

## DOCTORAL THESIS

### Multi-cue visual tracking: feature learning and fusion

Lan, Xiangyuan

*Date of Award:*  
2016

[Link to publication](#)

#### General rights

Copyright and intellectual property rights for the publications made accessible in HKBU Scholars are retained by the authors and/or other copyright owners. In addition to the restrictions prescribed by the Copyright Ordinance of Hong Kong, all users and readers must also observe the following terms of use:

- Users may download and print one copy of any publication from HKBU Scholars for the purpose of private study or research
- Users cannot further distribute the material or use it for any profit-making activity or commercial gain
- To share publications in HKBU Scholars with others, users are welcome to freely distribute the permanent URL assigned to the publication

# Abstract

As an important and active research topic in computer vision community, visual tracking is a key component in many applications ranging from video surveillance and robotics to human computer interaction. Although it has been extensively studied in last two decades and significant progress has been made in recent years, it still remains challenging mainly due to numerous appearance variations caused by illumination, pose, occlusion, shape deformation and so on. Therefore, how to design an effective object appearance model which is able to handle different appearance changes has become one of the key problems in visual tracking. Since a single visual cue (feature) may not be sufficient to account for different appearance variations, multiple visual cues which describe different characteristics of the tracked object, e.g. color, texture, shape can be exploited jointly to achieve a more robust appearance model. In this thesis, we propose new appearance models based on multiple visual cues and address several research issues in feature learning and fusion for visual tracking.

Feature extraction and feature fusion are two key modules to construct the appearance model for the tracked target with multiple visual cues. Feature extraction aims to extract informative features for visual representation of the tracked target, and many kinds of hand-crafted feature descriptors which capture different types of visual information have been developed. However, since large appearance variations, e.g. occlusion, illumination may occur during tracking, the target samples may be contaminated/corrupted. As such, the extracted raw features may not be able to capture the intrinsic properties of the target appearance. Besides, without ex-

Explicitly imposing the discriminability, the extracted features may potentially suffer background distraction problem. To extract uncontaminated discriminative features from multiple visual cues, this thesis proposes a novel robust joint discriminative feature learning framework which is capable of 1) simultaneously and optimally removing corrupted features and learning reliable classifiers, and 2) exploiting the consistent and feature-specific discriminative information of multiple feature. In this way, the features and classifiers learned from potentially corrupted tracking samples can be better utilized for target representation and foreground/background discrimination.

As shown by the Data Processing Inequality, information fusion in feature level contains more information than that in classifier level. In addition, not all visual cues/features are reliable, and thereby combining all the features may not achieve a better tracking performance. As such, it is more reasonable to dynamically select and fuse multiple visual cues for visual tracking. Based on aforementioned considerations, this thesis proposes a novel joint sparse representation model in which feature selection, fusion, and representation are performed optimally in a unified framework. By taking advantages of sparse representation, unreliable features are detected and removed while reliable features are fused on feature level for target representation. In order to capture the non-linear similarity of features, the model is further extended to perform feature fusion in kernel space. Experimental results demonstrate the effectiveness of the proposed model.

Since different visual cues extracted from the same object should share some commonalities in their representations and each feature should also have some diversities to reflect its complementarity in appearance modeling, another important problem in feature fusion is how to learn the commonality and diversity in the fused representations of multiple visual cues to enhance the tracking accuracy. Different from existing multi-cue sparse trackers which only consider the commonalities among the sparsity patterns of multiple visual cues, this thesis proposes a novel multiple sparse representation model for multi-cue visual tracking which jointly exploits

the underlying commonalities and diversities of different visual cues by decomposing multiple sparsity patterns. Moreover, this thesis introduces a novel online multiple metric learning to efficiently and adaptively incorporate the appearance proximity constraint, which ensures that the learned commonalities of multiple visual cues are more representative. Experimental results on tracking benchmark videos and other challenging videos show that the proposed tracker achieves better performance than the existing sparsity-based trackers and other state-of-the-art trackers.

In short, the major contributions of this thesis are summarized as follows.

- A robust feature learning model is proposed to learn uncontaminated and discriminative features from multiple visual cues for feature extraction in visual tracking.
- A robust feature fusion model based on joint sparse representation is proposed to dynamically select and fuse multiple visual cues on feature level for visual tracking.
- A multiple sparse representation framework is proposed to explicitly model the commonality and diversity in the representation of multiple visual cues, which achieves more accurate appearance modeling with multiple features in visual tracking.

**Keywords:** visual tracking, feature learning, feature fusion

# Table of Contents

|   |            |
|---|------------|
| <b>Declaration</b>                                      | <b>i</b>   |
| <b>Abstract</b>   | <b>ii</b>  |
| <b>Acknowledgements</b>                                 | <b>v</b>   |
| <b>Table of Contents</b>                                | <b>vii</b> |
| <b>List of Tables</b>                                   | <b>xi</b>  |
| <b>List of Figures</b>                                  | <b>xii</b> |
| <b>List of Abbreviations</b>                            | <b>xv</b>  |
| <b>Chapter 1 Introduction</b>                           | <b>1</b>   |
| 1.1 Background and Motivation . . . . .                 | 1          |
| 1.1.1 Background . . . . .                              | 1          |
| 1.1.2 Motivations of This Project . . . . .             | 3          |
| 1.2 Review of Related Visual Tracking Methods . . . . . | 5          |
| 1.2.1 Generative Trackers . . . . .                     | 6          |
| 1.2.2 Discriminative Trackers . . . . .                 | 7          |
| 1.2.3 Sparsity-based Trackers . . . . .                 | 8          |
| 1.2.4 Metric Learning-based Trackers . . . . .          | 10         |
| 1.2.5 Feature Learning-based Trackers . . . . .         | 10         |
| 1.3 Contributions of This Thesis . . . . .              | 11         |

