

Recognition and Prominence Ranking of Alphanumeric Number Sequences in Images

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 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 that has been accepted for a degree or diploma by any university or 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

Dedicated to Tom Fellowes.

Acknowledgements

I thank everyone at DSTIL for fostering a highly supportive and motivating environment to work in: you make every day enjoyable and are all my second family. I especially thank Simon Vajda and Rodney Pilgrim for their direction throughout my honours year, the exceptional research guidance provided by Raj Vasa and Nicola Pastorello as well as for our fruitful discussions and well-valued feedback on this thesis, and Jake Renzella for his much-needed support and friendship throughout the last two years. I also thank Andrew Cain for his extraordinary efforts over many years to teach thousands of students (myself included) and with whom I have developed a highly-valued mentorship with in guiding me throughout life. And lastly, I thank my loving and patient partner Tom Fellowes, without whom this thesis (indeed this year) would not be possible.

Contents

Abstract	iii
Declaration	v
Acknowledgements	ix
Contents	ix
List of Abbreviations	xii
List of Figures	xiv
List of Tables	xv
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 Detection Strategies	7
2.1.1 Connected Component-based techniques	8
2.1.2 Learning-based techniques	13
2.2 Recognition Strategies	16
2.3 Metrics	21
2.3.1 Precision and Recall	21
2.3.2 The f -score	22

3	Dataset	25
3.1	A Data-Capturing Architecture	25
3.1.1	Methodology	26
3.1.2	What to Capture?	26
3.1.3	Describing the Metamodel	30
3.1.4	The Data Capturing Process	31
3.1.5	Layering	31
4	Benchmarking	33
4.1	Open Source Tools	33
4.2	Existing Pipelines From Literature	33
4.3	Hermes Approach	33
5	Processing Pipeline	35
6	Findings	37
7	Discussion	39
8	Conclusions and Future Work	41
References		53
A	Ethics Clearance	55
B	Prominence Ranking Survey Results	57

List of Abbreviations

AI Artificial Intelligence. 33, 34

CC Connected Component. 3, 6, 7, 11–17, 21, 26, 29, 30

CNN Convolutional Neural Network. 7, 8, 14, 19, 20, 30

DSTIL Deakin Software and Technology Innovation Laboratory. 6

ETL Extract-Transform-Load. 40

HOG Histogram of Oriented Gradient. 18, 23, 29

ICDAR International Conference on Document Analysis and Recognition. 21, 27, 30

LBP Local Binary Patterns. 18, 29

LoP Likelihood of Purchase. 36, 38

LPR License Plate Recognition. 2, 6

MLP Multilayer Perceptron. 19, 22, 23

MSER Maximally Stable Extremal Region. 14, 22, 29

NN Neural Network. 3, 6–8, 11, 18, 22–24, 29, 30, 34

OCR Optical Character Recognition. 1, 3, 7, 8, 11, 21–24, 29, 30

PNN Probabilistic Neural Network. 23, 24

RBN Racing Bib Number. 3, 4, 6, 7, 9, 11, 21, 23–26, 30, 35, 38, 39

SIFT Scale Invariant Feature Transform. 18

SURF Speeded-Up Robust Features. 18

SVM Support Vector Machine. 18, 19, 29

SWT Stroke Width Transformation. 13, 14, 25, 26, 29

TSR Traffic Sign Recognition. 2, 22, 23, 30

List of Figures

1.1	Sample racing bib numbers	3
1.2	Alphanumeric sequences observed in literature	3
2.1	Stroke analysis from Subramanian et al. [109]	8
2.2	Stroke Width Transformation from Epshtain et al. [28]	9
2.3	Using contrast-enhanced MSERs to detect text	10
2.4	Text energy for connecting candidates back together	11
2.5	Using graph spectrum to cluster CCs	11
2.6	Skeletonisation process of CCs	12
2.7	The Mask R-CNN framework for instance segmentation	15
2.8	Use of CNNs for object detection	15
2.9	A NN designed to recognise speed limit signs	17
2.10	A PNN used to recognise license plate characters	18
2.11	Text recognition pipelines dependent on heuristics	19
2.12	Character processing for feeding an OCR engine	20
2.13	Drawbacks of heuristic-driven approaches for RBN detection	20
2.14	Overlapping areas of ground truth and estimated targets	23

List of Tables

2.1	Top-scoring word recognition results from ICDAR 2011–2015.	16
2.2	Survey of text extraction literature	23
3.1	Summary of annotations captured in the dataset	29

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 [105], 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 [54, 125]. 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 [35, 114], object detection in videos [102], building recognition [112], image feature extraction to improve visual-based search engines [4, 79], and translation services of American Sign Language gestures [50]. Nonetheless, embedded text within images contains indexable data on the image’s semantics [103]; if text extraction is therefore not robust, information extraction suffers.

Text detection robustness is a factor that severely limits a text recognition pipeline. Research in overcoming such limitations have been contested numerous competitions [46, 80, 81, 99], where robustness is the key focus in the image processing pipelines proposed. This focus was reiterated by Chen et al. [17], 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 [103], and solving this issue can assist in applications of image processing and data indexing of content within images [29] 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 [2, 12]; Traffic Sign Recognition (TSR) [27, 61, 70, 97]; and, street number recognition, specifically a study by Netzer et al. [89], 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 Character Recognition (OCR) engines (such as Tesseract²) [6], machine learning algorithms [61, 65, 89] or deep Neural Network (NN) to classify the detected regions [51, 70, 98].

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)³, 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.

³While 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 [2]. *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 [27].



(c) Localisation of digits found from varying street view house numbers using the worker described in Netzer et al. [89].

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 [3] 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 [17, 28, 101], performance is likely to degrade as image complexity increases [68]. 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 [6], 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. [128].

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 that 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 [6] and similar areas [27, 61, 113] 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 [17, 27, 68]. 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 a 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 Optical Character Recognition (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 - Dataset Describes the dataset 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

Text capturing within unstructured photos from the wild are typically achieved in two stages: detection and recognition. Detection techniques are classified as either CC- or learning-based. The recognition phase uses traditional OCR engines or, more recently, artificial NNs.

In this chapter, we survey different applications where RBN recognition (and related works) are investigated. Various detection and recognition techniques discussed in literature are detailed.

2.1 Detection Strategies

Text extraction strategies have seen continuing interest in the literature, with many comprehensive surveys assessing the state of the art [13, 53, 54, 71, 125]. It is widely demonstrated that if text within an unstructured scene is *detected* reliably, then existing OCR engines can suitably extract these characters [104] once they exist in a structured context; thus not every extraction pipeline needs to self-contain a recognition strategy if commercial OCR packages suffice. 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) [17, 28, 49, 57, 66, 68, 77, 78, 100, 101, 109, 110, 126, 127] for unsupervised extraction; and, (2) learning-based (or *texture-based*) approach, which uses unique texture properties to supervise extraction text from its background [19, 21, 25, 38, 40, 63, 67, 73, 85, 92, 93, 115, 120, 123]. Some authors have proposed methods that mix both supervised and unsupervised techniques [7, 82, 88].

2.1.1 Connected Component-based techniques

Connected Component (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 such as the one originally presented by Horn [44].

Previous work required the use of a scanning window [19, 52, 73] that is limited by a constant image scale and discrete orientations of the sliding (thereby preventing text strokes in non-linear directions). Subramanian et al. [109] overcame this limitation by implementing an algorithm to detect text strokes by scanning an image horizontally and looking for sudden changes in background intensity (Figure 2.1b). However, this algorithm assumes a darker text on a lighter background to find such intensity changes, and consequently there are numerous parameters that must be fine-tuned. Additionally, the algorithm is only able to detect horizontal text only, and detected strokes are not grouped into characters, words and sentences.

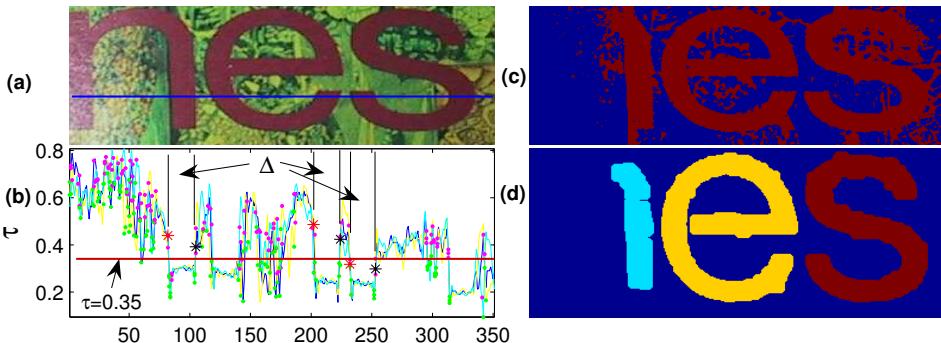


Figure 2.1: A study from Subramanian et al. [109] showed that stroke width could be determined from (a) the original image; (b) intensity plots of the image to determine stroke regions (τ is the intensity threshold and Δ is the stroke width); (c) the intensity at an optimal threshold; (d) the final thresholded image after morphological operations and CC analysis.

A study by Epshtain et al. [28] (and coincidentally Zhang and Kasturi [127]) built on the idea presented by Subramanian et al., and introduced the concept of Stroke Width Transformation (SWT), a local image operator that determines the most likely stroke of a given pixel by computing the per-pixel width. This was later expanded in Srivastav and Kumar [108]. The SWT approach overcame previous limitations by introducing a system that can detect text regardless of size, typeface, direction and language, making it one of the first widely cited multilingual text detection algorithms. Additionally, SWT overcame methods that required the use of an OCR filtering stage to reduce false positives [19, 20, 123]. A sample of SWT is shown in Figure 2.2.

