

Recognition and Prominence Ranking of Alphanumeric Number Sequences in Images

By Alex Cummaudo

BSc Swinburne

Supervised by Prof. Rajesh Vasa, Assoc. Prof. Andrew Cain

*A thesis submitted in partial fulfilment of the requirements for the
Bachelor of Information Technology (Honours)*



Deakin Software and Technology Innovation Laboratory
School of Information Technology
Deakin University, Australia

27 October 2017

Abstract

Text detection in natural images is a growing area with increasing applications, including traffic sign and license plate recognition, and text-based image search. Robustly detecting and recognising text is especially challenging when text is deformed, such as the photometric and geometric distortions of text worn by a moving subject in unstructured scenes. Existing methods of text detection in such cases are classified as learning-based or connected component (CC)-based, applying a mix of enhanced detection techniques—such as stroke width transformation (SWT), canny-edge detection and maximally stable extremal regions (MSERs)—and feeding candidates into optical character recognition (OCR) engines or neural networks to recognise the text. This study proposes applying a learning-based approach using deep-learning strategies to automate the recognition of racing bib numbers (RBNs) in a natural image dataset of various marathons, and then ranking detected subject’s photos in order of prominence. Experimental results showed that these deep-learning strategies performed favourably against other methods using a consistent dataset, prompting further investigation in the generality of the technique developed to other similar subject material.

Declarations

I certify that the the thesis entitled “Recognition and Prominence Ranking of Alphanumeric Number Sequences in Images” submitted for the degree of Bachelor of Information Technology (Honours) is the result of my own work and that where reference is made to the work of others, due acknowledgement is given. I also certify that any material in the thesis which has been accepted for a degree or diploma by any university of institution is identified in the text.

Alex Cummaudo, BSc *Swinburne*
27 October 2017

We certify that the thesis prepared by Alex Cummaudo entitled “Recognition and Prominence Ranking of Alphanumeric Number Sequences in Images” is prepared according to our expectations and that the honours coordinator can proceed to accept this submission for examination.

Prof. Rajesh Vasa
27 October 2017

Assoc. Prof. Andrew Cain
27 October 2017

Acknowledgements

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

I would also like to thank Andrew Cain for his extraordinary efforts over many years to teach hundreds of students (myself included) and who has developed a valued mentorship with me in guiding me throughout my academic life.

Contents

Abstract	iii
Declaration	v
Acknowledgements	vii
Contents	vii
List of Figures	x
List of Tables	xi
List of Abbreviations	xiii
1 Introduction	1
1.1 Background	2
1.2 Motivation	4
1.3 Research Goals	4
1.4 Thesis Organisation	6
2 Background	7
2.1 Applications of Image Recognition	7
2.2 RBN Recognition	7
2.3 Speed Limit Sign Recognition	7
2.4 License Plate Recognition	7
2.5 Detection Strategies	7
2.6 Recognition Strategies	7
2.7 Prominence Strategies	7

3	Data Set	9
4	Benchmarking	11
4.1	Open Source Tools	11
4.2	Existing Pipelines From Literature	11
4.3	Hermes Approach	11
5	Processing Pipeline	13
6	Findings	15
7	Discussion	17
8	Conclusions and Future Work	19
	References	23
A	Ethics Clearance	25
B	Prominence Ranking Survey Results	27

List of Figures

1.1	Sample racing bib numbers	3
1.2	Alphanumeric sequences observed in literature	3

List of Tables

List of Abbreviations

CC Connected Component.

CNN Convolutional Neural Network.

DSTIL Deakin Software and Technology Innovation Laboratory.

LPR Licence Plate Recognition.

NN Neural Network.

OCR Optical Character Recognition.

RBN Racing Bib Number.

TSR Traffic Sign Recognition.

Chapter 1

Introduction

Ever since the camera and phone were unified into smartphones, we have seen an increasing interest for image understanding (specifically to identify the content of an image) but text recognition still faces challenges within images of unstructured scenes. While successes in character recognition have a long history with Optical Character Recognition (OCR) engines [21], these are typically applied under strict conditions (e.g., flatbed scanners for documents without distracting backgrounds). Once applied within the context of a natural scene, real-world discrepancies pose serious shortcomings, such as illumination and viewpoint conditions, blur and glare variations, geometric and photometric distortion, and differences in font size and style. Overcoming these issues has motivated a variety of different techniques in order to realise potential applications that make use of text recognition at scale.

With the ubiquity of smartphone cameras, practical applications of natural image processing have increased. In the last two decades, we have seen the development of point-and-shoot product recognition [10, 24], object detection in videos [20], building recognition [22], image feature extraction to improve visual-based search engines [3, 15], and translation services of American Sign Language gestures [11]. Nonetheless, embedded text within images reveals the largest form of informative features about the image; if text extraction is therefore not robust, information extraction suffers.

