

Faculty of Engineering and Information Technology
University of Technology, Sydney

Detecting Text in Clutter Scene

A thesis submitted in partial fulfilment of
the requirements for the degree of
Doctor of Philosophy

by

Xia Cui

November 2014

CERTIFICATE OF ORIGINAL AUTHORSHIP

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Signature of Student:

Date:

Abstract

We often encounter cluttered visual scenes and need to identify objects correctly to navigate and interact with the world. As text takes the typical form of a human-designed informative visual object, retrieving texts in both indoor and outdoor environments is an important step towards providing contextual clues for a wide variety of vision tasks. Furthermore, it plays an invaluable role for multimedia retrieval and location based services.

Text detection from clutter background is nevertheless a challenging task because the text, being figures in image, can be presented in various ways with lots of room for uncertainty such as size, scale, font type, font texture and colour, unpredicted decorative elements put on the text, etc. The situation will be even more complicated if the text is presented in a clutter background where non-text objects possess similar low-level features to text. Further, all these objects are composed of distinct geometric shapes and they are similar with the essential composition elements of text objects. Pursuing a robust text feature descriptor is therefore always difficult because special feature descriptor is only a fragment of text existence. It needs the completely understanding of text.

Regarding the design, understanding, representation and calculating of text as one unitary process of text perceiving, we deal with the completely understanding and representation of text in image with many kinds of aspects in different levels. Without following the legend feature based solution, this research is motivated by perceptual image processing and the observation of painting masters. It will explore a brand new solution by investigating the spatial structure of text and the compositional complexity of the visual object (i.e. text) in image. The research will present the composition granularity indicator and expose novel discriminable attributes embedded inside text objects, which can successfully differentiate text regions and non-text regions on clutter backgrounds.

As figures in image with the clutter scene, it is merely the physical appearance of text which provides the perceptual content and plays a central role for text detection, i.e. location and coarse identification. During the view-construction of text, properties of individual character and textual organization of characters build up the physical appear-

ance. When observers see text appearance in clutter scene, they describe their feelings in terms of crowding effect and clutter. However, the appearance of text still has enough saliency to reveal an informative message. Accordingly, text not only has the characteristics of crowding effect and clutter but also follows the principles of saliency.

Significantly, the crowding effect of text is derived from the space regularity of in-built neighbouring letters which have commonalities beside their distinctiveness. In addition, low-level features of individual letters contribute to the commonalities and distinctiveness from the moment that the font is designed.

Therefore, the computational model of text appearance is built up to integrate the three-level properties, including features of individual characters (low-level features), properties for spatial regularity (i.e. neighbourhood, appearance similarity), and the crowding statistics property of space averaged over pooling regions.

In terms of image processing, if we consider the view construction of text, the features of individual characters in image processing are obtained on the basis of the properties of construction, including mean intensity, local RMS contrast, shape, pixel density, edge density, stroke width, straight line ratio, height to width ratio, stroke width to height ratio, etc.

For the purpose of calculating the properties of space regularity and the crowding space averaging property, the spatial elements and relations are quantified and these involve space granularity and composition rules.

If we examine the works of painters, especially impressionists, they use directional brushstroke or colour patches as space granularity to represent “formless” visual objects in space regularity instead of clear contour shape sketches. The space regularity of patches, i.e. repetitive patterns, can offer a compositional format to express an artist’s feelings about an object rather than to simply describe it. Secondly, it is the harmonious proportions among component parts that bridle component space patches into objects. If we consider the painter’s harmonious proportions, the component parts of an object can be said to react simultaneously so that they can be seen at one and the same time both together and separately.

Similarly, image is described by a set of grey space patches in multi-grey levels. In addition, each space patch groups pixels in position proximity and similarity, in just the same way as the colour patch is used by impressionists. The space organisation of them is also quantified as the measurement of space relations, especially in terms of the

neighbourhood and proportions among component parts. Moreover, the harmonious proportions among space patches are captured by the mathematical tool of geometric mean. Geometric mean (i.e., GM) is calculated over those space patches which possess the same grey level, and considered as the space granularity to form objects. Grey patches with the same GM are composed of GM regions, which are enlarged, extended kinds of pooling regions. Regions given by clusters which have resulted from similarity and neighbourhood are direct, compact pooling regions. Therefore, the statistical properties of space averaging are calculated over GM regions and image is represented as a set of GM regions over which text and other visual objects are analysed by GM indication.

Finally, the representation of an image and the three-level computational text model are put into practice to develop a new-brand algorithm on the public benchmark dataset and to design and implement an automatic processing system on the real big data of the bank cheque. The resulting performance of these tools/processes shows that they are highly competitive and effective.