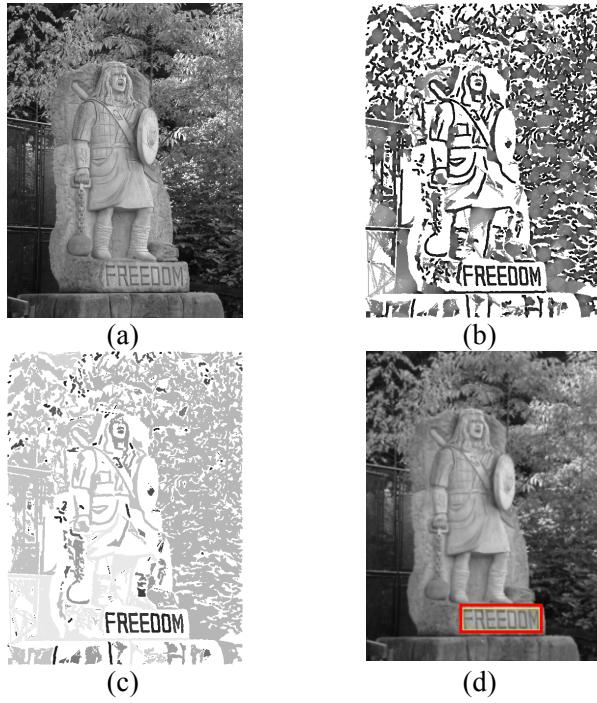


Figure 2.2: The Stroke Width Transformation (SWT) approach introduced in [28]. The original image (a) is converted to a binarised array with the most likely stroke width per-pixel (b), piping the information into geometric filtering (c) as text maintains fixed stroke width (excluding false positives such as foliage). The resulting detected text is shown in (d).

It is common to see edges computed from a raw image using the Canny-Edge Detection algorithm [11]. This was successfully applied in various CC-based studies [17, 28, 126]. While different researchers have exploited SWT and adapted it further [101, 127], when opposite edges are not parallel, the SWT forms candidates with holes appearing in stroke curves or joints. This is due to candidates formed by shooting rays from detected the edges along the gradient found, removing the rays if terminated by another edge pixel of a perpendicular gradient. Further limitations include undetected stronger highlights, blurry text, and text with a wide curvature.

An alternate approach that overcomes this limitation was introduced by Chen et al. [17], where the complimentary properties of Canny-Edges [11] and Maximally Stable Extremal Regions (MSERs) [83] were combined. MSER is a detection mechanism suited for region-based detection, is robust against varying viewpoints, scales and illuminations [84], and can be extracted from images efficiently [90]. A limitation of MSER is its sensitivity to image blur [84], but Chen et al. demonstrated that MSER can be edge-enhanced using Canny-Edges on a contrast-enhanced image (Figure 2.3), achieving comparable results to SWT presented by Epshtain et al. [28]. Multiple works have utilised MSERs in a wide range of applications, such as their use in teaching CNNs and real-time text extraction [39, 47, 61, 68, 124].

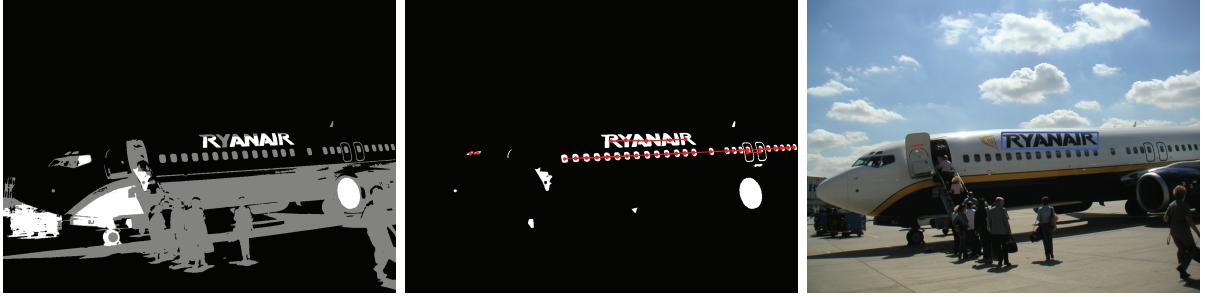


Figure 2.3: Extracting text from a natural image shown in Chen et al. [17]. *From left to right:* Detected Maximally Stable Extremal Regions (MSERs) of black-on-white objects; text candidates grouped to formed text lines after geometric and stroke width filtering; false positives rejected using text verification showing detected text in the blue box.

A significant requirement of all CC-based techniques are the requirements to cluster extracted components back together again. This, in turn, also helps to remove any false positives by removing properties that don't meet set criteria. Various proposals have been made:

- Epshtain et al. [28] use basic geometric filtering based on the stroke width detected and height ratios of candidates. Additionally, colours of candidates are averaged as it is expected that words be written in the same colour. These are then clustered into candidates pairs (of at least three letters), chained together if they share a similar direction.
- Zhang and Kasturi [127] investigate the spacial relationship and property similarity of two neighbouring candidates, computing their link energies to compute *text energy* (the probability a candidate is a true positive). The distance of the text energies are computed and, where beyond a set threshold, will be eliminated if not met. This is presented in Figure 2.4.
- Zhang and Kasturi [126] expand the use of graph spectrum, which has successfully been used in computer vision [96] to show image features in the form of a graph, use an adjacency matrix, then gather clusters of CCs based on the positive eigenvectors of the graph. This process is illustrated in Figure 2.5.
- Shivakumara et al. [101] propose the use of skeletal distance maps of a CC to remove small artefacts and reduce false positives. They define *simple* and *complex* CCs respectively as: (a) a single text string or false positive; (b) multiple text strings that are connected to each other. The skeletonisation process is shown in Figure 2.6.

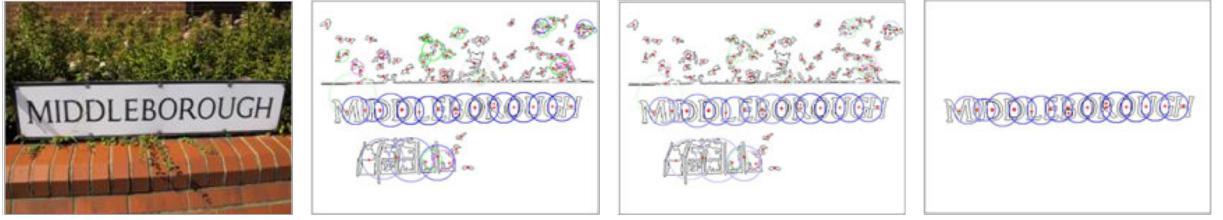


Figure 2.4: Using link and text energies for reconnecting character candidates shown in Zhang and Kasturi [127]. *From left to right:* original image; all link energies determined in a given image (note the false positives of background foliage); text energies calculated; all text energies greater than 0.5.

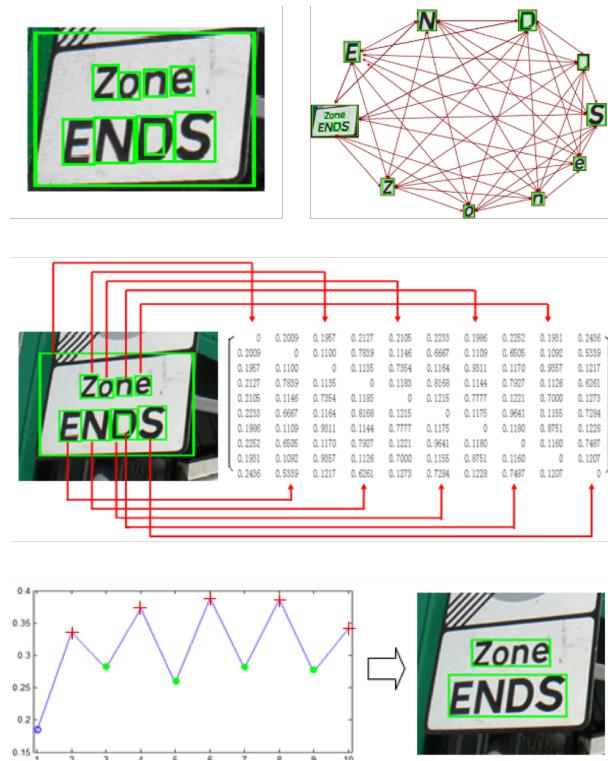


Figure 2.5: Process of grouping components via graph spectrum [126]. *Top-left:* 10 CCs detected. *Top-right:* generated graph from detected candidates. *Middle:* generated adjacency matrix. *Bottom-left:* Positive eigenvector resulting from the graph spectrum. *Bottom-right:* Resulting bounding boxes.



Figure 2.6: Developing a skeletal map using the process proposed by [101]. *Top row:* Original image is processed using Fourier-Laplacian filtering. A maximum distance map is developed and parsed through a morphological operation to remove smaller artefacts. *Middle row:* CC classification and further skeletonisation, showing five labelled subcomponents. *Bottom row:* The five sample subcomponents extracted from the skeletonisation process (in order). Note that subcomponents 4 and 5 are false positives.

2.1.2 Learning-based techniques

The learning-based approach has had a varied popularity over the years. While also referred to as a *texture*-based approach, it typically utilises learning-based methods to train a classifier using these textures, considering text as a special texture within the image. This is done by extracting certain features over a portion of the image (i.e., the texture) either via heuristics or machine learning. The texture is typically extracted by scanning at varying scales, then classifying areas of pixels based on certain features. A classifier uses these features to identify text from non-text to extract text from its background.