Text detection robustness is a factor which severely limits a text recognition pipeline. Research in overcoming such limitations were competed in Lucas et al. [16], where robustness was a key focus in the image processing pipelines proposed. This focus was reiterated by Chen et al. [6], who state the primary prerequisite for text-based recognition (especially within natural scenes) is the text location must be robustly located.

As with any data processing pipeline, false negatives increase where early stages of the

pipeline fail, and therefore detection of these potential candidates must be robust. We can reduce errors in a pipeline where: (1) there may be unwarranted stages of the pipeline, and therefore *excluding* unnecessary stages may also assist in reducing error cases, and (2) by piping through unmatched candidates to further pipelines, which can increase the detection. Both guide in improving robust detection, and this is therefore a key consideration made in our assessment of how useful that pipeline may be. Without the construct of robustness, we restrict these pipelines to very confined conditions, and its usefulness in products is not warranted.

1.1 Background

This study focuses on character recognition in unstructured scenes (Figure 1.1): specifically, short, alphanumeric number sequences. Previous works present methods to extract these sequences in various areas, namely: License Plate Recognition (LPR) systems [1, 5]; Traffic Sign Recognition (TSR) [7, 12, 14, 18]; and, street number recognition, specifically a study by Netzer et al. [17], using Google Street View¹ to determine the numerical value of street numbers. Figure 1.2 highlights typical usage of these sequences.

Different applications apply varying methods to parse short alphanumeric characters. There are typically two stages of any parsing method: *detection* and *recognition*. Detection refers to locating possible candidates and recognition refers to the representation of the text itself. Detection techniques usually are categorised as either Connected Component (CC)-based or learning or texture-based. CC-based detection will typically use a set of distinct properties on the image to detect relevant areas (such as width, stroke and colour) while learning-based feed images into a classifier that can distinguish candidates from false positives. The recognition phase can typically be achieved using optical image recognition OCR engines (such as Tesseract²), machine learning algorithms or deep neural networks (NNs) to classify the detected regions.

This study proposes the development of a learning-based detection and recognition pipeline using deep-learning neural networks within the context of unstructured photos, with a focus on marathon Racing Bib Numbers (RBNs), as shown in Figure 1.1.

¹<https://www.google.com/streetview/> last accessed 13 May 2017.

²<https://github.com/tesseract-ocr/tesseract> last accessed 14 May 2017.



Figure 1.1: Four RBNs in a sample marathon photo.



(a) Successful LPR character segmentation [1]. *Left to right*: original image; region segmentation; character segmentation after negation, height and orientation measurements.



(b) Successful recognition of speed sign digits shown in Eichner and Breckon [7].



(c) Localisation of digits found from varying street view house numbers using the worker described in Netzer et al. [17].

Figure 1.2: Various sample alphanumeric sequences observed in literature.

1.2 Motivation

Detection is harder when the photo is unstructured. Early investigations in Licence Plate Recognition (LPR) systems were systematic in the subject material assessed; a detailed survey by [2] showed that they work best with consistent lighting, specific colour and typeface detection, fixed detection regions, and non-noisy backgrounds. When applied in the context of images with unstructured backgrounds, these systematic approaches begin to have severe limitations as the text components cannot be easily determined.

While further investigations in the area utilise enhanced Connected Component (CC)-based detection [6, 8, 19], performance is likely to degrade as image complexity increases [13]. This is especially relevant when text is geometrically obfuscated, such as malformed Racing Bib Numbers (RBNs) as worn on a marathon runner's torso. Some studies have shown to overcome this by using facial recognition to find a more distinct candidate area [4], but nonetheless relies on such properties like a person's face to detect a number. Similarly, typical recognition techniques interpret text as segmented characters, rather than a single string, though there are exceptions such as in [25].

We also identify subject prominence ranking within natural scenes as an area that has little exploration within literature. (For example, the prominence of a *specific* marathon runner within a scene of many runners.) Prominence ranking is an important field in the context of RBN recognition: runners typically choose not to purchase photos where they have been recognised in an image but are not in the foreground. There are also varying factors which influence purchase likelihood, such as face visibility, eye contact with the camera, and blurriness. An assessment into how the prominence of a runner can be ordered in hundreds of identified photos (based from their recognised RBN) can be used by use of a Neural Network (NN).

This study forms part of an industry project under the Deakin Software and Technology Innovation Laboratory (DSTIL). As a part of the research project, access has been made to a labelled dataset of hundreds of thousands of marathon photos.

1.3 Research Goals

This study aims to develop a processing pipeline that both detects and recognises RBNs on a marathon runner, and then ranks the prominence of each runner detected in the photo. The intention is to explore the viability of artificial deep-learning NNs (such as Convolutional Neural

Networks or CNNs) in the pipeline. Previous studies in RBN recognition [4] and similar areas [7, 12, 23] were heavily heuristic and rule driven.

