# Sequence Modelling with Deep Learning for Visual Content Generation and Understanding

**by Zongxin Yang**

Thesis submitted in fulfilment of the requirements for the degree of

**Doctor of Philosophy**

under the supervision of Prof. Yi Yang

# CERTIFICATE OF ORIGINAL AUTHORSHIP

I, *Zongxin Yang* 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.

Production Note:

Signature: Signature removed prior to publication.

Date: 16^th Jul, 2021

# ABSTRACT

Although convolutional neural networks have proven to be effective and stable in image feature learning, sequence modelling is still critical for learning spatial and temporal context information. In an image scenario, different semantic structures can be regarded as a sequence arranged along the horizontal (or vertical) direction. Moreover, in a video scenario, temporal sequence modelling is necessary for understanding inter-frame relationships, such as object movement and occlusion. This thesis explores more effective spatial or temporal sequence modelling for image or video scenario understanding. For the former, an encoder-decoder framework is proposed to split an input scenario into a sequence of spatial features and reconstruct the input. By modelling spatial sequence information, the framework can even predict new scenes with very large scales in length while keeping a consistent style regarding the given input. For video understanding, the thesis processes temporal sequences in a recurrent manner (*i.e.*, frame by frame), which is more memory-efficient. In addition, the thesis proposes to implicitly impose the feature embedding of each target and relative background to be contrastive throughout the temporal sequence, promoting the results of downstream tasks accordingly. Besides, a novel transformation module is designed to model channel relationships for improving intra-frame representation ability. To validate proposed approaches and components, extensive experiments are conducted on image outpainting, instance segmentation, object detection, classification, video classification, and video object segmentation.

# ACKNOWLEDGMENTS

Firstly, I would like to thank my supervisor Professor Yi Yang, without whom I would not have been able to complete this thesis. I sincerely appreciate him for his selfless support and patience. He guided me throughout the research by providing tremendous help on his knowledge and insight into the related fields.

I would also like to thank my colleagues at University of Technology Sydney. I would like to thank Yunchao Wei, Linchao Zhu, Xin Yu, Ping Liu, Yu Wu, Qianyu Feng, Peike Li, Chen Liang, Jiaxu Miao, Yuhang Ding, and many others. I was really fortunate to work with them and participate in intellectual conversations with them.

Lastly, I would like to thanks my mother, Gexiu Liao, and my father, Jianlin Yang, for their support and love throughout the years.

# LIST OF PUBLICATIONS

1. **Z. YANG**, J. DONG, P. LIU, Y. YANG, AND S. YAN, *Very long natural scenery imageprediction by outpainting*, in ICCV, 2019, pp. 10561-10570.

2. **Z. YANG**, L. ZHU, Y. WU,AND Y. YANG, *Gated channel transformation for visual-recognition*, in CVPR, 2020, pp. 11794-11803.

3. **Z. YANG**, Y. WEI, AND Y. YANG, *Collaborative video object segmentation by foreground-background integration*, in ECCV, Springer, 2020, pp. 332-348.

# TABLE OF CONTENTS