# LARGE SCALE SEMANTIC CONCEPT DETECTION, FUSION, AND SELECTION FOR DOMAIN ADAPTIVE VIDEO SEARCH

JIANG YU-GANG

DOCTOR OF PHILOSOPHY

CITY UNIVERSITY OF HONG KONG

SEPTEMBER 2009

# CITY UNIVERSITY OF HONG KONG
# 香港城市大學

## Large Scale Semantic Concept Detection, Fusion, and Selection for Domain Adaptive Video Search
## 大規模語義概念的檢測、融合及選擇進行數據域自適應視頻檢索

Submitted to
Department of Computer Science
電腦科學系
in Partial Fulfillment of the Requirements
for the Degree of Doctor of Philosophy
哲學博士學位

by

JIANG Yu-Gang
姜育剛

September 2009
二零零九年九月

# ABSTRACT

This thesis investigates the problem of video search based on semantic concepts. We present approaches to handle three correlated issues that are critical to this problem: (1) *how to construct an effective feature representation for semantic concept detection,* (2) *how to exploit semantic context to improve the detection of these concepts,* and (3) *how to select the most suitable concept detectors to answer user queries.* In particular, as the target videos may come from different domains (genres or sources) with distinctive data characteristics, for each of the issues, we will need to cope with the domain changes.

Video frames are represented by bag-of-visual-words (BoW) derived from local keypoint features, which are invariant to rotation, scale and illumination. We first conduct a comprehensive study on the representation choices of BoW, including vocabulary size, weighting scheme, stop word removal, feature selection, spatial information, and visual bi-gram. The aim is to offer practical insights in how these choices will impact the performance of BoW for semantic concept detection. We also show how to further augment the BoW representation by exploring the linguistic and ontological aspects of visual words. A visual-word ontology is constructed to hierarchically specify their hyponym relationship, which is incorporated into BoW for improved video frame representation.

To exploit semantic context, we develop a novel and efficient domain adaptive semantic diffusion algorithm. Inter-concept relationship is modeled using a semantic graph, which treats concepts as nodes and the concept affinities as the weights of edges. It is then applied to refine the initial detection results through a function level graph diffusion process, aiming to recover the consistency and smoothness of the detection results over the graph. To handle the domain change

between training and test sets, our algorithm involves a graph adaptation process which iteratively refines the concept affinity based on the target domain data characteristics. This algorithm is efficient and scalable to large scale data sets.

For the selection of concept detectors, we focus on exploring heterogeneous knowledge sources for better measurement of query-detector similarity. Instead of using WordNet as in most existing works, we exploit the context information associated with Flickr images to estimate the similarity between queries and concept detectors. This similarity measure, named FCS, reflects the word correlation in images rather than text corpora. With an initial detector set selected by FCS for each query, we further propose a semantic context transfer algorithm that adapts the query-detector similarity to a target data set. The adaptation process is highly efficient, satisfying the critical requirement of online video search.

We evaluate all the proposed techniques on large scale video search benchmarks provided by TRECVID from years 2005 to 2008. Experimental evaluations demonstrate promising results of our techniques, and their potential to be applied to other applications such as visual object categorization and web scale image retrieval.

# CONTENTS

# LIST OF FIGURES

# LIST OF TABLES