# Accelerating Deep Convolutional Neural Networks via Filter Pruning

**by Yang He**

Thesis submitted in fulfilment of the requirements for the degree of

**Doctor of Philosophy**

under the supervision of Yi Yang

# Certificate of Authorship/Originality

I, Yang He, 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.

Signature: Yang He

Production Note:
Signature removed prior to publication.

Date: 1/July/2021

# Acknowledgements

First, I would like to express my tremendous gratitude to my supervisor, Professor Yi Yang, I really appreciate his mentoring and cultivating in my research and my life. Although I changed my major for my PhD study and had little knowledge in the area computer science when I entered the group, Yi has always been patient with me and given me quite a lot guidance. Every time I encountered difficulties, Yi is always here to help me. It is really my fortune to have Yi as my supervisor.

I am grateful for all my colleagues and teammates in Yi's group. I convey my gratitude to my co-supervisor, Prof. Liang Zheng, for his help and advises during this research. I feel so lucky to met Prof. Yanwei Fu, Guoliang Kang and Xuanyi Dong that they teached me quite a lot about doing research and writing code when I am a novice of the research field. It is so important to express my gratitude to Ping Liu, who is reliable and heart-warming, and gives me many kinds of help and advice. I am happy to collaborate with many creative teammates in our team and I really appreciate the kind and useful suggestions given by them. Moreover, I am thankful for Prof. Hanwang Zhang for his guidance and advises during my visiting to NTU. I have met quite a lot friends at NTU and thanks for your support. I also thank Qi Yao for helping me with my PhD study.

Finally, I want to express my special gratitude to my parents, who give me endless love, help, and encouragement in my life. How lucky to be your child.

Thanks for all the people that ever helped me and encouraged me.

<div align="right">

Yang He

Sydney, Australia, 2021.

</div>

# List of Publications

**Journal Papers**

J-1. **Yang He**, Xuanyi Dong, Guoliang Kang, Yanwei Fu, Chenggang Yan, and Yi Yang, "Asymptotic Soft Filter Pruning for Deep Convolutional Neural Networks," IEEE Transactions on Cybernetics, vol. 50, no. 8, pp. 3594–3604, 2019.

**Conference Papers**

C-1. **Yang He**, Guoliang Kang, Xuanyi Dong, Yanwei Fu, and Yi Yang, "Soft Filter Pruning for Accelerating Deep Convolutional Neural Networks," in International Joint Conference on Artificial Intelligence (IJCAI), pp. 2234–2240, 2018.

C-2. **Yang He**, Ping Liu, Ziwei Wang, Zhilan Hu, Yi Yang, "Filter Pruning via Geometric Median for Deep Convolutional Neural Networks Acceleration," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4340–4349, 2019.

C-3. **Yang He**, Yuhang Ding, Ping Liu, Linchao Zhu, Hanwang Zhang, Yi Yang, "Learning Filter Pruning Criteria for Deep Convolutional Neural Networks Acceleration," in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2009–2018, 2020.

# Contents

# List of Figures

# List of Tables

# ABSTRACT

## Accelerating Deep Convolutional Neural Networks via Filter Pruning

by

Yang He

The superior performance of deep Convolutional Neural Networks (CNNs) usually comes from the deeper and wider architectures, which cause the prohibitively expensive computation cost. To reduce the computational cost, works on model compression and acceleration have recently emerged. Among all the directions for this goal, filter pruning has attracted attention in recent studies due to its efficacy. For a better understanding of filter pruning, this thesis explores different aspects of filter pruning, including pruning mechanism, pruning ratio, pruning criteria, and automatic pruning. First, we improve the pruning mechanism with soft filter pruning so that the mistaken pruned filters can have a chance to be recovered. Second, we consider the asymptotic pruning rate to reduce the sudden information loss in the pruning process. Then we explore the pruning criteria to better measure the importance of filters. Finally, we propose the automatic pruning method to save human labor. Our methods lead to superior convolutional neural network acceleration results.

Dissertation directed by Professor Yi Yang

ReLER, Australian Artificial Intelligence Institute, School of Software