

# **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

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*  
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  
October 2017

---

Assoc. Prof. Andrew Cain  
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.





# Contents

<b>Abstract</b>	<b>iii</b>
<b>Declaration</b>	<b>v</b>
<b>Acknowledgements</b>	<b>vii</b>
<b>Contents</b>	<b>vii</b>
<b>List of Figures</b>	<b>x</b>
<b>List of Tables</b>	<b>xi</b>
<b>List of Abbreviations</b>	<b>xiii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background . . . . .	2
1.2 Motivation . . . . .	4
1.3 Research Goals . . . . .	4
1.4 Thesis Organisation . . . . .	6
<b>2 Background</b>	<b>7</b>
2.1 Detection Strategies . . . . .	7
2.1.1 CC-based techniques . . . . .	8
2.1.2 Learning-based techniques . . . . .	8
<b>3 Data Set</b>	<b>9</b>
<b>4 Benchmarking</b>	<b>11</b>
4.1 Open Source Tools . . . . .	11
4.2 Existing Pipelines From Literature . . . . .	11

4.3	Hermes Approach . . . . .	11
<b>5</b>	<b>Processing Pipeline</b>	<b>13</b>
<b>6</b>	<b>Findings</b>	<b>15</b>
<b>7</b>	<b>Discussion</b>	<b>17</b>
<b>8</b>	<b>Conclusions and Future Work</b>	<b>19</b>
	<b>References</b>	<b>26</b>
<b>A</b>	<b>Ethics Clearance</b>	<b>27</b>
<b>B</b>	<b>Prominence Ranking Survey Results</b>	<b>29</b>

# 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. 2, 4, 5, 7, 8

**CNN** Convolutional Neural Network. 5

**DSTIL** Deakin Software and Technology Innovation Laboratory. 4

**LPR** License Plate Recognition. 2, 4

**NN** Neural Network. 2, 4–7

**OCR** Optical Character Recognition. 1, 2, 5, 7

**RBN** Racing Bib Number. 2–6

**SWT** Stroke Width Transformation. 8

**TSR** Traffic Sign Recognition. 2





# 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 [50], 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 conditions, viewpoint and perspective differences, blur and glare variations, geometric and photometric distortion, and differences in font size and style [24, 59]. Overcoming these issues has motivated a variety of techniques 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 [15, 55], object detection in videos [47], building recognition [53], image feature extraction to improve visual-based search engines [3, 36], and translation services of American Sign Language gestures [20]. Nonetheless, embedded text within images contains indexable data on the image’s semantics [48]; 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 have been competed numerous competitions [18, 37, 38, 44], where robustness is the key focus in the image processing pipelines proposed. This focus was reiterated by Chen et al. [9], 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, and thus robustness, in a pipeline where: (1) there are unwarranted stages (*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.

Without the construct of robustness, we restrict these pipelines to very confined conditions, and its usefulness in products is not warranted. Therefore, the robustness of text extraction pipelines are imperative to gapping the semantic extraction of information from an image [48], and solving this issue can assist in applications of image processing and data indexing of content within images [14] of paramount proportions.

## 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, 6]; Traffic Sign Recognition (TSR) [12, 26, 31, 42]; and, street number recognition, specifically a study by Netzer et al. [41], using Google Street View<sup>1</sup> 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 Character Recognition (OCR) engines (such as Tesseract<sup>2</sup>) [4], machine learning algorithms [26, 28, 41] or deep Neural Network (NN) to classify the detected regions [21, 31, 43].

*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)<sup>3</sup>, as shown in Figure 1.1.*

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

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

<sup>3</sup>While referred to as numbers, some RBNs have alphabetic identifiers in them.



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 [12].



(c) Localisation of digits found from varying street view house numbers using the worker described in Netzer et al. [41].

Figure 1.2: Various sample alphanumeric sequences observed in literature.

## 1.2 Motivation

Detection is harder when the photo is unstructured. Early investigations in License 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 limitations as the text components cannot be easily determined.

While further investigations in the area utilise enhanced Connected Component (CC)-based detection [9, 13, 46], performance is likely to degrade as image complexity increases [30]. This is especially relevant when text is geometrically obfuscated, such as malformed Racing Bib Numbers (RBNs) as worn on a marathon runner's torso. Malformed, in this sense, is caused by non-flat bib sheets that tend to follow the runner's body shape, in addition to images that are taken in dynamical contexts. Some studies have shown to overcome this by using facial recognition to find a more distinct candidate area [4], but nonetheless rely on 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 Zhu et al. [62].

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 in-

tention is to explore the viability of artificial deep-learning NNs—such as Convolutional Neural Networks (CNNs)—in the pipeline. Previous studies in RBN recognition [4] and similar areas [12, 26, 54] 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 [9, 12, 30]. 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 approach outperform or is at parity with traditional OCR approaches with higher or equal recall and precision rates?

**RQ3)** Does a CNN-based OCR algorithm 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 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 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 NN.

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 NN in these contexts.

