

POLITECNICO DI TORINO

Department of Electronics and Telecommunications
Master of Science in ICT for Smart Societies

Master Degree Thesis

Optical Character Recognition



Supervisor:

Prof. Monica Visintin

Author:

Ani DEVER

April 2018

DEDICATION

I dedicate this dissertation work to my inspiring parents, for their endless support and encouragement during the challenges of my life.

ACKNOWLEDGEMENTS

I wish to present my special thanks to Autour team of Shared Reality Lab, McGill University. They have been a great team, whose help and motivation made my research successful.

I would like to express my gratefulness to Mathieu Bouchard, the lead developer of Autour, who taught me a lot about computer science and mobile application development.

I would like to pay my regards to Professor Monica Visintin, my advisor at Politecnico di Torino, who inspired my interest to machine learning and deep learning. Her encouragement made it possible to achieve my goals.

I would like to show my warm thank to Professor Jeremy Cooperstock, my advisor at McGill University, who gave me the chance being his intern and explore the research activities at his lab. I appreciate his trust in me, his moral and financial support.



ABSTRACT

Autour[5]¹ is an eyes-free, cartographical mobile application being developed by Shared Reality Laboratory of McGill University as a research project. Autour presently runs on iPhones starting from iPhone 4S and the corresponding Android version is being developed. The application is designed to give visually impaired individuals a better sense of their surroundings. It can be used hands-free by leaving the phone in a bag around the user's neck or handheld. It utilizes inbuilt capabilities and sensors of the smart phone such as accelerometer, gyroscope, compass and Global Positioning System (GPS), to determine the user's location and orientation. Then the relevant sound instructions are presented to the user either through a bone-conducting or open air headphones so as not to interfere with the ambient sounds. The application's aim is to use ambient audio to expose the sort of information that visual cues, e.g signs and markings, give to individuals with sight.

Autour has several modes: the tilt of the device allows users to choose between two of the modes: horizontal mode and vertical mode. By default, the vertical mode is Radar mode and the horizontal mode is Beam mode. In total, user has following modes: Shockwave, Browse, Tutorial, Menu and a mode that just waits for GPS lock. In any horizontal and vertical mode, user can tap the screen twice quickly to hear the address where the user is located, the status of the sensors and a summary of places around the user. While that is spoken, user may interrupt by tapping the screen once or tilt the device.

There are two different Sweep modes: Radar and Shockwave. Sweep

¹ <http://autour.mcgill.ca/en/>

modes are used when the device is being used horizontally. The selection of the sweep modes can be done in the settings. In Radar mode, a tap will start an automatic sweep when the device is pointed upwards with the screen facing towards the user. The user will hear a ticking sound indicating the progress of the scan, along with the names of the places around the user sorted by distance. In Shockwave mode sorting is done by direction, therefore the sweep is done as a circle starting from where the user is and grows until the specified maximum distance is reached.

Beam mode can be activated and deactivated in the same manner as in Sweep modes. When the Beam mode is activated the user will hear a tick, then places that are in that sector are enumerated in order of increasing distance.

Browse mode creates a list of places nearby and reads it to the user. The user can navigate through the list by swiping. Standby mode starts automatically when GPS does not work well. Autour stays in Standby mode until GPS data is available again.

Autour is a mobile application depending on external services such as GPS, FoursquareTM, Google PlacesTM or OpenStreetMap which can be unreliable. The application may fail if certain conditions are not satisfied e.g. if the user is traveling in a vehicle or if one the external services break unexpectedly. Autour does not recommends its user to rely on it for navigation and safety, it is not designed for mobilized user. It is not a replacement for a cane or a guide dog.

Shared Reality Lab is hosting the Autour team composed of a lead developer and students with different levels and backgrounds, a team of individuals with different interests and skills. Current objective is to add new features

to Autour (such as scene description, chat boxes and optical character recognition) in order to provide a broad visual feedback to its users. Ideally, the users will be able to know what kind of scene they are in, what is written in their point of view and will be able to ask questions about what is happening around.

The visual task of classifying an object once was a big challenge. Yet, recent advancements in artificial intelligence have sparked major progression in computer vision. Now machines are drastically successful at computer vision problems such as pattern recognition, object detection and classification. Currently, the most accomplished methods in the literature are deep learning based. Today it is possible build strong computer vision systems by the help of deep learning algorithms and frameworks.

This thesis documents the effort given to assemble an optical character recognition (scene text detection and recognition to be precise) system by gathering state-of-the-art methods. The primary goal of this thesis work is to build a system by pipelining two different open source projects provided by the authors & developers of the methods [69, 51]. Ultimately Autour will be able to read the texts out loud, present in the user's scene.



Autour

TABLE OF CONTENTS

DEDICATION	ii
ACKNOWLEDGEMENTS	iii
ABSTRACT	iv
LIST OF FIGURES	ix
1 Introduction	1
2 Background	3
2.1 OCR history	3
2.2 Methodologies	4
2.2.1 Text Detection	5
2.2.2 Text Recognition	10
3 Design and Implementation	13
3.1 Design	13
3.1.1 Text detection	14
3.1.2 Text recognition	17
3.2 Implementation	21
3.2.1 Text detection	21
3.2.2 Text recognition	22
3.2.3 Pipelining	23
3.2.4 Inference	27
3.2.5 Training	29
4 Results and Discussion	32
4.1 Results	32
4.2 Discussion	34
5 Benchmark Datasets	37
5.1 Text detection	37
5.2 Text recognition	39
6 Appendices	43
Appendix A - project layout	44

Appendix B - icdar.py	45
Appendix C - model.py	64
Appendix D - multigpu_train.py	68
Appendix E - eval.py	74
Appendix F - preproc.py	77
Appendix G - box_modifier.py	79
Appendix H - crnn_main.py	81
Appendix I - adaptor.cpp	87
Appendix J - recognize.py	89
Appendix K - crnn.py	91
Appendix L - utils.py	93
Appendix M - dataset.py	97
References	101
List of Abbreviations	108

LIST OF FIGURES

<u>Figure</u>	<u>page</u>
1–1 Scene text detection	2
2–1 Stepwise methodology	4
2–2 Integrated methodology	5
2–3 MSER detection of a sample image from ICDAR	7
2–4 Canny edge detection of a sample image from ICDAR	8
2–5 Text segmentation	11
3–1 Autour’s high-level data flow	13
3–2 Text detection high-level data flow	14
3–3 A generic FCN architecture	15
3–4 Text recognition high-level data flow	18
3–5 An RNN architecture	19
3–6 A generic LSTM architecture	20
3–7 An example of box manipulation	26
3–8 A sample image training image from ICDAR15	30
4–1 An instance of text detection	32
4–2 Cropped boxes from Fig 4–1	33
5–1 Typical images from COCO-text	37
5–2 Typical images from MSRA-TD500	38
5–3 Typical images from ICDAR incidental scene text dataset	39
5–4 Typical train images IIIT	40
5–5 Typical images from SVT	41
5–6 Typical character recognition images IC03	42
5–7 Typical images from IC13	42

6-1	Project layout	44
-----	--------------------------	----

CHAPTER 1

Introduction

According to World Health Organization (WHO) there are approximately 253 million people living with visual impairment: 36 million are blind and 217 million have moderate to severe vision impairment as of September 2017. Of these visually impaired people 90% are living in low- and middle-income countries. Furthermore WHO estimates that up to 80% of visual impairment and blindness in adults is preventable or treatable. Important progress is already being made by international communities such as WHO and World Blind Union in order to fight avoidable blindness. Besides that, the existence of projects like Autour is remarkably improving the life quality of blind community.

Being able to read visible texts in natural scenes e.g name tags at a door, street nameplates, store names, has a considerable importance in our daily life as they convey essential information about the environment. A technology allowing access to such information could provide a better independent travel experience, environmental awareness and improved self-confidence to blind users. This is the exact purpose of this thesis work, adding an end-to-end scene text detection and recognition engine to Autour in favor of exposing such opportunity.

Optical Character Recognition, abbreviated as OCR, can be defined as the identification and conversion of printed, typed or handwritten characters into machine encoded text. It is a process of transforming and importing already existing texts into machine environment. On the other hand, scene text detection and recognition is an open and more complicated task when compared

to standard OCR. Scene texts do not necessarily have uniform background or text alignment, they might appear in random sections of the image in different fonts and languages. These facts makes classic OCR techniques inadequate in Autour context. An example of scene text detection can be seen in Fig 1-1.



Figure 1-1: Scene text detection

figure taken from ICDAR's website: <http://rrc.cvc.uab.es/?ch=2&com=tasks>

CHAPTER 2

Background

2.1 OCR history

The idea of machines imitating human abilities is a long coming idea, and character recognition has always been a subset of this idea. The invention of character recognition is considered as Charles R. Carey’s retina scanner. Later on in 1954 The Reader’s Digest magazine installed an OCR mechanism to convert the sales reports into punched cards. As these advancements have triggered the progression, OCR started to have a bigger role in our life as in form of postal, passport and price tag scanner.

Nowadays, OCR is a domain of computer vision which obtains increasing popularity and significance. Top level conferences such as Conference on Computer Vision and Pattern Recognition (CVPR) and International Conference on Computer Vision (ICCV) have witnessed major discoveries in computer vision and pattern recognition. The International Conference on Document Analysis and Recognition (ICDAR) is a leading conference touching to character and text recognition along with camera and video based scene text analysis. The ICDAR hosts number of competitions: Arabic, Chinese Handwriting Recognition Competition, Writer Identification Contest. Nonetheless the most promising one for Autour’s case is the Robust Reading Competition. “Robust Reading”¹ refers to the research area dealing with the interpretation of written communication in unconstrained settings such as born-digital (such as the

¹ <http://rrc.cvc.uab.es/>

ones used in Web pages and email messages) and real scene images and videos. Typically Robust Reading is linked to the detection and recognition of textual information in scene images. The competition is organized around challenges that represent specific application domains for robust reading. Challenges are selected to cover a wide range of real-world situations such as Incidental Scene Text, Focused Scene Text [27].

2.2 Methodologies

Complete text detection and recognition systems can be analyzed with two commonly used methodologies: stepwise (see Fig 2-1) and integrated (see Fig 2-2). Integrated methodologies, have a goal of recognizing words where the detection and recognition procedures share information with character classification and/or use joint optimization strategies [64]. With an integrated methodology, character classification responses are considered the primary cues, and shared with detection and recognition modules. Stepwise methodologies, by contrast, have separated detection and recognition modules, and use a feed-forward pipeline to detect, segment and recognize text regions. They typically employ a coarse-to-fine strategy, which first localizes text candidates, and then verifies, segments, and recognizes them [61].

