# Feature Fusion, Feature Selection and Local N-ary Patterns for Object Recognition and Image Classification

*Author:*

Sheng WANG

*Supervisor:*

Qiang WU, Xiangjian HE

# 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.

Signature of Student:

_____

Date:

_____

*"Concern for man and his fate must always form the chief interest of all technical endeavors. Never forget this in the midst of your diagrams and equations."*

Albert Einstein

# *Abstract*

Doctor of Philosophy

**Feature Fusion, Feature Selection and Local N-ary Patterns for Object Recognition and Image Classification**

by Sheng WANG

Object recognition is one of the most fundamental topics in computer vision. During past years, it has been the interest for both academies working in computer science and professionals working in the information technology (IT) industry. The popularity of object recognition has been proven by its motivation of sophisticated theories in science and wide spread applications in the industry. Nowadays, with more powerful machine learning tools (both hardware and software) and the huge amount of information (data) readily available, higher expectations are imposed on object recognition. At its early stage in the 1990s, the task of object recognition can be as simple as to differentiate between object of interest and non-object of interest from a single still image. Currently, the task of object recognition may as well includes the segmentation and labeling of different image regions (i.e., to assign each segmented image region a meaningful label based on objects appear in those regions), and then using computer programs to infer the scene of the overall image based on those segmented regions. The original two-class classification problem is now getting more complex as it now evolves toward a multi-class classification problem. In this thesis, contributions on object recognition are made in two aspects. These are, improvements using feature fusion and improvements using feature selection. Three examples are given in this thesis to illustrate three different feature fusion methods, the descriptor concatenation (the low-level fusion), the confidence value escalation (the mid-level fusion) and the coarse-to-fine framework (the high-level fusion). Two examples are provided for feature selection to demonstrate its ideas, those are, optimal descriptor selection and improved classifier selection.

Feature extraction plays a key role in object recognition because it is the first and also the most important step. If we consider the overall object recognition process, machine learning tools are to serve the purpose of finding distinctive features from the visual data. Given distinctive features, object recognition is readily available (e.g., a simple threshold function can be used to classify feature descriptors). The proposal of Local N-ary Pattern (LNP) texture features contributes to both feature extraction and texture classification. The distinctive LNP feature generalizes the texture feature extraction process and improves texture classification. Concretely, the local binary pattern (LBP) is the special case of LNP with $n = 2$ and the texture spectrum is the special case of LNP with $n = 3$. The proposed LNP representation has been proven to outperform the popular LBP and one of the LBP's most successful extension - local ternary pattern (LTP) for texture classification.

# *Acknowledgements*

# Contents

# List of Figures

# List of Tables

# Abbreviations

| | |
|---|---|
| **a.k.a** | **a**lso **k**nown **a**s |
| **e.g.** | **e**xemplī **g**rātiā (for example) |
| **i.e.** | **i**d **e**st (that is) |
| **max** | **max**imum |
| **min** | **min**imum |
| **Q.E.D.** | **Q**uod **E**rat **Demonstrandum** (which had to be proven) |
| **s.t.** | **s**ubject **t**o |
| **vs.** | **v**er**s**us |

*To my paternal grandfather and maternal grandfather.*