

# **Exploring Visual Attention Mechanism for Scene Understanding in Image Captioning**

**by Zongjian Zhang**

Thesis submitted in fulfilment of the requirements for  
the degree of

**Doctor of Philosophy**

under the supervision of A/Prof. Qiang Wu,  
A/Prof. Yang Wang,  
Prof. Fang Chen

University of Technology Sydney  
Faculty of Engineering and Information Technology

September 2020

# Certificate of Original Authorship

I, **Zongjian Zhang** 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: 28 Sep 2020

# ABSTRACT

## **Exploring Visual Attention Mechanism for Scene Understanding in Image Captioning**

by

Zongjian Zhang

Supervisor: A/Prof. Qiang Wu

Co-supervisor: Prof. Fang CHEN, A/Prof. Yang WANG

Scene understanding is a high-level computer vision research task that requires multiple fundamental vision tasks on different visual elements. Image captioning is a typical scene understanding task that understands salient visual contents in a real-world scene of an image and then automatically describes its understandings via a natural language sentence. This thesis concentrates on exploring visual attention mechanism for the scene understanding in image captioning, aiming to achieve a comprehensive multimodal and transparent scene understanding ability. Specifically, four problems are studied to enhance the visual attention mechanism to fine-grained spatial level, to a comprehensive semantic level, and to a high-level ability of attending to visual relationships for interaction words.

Firstly, this thesis proposes a fine-grained and semantic-guided visual region attention model based on a novel Fully Convolutional Network (FCN)-Long Short Term Memory (LSTM) framework. It can attend to both object and “stuff” regions at a fine-grained grid-wise resolution and only focuses on the principal object information in each grid cell. In addition, grid-wise semantic labels are introduced to provide semantic guidance to ensure that related visual regions in different grid cells are correlated to each other. Moreover, an additional semantic context can be summarized from these textual semantic labels.

Secondly, this thesis proposes a novel high-resolution FCN encoder is used with

a residual attention structure to achieve a high-resolution attention supporting four high fine-grained resolutions (i.e.,  $27 \times 27$ ,  $40 \times 40$ ,  $60 \times 60$ ,  $80 \times 80$ ). A size-invariant “attention correctness” metric is further proposed for evaluating the attention accuracy under different resolutions. Based on the COCO-Stuff dataset, pixel-level evaluations are conducted on both object and “stuff” regions to analyse the performances of high-resolution fine-grained visual region attention.

Thirdly, this thesis concentrates on visual relationship attention exploring visual relationships between object regions under spatial constraints based on the parallel attention mechanism. The goal is to fully explore the relationships/interactions between visual regions for achieving an accurate description of interaction words in the caption. Moreover, it is trained implicitly through an unsupervised approach without using any explicit visual relationship annotations.

Last but not the least, this thesis takes a further step to achieve an adaptive attention module that can perform the role of both visual region attention or visual relationship attention adaptively according to the needs of language decoder. The dynamic linguistic context of the language decoder is fully leveraged for exploring and attending to related visual relationships for interactive words.

## Acknowledgements

First of all, I would like to express my deepest gratitude to my principal supervisor, A/Prof. Qiang Wu. Without his guidance, this thesis would not have been accomplished with high quality, and I would not have achieved qualified PhD research outcomes. With abundant knowledge and the insightful and rigorous way of thinking about research problems, he patiently guided me through many research challenges in this valuable journey of PhD study. Besides this degree itself, I've learned the methodology of exploring and solving a new problem, which I believe is more important for my future work and life. Moreover, he is also a great friend that cares about my personal life and offers helps when I face problems. I feel very lucky to finish my PhD under his supervision and can't imagine a better supervisor. Deeply appreciate it!

I am also profoundly grateful to my co-supervisors Prof. Fang Chen and A/Prof. Yang Wang for hosting me as a visiting student in Data 61 (formerly as NICTA) and providing great supports. Thanks to these, I can further extend my knowledge and techniques about data analytics, machine learning, and industrial projects. I believe these are the most valuable experiences for my future work career.

