國立臺灣科技大學

資 訊 工 程 系

# 碩 士 論 文

A Learning-based Algorithm for Natural Scene Recognition

研 究 生： Dang Duy Thang

學　　　號： M10215804

指導教授： 花凱龍博士

中 華 民 國 一 百 零 四 年 七 月 三 十 日

# Abstract

Scene recognition is an important problem in many application areas of image and video processing. Scene recognition has a wide range of applications, such as object recognition and detection, content-based image indexing and retrieval and intelligent vehicle and robot navigation. However, the natural scene images tend to be very complex and difficult to analyze due to changes of illumination and transformation. In this thesis, we will investigate into building a novel model to learn and recognize scenes in nature.

This study proposes a new approach that combines locality-constrained sparse coding (LCSP), Spatial Pyramid Pooling and linear SVM in end-to-end model. Firstly, interesting points of each image in the training set are extracted by a local descriptor as dense SIFT which represents local spatial information. These features known as codewords and each codeword is represented as part of a topic. Then we employs LCSP algorithm to learn the codeword distribution of those local features from the training dataset. Next, a modified Spatial Pyramid Pooling model is employed for encoding the spatial distribution of local features. Spatial Pyramid Pooling model has been remarkably successful in terms of both scene and object recognition. In the testing stage, a linear SVM will be used to classify local features which are encoded by Spatial Pyramid Pooling. The new system achieves very competitive results and leads to state-of-the-art performance on several benchmarks.

# Acknowledgements

Studying master of computer science at National Taiwan University of Science and Technology was one of the best decisions of my life. I would like to thank Prof. Kai-Lung Hua for the enthusiastic support and for many things that he has taught me. His kindness to meet at any time of day was crucial to my researching process and coming up with new advice.

Additionally, this work would not have been possible without the help of all the students at Video and Image Processing Laboratory (VIP lab) for the interesting discussions and research they have all done.

I would also like to thank my family for their continuing support in my studies and helping me with all the tough decisions along the way. I thank my parents, Dang Duy Quyet, Tran Thi Kim Thanh and my both of elder sisters, Thu Trang and Thanh Hoa who always encouraged me and gave me many helpful comments in my life, and also thank to my wife, Quynh Trang whose patience with my long time of work.

# Table of Contents

# List of Tables

# List of Illustrations