Early works in the area focused primarily on how to classify a region as either text or non-text [58, 67, 73, 106]. However, while it may be easy to procure training samples of text, it is often more difficult to procure non-text training samples [43, 111] due to the wide variance of ‘non-text’ samples, which need to be well-represented in the training set of the classifier.

A varied range of features have been utilised to mark textures. A recent feature made popular is the Histogram of Oriented Gradient (HOG), a object detection technique that utilises gradient orientation, first conceptualised in [1] and later coined by Freeman and Roth [30]. However, its popularity did not become widespread until [25], where Dalal and Triggs showed its applicability was shown on pedestrian detection. Later, HOG was shown to be useful in text detection by Hanif et al. [41]. Further features include: Local Binary Patterns (LBP) [91], which was shown to improve detection when combined with HOG [119]; wavelet energy [86, 117], which were applied to extract subtitles [38] and worked optimally with HOG in text extraction [92] (when compared to other features); Gabor filters [75]; Scale Invariant Feature Transform (SIFT) descriptors [79]; grey scale features [58]; Speeded-Up Robust Features (SURF) [4], edge map features [14]; and shape contexts [5].

Many of these features can be combined and fed into a single or multiple classifiers [38, 40, 41, 92, 115, 120, 123]. Example classifiers include Support Vector Machines (SVMs) [10, 23, 116], the Adaptive Boosting (AdaBoost) classifier [31]—and variants thereof ([32, 40, 107])—or NNs, such as Multilayer Perceptrons (MLPs), though Chen et al. [15] report that the use of SVMs showed better text texture verification than MLPs.

As shown in the wide range of features and classifiers, a main limitation of learning-based methods is the difficulty in selecting which combinations of features to use, the inability to detect sufficiently slanted text, as well as its reported high computational complexity due to the need to scan the image multiple time at different scales [28, 68]. Furthermore, it is typical for

these classifiers to required thousands of training images [16].

More recently though, deep-learning CNNs have been utilised for instance classification and per-pixel segmentation, as emphasised in Figure 2.8. While CNNs are a relatively old concept (see LeCun et al. [62]), they lost interest within image processing to Support Vector Machines (SVMs) and AdaBoost throughout the late 1990s and 2000s. In 2012, however, Krizhevsky et al. [59] sparked interest in their use by demonstrating far increased classification accuracy within images in the ImageNet Large Scale Visualisation Recognition Challenge [26], with only a few modifications made to the CNNs proposed by LeCun et al. more than a decade prior.

This has since led to the rise of enhanced CNNs such as FICS+++ [69], R-CNN [36], Fast/Faster R-CNN [37, 94], and Mask R-CNN [42]. A comparison of two of these networks using images in the Microsoft Common Objects in Context [74] dataset are shown in Figure 2.8, with the Mask R-CNN framework illustrated in Figure 2.7. The rise of these recent developments have sparked interest in multi-language machine learning libraries, such as MXNet¹, which are able to combine these academic works for use in large-scale industry-focused production code [18].

These methods are increasingly improving per-pixel classification of objects within natural scenes. However, applicability of these deep-learning networks in the sole context of *text extraction* is yet to be widely explored.

¹<http://mxnet.io> last accessed 30 June 2017.

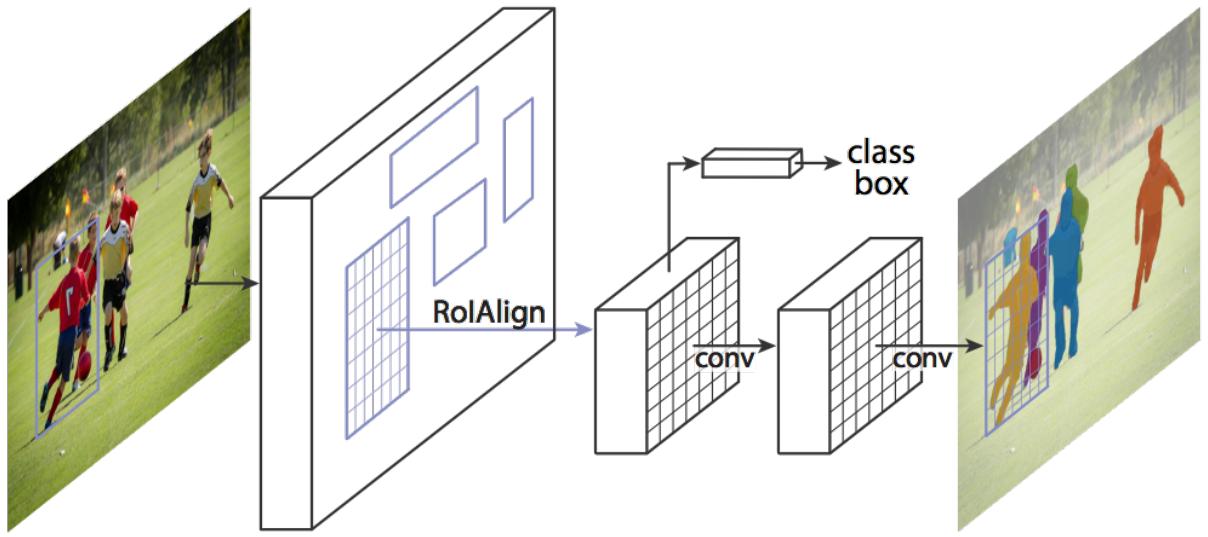


Figure 2.7: The Mask R-CNN framework for instance segmentation [42].

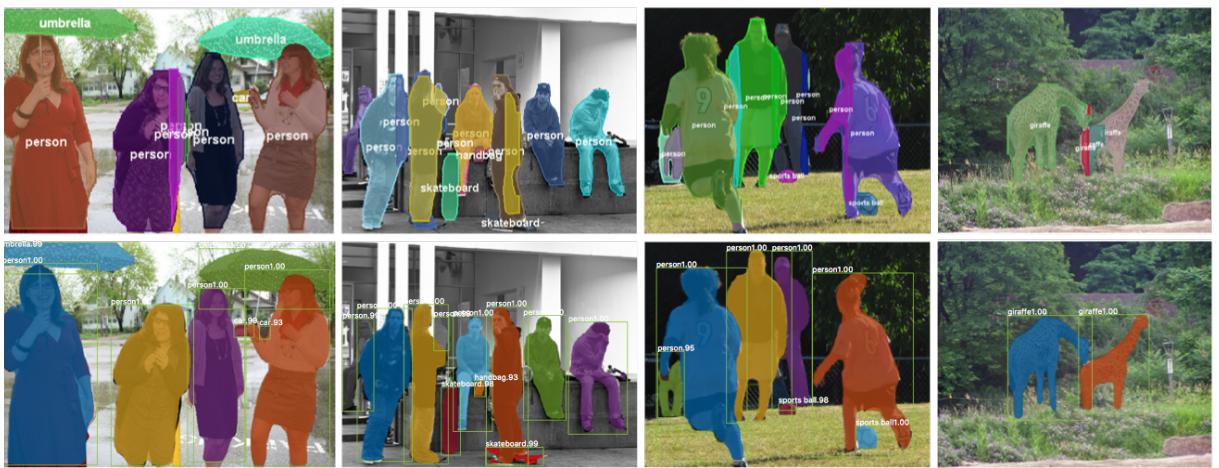


Figure 2.8: Use of CNNs have recently been shown to have accurate per-pixel object detection and classification. *Top row:* FCIS+++ [69]. *Bottom row:* Mask R-CNN [42]. Mask R-CNN outperforms FCIS+++ for overlapping instance segmentation.

2.2 Recognition Strategies

The International Conference on Document Analysis and Recognition (ICDAR) Robust Reading Competitions [55, 56, 80, 81, 99] broke down the issue of text extraction into two sub-problems: text locating and character recognition. Most of the literature discussed in Section 2.1 focused within the text locating sub-problem. For character recognition, the *Focused Scene Text* word recognition task² received three entries in 2011 and four in 2013. In 2015, the recognition challenge was redesigned³ using a new (and more challenging) *Incidental Scene Text* dataset of photos captured in-the-wild using Google Glass. The evaluation scheme uses the *Total Edit Distance* metric (described in [56]) and additionally the number of correctly recognised words for qualitative analysis. Table 2.1 summarises the top-scoring recognition rates from these competitions.

Table 2.1: Top-scoring word recognition results from ICDAR 2011–2015.

Year	Method	Dataset	Total Edit Distance	Correctly Recognised Words (%)
2011	TH-OCR System [76]	Focused	176.23	41.2
2013	PhotoOCR [8]	Focused	122.7	82.83
2015	MAPS [60]	Incidental	1128.0	32.93

It has been widely demonstrated that off-the-shelf commercial and open source OCR packages are able to correctly recognise text once the characters are extracted. In their conclusions of the ICDAR 2015 competition, Karatzas et al. note that the top performing methods will make use of commercial OCRs and conclude that conventional shape-based OCR engines can produce competitive results with pre-processed images.

This conclusion is emphasised in further works. The Open Source Tesseract OCR engine was used in Ben-ami et al. [6] to extract RBNs after preprocessed CC-based extraction. Similarly, leading commercial OCR engines used by Chen and Yuille [19] were able to achieve 93% recognition from binarised text regions after AdaBoost non-text classification, using ABBYY FineReader⁴, TOCR⁵ and Readiris Pro⁶. Similarly, Gatos et al. [34] used ABBYY FineReader and showed their extraction method showed a 50% improvement for indoor and outdoor scene images, an approach also applied by Chen and Yuille [19]. This is not to say that OCR engines

²See Challenge 2, Task 3 in [56, 99].

³See Challenge 4, Task 3 in [55].

⁴<https://www.abbyy.com> last accessed 3 July 2017.

⁵<http://www.transym.com> last accessed 3 July 2017.