I would also like to express my sincere gratitude to my mentor Dr. Zelin Li. I've gained lots of valuable suggestions and experiences of PhD research and research projects in Data 61. Moreover, he gave me lots of help that are very important to my personal life, which significantly makes my life in Australia easier.

Last but not least, I also want to thank my family. Thank you, my beloved wife, for accompanying me and fully supporting me when I am pursuing this PhD degree. Thank you, my selfless parents, for helping us taking care of our two lovely kids, so that we can work with dedication. Thank you, my kids, for trusting us and experiencing a harsh life transfer with us. You all are my motivations for doing good

life research, including this PhD.

Many thanks to UTS and Data 61 for providing scholarship and research platform to me, so I can accomplish myself with a PhD degree.

Zongjian Zhang  
Sydney, Australia, 2020.

# List of Publications

## Journal Papers

- J-1. **Z. Zhang**, Q. Wu, Y. Wang and F. Chen, “High-Quality Image Captioning With Fine-Grained and Semantic-Guided Visual Attention,” *IEEE Transactions on Multimedia*, vol. 21, no. 7, pp. 1681-1693, 2019.
- J-2. **Z. Zhang**, Q. Wu, Y. Wang and F. Chen, “Exploring Region Relationships Implicitly: Image Captioning with Visual Relationship Attention,” *Image and Vision Computing*, 2021.
- J-3. **Z. Zhang**, Q. Wu, Y. Wang and F. Chen, “Adaptively Exploring Visual Relationship through Linguistic Context,” *IEEE Transactions on Multimedia*, Under review.

## Conference Papers

- C-1. **Z. Zhang**, Q. Wu, Y. Wang and F. Chen, “Fine-Grained and Semantic-Guided Visual Attention for Image Captioning,” *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1709-1717, Mar. 12-15, 2018.
- C-2. **Z. Zhang**, Q. Wu, Y. Wang and F. Chen, “Size-Invariant Attention Accuracy Metric for Image Captioning with High-Resolution Residual Attention,” *Digital Image Computing: Techniques and Applications (DICTA)*, pp. 1-8, Dec. 10-13, 2018.
- C-3. **Z. Zhang**, Q. Wu, Y. Wang and F. Chen, “Visual Relationship Attention for Image Captioning,” *International Joint Conference on Neural Networks (IJCNN)*, pp. 1-8, Dec. 14-19, 2019.

# Contents

Abstract	iii
Acknowledgments	v
List of Publications	vii
List of Figures	xiv
<b>1 Introduction</b>	<b>1</b>
1.1 Thesis Overview . . . . .	1
1.2 Research Background . . . . .	4
1.3 Research Aims and Objectives . . . . .	6
1.4 Research Problems and Contributions . . . . .	8
1.4.1 Fine-Grained and Semantic-Guided Visual Region Attention . . . . .	8
1.4.2 High-Resolution Residual Visual Attention with Size-Invariant Attention Accuracy Evaluation Metric . . . . .	9
1.4.3 Visual Relationship Attention: Exploring Relationships between Adjacent Region . . . . .	10
1.4.4 Visual Relationship Attention: Adaptively Exploring Pairwise Relationships from Linguistic Context . . . . .	11
1.4.5 Relationships between four research problems . . . . .	12
1.5 Thesis Organization . . . . .	13
<b>2 Related Works</b>	<b>15</b>
2.1 Image Captioning on Different Content Level . . . . .	16