Acknowledgement

First and foremost I want to thank my supervisor Professor Longbing Cao. It has been an honour to be his Ph.D. student. He has taught me, both consciously and unconsciously, how good big data is done. I appreciate all his contributions of time, ideas, and funding to make my Ph.D. experience productive and stimulating. The joy and enthusiasm he has for his research was contagious and motivational for me, even during tough times in the Ph.D. pursuit. I am also thankful for the excellent example he has provided as a successful leader scientist and professor.

And I also want to thank my supervisor Dr. Qiang Wu. It has been a great gratitude to be his Ph.D student. He has taught me how good and joyful the thinking ability is done for research as a vigorous and rigorous scientist. I appreciate all his contributions of time and ideas to make my Ph.D full of creative thinking and enjoyment. The endless creative thinking and enthusiasm he has for his research encouraged me to make greater efforts. I am also thankful for the excellent example he has provided as a time manager.

The members of the AAI group have contributed immensely to my personal and professional time at UTS. The group has been a source of friendships as well as good advice and collaboration. I am especially grateful for the fun group of original AAI group members who stuck it out in grad school with me: Ziyue Zuo, Wei Wei, Jiahang Chen, Yuming Ou, Zhigang Zheng, XinHua Zhu, Chao Luo, Junfu Yin, Xuhui Fan, and SongYin. I would like to acknowledge honorary group member QiaoYu Sun who was here a couple of years ago. We worked together along with Qiang Wu on the i-Cheque system development, and I am very much appreciated her enthusiasm, intensity, and amazing ability to manipulate text images. And I am also grateful for these honorary group members who have come through the lab: Hongxiu Zhu, Wei Li, Xiaodong Yue. Other past and present group members that I have had the pleasure to work with or alongside of are grad students Jingjiu Li, Mu Li, Can Wang, She Zhong, Chunming Liu, DongYu, JinSong, Pochung Zhang, RenJing, Renhua Song and the numerous visitors and rotation students who have come through the lab.

When I started my PH.D from the cheque document image processing, I have appreciated the camaraderie and local expertise of bank as well as the AAI group early on. In regards to the system design and development, I have appreciated David for his collaboration and the impressive skills.

In my attempted measurements of the crowding structure in an image with clutter scene, I thank the following people for helpful discussions with us: JinSong Xu, Song Yin, DongYu. I thank for their inspirational discussions with us. And I would also like to acknowledge the AAI Facility.

For this dissertation I would like to thank my reading committee members for their time, interest, and helpful comments. I would also like to thank members of my DA defense committee, Jian Zhang, Jinyan Li, Weichang Ye, for their time and insightful questions.

I am grateful to our group's past and present administrative assistant: Lei Zheng, YangBo, Collin Wise and Rayee who kept us organized and was always ready to help. And I am especially great grateful for and touched by the kindness and very generous support from Dr. Dan Luo during tough times.

I gratefully acknowledge the funding sources that made my Ph.D. work possible. I was funded by the Australian IPRS Scholarship and APA for my first 3 years.

My time at UTS was made enjoyable in large part due to the many friends and groups that became a part of my life. I am grateful for time spent with roommates and friends, especially for elder friends as I made big decision for my life, and for many other people and memories. My time at UTS was also enriched by the AAI sport group, Omma's Church and Campsie Congregation.

Lastly, I would like to thank my family for all their love and encouragement. I am very grateful for my parents who raised me with a love of science and supported me in all my pursuits, especially for the presence of my parents here for three of my years here to take care of my kids. For my kids who safeguard my heart and restore my courage. And most of all for my loving, supportive, encouraging, and patient husband whose faithful support during the final stages of this Ph.D. is so appreciated. Thank you.

Xia Cui

The University of Technology, Sydney

June 2014

Table of Contents

CERTIFICATE OF ORIGINAL AUTHORSHIP	I
Abstract	III
Acknowledgement.....	VII
Chapter 1 Introduction	1
1.1 Previous work	1
1.1.1 OCR-based method for text detection.....	1
1.1.2 Feature-based method	2
1.2 Our motivation and aim	5
1.3 Methodology	7
1.4 The framework of our work.....	8
1.5 Organization of our work.....	12
1.6 Contribution	16
Chapter 2 Related Works	19
2.1 Guideline in Ergonomics	19
2.1.1 Legibility	20
2.1.2 Readability	20
2.1.3 Conspicuity	21
2.2 Crowding effect	22
2.2.1 Definition	22
2.2.2 Study objects	23
2.3 Saliency.....	24
2.3.1 General definition.....	24
2.3.2 Saliency map	25
2.3.3 Salient structure.....	26
2.3.4 Other computational schemes	26
2.4 Clutter scene.....	27
2.4.1 Definition	27
2.4.2 Clutter measurement	28