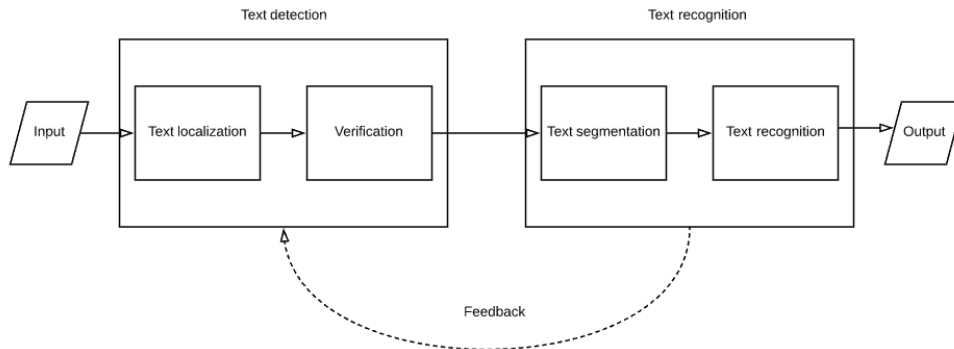


Figure 2–1: Stepwise methodology

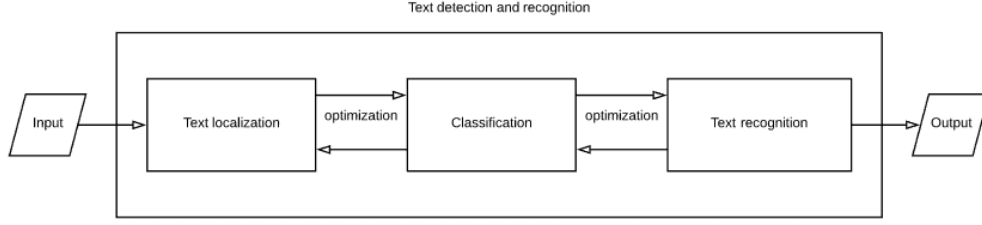


Figure 2-2: Integrated methodology

2.2.1 Text Detection

Existing methods for text detection can be roughly categorized into two major groups: connected component (CC) based methods [7, 11, 55], and region based methods (also called sliding window based methods) [4, 44, 8].

The region based methods use a multi-scale window to exhaustively scan through an image for candidate text regions. Then these regions are classified as text or non-text regions by a classifier which may exert intensity histogram, gradients or edges. Candidate text regions are first found with an edge map or gradient information. Subsequently, a refinement stage is conducted using heuristic rules or learned classifiers [66]. Region based methods commonly utilize an Adaboost classifier [31, 15, 8]. Recently deep learning approaches [59, 25, 9] are exploited, unsupervised learning techniques are used in [9] individually for detection and recognition, a linear SVM classifier is employed to classify the candidate regions. Although these methods can detect text effectively with high recall rate and robust to noise, their classification can be sensitive to false positives due to the large number of candidates. But due to heavy computations for intensive window scanning and advanced classification, these approaches are unsatisfactory to real-time applications.

The connected component(CC) based methods follows the connected component analysis theory, a case where subgroups of CCs are uniquely labeled

hinge on heuristics. The candidate text regions are generated based on extracted CCs of the region, which may differ by the used properties e.g spatial layout, color [65], edge [22], texture [28] or gradient [32] features. Color features are used under the assumption of "text is often produced in a consistent and distinguishable color so that it contrasts with the background" [33]. The family of edge/gradient-based approaches assumes that text exhibits a strong and symmetric gradient against its background. Thus, those pixels with large and symmetric gradient values could be regarded as text components. [64]. Then non-text regions are eliminated, a conditional random fields (CRFs) model is adopted in [46] for this purpose whereas in [70], a neural network is used. The advantage of the connected component methods is that their complexity typically does not depend on the properties of the text (range of scales, orientations, fonts) and that they also provide a segmentation which can be exploited in the OCR step. Their disadvantage is a sensitivity to clutter and occlusions that change connected component structure [44].

Especially, Maximally Stable Extremal Regions(MSER) [39] and Stroke Width Transform(SWT) [13] are two breakthrough techniques which already have achieved impressive performance in scene text detection. MSER is an affine feature region detection algorithm used for blob detection in images. The MSER algorithm extracts a number of co-variant regions called MSER which are connected components of the appropriately thresholded image. A new set of elements alleged extremal regions are introduced, these elements are distinguished regions possessing important properties: they are closed under the affine transformation of image coordinates and invariant to affine transformation of intensity. The word 'extremal' refers to the property that all pixels inside the MSER have either higher (bright extremal regions) or lower (dark extremal regions) intensity than all the pixels on its outer boundary. The

‘maximally stable’ in MSER describes the property optimized in the threshold selection process [40]. These properties makes MSER robust against scale, view point and lighting changes. This is why MSER is a widely adopted algorithm in context of text detection among the other region detectors e.g Harris-affine and Hessian-affine. An example of MSER method is depicted in Fig 2–3.



Figure 2–3: MSER detection of a sample image from ICDAR

The use of MSER based methods for text detection and recognition started with [11], where MSERs are classified using cross-correlation with training templates. However, its sensitivity against blur makes it unsatisfactory to utilize without an auxiliary technique; in [7], MSER and Canny edges [6] are combined to obviate MSER’s sensitivity against blur, which is achieved by removing the pixels outside the boundaries formed by Canny edges. In [23, 55] an intersection of Canny edges with MSER region is taken into consideration. Then a Stroke Width Transform step is placed after region detection by

Canny enhanced MSER. An example of MSER method is depicted in Fig 2–4 (OpenCV implementation of Canny [6]).



Figure 2–4: Canny edge detection of a sample image from ICDAR

Stroke Width Transform assumes that attached characters constructing text blocks are sharing similar attributes suchlike size, color and especially stroke width. SWT is a local image operator which computes per pixel the width of the most likely stroke containing the pixel. The output of the SWT is an image of size equal to the size of the input image where each element contains the width of the stroke associated with the pixel [13]. The biggest restraint of SWT method is its reliance on edge detection which may be unsuccessful in presence of blur and low-contrast. The approach defined in [23] uses SWT as filtering, text candidate regions where the stroke width varies largely are eliminated. Where [7] proposes a way of determining the stroke width based on distance transform, guaranteeing to provide stroke width information at every pixel of the original connected component with any stroke form. Different ways of enhancing SWT are introduced in [3] by edge orientation variance(EOV) and opposite edge pairs(OEP) or a combination of spatial-temporal analysis in [38].

However, these methods fall behind of those based on deep neural networks, in terms of both accuracy and adaptability, especially when dealing with challenging scenarios, such as low resolution and geometric distortion [69]. Deep learning based methods [68, 56, 35, 69, 26] achieved notable performance when compared to traditional methods. A Fully Convolutional Network (named Text-Block FCN) is used to generate a pixel-wise text/non-text salient map in [68], text blocks are detected via the FCN, followed by multi-line oriented text line extraction considering MSER components. The approach described in [56] introduces a novel Connectionist Text Proposal Network, a joint Convolutional Neural Network - Recurrent Neural Network model that directly localizes text sequences in convolutional layers. Authors of [56] have also developed a vertical anchor regression mechanism that jointly predicts vertical location and text/non-text score of each text proposal. The sequential proposals are naturally connected by a recurrent neural network, a Bi-directional Long Short Term Memory. Even though the method is reliable at horizontal and multi-scale text detection, it is not as competent as for text written at an angle. TextBoxes [35] is an end-to-end scene text detector inspired by SSD [36] and inherits the popular VGG-16 architecture [54]. Authors of [35] have proposed a single neural neural network to detect texts by directly predicting word bounding boxes. The approach adopts CRNN [51] as text recognizer in conjunction with TextBoxes, therefore it results in a simple pipeline and a single network to train. The work in [69] is named EAST, a simplistic scene text detection pipeline exploiting an FCN. The network has two stages only: an FCN which directly produces word or text-line level predictions, excluding redundant and slow intermediate steps such as candidate proposal, text region formation and word partition and NMS stage to obtain ultimate results. The method proposed in [26] is called Rotational Region

CNN based on Faster R-CNN [48] architecture. Authors take advantage of a Region Proposal Network to generate bounding boxes enclosing the candidate regions containing text. These proposals are classified and bounding boxes are refined before NMS is utilized to post-process the region candidates to yield the final results.

2.2.2 Text Recognition

Text recognition can be defined as the task of converting regions containing text into strings. Traditionally, text recognition has been focused on document images, where OCR techniques are well suited to digitize planar, paper-based documents. However, when applied to natural scene images, these document OCR techniques fail as they are tuned to the largely black-and-white, line-based environment of printed documents [25]. Despite the efforts and improvements, scene text recognition still remains as a challenging task considering the fact that the scene text regions might have complex background and non-uniform text patterns. In scene text detection-recognition case, detection modules generates bounding boxes around text candidate regions, after this point text recognition attempts to recognize the text represented within the box or refuse the box in case of a false positive detection.

The text recognition problem has been addressed in the literature on multiple levels: character recognition [9, 50], word recognition [45, 43] and text detection. The character recognition problem implicates constructing a recognizer to attain a probability distribution over all characters, when a character containing image is introduced. The character recognition problem is a classification problem that is generally addressed with the use of strong classifiers such as CNNs in [59], deformable parts models [53] or manually-engineered feature-extraction followed by a classifier [12]. The word recognition problem

is, much like phone recognition and handwriting recognition, a sequence recognition problem. Previous works have addressed this problem using CNNs [60], Conditional Random Fields (CRFs) [43, 45] and Pictorial Structures (PS) [2]. As for scene text (cropped word) recognition, the existing methods can be grouped into segmentation-based word recognition and holistic word recognition. The main method for scene text recognition is considered as segmentation-based word recognition. In general, segmentation-based word recognition methods integrate character segmentation and character recognition with language priors using optimization techniques, such as Markov models and CRF [66]. Methods of text segmentation are:

- text binarization: operates to extract text pixels and remove the background pixels
- text line segmentation: attempts to convert a region of multiple text lines into multiple sub-regions of single text lines
- character segmentation: separates a text region into multiple regions of single characters [64]

An instance of text segmentation is shown in Fig 2–5.



