoshui Huang, Jiayuan Zhang, **Lingxiang Yao** and Qiang Wu. "Kpsnet: Keypoint Detection and Feature Extraction for Point Cloud Registration." 2019 IEEE International Conference on Image Processing (ICIP) (2019): 2576-2580.

C-4. **Lingxiang Yao**, Worapan Kusakunniran, Qiang Wu, J. Zhang and Zhen-min Tang. "Robust CNN-based Gait Verification and Identification using Skeleton Gait Energy Image." 2018 Digital Image Computing: Techniques and Applications (DICTA) (2018): 1-7.

**Accepted Papers**

J-1. **Lingxiang Yao**, Worapan Kusakunniran, Qiang Wu, Jian Zhang and Jingsong Xu. "Collaborative Feature Learning for Gait Recognition under Cloth Changes," IEEE Transactions on Circuits and Systems for Video Technology, 2021.

J-2. **Lingxiang Yao**, Worapan Kusakunniran, Qiang Wu, Jingsong Xu and Jian Zhang. "Recognizing Gaits acrossWalking and Running Speedes", ACM Transactions on Multimedia Computing Communications and Applications, 2021.

# Contents

# List of Figures

# List of Tables

# ABSTRACT

## HUMAN GAIT RECOGNITION UNDER CHANGES OF WALKING CONDITIONS

by

Lingxiang Yao

Gait has been gathering extensive research interest for its non-fungible position in applications, *e.g.*, security surveillance and forensic identification. First, it is difficult to disguise one's gait, since walking is necessary for human mobility. Second, it works well in an unconstrained condition and can be attained at a distance without physical contact or proximal sensing. However, although recently different methods have been proposed for gait recognition, gait analysis is still in its infancy. Most methods enable to garner a remarkable recognition performance when the gallery and the probe are in a similar situation. However, when exterior factors affect a person's gait and changes occur in human appearances, a significant performance degradation happens.

Among these exterior factors, clothing variations and mode changes can be treated as the most influential factors for gait recognition. It is advisable to identify a person using gait, since each person exhibits his/her walking patterns in a sufficiently unique and fairly characteristic way. However, clothing variations can significantly influence available features to be used in the future recognition process, while walking/running modes can change human motions made by limbs and thus dramatically influence the instinct walking patterns of each person. Hence, in this thesis different methods have been proposed for gait recognition to handle the difficulties of clothing variations and walking/running mode changes.

First, given that model-based methods are less vulnerable to clothing variances, a more robust model-based gait feature, Skeleton Gait Energy Image (SGEI), is formed to handle this cloth-changing gait recognition problem. Then, since clothing changes

can cause different impacts to different body parts, a part-based collaborative spatio-temporal feature learning method is also proposed for cloth-changing gait recognition by concatenating features from the non/less affected body parts under the correlative $H-W$ and $T-W$ views. Based on the aforementioned two methods, another efficient network is proposed for cloth-changing gait recognition. This network consists of two sub-networks, aiming to produce part-based features from the non/less affected body parts and the estimated skeleton key-point regions. Moreover, in order to address the walking-vs-running problem in a cross-mode manner, a feasible hybrid method is also proposed in this thesis. Distinct from most cross-mode gait recognition methods, this method focuses on learning mode-invariant features for each person from their innate patterns between walking and running modes. Multi-task learning strategies are also used to enhance the efficiency of these learned features. Finally, given that the above-mentioned methods are all proposed based on sufficient input data, a complementary solution is given when only a few gait frames can be offered.

To sum up, the main objective of this thesis is to address the problems of clothing variations and walking/running mode changes for gait recognition, thus four different methods have been proposed in this thesis. Besides, related experiments have proved that these proposed methods can obtain a remarkable performance when tackling the cloth-changing and walking-vs-running gait recognition problems.