This primary aim is developed into three key objectives:

Goal 1: *Detect RBNs using a CNN*

Literature has shown that heuristic-based detection algorithms (that are CC-based) are able to detect text within photos [6, 7, 13]. We propose to apply these rule-based techniques to a large labelled dataset within the context of RBNs, and contrast them against a learning-based detection and recognition algorithms (using NNs). By benchmarking against existing libraries and open source tools, we explore if heuristic-based detection algorithms (focusing namely on CC-based detection) outperforms learning-based detection methods. For this goal the research question is framed as:

RQ1) Do CNNs detect RBNs with equal or higher recall and precision rates than CC-based methods?

Goal 2: *Design a CNN that can recognise RBNs*

Typically, traditional alphanumeric sequence parsing can be performed by character segmentation, and then piping those characters into OCR engines. In the context of marathon photos, we explore answers to the following:

RQ2) Does a CNN-based OCR outperform or is at parity with traditional OCR with higher or equal recall and precision rates?

RQ3) Does a CNN-based OCR perform *without* the use of character segmentation?

Goal 3: *Rank prominence of alphanumeric sequences*

Our research objective is aimed to compare if humans are always better at ranking the prominence of an RBN than that of a NN. We can therefore propose the followings research questions:

RQ4) Can a deep-learning NN be trained to rank marathon runners by prominence?, and if so

RQ5) Does a trained deep-learning NN rank prominence of a runner better or equal to that of a human?

1.4 Thesis Organisation

This thesis is organised into the chapters as outlined below. An appendix follows with additional supplementary material.

Chapter 2 - Background Provides an overview of prior studies broadly around the areas of number detection and recognition in image processing and artificial NNs.

Chapter 3 - Data Set Describes the data set to be used, data treatment steps, possible techniques in closer depth to develop a number recognition pipeline, and explores ways to develop prominence ranking techniques.

Chapter 4 - Benchmarking Collates results of a series of experiments using our dataset amongst existing open source tools and pipelines presented in previous work

Chapter 5 - Processing Pipeline Discusses the proposed processing pipeline developed that satisfies the aims of this study.

Chapter 6 - Findings Outlines the method used for validation and presentation of our results.

Chapter 7 - Discussion Presents implications that were found from the results of our findings and limitations.

Chapter 8 - Conclusions and Future Work Draws a number of conclusions and alleviates gaps in the findings of this work by presenting future studies.

Summary

In this chapter we identified some shortcomings in text recognition, developed the context of the study—namely RBN detection. We discussed the general stages that exist for text parsing within natural scenes, detection and recognition, and introduced typical techniques that are applied in this context. We outlined the research aims this study achieves, and how the thesis is organised. The following chapter will detail applications of image processing, using neural networks for image processing, and outline what techniques have been used in previous studies to achieve this.

Chapter 2

Background

We have introduced the context of image processing and neural networks in the related field, and discussed how text capturing within photos is typically achieved in two stages: detection and recognition. Detection techniques are generally classified as either CC- or learning-based. The recognition phases can be applied using traditional OCR engines or, more recently, using artificial NNs. This chapter surveys a number of broad applications where RBN recognition (and related works) are investigated using such phases. The various detection and recognition techniques discussed in the literature are also detailed. We also broadly define the applications of artificial deep-learning NNs in these contexts.

2.1 Applications of Image Recognition

2.2 RBN Recognition

[9]