Figure 2–5: Text segmentation

figure taken from ICDAR’s website: <http://rrc.cvc.uab.es/?ch=2&com=tasks>

Text segmentation step is used by approaches [4, 25] to achieve accurately bounded characters. These approaches treat isolated character classification

and subsequent word recognition separately. As they read each character independently they do not unleash the full potential of word context information in the recognition. Nonetheless their performance is severely harmed by the difficulty of character segmentation or separation. Importantly, recognizing each character independently discards meaningful context information of the words, significantly reducing its reliability and robustness [17]. Great portion of the work in this field relies on lexicon-dependent approaches. A lexicon is a set of label sequences that prediction is constrained to, *e.g.* a spell checking dictionary [51]. As it currently stands, it is difficult to recognize words with high accuracy without any language model due to the character confusion problem, therefore, all of the previous systems rely on lexicons to improve the results. However, since lexicons can be very large, authors of [2] make the distinction in our approach between where query time is linear in the size of the lexicon and those approaches where it is constant.

Given that texts are formed by sequentially ordered characters it is possible to see them as sequence-like objects. A system aspiring to recognize these objects are supposed to predict a series of labels, therefore this task can be cast as a sequence recognition problems where the length of the sequence may vary severely. RNNs are primarily designed to handle sequences and they are able to effectively learn continuous sequential features. Novel approaches [51, 17] are exploiting RNN’s capability of learning continuous sequential features successfully by combining it with convolutional layers in order to leverage both the advantages of CNN and RNN.

CHAPTER 3

Design and Implementation

3.1 Design

As described and mentioned before, the goal of this dissertation is to design, build and integrate a scene text detector-recognizer to Autour for a possible use case that can be seen in Fig 3–1. The process starts when user captures a photo, that photo is then uploaded to a server which contains the deep learning models (OCR, scene description and chat box). Once the photo is received by the server, it is inferred through the neural networks. Each model generates its output and these outputs can be combined (text regions with detected objects) if there is a relevant match. A stepwise approach (see Fig 2–1) is followed in order to construct OCR pipeline to be invoked upon query. The scene text detection-recognition module consists of two major modules as is evident from its name.

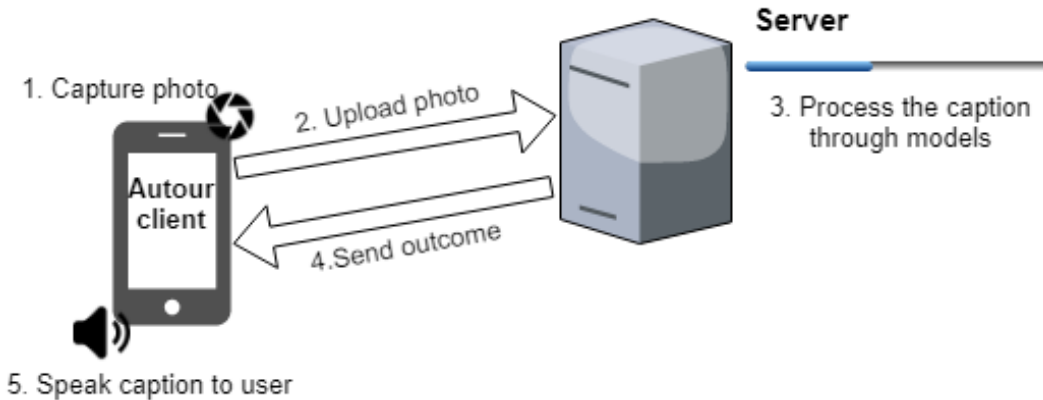


Figure 3–1: Autour’s high-level data flow

3.1.1 Text detection

Text detection module relies on EAST: An Efficient and Accurate Scene Text Detector [69] and its reimplementation¹ with TensorFlow [1]. A high-level overview of text detection module is depicted in Figure 3–2. The fundamental element of the proposed algorithm is a fully-convolutional neural network which follows the general principles of [21]. The proposed model abandons unnecessary intermediate components and steps, and allows for end-to-end training and optimization. The resultant system, equipped with a single, light-weighted neural network, surpasses all previous methods by an obvious margin in both performance and speed.

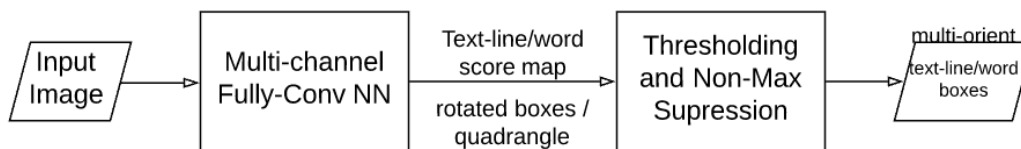


Figure 3–2: Text detection high-level data flow

Various parameters should be taken into consideration whilst designing neural networks aimed to detect text. Predicting bounding geometries of a small text region requires low level information in early stages, while large text regions need features from late-stage of a neural network. A network should be designed in a way that it should be able to utilize features from early and late stages of its structure, considering the fact that real scene text regions are not uniform. HyperNet [30] meets these conditions on features maps, however merging a large number of channels on large feature maps would significantly increase the computation overhead for later stages. In remedy of this, authors

¹ <https://github.com/argman/EAST>

of EAST have adopted the idea from U-shape [49] to merge feature maps gradually, while keeping the up-sampling branches small [69].

EAST’s model can be decomposed in to three parts: feature extractor *stem*, feature-merging *branch* and output layer. An image is passed into the FCN, starting from the feature extractor and multiple channels of pixel-level text score map and geometry are generated at the output layer. A generic FCN architecture can be seen in Fig 3–3. The *stem* can be a convolutional network with interleaving convolution and pooling layers. Four levels of feature maps are extracted from the stem, whose sizes are respectively $\frac{1}{32}$, $\frac{1}{16}$, $\frac{1}{8}$ and $\frac{1}{4}$ of the input image. As a base model of stem, three different networks are examined: VGG16 [54] which is a commonly used model in many tasks, PVANET [19] a light-weight substitute of the feature extractor Faster-RCNN [48] and its double-channelled version *PVANET2x*. A generic FCN architecture developed for classification is depicted in Fig 3–3.

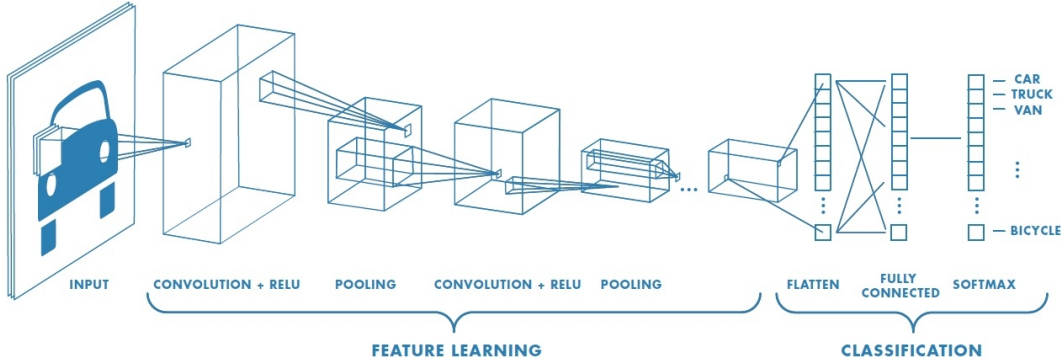


Figure 3–3: A generic FCN architecture

Image taken from from Matworks - <https://www.mathworks.com/discovery/convolutional-neural-network.html>

In the *branch* the features are gradually merged, in each merging stage, the feature map from the last stage is first fed to an unpooling layer to double its size, and then concatenated with the current feature map [69]. The final

output layer contains several 1×1 convolution operations to project 32 channels of feature maps into 1 channel of score map and a multi-channel geometry map. Authors of EAST have experimented two geometry shapes for text regions, rotated box (RBOX) and quadrangle (QUAD), the geometry output can be one of them. For RBOX, the geometry is represented by 4 channels of axis-aligned bounding box (AABB) \mathbf{R} and 1 channel rotation angle θ , \mathbf{R} is formulated in [21]. The authors of [21] have defined the left top and right bottom points of the target bounding box in output coordinate space as $p_t = (x_t, y_t)$ and as $p_b = (x_b, y_b)$ respectively, then each pixel i is located at (x_i, y_i) in the output feature map describes a bounding box with a 5-dimensional vector as $\hat{t}_i = i\{\hat{s}, \hat{d}x^t = x_i - x_t, \hat{d}y^t = y_i - y_t, \hat{d}x^b = x_i - x_t, \hat{d}y^b = y_i - y_t\}$, 4 channels $\hat{d}x^t, \hat{d}y^t, \hat{d}x^b, \hat{d}y^b$ denote 4 distances from the pixel location to the top, right, bottom, left boundaries of the rectangle respectively and \hat{s} is the confidence score of being an object in. In EAST, \hat{s} denotes the rotation angle θ . For QUAD, 8 numbers are used to denote the coordinate shift from four corner vertices (same as RBOX description) of the quadrangle to the pixel location. As each distance off-set contains two numbers $(\Delta x_i, \Delta y_i)$ the geometry output contains 8 channels [69].

Non-Maximum Suppression (NMS) stage has made its place in object detection methods as a post-processing step. As a popular pattern (object, text) detection approach the geometry scores are filtered by a predefined thresholding, valid geometries are then merged through NMS. A greedy NMS algorithm greedily selects high scoring detections and deletes close-by less confident neighbors since they are likely to cover the same object, a reasonable practice, since the actual goal is to generate exactly one detection per object. Yet the greedy NMS method makes hard decision by deleting detections and bases this decision on one fixed parameter that controls how wide the suppression is [20].

A naive NMS (running in $O(n^2)$, where number of candidate geometries is noted as n) would not be efficient, as the predictions are in great number. As an alternative a weighted merge algorithm is developed by the authors of EAST instead of a greedy NMS by assuming that closely located pixels are likely to be highly correlated. Developed algorithm merges the geometries row by row, and while merging geometries in the same row, it iteratively merges the geometry currently encountered with the last merged one. The coordinates of merged quadrangle are weight-averaged by the scores of two given quadrangle. To be clear, given two quadrangles g, p , the proposed method *WEIGHTEDMERGE* is defined as follows: $a = \text{WEIGHTEDMERGE}(g, p)$ then $a_i = V(g)g_i + V(p)p_i$ and $V(a) = V(g) + V(p)$ where a_i is one of the coordinates of a subscripted by i and $V(a)$ is the score of geometry a .