⁶<http://www.irislink.com/readiris> last accessed 3 July 2017.

are always needed.

Recent interest in developing novel general character recognition strategies have also been investigated. Wang et al. [118] compared ABBYY FineReader with their novel PLEX approach and demonstrated that their text recognition system can outperform traditional OCR engines without the use of a text detector. This inspired more recent works: in 2016, Lee and Osindero [64] developed a lexicon-free photo OCR frameworks. Furthermore, application-specific recognition has been investigated using variant techniques.

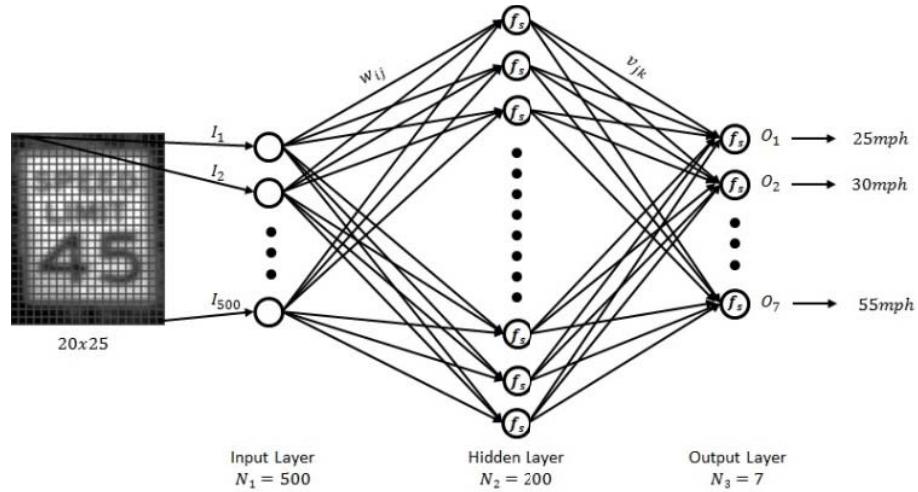


Figure 2.9: The artificial Multilayer Perceptron (MLP) NN designed in Kundu and Mackens [61] to recognise US-style speed limit signs.

A 2015 study into Traffic Sign Recognitions (TSRs) to detect US-style speed limit signs achieved recognition without the use of any OCR packages. In [61], Kundu and Mackens were able to extract a speed limit sign via the use of MSERs and template matching. The resulting detected signs were scaled to a grayscaled size of 20×25 pixels and fed into a Multilayer Perceptron (MLP) NN of 200 neurons in the hidden layer. The output layer of the network consisted of seven nodes, each representing the seven kinds of speed limit signs in US cities (25, 30, 35, 40, 45, 50, and 55 miles per hour). This architecture is shown in Figure 2.9. When trained with 13,289 images of text cases and 4,319 non-text cases, the results showed that their recognition classifier was able to correctly recognise speed limit signs with an accuracy of 98.04%. Similar results were achieved using a feed-forward MLP in [27], using UK/Poland style speed limits scaled to 20×20 pixels (grayscale) and 12 output layer neurons (10...100 kilometres per hour, the national speed limit sign, and non-sign neurons).

However, works in TSR systems that utilise networks are generally non-generalisable, and only work in a limited context (i.e., by classifying speed limit signs of known outputs). In our

context of RBN recognition, we have a known character output range of 36 possibilities: 0–9 and A–Z = 10 + 26. We do not consider lowercase letters.

Beyond TSR systems, however, we see the use of more generalisable networks: Netzer et al. [89] trained neural network to recognise street number characters from Google Street View with higher precision and recall than that of HOG and the Tesseract OCR engine, showing that the applicability of NNs for recognition can outperform traditional means. Anagnostopoulos et al. [2] used a Probabilistic Neural Network (PNN) to recognise single characters of the same 36 possibility range (i.e., uppercase alphanumeric characters) corresponding to the input grayscale vector of 9×12 pixels ($9 \times 12 = 108$ input neurons) for a single character. Figure 2.10 illustrates the PNN architecture used in this study. Furthermore, investigations in comparing different architectures of NNs for this context is given in Lee and Osindero [64].

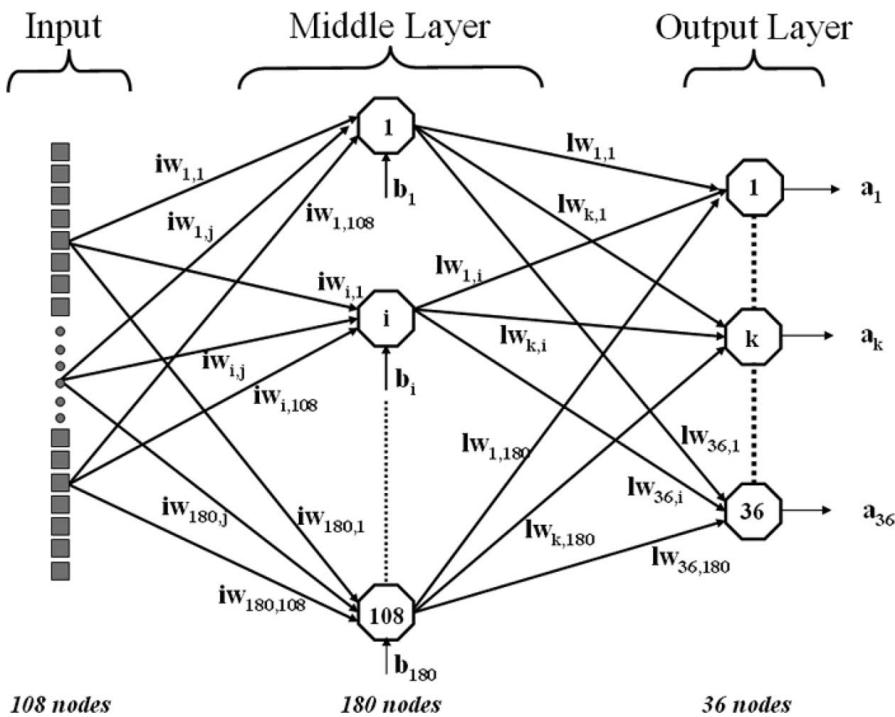


Figure 2.10: Anagnostopoulos et al. [2] developed the architecture of a PNN to determine a single character from within a license plate.

However, these pipelines are still entirely dependent on good detection strategies, and in the case where classifiers are used, quality training data must be supplied. The pipeline developed by Ben-ami et al. [6] for RBN detection (Figure 2.11) is dependent on quality facial detection: the OpenCV implementation [72] was used in this study. In order to detect the torso region—and thus detect the bib sheet itself—a heuristically-driven calculation is used to hypothesise

where the torso bounds ($T_b = T_h \times T_w$) are, given by the face height (F_h), face width (F_w):

$$T_b = (3 \times F_h) \times \left(\frac{7}{3}F_w\right)$$

The location of T_b is horizontally located at the centre of the face midpoint and vertically halfway below the face height ($0.5 \times F_w$). Hence, the dependency on these heuristics are heavily driven by a tolerant bounding box, which in itself depends on the face detection itself. Additionally, these heuristics are not always accurate—Figure 2.13 highlights a case where the face detection approach has missed the RBN entirely. As Fu et al. [33] have also noted, even if the face is detected, the bib sheet can be obfuscated (e.g., by another runner) or if the camera is capturing on the runner’s side, the OCR engine or face-detection algorithm may produce poor results. In our study, we ignore these limitations for the intention of capturing prominent runners.

Furthermore, the use of Tesseract in the study meant significant filtering, separation and character alignment is needed for the OCR engine to read characters correctly (Figure 2.12). This is yet another step that can possibly be avoided via the use of well-trained NNs, as demonstrated in previous works. Investigation into applying such networks in this context (i.e., alphanumeric sequences on *human* subjects) is largely lacking.



Figure 2.11: The RBN recognition pipeline by Ben-ami et al. [6]. *From left to right:* Input image; face detection results in red and the estimated RBN region hypothesis in green; SWT of the hypothesis region; tag region detection in blue; tag region after digit processing.



Figure 2.12: Character processing in [6]. *From left to right:* A detected tag; separation via SWT CC analysis; binarised characters; CC analysis and filtering; separation and alignment.