2.1.1	Image-Level Captioning . . . . .	17
2.1.2	Region-Level Captioning . . . . .	19
2.1.3	Relationship-Level Captioning . . . . .	21
2.2	Visual Attention Mechanism for Image Captioning . . . . .	21
2.2.1	Visual Attention based on Different Representations Provided by the Image Encoder . . . . .	22
2.2.2	Visual Attention based on different structures . . . . .	26
2.2.3	Quantitative Evaluation on Visual Attention . . . . .	28
2.2.4	Visual Attention Exploring Higher-Level Visual Relationships	29
2.3	Non-Factual Image Captioning . . . . .	30
2.4	Datasets . . . . .	33
2.4.1	Microsoft COCO . . . . .	33
2.4.2	Flickr . . . . .	33
2.4.3	COCO-Stuff . . . . .	34
2.4.4	Caption Entities . . . . .	34
2.5	Evaluation Metrics . . . . .	34
2.5.1	BLEU . . . . .	34
2.5.2	ROUGE . . . . .	35
2.5.3	METEOR . . . . .	35
2.5.4	CIDEr . . . . .	36
2.5.5	SPICE . . . . .	36
<b>3</b>	<b>Mathematical Model for Image Captioning</b>	<b>37</b>
3.1	Image Encoder . . . . .	37
3.2	Language Decoder . . . . .	38

<b>4 Fine-Grained and Semantic-Guided Visual Region Attention</b>	<b>40</b>
4.1 Introduction . . . . .	40
4.2 FCN-LSTM Framework for Image Captioning . . . . .	42
4.2.1 FCN Encoder . . . . .	44
4.2.2 LSTM Decoder . . . . .	44
4.3 Fine-Grained Grid-Wise Soft-Attention . . . . .	46
4.4 Semantic-Guided Attention . . . . .	48
4.4.1 Saliency Pooling Layer . . . . .	49
4.4.2 Attention Distribution Prediction Layer . . . . .	51
4.4.3 Joint Context Computation Layer . . . . .	54
4.5 Experiments . . . . .	55
4.5.1 Datasets and Metrics . . . . .	55
4.5.2 Experiment Settings . . . . .	56
4.5.3 Experiment 1 - Evaluation on the Performance Contributed by Fine-Grained Grid-Wise Attention . . . . .	59
4.5.4 Experiment 2 - Evaluation on the Performance Contributed by Semantic Guidance . . . . .	60
4.5.5 Experiment 3 - Evaluation in Terms of Caption Diversity . . . . .	63
4.5.6 Computational Costs . . . . .	65
4.6 Conclusion . . . . .	65
<b>5 High-Resolution Residual Visual Attention with Size-Invariant Attention Accuracy Evaluation Metric</b>	<b>67</b>
5.1 Introduction . . . . .	67

5.2	Image Captioning Model with High-Resolution Residual attention . . .	71
5.2.1	Image Captioning Model . . . . .	71
5.2.2	High-Resolution Residual Attention Model . . . . .	74
5.3	“Normalised Attention Correctness” Metric . . . . .	77
5.4	Experiments . . . . .	80
5.4.1	Datasets and Metrics . . . . .	80
5.4.2	Experiment Settings . . . . .	81
5.4.3	Experiment-1: Evaluation on Image Captioning Performance by Increasing Resolution and Using New Residual Attention Model . . . . .	82
5.4.4	Experiment-2: Evaluation on the Efficiency of the Proposed “normalized attention correctness” Metric for Measuring Attention Accuracy . . . . .	84
5.5	Conclusion . . . . .	86
<b>6</b>	<b>Visual Relationship Attention: Exploring Relationships between Adjacent Regions</b>	<b>87</b>
6.1	Introduction . . . . .	87
6.2	The Overall Module Structure for Image Captioning . . . . .	92
6.3	Visual Relationship Attention . . . . .	95
6.3.1	Visual Relationship Encoder . . . . .	96
6.3.2	Visual Relationship Selector . . . . .	99
6.4	Spatial Constraints . . . . .	99
6.4.1	Hard-way dropout . . . . .	100
6.4.2	Soft-way re-weighting . . . . .	101
6.5	Training Objectives . . . . .	102

