UTS
UNIVERSITY
OF TECHNOLOGY
SYDNEY

# Learning from Imperfect Supervision in Visual Pattern Classification and Localization

**by Fan Ma**

Thesis submitted in fulfilment of the requirements for the degree of

**Doctor of Philosophy**

under the supervision of Yi Yang

University of Technology Sydney

Faculty of Engineering and Information Technology

August 2022

# CERTIFICATE OF ORIGINAL AUTHORSHIP

I, Fan Ma declare that this thesis, is submitted in fulfilment of the requirements for the award of

*Doctor of Philosophy*, in the *Faculty of Engineering and Information Technology* at the University of

Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I

certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

This research is supported by the Australian Government Research Training Program.

Production Note:

**Signature:** Signature removed prior to publication.

Date: 24th August 2022

# ABSTRACT

Machine learning algorithms have achieved tremendous success on various computer vision tasks in past decades. Large-scale well-annotated data, such as ImageNet and ActivityNet, are necessary for learning a valuable model. However, high-quality training samples are often insufficient in practice, and it is labor-intensive and time-consuming to produce intense supervision for different learning tasks. Designing algorithms with imperfect training data thus becomes significant in the current data explosion era.

In this dissertation, imperfect supervision is categorized into three classes: 1) Limited supervision where only a small portion of training samples are annotated; 2) Noisy supervision where some labels of training samples are corrupted; 3) Weak supervision where the labels of training data are imprecise to provide expected outputs. Several models are developed to learn from the supervision of different data types. A self-paced co-training algorithm is proposed to improve the model performance when limited training samples are available. I have also proved that our algorithm can achieve a better model with diverse classifiers. Moreover, a self-reweighting mechanism based on online learned class centroids is introduced to prevent the model from deteriorating by noisy supervision. Experiments are conducted on several image recognition datasets demonstrating the superiority of our designed algorithms under both limited and noisy supervision. Furthermore, two practical applications of temporal localization are studied when weak supervision is available. The first task is the temporal action localization, where only a single frame is annotated for each action instance. The goal is to produce

precise temporal boundaries for action instances. An efficient frame expanding algorithm has been introduced to improve the temporal action localization performance. The other task uses query language to temporally localize moments in videos where only language-video pairs are available in the training data. The connections between the video clips and concepts in query sentences are formed by decoupling the core concepts in the query sentence.

This thesis demonstrates that our well-designed algorithms yield excellent results when only imperfect data are available in various vision tasks, ranging from image classification, object detection, and temporal localization in videos.

Dissertation directed by Professor Yi Yang

Australian Artificial Intelligence Institute, School of Computer Science, Faculty of Engineering and Information Technology, University of Technology Sydney

# ACKNOWLEDGMENTS

First and foremost, I want to acknowledge the help and guidance of my principal supervisor, Prof. Yi Yang. Yi is a brilliant person not only in research fields but also in social communication. In the past four years, I have learned a lot from him. He encourages and helps me to build industry connections to stay at the forefront of several computer vision tasks. Moreover, he always points out my problem quickly and instructs me to make life-long plans. "Think before doing" is the most vital spirit I have learned from him. His calm and rational advice has been essential to my survival of several years as a graduate student and to my development as a researcher.

I would also like to thank my co-supervisor Prof. Deyu Meng, who was also my master tutor. I got to know Yi because of his recommendation. He always encourages me to stay focused on the research and to do some exciting work. His passion and optimism for research and life also inspire me to enjoy the process of doing academic research. Many thanks to my co-authors, Linchao Zhu, Zheng Shou, Xin Yu, Zhicheng Yan, Yu Wu. They have provided many insightful ideas and contributed a lot to the paper writing. I also want to thank my colleagues, Xiaohan Wang, Ruijie Quan, and Zhun Zhong, who are also my roommates and friends. We have had a wonderful time in the past four years.

I also want to thank my group members and colleagues, Xiaojun Chang, Jiaxu Miao, Hehe Fan, Pingbo Pan, Yanbin Liu, Zongxin Yang, Qingji Guan, Yawei Luo, Qianyu Feng, Yunqiu Xu, Guang Li, Youjiang Xu, Xuanmeng Zhang, Yuhang Ding, Guangrui Li, Yang He, Ping Liu, Yutian Lin, Peike Li, Xuanyi Dong, Tianqi Tang, Bingwen Hu, Minfeng Zhu, Fengda Zhu, Hu Zhang, Zhedong Zheng, Liang Zheng, Xiaolin Zhang, Aming Wu,

and many others. I was very fortunate to collaborate with or discuss with them. These discussions inspired and motivated many of my research works.

Lastly, I would like to thank my family and my girlfriend Yue Yang for their selfless support and love.

<div align="right">

Fan Ma

Sydney, Australia, 2022

</div>

# LIST OF PUBLICATIONS

**CONFERENCE PAPERS:**

1. **Fan Ma**, Linchao Zhu, Yi Yang, Shengxin Zha, Gourab Kundu, Matt Feiszli, Zheng Shou. "SF-Net: Single-Frame Supervision for Temporal Action Localization", European Conference on Computer Vision (ECCV 2020)

2. **Fan Ma**, Zheng Shou, Linchao Zhu, Haoqi Fan, Yilei Xu, Yi Yang, Zhicheng Yan. "Unified Transformer Tracker for Object Tracking", IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2022)

**JOURNAL PAPERS :**

3. **Fan Ma**, Deyu Meng, Xuanyi Dong, Yi Yang, "Self-Paced Multi-view Co-training", Journal of Machine Learning Research (JMLR)

4. **Fan Ma**, Yu Wu, Xin Yu, Yi Yang, "Learning With Noisy Labels via Self-Reweighting From Class Centroids", IEEE Transactions on Neural Networks and Learning Systems (TNNLS)

5. **Fan Ma**, Linchao Zhu, Yi Yang, "Weakly Supervised Moment Localization with Decoupled Consistent Concept Prediction", International Journal of Computer Vision (IJCV)

# TABLE OF CONTENTS

# LIST OF FIGURES