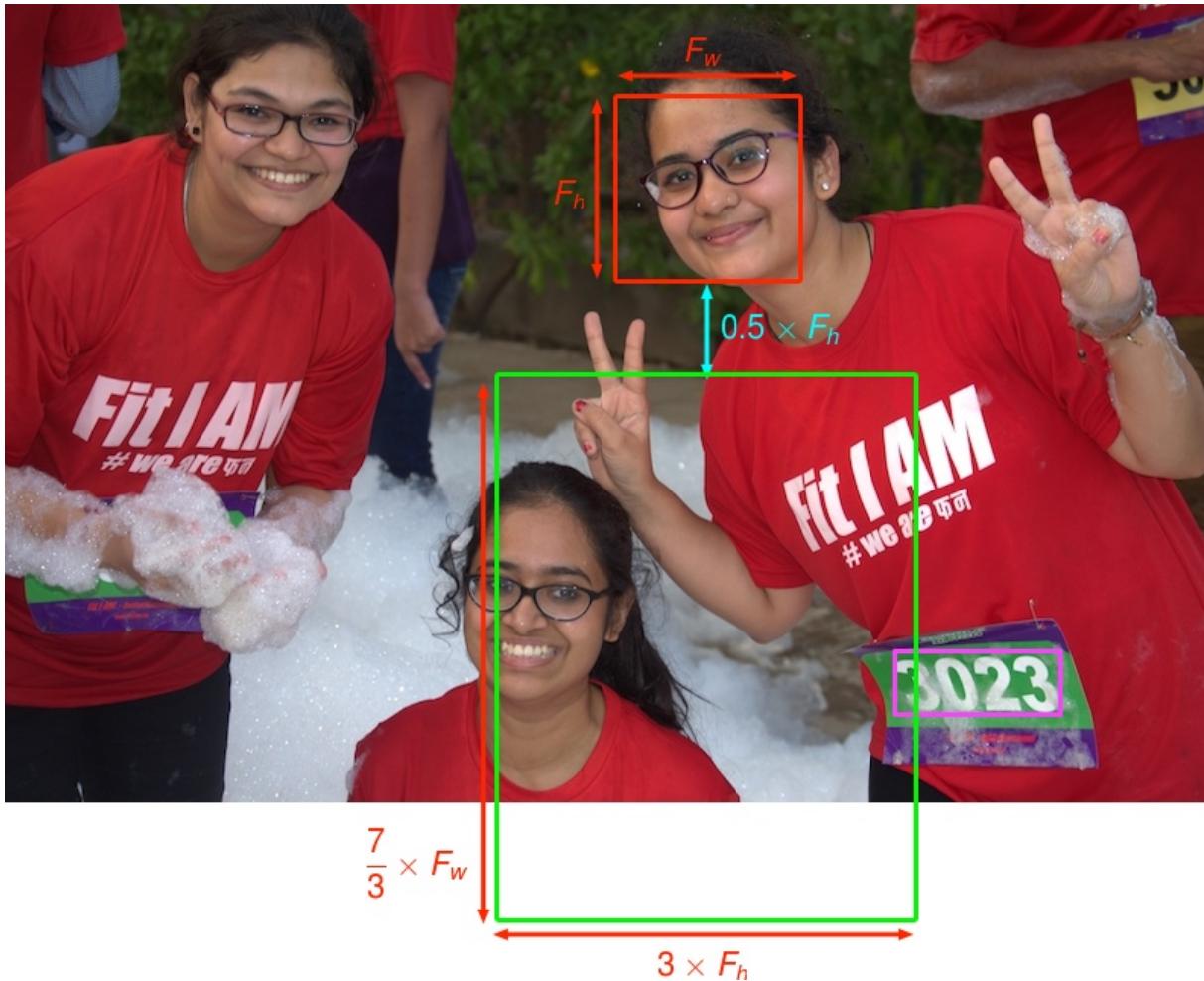


Figure 2.13: Cases where false negatives occur with the heuristic-based approach proposed in [6]. The expected RBN region (shown in magenta) is missed as the subject is leaning in the photo. Face detection in this photo uses the same OpenCV implementation [72] that was used in the study.

2.3 Metrics

Throughout our survey, we note the standard the evaluation scheme for text extraction first proposed for use in image processing in the International Conference on Document Analysis and Recognition (ICDAR) competitions [55, 56, 80, 81, 99]. This scheme was designed to be easy to understand and compute, reward text extraction useful for natural scenes, and heavily punish trivial solutions. The intention behind these metrics were to develop a measure of ‘robustness’ a text extraction pipeline can achieve.

2.3.1 Precision and Recall

Generally in information retrieval, the precision (p) and recall (r) metrics are used, first defined in the six evaluation criteria for information retrieval systems by Cleverdon et al. [22]. Precision refers to the proportion of relevant matches actually retrieved in the results, while recall refers to the proportion of relevant matches retrieved in total relevant instances. We use recall and precision metrics to assess the *effectiveness* of an information retrieval system [95].

In the context of image processing, systems that over-estimate are punished with a low precision score, while systems that under-estimate are punished with a low recall score [81]. Therefore, precision is the number of correct candidates (c) divided by the number of total estimates found (E):

$$p = \frac{c}{|E|}$$

And recall is defined as the number of correct estimates divided by the total number of ground-set truth targets (T):

$$r = \frac{c}{|T|}$$

However, it is not realistic for a given text extraction pipeline to *exactly* agree with the rectangle bounds manually tagged by a human. Lucas et al. [81] first proposed changes to these calculations to better suit their usage in the context of information extraction from within images. They adopt a more flexible notion of what a ‘match’ is. They define a new match measure (m_p) between two rectangles (i.e., the ground truth and the system’s detected candidate) as “the area of intersection of both rectangles divided by the area of the minimum bounding box containing both rectangles” [81]. This allows for a match value of one when the candidate is identical to the ground truth, and zero where the candidate has no intersection at all to the ground truth.

Therefore, the best match, $m(r, R)$, of a rectangle r in a set of rectangles R is:

$$m(r, R) = \max m_p(r, r') \mid r' \in R$$

Lastly, we can redefine the recall and precision metrics to be more forgiving in the image extraction context:

$$p' = \frac{\sum_{r_e \in E} m(r_e, T)}{|E|}$$

$$r' = \frac{\sum_{r_t \in T} m(r_t, E)}{|T|}$$

One issue with the calculation is that the metric assumes a rectangular bounding box. This limitation means that words of non-rectangular nature, such as curved or rhomboidal text, are not well-represented.

2.3.2 The f -score

A common metric used when developing text extraction pipelines is the f -score, a single measure of quality that combines both precision and recall values. We are able to compute this metric using the standard measure across many studies, as contrasted in Table 2.2.

The f -score algorithm is given in the context of image processing in Lucas et al. [81]. Relative weights controlled by an α value of 0.5 give equal weight to both precision and recall metrics:

$$f = \frac{1}{\frac{\alpha}{p'} + \frac{1-\alpha}{r'}}$$

As Chen et al. [17] report, even when all text is correctly localised, it is likely that the f -score will vary between 0.8–1.0. This is because E boundaries are unlikely to match *exactly* with the manually labelled T boundaries. An illustration of this phenomena is shown in Figure 2.14.

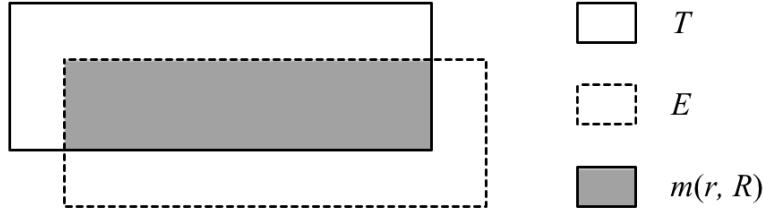


Figure 2.14: Overlapping areas of the ground truth targets, T , the estimated target boundaries E and the best match $m(r, R)$. (Adapted from [123].)

Table 2.2: A survey of text extraction literature, separated into CC- and learning-based detection methods.

Ref.	Year	Scientific Background		Precision	Recall	f -score	Dataset(s)	Performance	
		Detection	Recognition					Platform	Time (s)
[6]	2012	<i>CC-based:</i> SWT; Canny-Edges; Binary Conversion; Geometric Filtering; CC-Alignment	OCR (Tesseract)	0.65	0.62	0.63	N/A	N/S	N/S
[17]	2011	<i>CC-based:</i> Canny-Edges; MSER; SWT, Geometric, Template Matching	OCR (N/S [†])	0.73	0.60	0.66	[80, 81]	2.5 GHz	0.20
[68]	2012	<i>CC-based:</i> MSER; geometric and SWT filters; skeletal distance mapping	N/S	0.59 [81] 0.59 [99]	0.59 [81] 0.62 [99]	0.59 [81] 0.61 [99]	[80, 99]	N/S	N/S
[127]	2011	<i>CC-based:</i> character detection with HOG and SWT; link energies via spacial relationships between characters	N/S	0.73	0.62	0.67	[81]	N/S	N/S
[101]	2011	<i>CC-based:</i> Fourier-Laplacian Filtering; clustering based on maximum distance via K-Means; skeletal distance mapping	N/S	0.76 [81] 0.81 [46]	0.86 [81] 0.93 [46]	0.81 [81] 0.87 [46]	[46, 81]	2.0 GHz	7.80
[28]	2010	<i>CC-based:</i> Canny-Edges; SWT; Modified CC algorithm; Geometric Filtering	OCR (N/S)	0.73	0.60	0.66	[80, 81]	N/S	0.94
[126]	2010	<i>CC-based:</i> Canny-Edges; HOG; geometric filtering; graph spectrum	N/S	0.67	0.46	0.55	[81]	N/S	N/S
[123]	2005	<i>Learning-based:</i> SVM to reject non-text; wavelet movement; HOG; OCR filtering	N/S	N/S	0.97	N/S	[46]	1.6 GHz	8.30
[92]	2010	<i>Learning-based:</i> Waldboost Classifier; HOG; LBP; Gabor Wavelets	N/S	0.56	0.70	0.68	[81]	3.4 GHz	0.37
[38]	2004	<i>Learning-based:</i> Wavelet transformation; feature estimation; pixel-block classification	N/S	0.87	0.90	0.88	[45]	N/S	N/S
[40]	2009	<i>Learning-based:</i> Mean Difference, Standard Deviation and HOG features; AdaBoost and CAdaBoost classifiers; NN-based localiser	N/S	0.25	0.35	0.35	[81]	2.2 GHz	N/S

[†]Not Specified

Conclusion

In this chapter, we have presented a survey of literature in various application contexts: RBN and TSR recognition, recognition of alphanumeric sequences ‘in the wild’, and additionally object instance segmentation. We also present the varied range of techniques used to both detect and recognise the text, using both heuristic-based and NN-based approaches. This said, we acknowledge that datasets within the survey are not always recent (ranging as far back as 2003) due to tendencies for researchers to use the more popular (though aged) ICDAR datasets.

The state of the art of learning-based detection approaches such as CNNs for have gained wide popularity, albeit for object segmentation. These approaches are yet to be applied within the context of alphanumeric sequences. Recent years have had a heavier focus on heuristic-based detection strategies using CC-based methods, while a majority of learning-based detection methods have had far fewer recent investigations. Furthermore, recognition of characters using NNs are not yet widely used, and off-the-shelf OCR packages are still standard.