6.6	Experiments . . . . .	103
6.6.1	Datasets and Metrics . . . . .	103
6.6.2	Image Features . . . . .	104
6.6.3	Implementation Details . . . . .	104
6.6.4	Quantitative Analysis . . . . .	105
6.6.5	Qualitative Evaluation on Visual Relationship Attention . . . . .	110
6.6.6	Ablated Study . . . . .	113
6.7	Conclusion . . . . .	119
<b>7</b>	<b>Visual Relationship Attention: Adaptively Exploring Pairwise Relationships from Linguistic Context</b>	<b>120</b>
7.1	Introduction . . . . .	120
7.2	Brief Structure of Image Captioning Model . . . . .	126
7.3	Faster R-CNN Image Encoder . . . . .	127
7.4	2-layer LSTM Language Encoder . . . . .	128
7.5	Adaptive Attention Module . . . . .	130
7.5.1	Visual Relationship Attention . . . . .	130
7.5.2	Visual Region Attention . . . . .	134
7.5.3	Attention Modulator . . . . .	134
7.6	Loss Functions . . . . .	136
7.6.1	Caption Loss . . . . .	136
7.6.2	Attention Loss . . . . .	137
7.7	Interaction-focused Dataset . . . . .	138
7.8	Experiments . . . . .	139
7.8.1	Datasets and Metrics . . . . .	139

7.8.2	Model Parameter Details . . . . .	140
7.8.3	Ablated Models Revised based on Adaptive Attention Module (AdpAtt) . . . . .	141
7.8.4	Quantitative Evaluation of Adaptive Attention Module via Captioning performance . . . . .	142
7.8.5	Qualitative Evaluation of Adaptive Attention Module . . . . .	146
7.9	Conclusion . . . . .	147
<b>8</b>	<b>Conclusion and Future Work</b>	<b>148</b>
8.1	Conclusion . . . . .	148
8.2	Future Work . . . . .	153
8.2.1	Exploring visual relationship between object and “stuff” . . . . .	154
8.2.2	Extend the ability of visual attention mechanism to adjective words . . . . .	154
8.2.3	Extend the visual relationship attention to the temporal dimension . . . . .	155
8.2.4	Optimization on SPICE metric in the second-stage Reinforcement Learning based training . . . . .	155
	<b>Bibliography</b>	<b>156</b>

# List of Figures

2.1	Two basic structures for visual attention. . . . .	27
4.1	The overview of the proposed framework. . . . .	43
4.2	The detailed structure of the proposed fine-grained and semantic-guided attention model. . . . .	48
4.3	An illustration of the saliency pooling layer for single field (a) and entire image (b). . . . .	50
4.4	An illustration of the attention distribution prediction layer (a) and joint context computation layer (b). . . . .	52
5.1	The overview of the proposed framework. . . . .	68
5.2	The detailed structure of high-resolution residual attention model. . .	75
5.3	Illustration for two attention evaluation metrics. . . . .	78
5.4	Qualitative Analysis of the Advantages Provided by Higher-Resolution Residual Attention. . . . .	83
5.5	The comparisons of captioning performance between soft-attention model and the residual attention model across four fine-grained resolutions (i.e., $27 \times 27$ , $40 \times 40$ , $60 \times 60$ , $80 \times 80$ ). . . . .	84
5.6	The evaluation effects of “normalized attention correctness” and “attention correctness” between soft-attention model and the residual attention model across four fine-grained resolutions (i.e., $27 \times 27$ , $40 \times 40$ , $60 \times 60$ , $80 \times 80$ ). . . . .	84

5.7	The comparisons of attention accuracy between soft-attention model and the residual attention model across four fine-grained resolutions (i.e., $27 \times 27$ , $40 \times 40$ , $60 \times 60$ , $80 \times 80$ ).	85
6.1	The detailed framework structure of proposed image captioning model.	89
6.2	The detailed structure of parallel attention.	97
6.3	Qualitative analysis of differences between visual relationship attention (VRAtt-Soft) and visual region attention (RegionAtt) for image captioning.	111
6.4	Qualitative analysis of differences between visual relationship attentions under Soft-Way Spatial Constraints and Hard-Way Spatial Constraints.	118
7.1	Illustration of image captioning framework with adaptive attention module.	121
7.2	The detailed structure of image captioning model with proposed adaptive attention module.	126
7.3	The process of pairwise relationship attention map generation.	131
7.4	The process of pairwise visual relationship feature generation.	133
7.5	Qualitative analysis of our adaptive attention module.	146