# Robust Face Recognition

Changxing Ding

Faculty of Engineering and Information Technology

University of Technology Sydney

A thesis submitted for the degree of

*Doctor of Philosophy*

July, 2016

# Certificate of Original Authorship

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Student: Changxing Ding

Date: 06/07/2016

I would like to dedicate this thesis to my loving wife and parents.

# Acknowledgements

Many people have significantly influenced me - both academically and personally - during my time at University of Technology Sydney (UTS). I would like to take this opportunity to express my sincere gratitude to them for their support during my PhD study.

My deepest gratitude goes to Prof. Dacheng Tao, who has been not just a great supervisor but also a sincere friend. Prof. Tao has taught me how to perform research from scratch and has guided me into academia. He has been a constant source of suggestions and inspiration, and I have benefitted significantly from our discussions and his mentorship. I would not have been able to publish scientific papers in top journals without his insightful instructions, high scientific standards, patient guidance, and consistent encouragement. Also, many thanks to his guidance and suggestions for my future career during the final stage of my PhD study.

I owe special thanks to Prof. Larry S. Davis at the University of Maryland, College Park (UMD) for a number of insightful discussions and suggestions on my first work on face recognition. His careful criticism bridging theory and application broadened my thinking and helped me to improve the work. I would also like to thank Dr. Jonghyun Choi from UMD for his collaboration on experiments and developing the first work. I also thank Dr. Chang Xu, Beijing University, who answered some of my general questions about machine learning and helped me develop the model for pose-invariant face recognition.

I have been so fortunate to work in the Centre for Quantum Computation and Intelligent Systems (QCIS). I appreciate the support of the center director, Prof. Chengqi Zhang, who has led QCIS to be a leading research center where I had the chance to meet and get to know many world-famous

# Abstract

Face recognition is one of the most important and promising biometric techniques. In face recognition, a similarity score is automatically calculated between face images to further decide their identity. Due to its non-invasive characteristics and ease of use, it has shown great potential in many real-world applications, e.g., video surveillance, access control systems, forensics and security, and social networks. This thesis addresses key challenges inherent in real-world face recognition systems including pose and illumination variations, occlusion, and image blur. To tackle these challenges, a series of robust face recognition algorithms are proposed. These can be summarized as follows:

In Chapter 2, we present a novel, manually designed face image descriptor named "Dual-Cross Patterns" (DCP). DCP efficiently encodes the seconder-order statistics of facial textures in the most informative directions within a face image. It proves to be more descriptive and discriminative than previous descriptors. We further extend DCP into a comprehensive face representation scheme named "Multi-Directional Multi-Level Dual-Cross Patterns" (MDML-DCPs). MDML-DCPs efficiently encodes the invariant characteristics of a face image from multiple levels into patterns that are highly discriminative of inter-personal differences but robust to intra-personal variations. MDML-DCPs achieves the best performance on the challenging FERET, FRGC 2.0, CAS-PEAL-R1, and LFW databases.

In Chapter 3, we develop a deep learning-based face image descriptor named "Multimodal Deep Face Representation" (MM-DFR) to automatically learn face representations from multimodal image data. In brief, convolutional neural networks (CNNs) are designed to extract

complementary information from the original holistic face image, the frontal pose image rendered by 3D modeling, and uniformly sampled image patches. The recognition ability of each CNN is optimized by carefully integrating a number of published or newly developed tricks. A feature level fusion approach using stacked auto-encoders is designed to fuse the features extracted from the set of CNNs, which is advantageous for non-linear dimension reduction. MM-DFR achieves over 99% recognition rate on LFW using publicly available training data.

In Chapter 4, based on our research on handcrafted face image descriptors, we propose a powerful pose-invariant face recognition (PIFR) framework capable of handling the full range of pose variations within $\pm 90°$ of yaw. The framework has two parts: the first is Patch-based Partial Representation (PBPR), and the second is Multi-task Feature Transformation Learning (MtFTL). PBPR transforms the original PIFR problem into a partial frontal face recognition problem. A robust patch-based face representation scheme is developed to represent the synthesized partial frontal faces. For each patch, a transformation dictionary is learnt under the MtFTL scheme. The transformation dictionary transforms the features of different poses into a discriminative subspace in which face matching is performed. The PBPR-MtFTL framework outperforms previous state-of-the-art PIFR methods on the FERET, CMU-PIE, and Multi-PIE databases.

In Chapter 5, based on our research on deep learning-based face image descriptors, we design a novel framework named Trunk-Branch Ensemble CNN (TBE-CNN) to handle challenges in video-based face recognition (VFR) under surveillance circumstances. Three major challenges are considered: image blur, occlusion, and pose variation. First, to learn blur-robust face representations, we artificially blur training data composed of clear still images to account for a shortfall in real-world video training data. Second, to enhance the robustness of CNN features to pose variations and occlusion, we propose the TBE-CNN architecture, which efficiently extracts complementary information from holistic face images and patches cropped around facial components. Third, to further promote

the discriminative power of the representations learnt by TBE-CNN, we propose an improved triplet loss function. With the proposed techniques, TBE-CNN achieves state-of-the-art performance on three popular video face databases: PaSC, COX Face, and YouTube Faces.

# Contents

# List of Figures

# List of Tables