Chapter 3

Dataset

How do we capture large and annotated datasets to train an Artificial Intelligence (AI) model at scale? What is an annotation, and how can we constrain them? These are important questions when supplying data to train graphics-based AIs. In this chapter, we discuss a theoretical data tagging metamodel that enables us to concisely describe the architecture of how quality labelled data is captured for these purposes. Additionally, we discuss a partial implementation of this metamodel which we have used to gather our training dataset.

3.1 A Data-Capturing Architecture

Why should we care about the architecture of our dataset capturing, so long as the AI is well-trained? This question fundamentally leads to the *provenance* or *lineage* of our data: not all training data is initially perfect. At times, we will need to modify our training data, or in the cases of some AI systems, training data is re-fed back into the system from the AI itself.

Hence, it is imperative to understand how the transition and flow of data occurs [9, 24, 48] in these AI systems. Where does it comes from? How is it derived and updated, and how might this affect our trained AI models over time? A lack of understanding in training data provenance for large-scale AIs implies difficulty in sourcing the cause of errors due to these mutations. It may be a difference in the data training format, or a mismatch of values on the same training data. Overcoming this without systematic record-keeping is challenging.

We therefore propose a grammar to describe layered-approach of a data-capturing architecture. As such a grammar is textual, it can be version controlled: therefore data provenance is recorded, which allows us to source how the mutations in training data can occur and what effect this has our AI systems. From an exploratory analysis of developing this architecture (Section 3.1.1), we propose a theoretical and generalisable metamodel that can be used to

capture any form image data for training purposes (e.g., frames of a video, images of natural scenes). This metamodel is largely inspired by the works of Wickham [121, 122] and Moody [87].

3.1.1 Methodology

A partial implementation of this metamodel, which we call *Argus*^{1,2}, was developed for the purposes of capturing our dataset. Refer to Section ?? for more about this implementation.

3.1.2 What to Capture?

What features do we want to capture from our dataset? To answer this question, we need to consider what information we see as relevant in a standard marathon photo (such as Figure 1.1). There are two annotations to describe features of the entire photo:

1. whether or not the photo is considered as **crowded** (TRUE or FALSE), and
2. an *optional* collection of **runners** in the photo, given that the photo is not crowded.

We train a NN based on photos that are not crowded; photos that are considered crowded are typically not desirable as runners prefer photos where they are the key subject. We therefore discard these photos to remove any bias on the network. We can only tag runners on the condition that the photo is not marked as crowded. Additionally, in some photos, all RBNs are missing within the photo (e.g., hidden behind other runners, cropped out of view) and some photos contain no runners at all. We therefore identify this as an *optional* collection as we must associate an RBN to a runner—the photo is not crowded, but there are no runners to tag, thereby making the annotation optional. Thus, an *implicit attribute* exists with such a collection: the *count* of the runners in a photo is something we can automatically count.

We identify the following two annotations of what we would label for a given runner’s bib sheet, a feature important to identify RBNs:

1. a polygon around the runner’s **bib sheet** that contains the *x* and *y* coordinates, given that there are exactly four vertices of this polygon
2. a string with the runner’s **RBN**, once the bib sheet is tagged (exists).

¹<http://www.deakin.edu.au/~ca/argus> last accessed 5 July 2017.

²Inspired by the *all-seeing* giant of the same name from Greek mythology: “With his multiple sets of eyes, Argus could see nearly everything in his vicinity”. See <http://www.loggia.com/myth/argus.html>.

A further feature that is important is the runner's face region, which could compromise of five annotations:

1. a rectangle around the runner's **face** that contains the two x and y coordinates of the two opposite vertices, given that the *bottom* of this rectangle are above the *top* of the bib sheet,

and once this face bounds has been tagged, more annotations can be extracted:

2. whether or not the runner is **wearing a hat** (TRUE or FALSE),
3. whether or not the runner is **wearing (sun)glasses** (TRUE or FALSE), and
4. the **gender** of the runner (MALE, FEMALE, or UNSURE).

The bib sheet polygon would contain four *explicit attributes* required as input by the tagger: the coordinates of its four vertices, $\{(x_1, y_1), (x_2, y_2), (x_3, y_3), (x_4, y_4)\}$. Similarly, there are two explicit attributes required for the runner's face rectangle: the coordinates of the top left and the bottom right of the rectangle (i.e., the two opposite vertices) of the runner's face bounds, $\{(x_1, y_1), (x_2, y_2)\}$. As both of these are shapes, we can automatically calculate the implicit attributes from these coordinates: $\{top, left, bottom, right, width, height\}$.

We can now identify another feature, the prominence of the runner within the photo. These would consist of three annotations:

1. whether or not the runner's **face** is **visible** (TRUE or FALSE),
2. whether or not the runner is **blurry** (TRUE or FALSE),
3. the **Likelihood of Purchase (LoP)** that the runner buys the photo (NO, MAYBE, YES),

A runner's face is considered 'invisible' if it is obstructed by another object, cropped out of the image, or if the runner is looking down. Runners are more likely purchase a photo if they are looking at the camera, and likewise if they are blurry in the image; these therefore have a significant impact on their prominence. We also define the Likelihood of Purchase (LoP) value as a qualitative metric. We ask the data tagger to 'picture' themselves as runner, and if so, we ask if they would purchase it. We then use this metric to train positive samples of good photos (LoP = YES) and samples of bad photos (LoP = NO). Where MAYBE is provided, the runner is ignored to prevent training the network with indeterminate samples.

Another feature we can capture is a way to identify runners by the colour of their clothing. (Within our dataset, various clothing items of distinct colours are given to runners for particular races.) These features comprise of four annotations:

1. an *optional* **colour** of the runner's **hat**, given that they were annotated as wearing a hat,
2. the **colour** of the runner's **shirt**,
3. an *optional* **colour** of the runner's **shorts**, and
4. an *optional* **colour** of the runner's **shoes**.

The colour of a runner's shirt is required, as we expect that a bib would be detected on the shirt of a runner (which would be tagged). The other colour annotations are all optional, as it is likely that some of these clothing items may be cropped out of the photo or is not visible. Furthermore, the hat colour can only be specified on the condition that the tagger has marked the runner has wearing a hat.

We can segment groups of the annotations we extract into features of two categories: annotations that feature at the *image*-level and those that do not, which we call *segment*-level features. The image-level features are those which apply to the entire image (i.e., *PhotoCrowded*, *Runners*). Those that apply at the segment level can be grouped into what types of features we are extracting: the runner's bib, face, their prominence and their colour identification.

In summary, we have identified a total of 5 features of 15 annotations, which are summarised within Table 3.1. Each of these annotations can be classified with a name and type, conditions for the annotation to be valid, dependencies on other annotations to exist before the annotation can be made, explicit attributes which the tagger must specify, implicit attributes which we can automatically compute, possible values which limit the range of data for that annotation, and whether or not the annotation is optional.

Table 3.1: A summary of the annotations we wish to capture from our dataset. Image and segment-level features are separated using the double line.

Feature	Name	Type	Conditions	Dependencies	Explicit Attributes	Implicit Attributes	Possible Values	Optional
Image-Level	<i>PhotoCrowded</i>	Boolean	None	None	None	None	{TRUE, FALSE}	No
	<i>Runners</i>	Collection	{ <i>PhotoCrowded</i> = FALSE}	None	None	{count}	N/A	Yes
Bib	<i>BibSheet</i>	Polygon	{vertices = 4}	None	{x ₁ ...x ₄ , y ₁ ...y ₄ }	{top, left, bottom, right, width, height}	N/A	No
	<i>RBN</i>	String	None	<i>BibSheet</i>	None	None	N/A	No
Face	<i>FaceBounds</i>	Rectangle	{bottom > BibSheet _{top} }	<i>BibSheet</i>	{x ₁ , y ₁ , x ₂ , y ₂ }	{top, left, bottom, right, width, height}	N/A	No
	<i>FaceVisible</i>	Boolean	None	<i>FaceBounds</i>	None	None	{TRUE, FALSE}	No
	<i>WearingHat</i>	Boolean	None	<i>FaceBounds</i>	None	None	{TRUE, FALSE}	No
	<i>WearingGlasses</i>	Boolean	None	<i>FaceBounds</i>	None	None	{TRUE, FALSE}	No
	<i>Gender</i>	Category	None	<i>FaceBounds</i>	None	None	{MALE, FEMALE}	No
Prominence	<i>LoP</i>	Category	None	<i>FaceBounds</i>	None	None	{NO, MAYBE, YES}	No
	<i>Blurry</i>	Boolean	None	<i>FaceBounds</i>	None	None	{TRUE, FALSE}	No
Colours	<i>ShirtColor</i>	Color	None	<i>FaceBounds</i>	{red, green, blue}	None	N/A	No
	<i>ShoeColor</i>	Color	None	<i>FaceBounds</i>	{red, green, blue}	None	N/A	Yes
	<i>ShortsColor</i>	Color	None	<i>FaceBounds</i>	{red, green, blue}	None	N/A	Yes
	<i>HatColor</i>	Color	{ <i>WearingHat</i> = TRUE}	<i>FaceBounds</i>	{red, green, blue}	None	N/A	Yes

3.1.3 Describing the Metamodel

In the example above, we have seen some components that make up an annotation. We can now develop our metamodel using the terminology described: types, conditions and dependencies, explicit and implicit attributes, possible values, and optionality.