|                  |   |           |
|------------------|---|-----------|
| 1.4              | Overview of This Thesis . . . . .   | 12        |
| <br>             |   |           |
| <b>Chapter 2</b> | <b>Robust Joint Discriminative Feature Learning for Multi-Cue Visual Tracking</b>                                     | <b>16</b> |
| 2.1              | Introduction . . . . .  | 16        |
| 2.2              | Related Work . . . . .  | 19        |
| 2.2.1            | Shared and Feature-Specific Information among Multiple Features/Modalities/Views for Pattern Classification . . . . . | 20        |
| 2.3              | Proposed Model . . . . .  | 20        |
| 2.3.1            | Robust Joint Discriminative Feature Learning . . . . .  | 20        |
| 2.3.2            | Optimization . . . . .  | 23        |
| 2.4              | Implementation Details . . . . .  | 27        |
| 2.4.1            | Target Representation . . . . .   | 27        |
| 2.4.2            | Observation Likelihood for Particle Filtering . . . . .   | 27        |
| 2.5              | Experiments . . . . .   | 28        |
| 2.5.1            | Experimental Setting . . . . .  | 28        |
| 2.5.2            | Experimental Results . . . . .  | 29        |
| 2.6              | Conclusion . . . . .  | 32        |
| <br>             |   |           |
| <b>Chapter 3</b> | <b>Joint Sparse Representation and Robust Feature-Level Fusion for Multi-Cue Visual Tracking</b>                      | <b>34</b> |
| 3.1              | Introduction . . . . .  | 34        |
| 3.2              | Related Work . . . . .  | 36        |
| 3.2.1            | Multi-Task Joint Sparse Representation . . . . .  | 36        |
| 3.3              | Robust Feature-Level Fusion for Multi-Cue Tracking . . . . .  | 37        |
| 3.3.1            | Particle Filter . . . . .   | 37        |
| 3.3.2            | Feature-Level Fusion Based on Joint Sparse Representation . . . . .   | 38        |
| 3.3.3            | Detecting Unreliable Visual Cues for Robust Fusion . . . . .  | 39        |
| 3.3.4            | Optimization Procedure . . . . .  | 42        |
| 3.3.5            | Kernelized Framework for Robust Feature-Level Fusion . . . . .  | 45        |

|       |  |    |
|-------|--|----|
| 3.3.6 | Computational Complexity . . . . .                       | 46 |
| 3.4   | Implementation Details . . . . .                         | 47 |
| 3.4.1 | Template Update Scheme . . . . .                         | 47 |
| 3.4.2 | Tight Lipschitz Constant Estimation . . . . .            | 48 |
| 3.5   | Experiments . . . . .                                    | 50 |
| 3.5.1 | Unreliable Feature Detection on Synthetic Data . . . . . | 50 |
| 3.5.2 | Visual Tracking Experiments . . . . .                    | 51 |
| 3.6   | Conclusion . . . . .                                     | 69 |

**Chapter 4 Multiple Sparse Representations with Commonality and Diversity Modeling for Multi-Cue Visual Tracking 70**

|       |   |    |
|-------|---|----|
| 4.1   | Introduction . . . . .  | 70 |
| 4.2   | Related Work . . . . .  | 74 |
| 4.2.1 | Commonality and Diversity Modeling in Pattern Classification and Recognition . . . . .                    | 74 |
| 4.3   | Proposed Tracking Algorithm . . . . .   | 75 |
| 4.3.1 | Particle Filter . . . . .   | 75 |
| 4.3.2 | Fusion of Multiple Sparse Representations with the Proximity Constraint . . . . .                         | 75 |
| 4.3.3 | Learning Adaptive Proximity Constraint Using Online LogDet Regularized Multiple Metric Learning . . . . . | 78 |
| 4.3.4 | Optimization Procedure . . . . .  | 83 |
| 4.3.5 | Discussion . . . . .  | 84 |
| 4.4   | Implementation Details . . . . .  | 87 |
| 4.4.1 | Model Update Scheme . . . . .   | 87 |
| 4.5   | Experiments . . . . .   | 87 |
| 4.5.1 | Experiment Setting . . . . .  | 87 |
| 4.5.2 | Evaluation on Publicly Available Sequences . . . . .  | 89 |
| 4.5.3 | Evaluation on Visual Tracking Benchmark . . . . .   | 92 |
| 4.6   | Conclusion . . . . .  | 96 |

|                             |            |
|-----------------------------|------------|
| <b>Chapter 5 Conclusion</b> | <b>97</b>  |
| 5.1 Summary . . . . .       | 97         |
| 5.2 Future work . . . . .   | 99         |
| <b>Appendices</b>           | <b>101</b> |
| <b>Bibliography</b>         | <b>103</b> |
| <b>Curriculum Vitae</b>     | <b>120</b> |