3.1.2 Text recognition

Text recognition module is built on An End-to-End Trainable Neural Network for Image-based Sequence Recognition and Its Application to Scene Text Recognition [51] which originally is implemented in Torch [10] but for this project its PyTorch [47] ² port ³ is chosen. Convolutional Recurrent Neural Network is explicitly designed for recognizing sequence-like objects in images. CRNN is a combination of Deep Convolutional Neural Network (DCNN) and Recurrent Neural Network (RNN). This combination possesses desired properties from both architecture and it is not confined by their constraints. Like a DCNN, CRNN is able to learn informative representations directly from image data but it has much less parameters than a standard DCNN model. CRNN is capable of producing sequence of labels still it is not constrained by the length

² <http://pytorch.org/>

³ <https://github.com/mejjieru/crnn.pytorch>

of the sequence-like object. A high-level overview of text detection module is depicted in Fig 3-4.

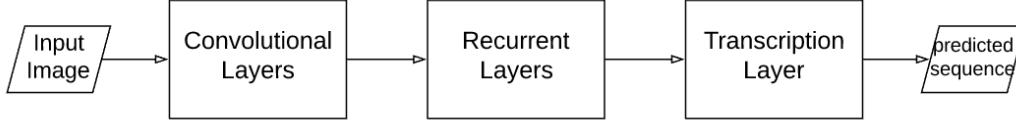


Figure 3-4: Text recognition high-level data flow

The convolutional layers automatically extract a feature sequence from each input image. These layers are adopted from a standard CNN by taking its convolutional and max-pooling layers and removing fully-connected layers so that it can be used to extract a sequential feature representation from an input image. Then a sequence of feature vectors is extracted from the feature maps produced by the convolutional layers and fed into recurrent layers. Convolutional layer architecture is based on well-known VGG [54] design. A deep bidirectional Recurrent Neural Network, as recurrent layers, is built on the top of the convolutional layers, in order to make prediction for each frame of the feature sequence. An RNN is favored in this context because of several reasons: its strong capability of capturing contextual information within a sequence, its ability to back-propagate error differentials to its input, its strength to operate on sequences of arbitrary lengths, traversing from starts to ends. However, traditional RNN units deteriorates from vanishing gradient problem which prevents the gradients *i.e.* weights of the network, from changing its value. Vanishing gradient problem confines the range of context the unit can store and causes difficulty to train the network. The basic structure of a recurrent neuron is shown in the left hand side of Fig 3-5, the unfolding of three time steps of the RNN is depicted in the right-hand side.

Long Short Term Memory (LSTM) [18] is specifically designed to address

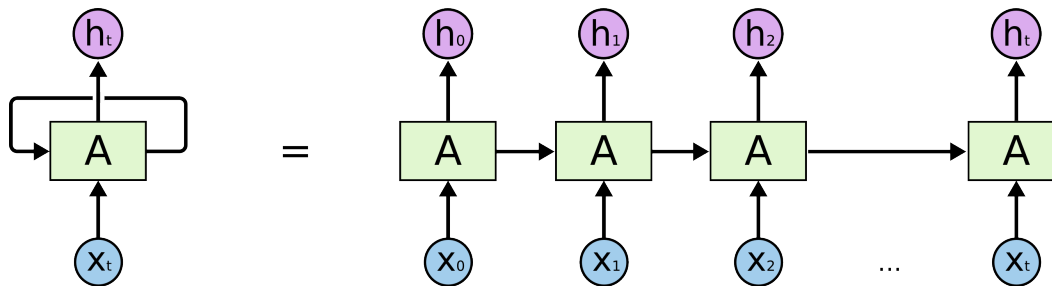


Figure 3-5: An RNN architecture

Image taken from from <https://wiki.tum.de/display/lfdv/Recurrent+Neural+Networks+-+Combination+of+RNN+and+CNN>

this problem. A common LSTM composition consists of a memory cell, an input gate, an output gate and a forget gate. Essentially, the memory cell stores past values *e.g.* states, and the input and output gates allow the cell to store values for either long or short time periods. This time interval is defined by designing the forget gate's activation, namely, the use of linear identity function (which has 1 as its derivative) will avoid the gradient to vanish. Nevertheless, in image-based sequences, information from both forward and backward directions are valuable and complementary to each other yet traditional LSTM's are unidirectional, allowing the use of only past information. To overcome this issue a bidirectional LSTM is designed by merging a forward and a backward LSTM. A generic LSTM is depicted in Fig 3-6 with its gates. The implementation can be seen on **crnn.py** on page 91 as class named *BidirectionalLSTM*.

The transcription layer, as the final component of CRNN is adopted to translate the per-frame predictions by the recurrent layers into a label sequence. For this reason conditional probability defined by Connectionist Temporal Classification (CTC) layer proposed in [14] is adopted . CTC is proposed to solve a fundamental issue of RNNs. RNNs can only be trained to make a series of independent label classifications. This means that the training data must be pre-segmented, and that the network outputs must be post-processed

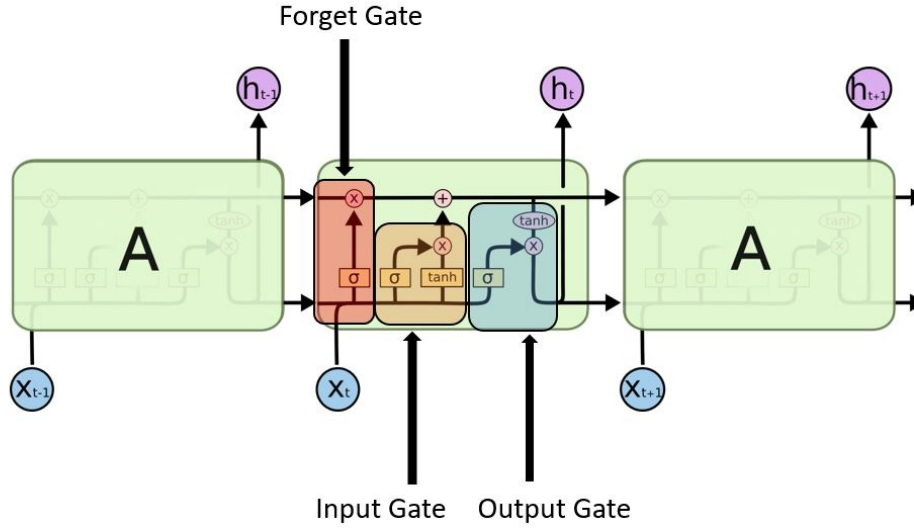


Figure 3–6: A generic LSTM architecture

Image taken from: <https://wiki.tum.de/display/lfdv/Recurrent+Neural+Networks+-+Combination+of+RNN+and+CNN>

to give the final label sequence. CTC models all aspects of the sequence within a single network architecture and trains the network to label the entire input sequence at once [14].

Both lexicon-based and lexicon-free transcription are examined. However, for extensive lexicons, such as Hunspell spell-checking dictionary⁴, it would be highly time consuming to execute an exhaustive search over the lexicon. The authors of [51] have discovered that label sequences predicted via lexicon-free transcription are generally close to the ground-truth under the edit distance metric. Exploiting this fact, CRNN authors were able to limit the search to the nearest neighbor candidates in order to cope with long search issue.

The objective function (negative log-likelihood of conditional probability of ground truth) which calculates a cost value directly from the image and its

⁴ <https://hunspell.github.io>

ground truth label sequence is defined. Stochastic gradient algorithm is used to train the network, where gradients are calculated by the back-propagation algorithm [51].

3.2 Implementation

3.2.1 Text detection

The text detection module is implemented using Google’s open-source machine learning framework TensorFlow. Main programming language of the project is Python. However, locality aware NMS module is written in C++ by the authors of [51] which provides drastic time savings. The project depends on several Python packages. SciPy, is a Python ecosystem consists of packages like NumPy and Matplotlib which are commonly used in mathematics, science, and engineering applications. Especially NumPy is a preferred Python library especially for scientific applications like neural networks, providing support for large, multi-dimensional arrays and matrices, along with a broad set of high-level mathematical functions to perform on these arrays. Polygon function of Shapely, a Python package for manipulation and analysis of planar geometric objects, ease area calculations of polygons, use can be seen on page 45 of **icdar.py**. Besides that there are several architectural differences between the implementation and the published paper: paper adopts three different base networks VGG16 [54], PVANET [19] and *PVANET2x* for its feature extractor *stem*, for the Tensorflow implementation of ResNet50 [16] is selected. ResNet (Residual Network) introduces residual learning, which eases the training of the network and provides increased depth. The variants of ResNet (ResNet50, ResNet101, ResNet152) have proved their success and accuracy on image classification. As ResNet V1 50 implementation ⁵ of TensorFlow’s TF-Slim library,

⁵ <https://github.com/tensorflow/models/blob/master/research/slim/nets/>

a lightweight package for defining, training and evaluating models, is utilized. The paper experiments with two geometry shapes for text regions: rotated box (RBOX) and quadrangle (QUAD). EAST authors have adopted balanced cross entropy loss inspired by [63] to facilitate a simpler training procedure. The loss function is defined as $-\beta Y^* \log \hat{Y} - (1 - \beta)(1 - Y^*) \log(1 - \hat{Y})$ where \hat{Y} is prediction of the score map (named *F_score* in the code), Y^* is the ground truth and β is the balancing factor between positive and negative samples. Alternatively the implementation uses dice loss (optimize Intersection over Union of segmentation) for both score map prediction and RBOX regression. Briefly, dice loss is a measure of overlap, proposed in [41] as a loss function. Authors of [41] states that dice coefficient is a quantity ranging between 0 and 1 aimed to be maximized. The dice coefficient D between two binary volumes is be written as:

$$D = \frac{2 \sum_i^N p_i g_i}{\sum_i^N p_i^2 + \sum_i^N g_i^2}$$

where the sums run over the N voxels, of the predicted binary segmentation volume $p_i \in P$ and the ground truth binary volume $g_i \in G$. The code definition of loss function and dice coefficients can be seen on page 64, defined in **model.py** as *loss* and *dice_coefficient* functions respectively.

3.2.2 Text recognition

The text recognition module is originally implemented ⁶ by the authors of [51], using Torch [10] an open source machine learning library implemented in C with a wrapper in LuaJIT scripting language. The original implementation is a good reflection of the published paper with custom implementations

⁶ <https://github.com/bgshih/crnn>

for the LSTM units (in Torch7/CUDA), the transcription layer (in C++) and the BK-tree data structure (in C++). However for Autour’s OCR PyTorch port is selected, which is slightly different than the original. The original implementation contains two modes of transcription, namely the lexicon-free and lexicon-based transcriptions. In lexicon-free mode, predictions are made without any lexicon. In lexicon-based mode, predictions are made by choosing the label sequence that has the highest probability. In the paper, ADADELTA [67] is used for optimization to automatically calculate per-dimension learning rates. The PyTorch port is built on PyTorch, an open-source deep learning framework mainly developed by Facebook’s artificial intelligence research team. It wraps the core Torch binaries in Python. The project has two dependencies: lmdb and warp-ctc. Lmdb is a universal Python binding for the LMDB Lightning Database, it is used to create a tiny database to manipulate (resize, label, batch) the images. Warp-CTC is a CTC implementation developed by Baidu’s AI Lab. Since the project is implemented in PyTorch, Warp-CTC’s PyTorch binding⁷ is used. Unlike the original implementation PyTorch port does not provide lexicon-based transcriptions despite this it enables to choose between ADAM, ADADELTA and RMSprop as optimizer.