Types We have so far identified *Boolean*, which is just a restricted *Category* consisting of the possible values {TRUE, FALSE}. We have *Polygons* and *Rectangles*, and therefore we must conceptualise some form of *Shape*. We also have *Colours*, which can be represented as an RGB represented hexadecimal *String*. This only leaves our *Collection* of runners: a data type which encompasses multiple annotations. We therefore a data schema that is generalisable from *Strings*, *Shapes*, *Collections* and *Categories*.

Conditions A condition is simply a construction rule which disallows a data tagger to create an annotation until a condition is met by the annotator. For instance, we can never have a *BibSheet* above a *FaceBounds*, so when we can specify this as a condition. These are not bound to just geometric constraints; we also see that we cannot tag *Runners* if the *PhotoCrowded* annotation is marked as TRUE.

Dependencies Some annotations are dependent on others existing. These are construction rules which we label as dependencies. We cannot annotate an RBN if there is no *BibSheet* annotated. Likewise, a *FaceBounds* is dependent on where the *BibSheet* is, and so the *BibSheet* must be marked up first. These dependencies add order in which features are being extracted, which is important to the tagger marking up the photo.

Possible Values In the instance where a data type must be restricted to a set of values (i.e., *Categorical* data), we set what those restrictions are as a construction rules of our annotations.

Optionality Not all annotations are needed—in these cases the annotation is considered optional, but by default is required.

Attributes Attributes exist to as a means capture further information about an annotation. We represent attributes in our metamodel as either implicit or explicit. As stated in Section 3.1.2, those attributes which are required for the tagger to manually markup are explicit (i.e., manually

must be added) whereas those that can be computed by the system are implicit (i.e., inferred from explicit attributes).

These begin to form layers of data within the images we annotate. This layer can then be represented in a hierarchical data format (e.g., a tree) to describe the layering of a fully annotated photo.

We therefore highlight some components which allow us to describe this architecture in a more generalisable way.

Runners: [Collection]

1. Annotation, and further, an optional annotation
2. A condition

3.1.4 The Data Capturing Process

Now that we have defined which annotations exist, we describe the process by which we are able to capture this data. At its core, this is an Extract-Transform-Load (ETL) process: we extract the data from our

3.1.5 Layering

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] Method of and apparatus for pattern recognition.
- [2] 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*.
- [3] 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*.
- [4] Bay, H., A. Ess, T. Tuytelaars, and L. J. Van Gool (2008). Speeded-Up Robust Features (SURF). *Computer Vision and Image Understanding*.
- [5] Belongie, S. J., J. Malik, and J. Puzicha (2001). Matching Shapes. *ICCV*.
- [6] Ben-ami, I., T. Basha, and S. Avidan (2012). Racing Bib Numbers Recognition. In *British Machine Vision Conference 2012*.
- [7] Bengio, Y., P. Lamblin, D. Popovici, and H. Larochelle (2006). Greedy Layer-Wise Training of Deep Networks. *NIPS*.
- [8] Bissacco, A., M. Cummins, Y. Netzer, and H. Neven (2013). PhotoOCR - Reading Text in Uncontrolled Conditions. *ICCV*.
- [9] Buneman, P., S. Khanna, and W.-C. Tan (2000). Data Provenance: Some Basic Issues. In *Computer Vision – ACCV 2010*.
- [10] Burges, C. J. C. (1998). A Tutorial on Support Vector Machines for Pattern Recognition. *Data Min. Knowl. Discov.*
- [11] Canny, J. F. (1986). A Computational Approach to Edge Detection. *IEEE Trans. Pattern Anal. Mach. Intell..*

- [12] Cano-Perez, J. and J. C. Pérez-Cortes (2003). Vehicle License Plate Segmentation in Natural Images. *IbPRIA*.
- [13] Chen, D. and J. Luettin (2000). A survey of text detection and recognition in images and videos.
- [14] Chen, D., J.-M. Odobez, and H. Bourlard (2004a). Text detection, recognition in images and video frames. *Pattern Recognition*.
- [15] Chen, D., J.-M. Odobez, and H. Bourlard (2004b). Text detection, recognition in images and video frames. *Pattern Recognition*.
- [16] 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..*
- [17] 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*.
- [18] Chen, T., M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang (2015). MXNet: A Flexible and Efficient Machine Learning Library for Heterogeneous Distributed Systems. *arXiv.org*.
- [19] Chen, X. and A. L. Yuille (2004a). 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*.
- [20] Chen, X. and A. L. Yuille (2004b). Detecting and Reading Text in Natural Scenes. *CVPR*.
- [21] Chen, X. and A. L. Yuille (2005). A Time-Efficient Cascade for Real-Time Object Detection - With applications for the visually impaired. *CVPR Workshops*.
- [22] Cleverdon, C., J. Mills, and M. Keen (1966). Factors Determining the Performance of Indexing Systems. Technical report, ASLIB Cranfield Research Project.
- [23] Cortes, C. and V. Vapnik (1995). Support-Vector Networks. *Machine Learning*.
- [24] Cui, Y. and J. Widom (2003). Lineage tracing for general data warehouse transformations. *The VLDB Journal The International Journal on Very Large Data Bases*.

- [25] Dalal, N. and B. Triggs (2005). Histograms of Oriented Gradients for Human Detection. *CVPR*.
- [26] Deng, J., W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei (2009). ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [27] 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)*.
- [28] Epshtain, B., E. Ofek, and Y. Wexler (2010). Detecting text in natural scenes with stroke width transform. *CVPR*.
- [29] 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*.
- [30] Freeman, W. T. and M. Roth (1995). Orientation histograms for hand gesture recognition. *International workshop on automatic*
- [31] Freund, Y. and R. E. Schapire (1996). Experiments with a New Boosting Algorithm. *ICML*.
- [32] Friedman, J., R. Tibshirani, and T. Hastie (2000). Additive logistic regression: a statistical view of boosting (With discussion and a rejoinder by the authors). *The annals of statistics*.
- [33] 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)*.
- [34] Gatos, B., I. Pratikakis, and K. Kepene (2005). Text detection in indoor/outdoor scene images. *Proc First Workshop of*
- [35] 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*.
- [36] Girshick, R., J. Donahue, T. Darrell, and J. Malik (2014). Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

- [37] Girshick, R. B. (2015). Fast R-CNN. *ICCV*.
- [38] Gllavata, J., R. Ewerth, and B. Freisleben (2004). Text Detection in Images Based on Unsupervised Classification of High-Frequency Wavelet Coefficients. *ICPR*.
- [39] Gonzalez, Á., L. M. Bergasa, J. J. Y. Torres, and S. Bronte (2012). Text location in complex images. *ICPR*.
- [40] Hanif, S. M. and L. Prevost (2009). Text Detection and Localization in Complex Scene Images using Constrained AdaBoost Algorithm. *ICDAR*.
- [41] Hanif, S. M., L. Prevost, and P. Negri (2008). A cascade detector for text detection in natural scene images. *ICPR*.
- [42] He, K., G. Gkioxari, P. Dollár, and R. B. Girshick (2017). Mask R-CNN. *CoRR*.
- [43] Heisele, B., T. Serre, S. Mukherjee, and T. A. Poggio (2001). Feature Reduction and Hierarchy of Classifiers for Fast Object Detection in Video Images. *CVPR*.
- [44] Horn, B. (1986). *Robot Vision*. MIT Press.
- [45] Hua, X.-S., L. Wenyin, and H. Zhang (2001). Automatic Performance Evaluation for Video Text Detection. *ICDAR*.
- [46] 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.*.
- [47] Huang, X., T. Shen, R. Wang, and C. Gao (2015). Text detection and recognition in natural scene images. In *2015 International Conference on Estimation, Detection and Information Fusion (ICEDIF)*.
- [48] Ikeda, R. and J. Widom (2009). Data lineage: A survey.
- [49] Jain, A. K. and B. Y. 0002 (1998). Automatic text location in images and video frames. *ICPR*.
- [50] 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 (TENSYMP)*.

- [51] Jin, J., K. Fu, and C. Zhang (2014). Traffic Sign Recognition With Hinge Loss Trained Convolutional Neural Networks. *IEEE Transactions on Intelligent Transportation Systems*.
- [52] Jung, C., Q. Liu, and J. Kim (2009). A stroke filter and its application to text localization. *Pattern Recognition Letters*.
- [53] Jung, K., K. In Kim, and A. K Jain (2004). Text information extraction in images and video: a survey. *Pattern Recognition*.
- [54] Jung, K., K. I. Kim, and A. K. Jain (2004). Text information extraction in images and video - a survey. *Pattern Recognition*.
- [55] Karatzas, D., L. Gomez-Bigorda, A. Nicolaou, S. K. Ghosh, A. D. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu, F. Shafait, S. Uchida, and E. Valveny (2015). ICDAR 2015 competition on Robust Reading. *ICDAR*.
- [56] Karatzas, D., F. Shafait, S. Uchida, M. Iwamura, L. G. i. Bigorda, S. R. Mestre, J. Mas, D. F. Mota, J. A. Almazan, and L. P. de las Heras (2013). ICDAR 2013 Robust Reading Competition. In *2013 12th International Conference on Document Analysis and Recognition (ICDAR)*.
- [57] Kim, H.-K. (1996). Efficient Automatic Text Location Method and Content-Based Indexing and Structuring of Video Database. *J. Visual Communication and Image Representation*.
- [58] Kim, K. I., K. Jung, and J. H. Kim (2003). Texture-Based Approach for Text Detection in Images Using Support Vector Machines and Continuously Adaptive Mean Shift Algorithm. *IEEE Trans. Pattern Anal. Mach. Intell.*.
- [59] Krizhevsky, A., I. Sutskever, and G. E. Hinton (2012). ImageNet Classification with Deep Convolutional Neural Networks. *NIPS*.
- [60] Kumar, D., M. N. A. Prasad, and A. G. Ramakrishnan (2012). *MAPS: midline analysis and propagation of segmentation*. ACM.
- [61] Kundu, S. K. and P. Mackens (2015). Speed Limit Sign Recognition Using MSER and Artificial Neural Networks. *ITSC*.
- [62] LeCun, Y., L. Bottou, Y. Bengio, and P. Haffner (1998). Gradient-based learning applied to document recognition. *Proceedings of the IEEE*.