2.3 Speed Limit Sign Recognition

2.4 License Plate Recognition

2.5 Detection Strategies

2.6 Recognition Strategies

2.7 Prominence Strategies

Chapter 3

Data Set

Chapter 4

Benchmarking

4.1 Open Source Tools

4.2 Existing Pipelines From Literature

4.3 Hermes Approach

Chapter 5

Processing Pipeline

Chapter 6

Findings

Chapter 7

Discussion

Chapter 8

Conclusions and Future Work

References

- [1] Anagnostopoulos, C.-N., I. Anagnostopoulos, V. Loumos, and E. Kayafas (2006). A License Plate-Recognition Algorithm for Intelligent Transportation System Applications. *IEEE Trans. Intelligent Transportation Systems*.
- [2] Anagnostopoulos, C.-N., I. Anagnostopoulos, I. D. Psoroulas, V. Loumos, and E. Kayafas (2008). License Plate Recognition From Still Images and Video Sequences - A Survey. *IEEE Trans. Intelligent Transportation Systems*.
- [3] Bay, H., A. Ess, T. Tuytelaars, and L. J. Van Gool (2008). Speeded-Up Robust Features (SURF). *Computer Vision and Image Understanding*.
- [4] Ben-ami, I., T. Basha, and S. Avidan (2012). Racing Bib Numbers Recognition. In *British Machine Vision Conference 2012*, pp. 19.1–19.10. British Machine Vision Association.
- [5] Cano-Perez, J. and J. C. Pérez-Cortes (2003). Vehicle License Plate Segmentation in Natural Images. *IbPRIA 2652*(Chapter 17), 142–149.
- [6] Chen, H., S. S. Tsai, G. Schroth, D. M. Chen, R. Grzeszczuk, and B. Girod (2011). Robust text detection in natural images with edge-enhanced Maximally Stable Extremal Regions. *ICIP*.
- [7] Eichner, M. L. and T. P. Breckon (2008). Integrated speed limit detection and recognition from real-time video. In *2008 IEEE Intelligent Vehicles Symposium (IV)*, pp. 626–631. IEEE.
- [8] Epshtein, B., E. Ofek, and Y. Wexler (2010). Detecting text in natural scenes with stroke width transform. *CVPR*.
- [9] Fu, C., C.-W. Cheng, W.-H. Shen, Y.-L. Wei, and H.-M. Tsai (2015). LightBib: Marathoner Recognition System with Visible Light Communications. In *2015 IEEE International Conference on Data Science and Data Intensive Systems (DSDIS)*, pp. 572–578. IEEE.

- [10] Girod, B., V. Chandrasekhar, D. Chen, N.-M. Cheung, R. Grzeszczuk, Y. Reznik, G. Takacs, S. Tsai, and R. Vedantham (2011). Mobile Visual Search. *IEEE Signal Processing Magazine* 28(4), 61–76.
- [11] Jin, C. M., Z. Omar, and M. H. Jaward (2016). A mobile application of American sign language translation via image processing algorithms. In *2016 IEEE Region 10 Symposium (TENSYP)*, pp. 104–109. IEEE.
- [12] Kundu, S. K. and P. Mackens (2015). Speed Limit Sign Recognition Using MSER and Artificial Neural Networks. *ITSC*.
- [13] Li, Y. and H. Lu (2012). Scene text detection via stroke width. *ICPR*.
- [14] Lian, Z., X. Jing, S. Sun, and H. Huang (2016). Frequency Selective Convolutional Neural Networks for Traffic Sign Recognition. In *2016 IEEE 83rd Vehicular Technology Conference (VTC Spring)*, pp. 1–5. IEEE.
- [15] Lowe, D. G. (2004). Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision* 60(2), 91–110.
- [16] Lucas, S. M., A. Panaretos, L. Sosa, A. Tang, S. Wong, and R. Young (2003). ICDAR 2003 robust reading competitions. In *Seventh International Conference on Document Analysis and Recognition*, pp. 682–687. IEEE Comput. Soc.
- [17] Netzer, Y., T. Wang, and A. Coates (2011). Reading digits in natural images with unsupervised feature learning. *NIPS workshop on . . .*
- [18] Seo, Y.-W., J. Lee, W. Zhang, and D. Wettergreen (2015). Recognition of Highway Workzones for Reliable Autonomous Driving. *IEEE Transactions on Intelligent Transportation Systems* 16(2), 1–11.
- [19] Shivakumara, P., T. Q. Phan, and C. L. Tan (2011). A Laplacian Approach to Multi-Oriented Text Detection in Video. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33(2), 412–419.
- [20] Sivic, J. and A. Zisserman (2003). Video Google - A Text Retrieval Approach to Object Matching in Videos. *ICCV*.

- [21] Smith, R. W. (1987). *The Extraction and Recognition of Text from Multimedia Document Images*. Ph. D. thesis, University of Bristol.
- [22] Takacs, G., V. Chandrasekhar, N. Gelfand, Y. Xiong, W.-C. Chen, T. Bismpiagiannis, R. Grzeszczuk, K. Pulli, and B. Girod (2008). Outdoors augmented reality on mobile phone using loxel-based visual feature organization. *Multimedia Information Retrieval*.
- [23] Torresen, J., J. W. Bakke, and L. Sekanina (2004). Efficient recognition of speed limit signs. In *The 7th International IEEE Conference on Intelligent Transportation Systems*, pp. 652–656. IEEE.
- [24] Tsai, S. S., D. M. Chen, V. Chandrasekhar, G. Takacs, N.-M. Cheung, R. Vedantham, R. Grzeszczuk, and B. Girod (2010). Mobile product recognition. *ACM Multimedia*.
- [25] Zhu, L., C.-S. Yang, and J.-S. Pan (2016). Detection and Recognition of Speed Limit Sign from Video. *ACIIDS*.

Appendix A

Ethics Clearance

Appendix B

Prominence Ranking Survey Results