3.2.3 Pipelining

Testings of this project is started in McGill University’s Shared Reality Lab computer with a dedicated NVIDIA Tesla K40c GPU with CUDA 7.5 installed and a 4 core Intel Core i5 Haswell CPU at 3.50GHz. Tesla series GPUs are targeted for high-performance computing applications such as stream processing, deep learning(both training and inference). Tesla K40c model has 2880 processor cores, 12GB of memory with type DDR5 SDRAM.

⁷ <https://github.com/SeanNaren/warp-ctc>

GPUs differently from CPUs are designed to compute the same instructions in parallel, they might have thousands of cores, more computational units and a higher bandwidth to retrieve from memory. On the other hand, DNN's are engineered in a manner such that at each layer of the network thousands of identical artificial neurons perform the same computation (especially huge amounts of matrix multiplications). As a consequence, the structure of a DNN fits greatly with the sorts of computation that a GPU can efficiently (by parallelizing) perform. Exclusively NVIDIA is pioneering artificial intelligence with their specially designed hardware and their SDKs such as cuDNN, cuSPARSE and cuBLAS to support main deep learning frameworks *e.g.* TensorFlow, Keras and PyTorch. CUDA Toolkit of NVIDIA supplies a development environment for building high performance GPU-accelerated applications. With CUDA, developers are able to greatly fasten computing applications by exploiting the power of GPUs. In GPU-accelerated applications, CPU handles the sequential part of the workload since it is optimized for single-threaded performance, while the computationally demanding fragment of the application runs on GPU cores in parallel. As the modules are implemented in different frameworks, it is critical to make sure that both frameworks will be able to access GPU for computations by installing their GPU-enabled versions.

The main contribution of this thesis is to combine the two aforementioned modules in a pipeline. By default, EAST's implementation is generating axis-aligned minimum bounding boxes (AABB) are wrapped around the text region. Namely, the text detection module produces bounding boxes in a following manner $(x_1, y_1, x_2, y_2, x_3, y_3, x_4, y_4)$. Meanwhile CRNN's implementation uses Python Imaging Library (PIL) which accepts rectangles as input, in this case it was not possible to use text detection's output without modifying.

The solution is to produce bounding boxes using its ($\min(x)$, $\max(x)$, $\min(y)$, $\max(y)$) coordinates which encloses the AABB generated by text detection. These bounding boxes are very likely to have bigger size as they have excess pixels with no text content, modified to fit the conditions. In other words, a box with coordinates (377, 117, 463, 120, 465, 130, 378, 150) is enlarged to (377, 465, 117, 150). Finally these boxes are cropped from the image, represented as NumPy ndarrays and passed to text recognizer (**recognize.py**), starting from line 87 of **eval.py** on page 74. However, text detection module intrinsically produces tightly wrapped bounding boxes around the text region, this causes separation of words which are meant to be read and/or pronounced together.

A straightforward search algorithm is developed to predict such text regions (to be combined), so that each cropped box can be passed in recognition module individually. This algorithm is called *match_check* (defined in **box_modifier.py** on page 74) and it basically looks for an overlap of bounding boxes over all possible permutations of two boxes. Since the boxes are enlarged, the probability of overlap of relevant boxes is naturally increased. This approach lowers the probability of false transcriptions by breaking word sequences. An image is shown in Fig 3-7 to provide better understanding of **box_modifier.py**. Green boxes are created by the text detection module, yet red boxes are cropped to pass into text recognizer then as the overlapping boxes (the ones at top right-hand corner and the ones at mid right-hand side) are combined together (after recognition) to form a bigger box.

Recognition module then transcribes each cropped box into a string. In order to have accurate and relevant results a spellchecking library PyEnchant is placed after recognition. Each string is run through the dictionary to check

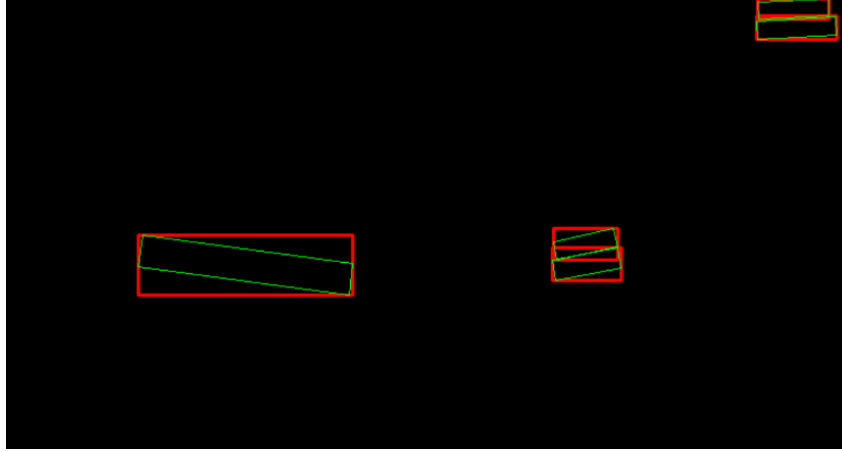


Figure 3-7: An example of box manipulation

whether the word is correctly spelled, in case the word is misspelt the first element from a suggestion list (where the suggested words are ordered from most likely replacement to least likely) is returned. It is an essential step, especially for a city like Montreal where both English and French is commonly used in daily life.

To sum up, this thesis work gathers two independent open-source work and builds a bridge between them. Output of the text detection module is manipulated and provided to text recognition, then the results are improved with the help of a spellchecking dictionary. Exploiting the location information of the boxes, the strings ‘likely to be related’ are combined with a developed algorithm and their box coordinates are merged. As the developed and provided projects of the authors of EAST [69] and CRNN [51] are performing satisfactory, no modifications are done on top of values selected by them. Python scripts of the original projects are divided into smaller parts in order to have a clear view and understanding e.g **preprocess.py**, **detect.py** and **recognize.py**. Project layout in an Integrated Development Environment can be seen on page 44. By default trained CRNN model is stored in

directory *model_CRNN*, trained text detection (EAST) model is stored in directory *east_icdar2015_resnet_v1_50_rbox* and the input images are taken from directory *images*.

3.2.4 Inference

The inference can be initiated by invoking **eval.py** found on page 79. Multiple flags are set with *tf.app.flags.DEFINE_string* to define variables such as input data path, model path and GPU list. By its design Tensorflow has two main steps: constructing a dataflow graph (*tf.Graph*) and executing the graph in a session (*tf.Session*). Generally most Tensorflow programs begins with a dataflow construction stage, where nodes (*tf.Operation*) and edges (*tf.Tensor*) are specified in an abstract way. Then a *tf.Session* is created, running the session allows the graph to be executed on resources like CPU or GPU. This approach brings few favorable attributes such as parallelism, distributed execution, compilation and portability. In **eval.py** they are defined as *tf.get_default_graph().as_default()* and *tf.Session(config=tf.ConfigProto(allow_soft_placement=True)) as sess* respectively, *allow_soft_placement* flag ignores *tf.device* annotations that attempt to place CPU-only operations on a GPU device. When the session is created the trained model is restored from *checkpoint_path* defined with a flag.

The function *get_images* (defined at **preproc.py**, page 77), is used to load the images from *test_data_path*. Then the images are resized by the function *resize_image* (there is no certain input size, however by design it has to be multiple of 32 and less than 2400 pixels per vertex in order to confine GPU memory usage).

In line 65 of **eval.py** the session is run, the input images are fed into model. Within the function *model* (defined in **model.py**) the model is easily built through the use of argument scoping, *arg_scope*, provided by *tf.slim*. The

function *model* takes images as input and generates score (F_score) for each geometry ($F_geometry$), where $F_geometry$ is represented by 4 channel of axis aligned bounding box and 1 channel rotation angle.

The obtained scores and geometries are passed to function *detect* (defined at **detect.py**, on page) as parameters: *score_map* and *geo_map*. Besides these, *detect* has three different constant threshold parameters: *score_map_thresh*, *box_thresh*=0.2, *nms_thres*=0.6. They are thresholds for respectively: score map, box, non-maximum suppression. Score maps are filtered according to threshold value and detected boxes are sorted via the y-axis afterwards. The sorted boxes are then restored with function *restore_rectangle* (defined in **icdar.py**, on page 45) before NMS.

NMS stage (implemented in C++ and ported to project) is initiated at line 37 of **detect.py** by invoking the function *merge_quadrangle_n9*, alternatively there is an unused Python version of NMS is included within the project under the name **locality_aware_nms.py**. The script **adaptor.cpp** (can be seen on page 87) describes the bindings. The obtained bounding boxes with low scores are filtered by the average score map, differently from paper.

Acquired bounding box coordinates are checked for their coherence, starting at line 83 of **eval.py**. In the sequel, these boxes are cropped as mentioned in section 3.2.3 and fed into recognizer.

Each step before this point was a component of text detection, by calling **recognize.py** system switches to text recognition. When the function *recognize* is invoked an instance of CRNN model (defined in **crnn.py**, on page 91) is created depending on CUDA availability. Then the trained model is loaded from *model_path*. The defined alphabet (at line 13 of **recognize.py**) is then converted from string to labels by function *decode* in class *strLabelconverter* defined in **utils.py** (can be found on page 93). The cropped bounding

box image is loaded, resized (to 100×32) and passed into a PyTorch variable with the same name (*image*). Raw predictions are attained after an inference through the network (see line 45 of **recognize.py**).

The inference of an high definition image through the text detection network takes less than 1.5 seconds and less than 50 millisecond at NMS stage (around 400 milliseconds when NMS Python is used), an inference through the whole pipeline takes less than 3 seconds on the mentioned computer of Shared Reality Lab. On a mid-end computer with no CUDA support and an Intel i5 CPU at 2.50GHz (model runs on CPU), text detection inference takes about 4 seconds and the whole pipeline inference runs around 6 seconds. Bearing in mind that, as GPU limitations are not explicitly set, deep learning frameworks allocates all GPU memory for the process whilst they run. The size of the input image and the number of text regions in the image plays a crucial role on inference timing as well as GPU performance .