- [63] Lee, C. W., K. Jung, and H. J. Kim (2003). Automatic text detection and removal in video sequences. *Pattern Recognition Letters*.
- [64] Lee, C.-Y. and S. Osindero (2016). Recursive Recurrent Nets with Attention Modeling for OCR in the Wild.
- [65] Lee, E. R., P. K. Kim, and H. J. Kim (1994). Automatic Recognition of a Car License Plate using Color Image Processing. *ICIP*.
- [66] Lee, S., M. S. Cho, K. Jung, and J. H. Kim (2010). Scene Text Extraction with Edge Constraint and Text Collinearity. *ICPR*.
- [67] Li, H., D. S. Doermann, and O. E. Kia (2000). Automatic text detection and tracking in digital video. *IEEE Trans. Image Processing*.
- [68] Li, Y. and H. Lu (2012). Scene text detection via stroke width. *ICPR*.
- [69] Li, Y., H. Qi, J. Dai, X. Ji, and Y. Wei (2016). Fully convolutional instance-aware semantic segmentation. *arXiv.org*.
- [70] 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)*.
- [71] Liang, J., D. S. Doermann, and H. Li (2005). Camera-based analysis of text and documents - a survey. *IJDAR*.
- [72] Lienhart, R. and J. Maydt (2002). An extended set of Haar-like features for rapid object detection. *ICIP*.
- [73] Lienhart, R. and A. Wernicke (2002). Localizing and segmenting text in images and videos. *IEEE Trans. Circuits Syst. Video Techn.*.
- [74] Lin, T.-Y., M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick (2014). Microsoft COCO - Common Objects in Context. *ECCV*.
- [75] Liu, C. L., M. Koga, and H. Fujisawa (2005). Gabor feature extraction for character recognition: comparison with gradient feature. In *Eighth International Conference on Document Analysis and Recognition (ICDAR'05)*.

- [76] Liu, H. and X. Ding (2005). Handwritten Character Recognition Using Gradient Feature and Quadratic Classifier with Multiple Discrimination Schemes. *ICDAR*.
- [77] Liu, Y., S. Goto, and T. Ikenaga (2006). A Contour-Based Robust Algorithm for Text Detection in Color Images. *IEICE Transactions*.
- [78] Liu, Z. and S. Sarkar (2008). Robust outdoor text detection using text intensity and shape features. *ICPR*.
- [79] Lowe, D. G. (2004). Distinctive Image Features from Scale-Invariant Keypoints. *International Journal of Computer Vision*.
- [80] Lucas, S. M. (2005). ICDAR 2005 text locating competition results. In *Eighth International Conference on Document Analysis and Recognition (ICDAR'05)*.
- [81] 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*.
- [82] Mairal, J., F. R. Bach, J. Ponce, G. Sapiro, and A. Zisserman (2008). Discriminative learned dictionaries for local image analysis. *CVPR*.
- [83] Matas, J., O. Chum, M. Urban, and T. Pajdla (2002). Robust Wide Baseline Stereo from Maximally Stable Extremal Regions. *BMVC*.
- [84] Mikolajczyk, K., T. Tuytelaars, C. Schmid, A. Zisserman, J. Matas, F. Schaffalitzky, T. Kadir, and L. J. Van Gool (2005). A Comparison of Affine Region Detectors. *International Journal of Computer Vision*.
- [85] Minetto, R., N. Thome, M. Cord, J. Fabrizio, and B. Marcotegui (2010). SnooperText - A multiresolution system for text detection in complex visual scenes. *ICIP*.
- [86] Mohan, A., C. Papageorgiou, and T. A. Poggio (2001). Example-Based Object Detection in Images by Components. *IEEE Trans. Pattern Anal. Mach. Intell.*.
- [87] Moody, D. L. (2009). The “Physics” of Notations - Toward a Scientific Basis for Constructing Visual Notations in Software Engineering. *IEEE Trans. Software Eng.*.
- [88] Mutch, J. and D. G. Lowe (2006). Multiclass Object Recognition with Sparse, Localized Features. *CVPR*.

- [89] Netzer, Y., T. Wang, and A. Coates (2011). Reading digits in natural images with unsupervised feature learning. *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*.
- [90] Nistér, D. and H. Stewénius (2008). Linear Time Maximally Stable Extremal Regions. *ECCV*.
- [91] Ojala, T., M. Pietikainen, and D. Harwood (1994). Performance evaluation of texture measures with classification based on Kullback discrimination of distributions. In *12th International Conference on Pattern Recognition*.
- [92] Pan, Y.-F., C.-L. Liu, and X. Hou (2010). Fast scene text localization by learning-based filtering and verification. In *2010 17th IEEE International Conference on Image Processing (ICIP 2010)*.
- [93] Phan, T. Q., P. Shivakumara, and C. L. Tan (2009). A Laplacian Method for Video Text Detection. In *2009 10th International Conference on Document Analysis and Recognition*.
- [94] Ren, S., K. He, R. B. Girshick, and J. S. 0001 (2017). Faster R-CNN - Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Trans. Pattern Anal. Mach. Intell.*.
- [95] Rijsbergen, C. J. V. (1979). *Information Retrieval*. Butterworth-Heinemann.
- [96] Sarkar, S. and K. L. Boyer (1996). Quantitative Measures of Change based on Feature Organization - Eigenvalues and Eigenvectors. *CVPR*.
- [97] Seo, Y.-W., J. Lee, W. Zhang, and D. Wettergreen (2015). Recognition of Highway Work-zones for Reliable Autonomous Driving. *IEEE Transactions on Intelligent Transportation Systems*.
- [98] Sermanet, P. and Y. LeCun (2011). Traffic sign recognition with multi-scale Convolutional Networks. *IJCNN*.
- [99] 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)*.
- [100] 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*.

- [101] 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.*
- [102] Sivic, J. and A. Zisserman (2003). Video Google - A Text Retrieval Approach to Object Matching in Videos. *ICCV.*
- [103] 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..*
- [104] Smith, R. (2007). An Overview of the Tesseract OCR Engine. In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007) Vol 2.*
- [105] Smith, R. W. (1987). *The Extraction and Recognition of Text from Multimedia Document Images.* Ph. D. thesis.
- [106] Sobottka, K., H. Bunke, and H. Kronenberg (1999). Identification of Text on Colored Book and Journal Covers. *ICDAR.*
- [107] Sochman, J. and J. Matas (2005). WaldBoost ? Learning for Time Constrained Sequential Detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05).*
- [108] Srivastav, A. and J. Kumar (2008). Text detection in scene images using stroke width and nearest-neighbor constraints. In *TENCON 2008 - 2008 IEEE Region 10 Conference (TENCON).*
- [109] Subramanian, K., P. Natarajan, M. Decerbo, and D. A. Casta n  (2007). Character-Stroke Detection for Text-Localization and Extraction. *ICDAR.*
- [110] Sun, Q., Y. Lu, and S. Sun (2010). A Visual Attention Based Approach to Text Extraction. *ICPR.*
- [111] Sung, K. K. and T. A. Poggio (1998). Example-Based Learning for View-Based Human Face Detection. *IEEE Trans. Pattern Anal. Mach. Intell..*
- [112] Takacs, G., V. Chandrasekhar, N. Gelfand, Y. Xiong, W.-C. Chen, T. Bismigiani, R. Grzeszczuk, K. Pulli, and B. Girod (2008). Outdoors augmented reality on mobile phone using loxel-based visual feature organization. *Multimedia Information Retrieval.*

- [113] 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*.
- [114] 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*.
- [115] Tu, Z., X. Chen, A. L. Yuille, and S. C. Zhu (2003). Image Parsing - Unifying Segmentation, Detection, and Recognition. *ICCV*.
- [116] Vapnik, V. (1999). *The Nature of Statistical Learning Theory*. Springer Science & Business Media.
- [117] Viola, P. A., M. J. Jones, and D. Snow (2003). Detecting Pedestrians Using Patterns of Motion and Appearance. *ICCV*.
- [118] Wang, K., B. Babenko, and S. J. Belongie (2011). End-to-end scene text recognition. *ICCV*.
- [119] Wang, X., T. X. Han, and S. Yan (2009). An HOG-LBP human detector with partial occlusion handling. *ICCV*.
- [120] Wang, X., L. Huang, and C. Liu (2009). A New Block Partitioned Text Feature for Text Verification. In *2009 10th International Conference on Document Analysis and Recognition*.
- [121] Wickham, H. (2007). Reshaping data with the reshape package. *Journal of Statistical Software*.
- [122] Wickham, H. (2010). A Layered Grammar of Graphics. *Journal of Computational and Graphical Statistics*.
- [123] Ye, Q., Q. Huang, W. G. 0001, and D. Zhao (2005). Fast and robust text detection in images and video frames. *Image Vision Comput..*
- [124] Zagoruyko, S. and N. Komodakis (2015). Learning to compare image patches via convolutional neural networks. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [125] 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)*.

- [126] Zhang, J. and R. Kasturi (2010). Text Detection Using Edge Gradient and Graph Spectrum. *ICPR*.
- [127] Zhang, J. and R. Kasturi (2011). Character Energy and Link Energy-Based Text Extraction in Scene Images. In *Computer Vision – ACCV 2010*.
- [128] 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