UNIVERSITY OF TECHNOLOGY SYDNEY

Faculty of Engineering and Information Technology

# Context-aware Image Semantic Segmentation

by

**Ye Huang**

A THESIS SUBMITTED
IN FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE

**Doctor of Philosophy**

Sydney, Australia

2022

## CERTIFICATE OF ORIGINAL AUTHORSHIP

I, Ye Huang, declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the School of Electrical and Data Engineering, 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.

Signature: Production Note:
Signature removed prior to publication.

Date: May-06-2022

# ABSTRACT

**Context-aware Image Semantic Segmentation**

by

Ye Huang

Semantic segmentation is a fundamental task for computer vision applications. However, the existing solutions have many issues when handling difficult cases. This thesis develops three novel approaches which have improved the generalization ability of the existing solutions at significantly reduced computation costs. Extensive experiments conducted on multiple benchmark datasets have demonstrated the superior performance of the proposed approaches.

**Scale-invariant:** The state-of-the-art semantic segmentation solutions usually leverage different receptive fields via multiple parallel branches to handle objects of different sizes. However, employing separate kernels for individual branches degrades the generalization of the network to objects with different scales, and the computational cost increases with the increase of the number of branches. In this thesis, a novel network structure, namely *Kernel-Sharing Atrous Convolution (KSAC)*, is proposed, where branches with different receptive fields share the same kernel, i.e., letting a single kernel "see" the input feature maps more than once with different receptive fields.

**Seamless dual attention:** Spatial and channel attentions, modelling the semantic inter-dependencies in spatial and channel dimensions respectively, have recently been widely used for semantic segmentation. However, computing spatial attention and channel attention separately sometimes causes errors, especially in those difficult cases. In this research, a Channelized Axial Attention (CAA) is developed to seamlessly *integrate* channel attention and spatial attention into a single operation with negligible computation overhead. Furthermore, a novel grouped vec-

torization approach is developed to allow the proposed model to run with very little memory consumption without slowing down the computation.

**Class-aware regularization:** Recent segmentation methods utilizing class-level information in addition to pixel features have achieved notable success in boosting the accuracy of existing network models. However, the extracted class-level information was simply concatenated to pixel features, without being explicitly exploited to learn better pixel representation. Moreover, these approaches learn soft class centers based on coarse mask prediction, which is prone to error accumulation. Motivated by the fact that humans can recognize an object by itself no matter which other objects it appears with and aiming to use class-level information more effectively, a universal Class-Aware Regularization (CAR) approach is proposed to optimize the intra-class variance and inter-class distance during feature learning. Furthermore, the class center in the proposed approach is directly generated from ground truth instead of from the error-prone coarse prediction. The proposed CAR can be easily applied to most existing segmentation models and can largely improve their accuracy at no additional inference overhead.

Dissertation directed by Prof. Xiangjian He and Dr Wenjing Jia
School of Electrical and Data Engineering

# Dedication

Dedicated to my family. Dedicated to the world peace.

# Acknowledgements

# List of Publications

**Journal Papers**

1. **Ye Huang**, Qingqing Wang, Wenjing Jia and Xiangjian He, See More Than Once–Kernel-Sharing Atrous Convolution for Semantic Segmentation, in Neurocomputing, Volume 443, 5 July 2021, Pages 26-34.

2. Qingqing Wang, **Ye Huang**, Wenjing Jia, Xiangjian He, Michael Blumenstein, Shujing Lyu and Yue Lu, FACLSTM: ConvLSTM with focused attention for scene text recognition, Science China Information Sciences 2020

**Conference Papers**

1. **Ye Huang**, Di Kang, Liang Chen, Xuefei Zhe, Wenjing Jia, Xiangjian He, Linchao Bao, CAR: Class-aware Regularizations for Semantic Segmentation, accepted by ECCV 2022.

2. **Ye Huang**, Di Kang, Wenjing Jia, Xiangjian He and Liu Liu, Channelized Axial Attention – Considering Channel Relation within Spatial Attention for Semantic Segmentation, Proceedings of the AAAI Conference on Artificial Intelligence, 36(1), 1016-1025.

3. Qingqing Wang, Wenjing Jia, Xiangjian He, Yue Lu, Michael Blumenstein, **Ye Huang** and Shujing Lyu, DeepText: Detecting text from the wild with multi-ASPP-assembled deeplab, Proceedings of 2019 International Conference on Document Analysis and Recognition

# Contents

# List of Figures

# List of Tables

# Abbreviation

- H - Height

- W - Width

- Channels - The size of last dimension of the 4D feature map.

- Encoding - Nerual network encoded 3 channels image input to the faeture map with multiple channels.

- mIOU - Mean Intersection over Union

- KSAC - Kernel Sharing Astrous Convolution

- CAA - Channelized Axial Attention

- CAR - Class-aware Regularization

- OS - Output stride [3]

- FCN - Fully Convolutional Networks [37]

- ASPP - Atrous Spatial Pyramid Pooling [4]

- PPM - Pyramid Pooling module [69]

- FPN - Feature Pyramid Networks [32]

- SA - Self-attention [53]

- ACFNet - Attentional Class Feature Network [64]

- OCR - Object-Contextual Representations [61]

- CPNet - Context Prior Network [60]