3.2.5 Training

Since there are two different modules implemented on two different frameworks, they have to be trained separately. Text detection model can be trained by running **multigpu_train.py**. Several flags are set in the beginning of the script e.g. *batch_size_per_gpu*, *learning_rate*, *save_checkpoint_steps* and *save_summary_steps*. At the function *tower_loss* (starting at line 28 of *multigpu_train.py*) the inference graph is built by calling the model defined at **model.py**. *Model_loss* is computed with *loss* function defined in **model.py** and by adding regularization losses to *model_loss*, *total_loss* is computed (see line 35 of **multigpu_train.py**) In the *main* function, placeholders for input images and geometry maps are created global step, learning rate and optimizer is defined. The session *tf.Session* starts at line 137, and runs until loss is diverged, if not until the defined step number is reached.

To train the text detection model benchmark datasets (can be found in section 5.1) can be used. For training, a separate text file should be provided for each image in dataset. For instance, if the scene image is represented as `img_1.jpg` the corresponding text file should be named as `img_1.txt`. Each text region should be noted with its coordinates in format $(x_1, y_1, x_2, y_2, x_3, y_3, x_4, y_4)$ followed by its text representation. The annotations are loaded by using *load_annotation* function described in **icdar.py**. A sample training image is depicted in Fig 3–8, and corresponding text file should be as follows:

```
377,117,463,117,465,130,378,130,Genaxis Theatre
493,115,519,115,519,131,493,131,[06]
374,155,409,155,409,170,374,170,###
492,151,551,151,551,170,492,170,62-03
376,198,422,198,422,212,376,212,Carpark
494,190,539,189,539,205,494,206,###
374,1,494,0,492,85,372,86,###
```



Figure 3–8: A sample image training image from ICDAR15

The network is trained end-to-end using only training images from ICDAR 2015 and ICDAR 2013 with ADAM [29] optimizer. Learning rate of ADAM starts from $1e-3$, decays to one-tenth every 27300 mini batches, and stops at $1e-5$ (staged learning rate decay), implementation instead accepts linear learning rate decay [69], can be seen on page 68. Such trained model

achieves 77.32% recall and 84.66% precision rate on ICDAR 2015 Incidental Scene Text Detection Challenge. It would be beneficial to keep in mind that the script **multigpu_train.py** is developed considering ICDAR images as input (**icdar.py** is developed specifically to adapt to **multigpu_train.py**).

Text recognition module can be trained by running **crnn_main.py** (can be seen on page 81). The script starts with flags e.g. *ngpu* (number of GPUs to use), *batchSize*, *niter* (number of epochs to train for) and optimizer option (ADAM, ADADELTA or RMSprop). The function *trainBatch* takes the CNN model, optimizer (set up starts at line 114 of **crnn_main.py**) and CTCLoss from warp-ctc as criterion, trains a batch and returns *cost* (defined at line 187) until the *niter* is reached. To train the text recognition model, benchmark datasets (can be found in section 5.2) can be used. An LMDB Lightning Database should be created by using function **dataset.py**, any kind of preprocessing e.g. transformations, tensorization and batching is handled by class *lmdbDataset* defined in the same script. The training of CRNN model takes little more than 2 days using MJSynth synthetic dataset introduced in [24].

CHAPTER 4

Results and Discussion

4.1 Results

This section aims to demonstrate how the pipeline works by showing its results. A photo of an hotel taken in downtown Montreal is used. The implementation of EAST authors generates cyan colored bounding boxes, yellow colored bounding boxes are cropped from the image (see Fig 4-1).



Figure 4-1: An instance of text detection



Figure 4-2: Cropped boxes from Fig 4-1

The raw output of the text recognition module is as following:

```

1-----2----1----6----- => 1216
s-----c--o--m-f-o-r-t--- => scomfort
c-----o--m--f--o--r--t--- => comfort
s-----u---i--t--e---s--- => suites
s-----u---i--t--e---s----- => suites
h-----o-----t---e---l--- => hotel
e-----t--o----h---- => etoh
p-----e---r--i-o---d--ee-- => periode
i--n--t-e--r---d-iitt-e--- => interdire
2-----z----- => 2z
s-----o---u---n---d--- => sound

```

Note that '*scomfort*' is the output of bigger box containing 'comfort' string found at top left-hand corner of the image, '*etoh*' is the output of small box containing reflected 'hotel' string and '2z' is the output of graffiti written on pole. After spellchecking filtering, accent removal of French words and box merging the outcome is [*'comfort suites', 'hotel', 'suites comfort', 'periode interdite', 'comfort suites', '1216'*], followed by the box centre coordinate in percentage with regards to the image. To be more clear, if the box centre is

at (120, 200) where the image size is 480×600 the returned location will be (0.25, 0.2).

The project is maintained in this ¹ repository, the trained ResNet model can be downloaded from this ² link, a complete trained EAST model can be found in this ³ Google drive folder and trained CRNN model can be found in this ⁴ Dropbox folder.

4.2 Discussion

Currently, this OCR engine has its role within an existing mobile application. Depending on external factors, the application is able to give detailed visual description of users' scene with a tolerable latency in real-life scenarios. However, it is necessary to be aware of weaknesses of the system. It is clear that today's mobile platforms are not able to run such sized models (scene description and OCR in Autour's case) with their resource limitations. As an initial phase, Autour server runs these models sequentially. Namely, the scene description model is run then it is followed by OCR model. Even though it is a minor setback (which can be solved by assigning multiple resources or switching to cloud computing platforms enabling dynamic scaling), it introduces noticeable latency.

The current string combinations are made only in pairs. Namely, the system will fail to combine if there are three consecutive words present in the scene. Besides that, the filtering method, as it stands, relies on a spellchecking dictionary and the system has no apriori information about the language of

¹ <https://github.com/heildever/AutourOCR>

² download.tensorflow.org/models/resnet_v1_50_2016_08_28.tar.gz

³ <https://drive.google.com/file/d/0B3APw5BZJ67ETHNPau9xUkVoV0U/>

⁴ <https://www.dropbox.com/s/dboqjk20qjkpta3/crnn.pth?dl=0>

the text (options may vary in multi-cultural cities) which means if the word is misspelt the system suggests French words only. One of the intended improvements is to make use of location information coming from FoursquareTM, Google PlacesTM or OpenStreetMap look for possible matches i.e. in case of Fig 4-1 the hotel name ‘Comfort Suites’ can be confirmed by location information. As blind users of Autour indicated that they are not interested in every text in their scene.

The text detection module holds a prestigious ranking among published methods, and its implementation is performing satisfactory for Autour’s use. However, ICDAR 2015 dataset is constructed by images containing mostly horizontally oriented text regions. Since the model is trained on such dataset, this may cause miss or imprecise predictions in case system encounters a vertically text instance.

CRNN inherently is able to recognize text sequences written with a straight orientation with similar characters. In the other words, system will fail to recognize if the text is written upside down or slightly angled. On the other hand, authors of CRNN [51] made a tweak to recognize English texts by adopting 1×2 sized rectangular pooling windows instead of the conventional squared one in the 3rd and 4th max-pooling layers. Assuming that typically a feature sequence of 25 frames can be generated from an image sized 100×32 and an image with that size contains up to 10 characters which exceeds the length of most English words. Meaning that system might misrecognize the words longer than 10 characters.

Nevertheless CRNN [51] is still accepted as state-of-the-art and used as a recognition module such as TextBoxes [35] and TextBoxes++ [34]. In addition to this the authors of CRNN have published a model which is able to recognize several types of irregular text, including perspective text and curved text [52].

Differently, a unified network is proposed in [37] for simultaneous detection and recognition. A new differentiable operator is introduced in order to share convolutional features between detection and recognition. Leveraging from convolution sharing strategy and the joint training method (enabling to learn more generic features), the method outperforms almost all state-of-the-art methods in both text localization and end-to-end tasks of ICDAR Incidental Scene Text challenge. Rankings can be seen on ICDAR’s website ⁵ .

⁵ <http://rrc.cvc.uab.es/?ch=4&com=evaluation&task=1>

CHAPTER 5

Benchmark Datasets

5.1 Text detection

COCO-Text [57] is a large scale dataset for text detection and recognition in natural images. The dataset is based on the MS COCO dataset, which contains images of complex everyday scenes. The dataset is organized around three tasks: Text localization, Cropped Word Recognition and End-to-End Recognition. The images were not collected with text in mind and thus contain a broad variety of text instances. The dataset contains 63,686 images with 173,589 labeled text regions in which 43,686 are chosen to be the training set and the rest 20,000 for testing. Three images from COCO-Text are shown in Fig 5–1.



Figure 5–1: Typical images from COCO-text

MSRA-TD500 [62] The MSRA Text Detection 500 Database contains 500 natural images, which are taken from indoor (office and mall) and outdoor (street) scenes using a pocket camera. The indoor images are mainly signs, doorplates and caution plates while the outdoor images are mostly guide boards and billboards in complex background. The dataset is divided into two

parts: training set and test set. The training set contains 300 images randomly selected from the original dataset and the remaining 200 images constitute the test set. All the images in this dataset are fully annotated. Several images from MSRA-TD500 are depicted in Fig 5–2, text regions are emphasized with bounding boxes.



Figure 5–2: Typical images from MSRA-TD500

ICDAR 2015 [27] Competition on Robust Reading is structured in four challenges addressing text extraction in different application domains, namely born-digital images, focused scene images, incidental scene text and video text. Incidental Scene Text refers to text that appears in the scene without the user having taken any prior action to cause its appearance in the field of view, or improve its positioning or quality in the frame. While focused scene text is the expected input for applications such as translation on demand, incidental scene text covers another wide range of applications linked to wearable cameras or massive urban captures where the capture is difficult or undesirable to control. For incidental scene text a new dataset is introduced, a dataset of 1,670 images (17,548 annotated regions) acquired using the Google Glass.

Authors of [27] have used 1000 of images are used for training and the remaining are for testing. Few images from ICDAR’S incidental scene text dataset are shown in Fig 5–3.