### 2.1 Detection Strategies

Text extraction strategies have seen continuing interest in the literature, with many comprehensive surveys assessing the state of the art [7, 23, 24, 32, 59]. It is widely demonstrated that if text within an unstructured scene is *detected* reliably, then existing OCR engines can suitably extract these characters [49] once they exist in a structured context; thus not every . A survey into the two prominent detection strategies is given in Sections 2.1.1 to 2.1.2.

These two prominent strategies have a varied nomenclature: (1) the CC-based (or *region*-based) approach, that utilise different region properties (e.g., colour, edges, CCs) [9, 13, 19, 25, 29, 30, 34, 35, 45, 46, 51, 52, 60, 61] for unsupervised extraction; and, (2) learning-based (or *texture*-based) approach, which uses unique texture properties to supervise extraction text from its background [10, 11, 16, 27, 56–58]. Additionally, some authors have proposed methods to combine these unsupervised and supervised techniques [5, 39, 40].

### 2.1.1 CC-based techniques

CC-based approaches generate separated CCs using properties such as stroke width, pixel colour and edges, typically applying geometric and texture filters to reduce false positives. Neighbouring pixels are then ‘grouped’ using an algorithm originally presented by Horn [17].

Previous work required the use of a scanning window [10, 22, 33] which is limited by a constant image scale and discrete orientations of the sliding (thereby preventing text strokes in non-linear directions). However, a study by Epshtein et al. [13] (and coincidentally Zhang and Kasturi [61]) introduced the concept of Stroke Width Transformation (SWT).

### 2.1.2 Learning-based techniques

typically referred to as a *learning*-based approach, due to the common use of machine learning methods utilised

Typically, texture-based approaches utilise supervised learning methods, though it is typical for these classifiers to required thousands of training images [8]. Additionally, these methods

## Summary



# **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] Bengio, Y., P. Lamblin, D. Popovici, and H. Larochelle (2006). Greedy Layer-Wise Training of Deep Networks. *NIPS*.
- [6] Cano-Perez, J. and J. C. Pérez-Cortes (2003). Vehicle License Plate Segmentation in Natural Images. *IbPRIA 2652*(Chapter 17), 142–149.
- [7] Chen, D. and J. Luettin (2000). A survey of text detection and recognition in images and videos.
- [8] Chen, D., J.-M. Odobez, and J.-P. Thiran (2004). A localization/verification scheme for finding text in images and video frames based on contrast independent features and machine learning methods. *Sig. Proc. - Image Comm.*.
- [9] 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*.

- [10] Chen, X. and A. L. Yuille (2004). Detecting and reading text in natural scenes. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004*, pp. 366–373. IEEE.
- [11] Chen, X. and A. L. Yuille (2005). A Time-Efficient Cascade for Real-Time Object Detection - With applications for the visually impaired. *CVPR Workshops*.
- [12] 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.
- [13] Epshtein, B., E. Ofek, and Y. Wexler (2010). Detecting text in natural scenes with stroke width transform. *CVPR*.
- [14] Faloutsos, C., R. Barber, M. Flickner, J. Hafner, W. Niblack, D. Petkovic, and W. Equitz (1994). Efficient and effective Querying by Image Content. *Journal of Intelligent Information Systems* 3(3-4), 231–262.
- [15] 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.
- [16] Hanif, S. M. and L. Prevost (2009). Text Detection and Localization in Complex Scene Images using Constrained AdaBoost Algorithm. *ICDAR*.
- [17] Horn, B. (1986, January). *Robot Vision*. MIT Press.
- [18] Hua, X.-S., L. Wenyin, and H. Zhang (2004). An automatic performance evaluation protocol for video text detection algorithms. *IEEE Trans. Circuits Syst. Video Techn.*.
- [19] Jain, A. K. and B. Y. 0002 (1998). Automatic text location in images and video frames. *ICPR*.
- [20] 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.
- [21] Jin, J., K. Fu, and C. Zhang (2014, September). Traffic Sign Recognition With Hinge Loss Trained Convolutional Neural Networks. *IEEE Transactions on Intelligent Transportation Systems* 15(5), 1991–2000.

- [22] Jung, C., Q. Liu, and J. Kim (2009, January). A stroke filter and its application to text localization. *Pattern Recognition Letters* 30(2), 114–122.
- [23] Jung, K., K. In Kim, and A. K Jain (2004, May). Text information extraction in images and video: a survey. *Pattern Recognition* 37(5), 977–997.
- [24] Jung, K., K. I. Kim, and A. K. Jain (2004). Text information extraction in images and video - a survey. *Pattern Recognition*.
- [25] Kim, H.-K. (1996). Efficient Automatic Text Location Method and Content-Based Indexing and Structuring of Video Database. *J. Visual Communication and Image Representation*.
- [26] Kundu, S. K. and P. Mackens (2015). Speed Limit Sign Recognition Using MSER and Artificial Neural Networks. *ITSC*.
- [27] Lee, C. W., K. Jung, and H. J. Kim (2003, November). Automatic text detection and removal in video sequences. *Pattern Recognition Letters* 24(15), 2607–2623.
- [28] Lee, E. R., P. K. Kim, and H. J. Kim (1994). Automatic Recognition of a Car License Plate using Color Image Processing. *ICIP* 2, 301–305.
- [29] Lee, S., M. S. Cho, K. Jung, and J. H. Kim (2010). Scene Text Extraction with Edge Constraint and Text Collinearity. *ICPR*.
- [30] Li, Y. and H. Lu (2012). Scene text detection via stroke width. *ICPR*.
- [31] 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.
- [32] Liang, J., D. S. Doermann, and H. Li (2005). Camera-based analysis of text and documents - a survey. *IJDAR*.
- [33] Lienhart, R. and A. Wernicke (2002). Localizing and segmenting text in images and videos. *IEEE Trans. Circuits Syst. Video Techn.*.
- [34] Liu, Y., S. Goto, and T. Ikenaga (2006). A Contour-Based Robust Algorithm for Text Detection in Color Images. *IEICE Transactions*.