2.5 Basic edge operators	34
2.5.1 Kirsch operator	34
2.5.2 Edge detection with embedded confidence	37
2.6 Summary	43
Chapter 3 Global Properties of Text Appearance	44
3.1 Characteristics of Crowding	44
3.1.1 Crowding, eccentricity and space density of objects ..	45
3.1.2 Anisotropy	46
3.1.3 Asymmetry	46
3.1.4 Crowding depends strongly on target/flanker similarity	47
3.1.5 Statistical properties-average	48
3.2 Theories of Crowding	48
3.2.1 Optical proposals	49
3.2.2 Neuronal proposals	49
3.2.3 Attention proposals	51
3.2.4 Computational proposals	52
3.3 Models of crowding	55
3.4 Breaking crowding	56
3.5 Correlates among crowding, clutter and saliency	56
3.5.1 Saliency – “pop-up” in crowding	56
3.5.2 Crowding in clutter	57
3.5.3 Four basic psychological principals	59
3.6 Summary	61
Chapter 4 the Properties of Individual Character	63
4.1 Aspects of text in typography	64
4.2 The anatomy of type of text	65
4.3 General shape	66
4.3.1 Simple shapes	66
4.3.2 Representation of shape	67
4.4 Width and relative dimension of character	69
4.4.1 Height	69
4.4.2 Height and mean intensity, local contrast	71

4.4.3 Height-to-width ratio.....	76
4.5 Stroke width and contrast.....	78
4.5.1 Stroke width-to-height ratio	78
4.5.2 Background colour, luminance contrast.....	79
4.6 Weight.....	80
4.6.1 Intra-class difference: different type font.....	81
4.6.2 Weight and stroke width	83
4.7 Straight line.....	84
4.8 Size and its related proportion	85
4.9 Summary	87
Chapter 5 Properties of Local Spatial Organization.....	89
5.1 Textual organization and space.....	90
5.2 Letter spacing and word spacing.....	91
5.2.1 Letter spacing	92
5.2.2 Word spacing	93
5.2.3 Setting in space organization calculation.....	93
5.3 Text line	94
5.3.1 Alignment and the importance of neighbourhood	94
5.3.2 Line width	95
5.3.3 Inter-line spacing.....	96
5.4 Keeping balance among spacing.....	97
5.5 Summary	98
Chapter 6 Representation of Image and Text in Clutter Scene	99
6.1 Space-averaged image representation.....	99
6.1.1 The size of pooling region.....	100
6.1.2 Pooling region generating	101
6.1.3 Statistics feature over GM regions.....	105
6.2 Computational model of text	109
6.2.1 Feature of letters.....	110
6.2.2 Letter-centred features	113
6.2.3 Word-centred description.....	114
6.3 Summary	115

Chapter 7 Text Detection Algorithm Based on the Space Averaged Crowding Model	119
7.1 Methodology	119
7.2 Image partition	121
7.2.1 Multi-grey connected component (MGCC)	122
7.2.2 Gradient-based partition	123
7.3 Features extraction	126
7.3.1 Basic features of physical appearance.....	126
7.3.2 The features of space relations among regions	130
7.3.3 Features over GM regions	131
7.4 Clusters analysis over GM and CRs	133
7.4.1 The composition of visual object over GM regions..	133
7.4.2 Analysis of GM regions	136
7.4.3 Cross validation among GM regions and CRs	138
7.5 Experiment	138
7.5.1 Data set.....	138
7.5.2 Evaluation.....	138
7.5.3 Results and discussion.....	139
7.6 Summary	142
Chapter 8 Automatic Processing of Bank Cheques.....	144
8.1 Signature extraction	146
8.1.1 The algorithm framework.....	148
8.1.2 Image partition	149
8.1.3 Context-aware saliency computation model (CSCM)	152
8.1.4 VO analysis	159
8.1.5 Signature experiment.....	160
8.1.6 Extended application of handwritten cheque selection	163
8.2 Automatic extraction of legal amount, payee name.....	165
8.2.1 The flow chart	165
8.2.2 Inferring.....	166
8.3 System.....	171
8.3.1 System block diagram	171
8.3.2 System evaluation	172

8.4 Summary	175
Chapter 9 Conclusion	177
Appendix A.....	183
BIBLIOGRAPHY.....	185