- [35] Liu, Z. and S. Sarkar (2008). Robust outdoor text detection using text intensity and shape features. *ICPR*.
- [36] Lowe, D. G. (2004). Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision* 60(2), 91–110.
- [37] Lucas, S. M. (2005). ICDAR 2005 text locating competition results. In *Eighth International Conference on Document Analysis and Recognition (ICDAR'05)*, pp. 80–84 Vol. 1. IEEE.
- [38] 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.
- [39] Mairal, J., F. R. Bach, J. Ponce, G. Sapiro, and A. Zisserman (2008). Discriminative learned dictionaries for local image analysis. *CVPR*, 10.
- [40] Mutch, J. and D. G. Lowe (2006). Multiclass Object Recognition with Sparse, Localized Features. *CVPR*.
- [41] Netzer, Y., T. Wang, and A. Coates (2011). Reading digits in natural images with unsupervised feature learning. *NIPS workshop on . . .*
- [42] 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.
- [43] Sermanet, P. and Y. LeCun (2011). Traffic sign recognition with multi-scale Convolutional Networks. *IJCNN*.
- [44] Shahab, A., F. Shafait, and A. Dengel (2011). ICDAR 2011 Robust Reading Competition Challenge 2: Reading Text in Scene Images. In *2011 International Conference on Document Analysis and Recognition (ICDAR)*, pp. 1491–1496. IEEE.
- [45] Shivakumara, P., W. Huang, T. Q. Phan, and C. L. Tan (2010). Accurate video text detection through classification of low and high contrast images. *Pattern Recognition*.
- [46] 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.



- [47] Sivic, J. and A. Zisserman (2003). Video Google - A Text Retrieval Approach to Object Matching in Videos. *ICCV*.
- [48] Smeulders, A. W. M., M. Worring, S. Santini, A. Gupta, and R. C. Jain (2000). Content-Based Image Retrieval at the End of the Early Years. *IEEE Trans. Pattern Anal. Mach. Intell.*.
- [49] Smith, R. (2007). An Overview of the Tesseract OCR Engine. In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007) Vol 2*, pp. 629–633. IEEE.
- [50] Smith, R. W. (1987). *The Extraction and Recognition of Text from Multimedia Document Images*. Ph. D. thesis, University of Bristol.
- [51] Subramanian, K., P. Natarajan, M. Decerbo, and D. A. Castañón (2007). Character-Stroke Detection for Text-Localization and Extraction. *ICDAR*.
- [52] Sun, Q., Y. Lu, and S. Sun (2010). A Visual Attention Based Approach to Text Extraction. *ICPR*.
- [53] Takacs, G., V. Chandrasekhar, N. Gelfand, Y. Xiong, W.-C. Chen, T. Bismpigiannis, R. Grzeszczuk, K. Pulli, and B. Girod (2008). Outdoors augmented reality on mobile phone using loxel-based visual feature organization. *Multimedia Information Retrieval*.
- [54] 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.
- [55] 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*.
- [56] Tu, Z., X. Chen, A. L. Yuille, and S. C. Zhu (2003). Image Parsing - Unifying Segmentation, Detection, and Recognition. *ICCV*.
- [57] Wang, X., L. Huang, and C. Liu (2009). A New Block Partitioned Text Feature for Text Verification. *ICDAR*.
- [58] Ye, Q., Q. Huang, W. G. 0001, and D. Zhao (2005). Fast and robust text detection in images and video frames. *Image Vision Comput.*.

- [59] Zhang, J. and R. Kasturi (2008). Extraction of Text Objects in Video Documents: Recent Progress. In *2008 The Eighth IAPR International Workshop on Document Analysis Systems (DAS)*, pp. 5–17. IEEE.
- [60] Zhang, J. and R. Kasturi (2010). Text Detection Using Edge Gradient and Graph Spectrum. *ICPR*.
- [61] Zhang, J. and R. Kasturi (2011). Character Energy and Link Energy-Based Text Extraction in Scene Images. In *Computer Vision – ACCV 2010*, pp. 308–320. Berlin, Heidelberg: Springer Berlin Heidelberg.
- [62] 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**