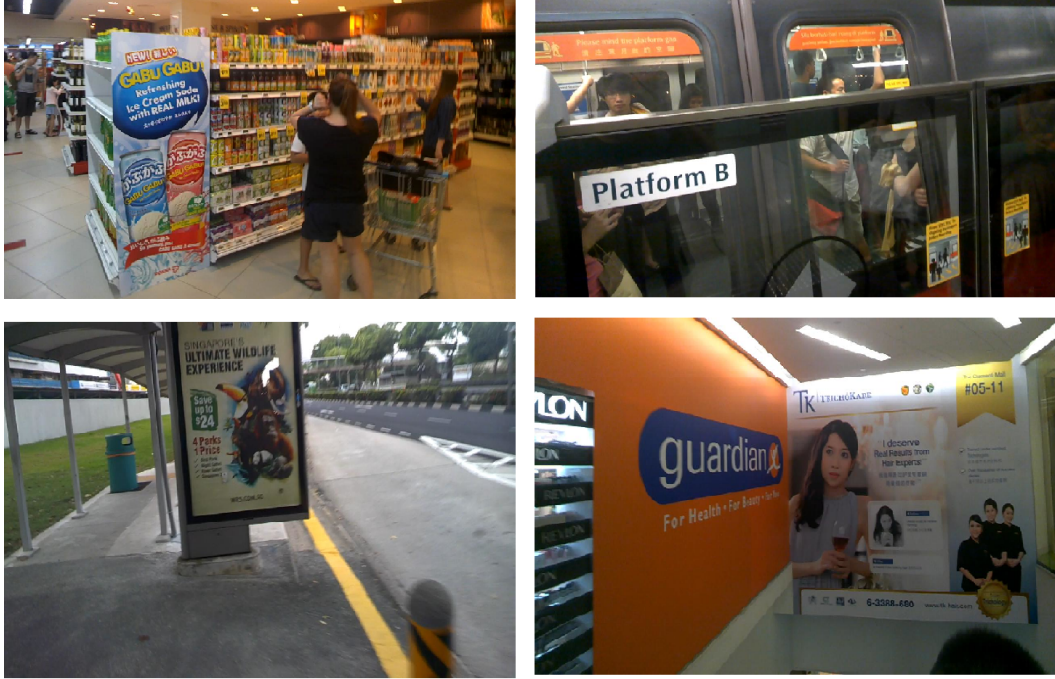


Figure 5–3: Typical images from ICDAR incidental scene text dataset

5.2 Text recognition

IIIT 5K [42] is harvested from Google image search. Query words like billboards, signboard, house numbers, house name plates, movie posters were used to collect images. The dataset contains 5000 cropped word images from Scene Texts and born-digital images. The dataset is divided into train and test parts. Three training images (above) and three test images (below) of IIIT5K are depicted in Fig 5–4.

SVT [58] The Street View Text dataset was harvested from Google Street View. The annotators were asked to find a place of business with its signboard, then have a clear point of view to minimize the skew of the text before taking a screen shot. Therefore the dataset includes business names and business



Figure 5-4: Typical train images IIIT

signs. Each business (hotel, restaurant and etc) and its sign is associated with representative text. Precisely; the annotators have generated annotations of horizontal bounding boxes and a non case-sensitive transcription for each text region. The word annotations are adopted to produce a dataset of cropped words called SVT-50. In total, SVT dataset contains 100 training and 250 testing images collected from 20 different cities. Four different SVT images can be seen in Fig 5-5.

IC03 datasets are created by the images that have been extracted from natural scenes for the ICDAR 2003 Robust Reading competitions. Four independent competitions were organized: Robust Reading, Robust Word Recognition, Robust Character Recognition and Text Locating. Recognition datasets are word recognition (made of 1157 words training set, 1111 words of testing set) and character recognition (made of 6185 characters training set, 5430 characters testing set). Few cropped character and cropped word images are depicted in Fig 5-6.

IC11 and IC13, IC11 is an extension of datasets used in earlier Robust Reading Competitions organized in ICDAR 2003 and 2005. Images of IC03 are taken as a base, the images without text are removed and around 100 images are added which are captured with a digital camera using auto focus and natural lighting. The final dataset consisted of 485 images containing



Figure 5–5: Typical images from SVT

text in a variety of colors and fonts on many different backgrounds and in various orientations. IC13 is an image dataset used for the ICDAR2013 Robust Reading Competition is almost the same as the dataset IC11. The difference from the ICDAR2011 dataset is revision of ground-truth texts at several images. In addition, a small number of images duplicated over training and test sets were excluded. Accordingly, ICDAR2013(IC13) dataset is a subset of ICDAR2011(IC11) dataset. The number of images of ICDAR2013 dataset is 462, which is comprised of 229 images for the training set and 233 images for the test set. Cropped characters (above) and cropped words (below) are shown in Fig 5–6.



Figure 5-6: Typical character recognition images IC03



Figure 5-7: Typical images from IC13

CHAPTER 6

Appendices

Appendix A - project layout

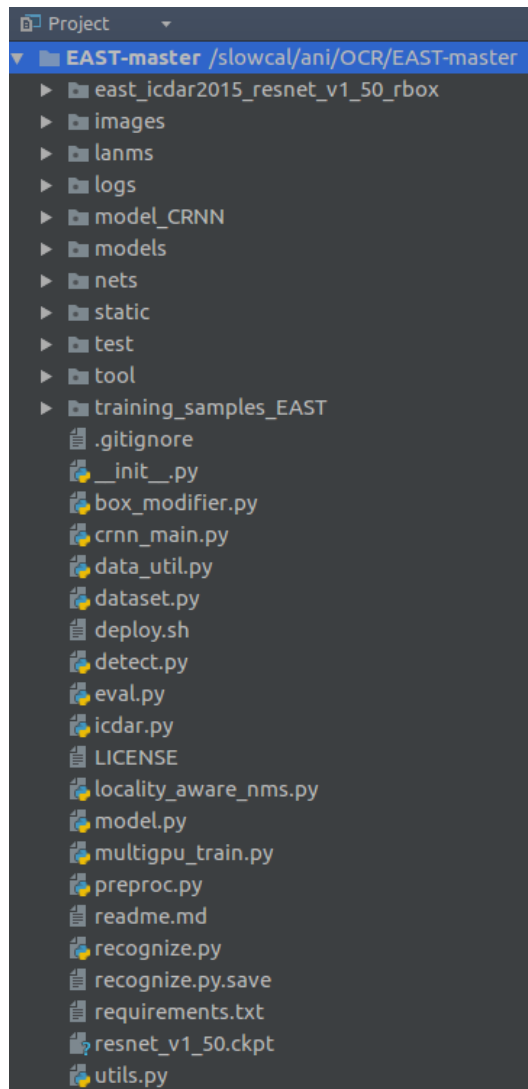


Figure 6-1: Project layout

Appendix B - icdar.py

```
1  # coding:utf-8
2  import glob
3  import csv
4  import cv2
5  import time
6  import os
7  import numpy as np
8  import scipy.optimize
9  import matplotlib.pyplot as plt
10 import matplotlib.patches as Patches
11 from shapely.geometry import Polygon
12
13 import tensorflow as tf
14
15 from data_util import GeneratorEnqueuer
16
17 tf.app.flags.DEFINE_string('training_data_path', '/data/ocr/icdar2015/',
18                           'training dataset to use')
19 tf.app.flags.DEFINE_integer('max_image_large_side', 1280,
20                             'max image size of training')
21 tf.app.flags.DEFINE_integer('max_text_size', 800,
22                             'if the text in the input image is bigger than this, then we
23                             ↳ resize'
24                             'the image according to this')
25 tf.app.flags.DEFINE_integer('min_text_size', 10,
26                             'if the text size is smaller than this, we ignore it during
27                             ↳ training')
28 tf.app.flags.DEFINE_float('min_crop_side_ratio', 0.1,
29                           'when doing random crop from input image, the'
30                           'min length of min(H, W)')
31 tf.app.flags.DEFINE_string('geometry', 'RBOX',
32                             'which geometry to generate, RBOX or QUAD')
33
34 FLAGS = tf.app.flags.FLAGS
35
36 def get_images():
```

```

37     files = []
38     for ext in ['jpg', 'png', 'jpeg', 'JPG']:
39         files.extend(glob.glob(
40             os.path.join(FLAGS.training_data_path, '*.{0}'.format(ext))))
41     return files
42
43
44 def load_annotation(p):
45     '''
46     load annotation from the text file
47     :param p:
48     :return:
49     '''
50     text_polys = []
51     text_tags = []
52     if not os.path.exists(p):
53         return np.array(text_polys, dtype=np.float32)
54     with open(p, 'r') as f:
55         reader = csv.reader(f)
56         for line in reader:
57             label = line[-1]
58             # strip BOM. \ufeff for python3, \xef\xbb\xbf for python2
59             line = [i.strip('\ufeff').strip('\xef\xbb\xbf') for i in line]
60
61             x1, y1, x2, y2, x3, y3, x4, y4 = list(map(float, line[:8]))
62             text_polys.append([[x1, y1], [x2, y2], [x3, y3], [x4, y4]])
63             if label == '*' or label == '###':
64                 text_tags.append(True)
65             else:
66                 text_tags.append(False)
67         return np.array(text_polys, dtype=np.float32), np.array(text_tags, dtype=np.bool)
68
69
70 def polygon_area(poly):
71     '''
72     compute area of a polygon
73     :param poly:
74     :return:
75     '''
76     edge = [
77         (poly[1][0] - poly[0][0]) * (poly[1][1] + poly[0][1]),

```

```

78         (poly[2][0] - poly[1][0]) * (poly[2][1] + poly[1][1]),
79         (poly[3][0] - poly[2][0]) * (poly[3][1] + poly[2][1]),
80         (poly[0][0] - poly[3][0]) * (poly[0][1] + poly[3][1])
81     ]
82     return np.sum(edge)/2.
83
84
85 def check_and_validate_polys(polys, tags, xxx_todo_changeme):
86     '''
87     check so that the text poly is in the same direction,
88     and also filter some invalid polygons
89     :param polys:
90     :param tags:
91     :return:
92     '''
93     (h, w) = xxx_todo_changeme
94     if polys.shape[0] == 0:
95         return polys
96     polys[:, :, 0] = np.clip(polys[:, :, 0], 0, w-1)
97     polys[:, :, 1] = np.clip(polys[:, :, 1], 0, h-1)
98
99     validated_polys = []
100    validated_tags = []
101    for poly, tag in zip(polys, tags):
102        p_area = polygon_area(poly)
103        if abs(p_area) < 1:
104            # print poly
105            print('invalid poly')
106            continue
107        if p_area > 0:
108            print('poly in wrong direction')
109            poly = poly[(0, 3, 2, 1), :]
110            validated_polys.append(poly)
111            validated_tags.append(tag)
112    return np.array(validated_polys), np.array(validated_tags)
113
114
115 def crop_area(im, polys, tags, crop_background=False, max_tries=50):
116     '''
117     make random crop from the input image
118     :param im:

```

```

119     :param polys:
120     :param tags:
121     :param crop_background:
122     :param max_tries:
123     :return:
124     '''
125     h, w, _ = im.shape
126     pad_h = h//10
127     pad_w = w//10
128     h_array = np.zeros((h + pad_h*2), dtype=np.int32)
129     w_array = np.zeros((w + pad_w*2), dtype=np.int32)
130     for poly in polys:
131         poly = np.round(poly, decimals=0).astype(np.int32)
132         minx = np.min(poly[:, 0])
133         maxx = np.max(poly[:, 0])
134         w_array[minx+pad_w:maxx+pad_w] = 1
135         miny = np.min(poly[:, 1])
136         maxy = np.max(poly[:, 1])
137         h_array[miny+pad_h:maxy+pad_h] = 1
138     # ensure the cropped area not across a text
139     h_axis = np.where(h_array == 0)[0]
140     w_axis = np.where(w_array == 0)[0]
141     if len(h_axis) == 0 or len(w_axis) == 0:
142         return im, polys, tags
143     for i in range(max_tries):
144         xx = np.random.choice(w_axis, size=2)
145         xmin = np.min(xx) - pad_w
146         xmax = np.max(xx) - pad_w
147         xmin = np.clip(xmin, 0, w-1)
148         xmax = np.clip(xmax, 0, w-1)
149         yy = np.random.choice(h_axis, size=2)
150         ymin = np.min(yy) - pad_h
151         ymax = np.max(yy) - pad_h
152         ymin = np.clip(ymin, 0, h-1)
153         ymax = np.clip(ymax, 0, h-1)
154         if xmax - xmin < FLAGS.min_crop_side_ratio*w or ymax - ymin <
            ↳ FLAGS.min_crop_side_ratio*h:
155             # area too small
156             continue
157         if polys.shape[0] != 0:
158             poly_axis_in_area = (polys[:, :, 0] >= xmin) & (polys[:, :, 0] <= xmax) \

```

```

159             & (polys[:, :, 1] >= ymin) & (polys[:, :, 1] <= ymax)
160         selected_polys = np.where(np.sum(poly_axis_in_area, axis=1) == 4)[0]
161     else:
162         selected_polys = []
163     if len(selected_polys) == 0:
164         # no text in this area
165         if crop_background:
166             return im[ymin:ymax+1, xmin:xmax+1, :], polys[selected_polys],
167                 ↪ tags[selected_polys]
168         else:
169             continue
170     im = im[ymin:ymax+1, xmin:xmax+1, :]
171     polys = polys[selected_polys]
172     tags = tags[selected_polys]
173     polys[:, :, 0] -= xmin
174     polys[:, :, 1] -= ymin
175     return im, polys, tags
176
177
178
179 def shrink_poly(poly, r):
180     '''
181     fit a poly inside the origin poly, maybe bugs here...
182     used for generate the score map
183     :param poly: the text poly
184     :param r: r in the paper
185     :return: the shrinked poly
186     '''
187     # shrink ratio
188     R = 0.3
189     # find the longer pair
190     if np.linalg.norm(poly[0] - poly[1]) + np.linalg.norm(poly[2] - poly[3]) > \
191        np.linalg.norm(poly[0] - poly[3]) + np.linalg.norm(poly[1] - poly[2]):
192         # first move (p0, p1), (p2, p3), then (p0, p3), (p1, p2)
193         ## p0, p1
194         theta = np.arctan2((poly[1][1] - poly[0][1]), (poly[1][0] - poly[0][0]))
195         poly[0][0] += R * r[0] * np.cos(theta)
196         poly[0][1] += R * r[0] * np.sin(theta)
197         poly[1][0] -= R * r[1] * np.cos(theta)
198         poly[1][1] -= R * r[1] * np.sin(theta)

```

```

199     ## p2, p3
200     theta = np.arctan2((poly[2][1] - poly[3][1]), (poly[2][0] - poly[3][0]))
201     poly[3][0] += R * r[3] * np.cos(theta)
202     poly[3][1] += R * r[3] * np.sin(theta)
203     poly[2][0] -= R * r[2] * np.cos(theta)
204     poly[2][1] -= R * r[2] * np.sin(theta)
205     ## p0, p3
206     theta = np.arctan2((poly[3][0] - poly[0][0]), (poly[3][1] - poly[0][1]))
207     poly[0][0] += R * r[0] * np.sin(theta)
208     poly[0][1] += R * r[0] * np.cos(theta)
209     poly[3][0] -= R * r[3] * np.sin(theta)
210     poly[3][1] -= R * r[3] * np.cos(theta)
211     ## p1, p2
212     theta = np.arctan2((poly[2][0] - poly[1][0]), (poly[2][1] - poly[1][1]))
213     poly[1][0] += R * r[1] * np.sin(theta)
214     poly[1][1] += R * r[1] * np.cos(theta)
215     poly[2][0] -= R * r[2] * np.sin(theta)
216     poly[2][1] -= R * r[2] * np.cos(theta)
217     else:
218         ## p0, p3
219         # print poly
220         theta = np.arctan2((poly[3][0] - poly[0][0]), (poly[3][1] - poly[0][1]))
221         poly[0][0] += R * r[0] * np.sin(theta)
222         poly[0][1] += R * r[0] * np.cos(theta)
223         poly[3][0] -= R * r[3] * np.sin(theta)
224         poly[3][1] -= R * r[3] * np.cos(theta)
225         ## p1, p2
226         theta = np.arctan2((poly[2][0] - poly[1][0]), (poly[2][1] - poly[1][1]))
227         poly[1][0] += R * r[1] * np.sin(theta)
228         poly[1][1] += R * r[1] * np.cos(theta)
229         poly[2][0] -= R * r[2] * np.sin(theta)
230         poly[2][1] -= R * r[2] * np.cos(theta)
231         ## p0, p1
232         theta = np.arctan2((poly[1][1] - poly[0][1]), (poly[1][0] - poly[0][0]))
233         poly[0][0] += R * r[0] * np.cos(theta)
234         poly[0][1] += R * r[0] * np.sin(theta)
235         poly[1][0] -= R * r[1] * np.cos(theta)
236         poly[1][1] -= R * r[1] * np.sin(theta)
237         ## p2, p3
238         theta = np.arctan2((poly[2][1] - poly[3][1]), (poly[2][0] - poly[3][0]))
239         poly[3][0] += R * r[3] * np.cos(theta)

```



```

240         poly[3][1] += R * r[3] * np.sin(theta)
241         poly[2][0] -= R * r[2] * np.cos(theta)
242         poly[2][1] -= R * r[2] * np.sin(theta)
243     return poly
244
245
246 def point_dist_to_line(p1, p2, p3):
247     # compute the distance from p3 to p1-p2
248     return np.linalg.norm(np.cross(p2 - p1, p1 - p3)) / np.linalg.norm(p2 - p1)
249
250
251 def fit_line(p1, p2):
252     # fit a line ax+by+c = 0
253     if p1[0] == p1[1]:
254         return [1., 0., -p1[0]]
255     else:
256         [k, b] = np.polyfit(p1, p2, deg=1)
257         return [k, -1., b]
258
259
260 def line_cross_point(line1, line2):
261     # line1 0= ax+by+c, compute the cross point of line1 and line2
262     if line1[0] != 0 and line1[0] == line2[0]:
263         print('Cross point does not exist')
264         return None
265     if line1[0] == 0 and line2[0] == 0:
266         print('Cross point does not exist')
267         return None
268     if line1[1] == 0:
269         x = -line1[2]
270         y = line2[0] * x + line2[2]
271     elif line2[1] == 0:
272         x = -line2[2]
273         y = line1[0] * x + line1[2]
274     else:
275         k1, _, b1 = line1
276         k2, _, b2 = line2
277         x = -(b1-b2)/(k1-k2)
278         y = k1*x + b1
279     return np.array([x, y], dtype=np.float32)
280

```

```

281
282 def line_verticle(line, point):
283     # get the verticle line from line across point
284     if line[1] == 0:
285         verticle = [0, -1, point[1]]
286     else:
287         if line[0] == 0:
288             verticle = [1, 0, -point[0]]
289         else:
290             verticle = [-1./line[0], -1, point[1] - (-1/line[0] * point[0])]
291     return verticle
292
293
294 def rectangle_from_parallelogram(poly):
295     '''
296     fit a rectangle from a parallelogram
297     :param poly:
298     :return:
299     '''
300     p0, p1, p2, p3 = poly
301     angle_p0 = np.arccos(np.dot(p1-p0, p3-p0)/(np.linalg.norm(p0-p1) *
302     ↪ np.linalg.norm(p3-p0)))
303     if angle_p0 < 0.5 * np.pi:
304         if np.linalg.norm(p0 - p1) > np.linalg.norm(p0-p3):
305             # p0 and p2
306             ## p0
307             p2p3 = fit_line([p2[0], p3[0]], [p2[1], p3[1]])
308             p2p3_verticle = line_verticle(p2p3, p0)
309
310             new_p3 = line_cross_point(p2p3, p2p3_verticle)
311             ## p2
312             p0p1 = fit_line([p0[0], p1[0]], [p0[1], p1[1]])
313             p0p1_verticle = line_verticle(p0p1, p2)
314
315             new_p1 = line_cross_point(p0p1, p0p1_verticle)
316             return np.array([p0, new_p1, p2, new_p3], dtype=np.float32)
317         else:
318             p1p2 = fit_line([p1[0], p2[0]], [p1[1], p2[1]])
319             p1p2_verticle = line_verticle(p1p2, p0)
320
321             new_p1 = line_cross_point(p1p2, p1p2_verticle)

```

```

321         p0p3 = fit_line([p0[0], p3[0]], [p0[1], p3[1]])
322         p0p3_verticle = line_verticle(p0p3, p2)
323
324         new_p3 = line_cross_point(p0p3, p0p3_verticle)
325         return np.array([p0, new_p1, p2, new_p3], dtype=np.float32)
326     else:
327         if np.linalg.norm(p0-p1) > np.linalg.norm(p0-p3):
328             # p1 and p3
329             ## p1
330             p2p3 = fit_line([p2[0], p3[0]], [p2[1], p3[1]])
331             p2p3_verticle = line_verticle(p2p3, p1)
332
333             new_p2 = line_cross_point(p2p3, p2p3_verticle)
334             ## p3
335             p0p1 = fit_line([p0[0], p1[0]], [p0[1], p1[1]])
336             p0p1_verticle = line_verticle(p0p1, p3)
337
338             new_p0 = line_cross_point(p0p1, p0p1_verticle)
339             return np.array([new_p0, p1, new_p2, p3], dtype=np.float32)
340         else:
341             p0p3 = fit_line([p0[0], p3[0]], [p0[1], p3[1]])
342             p0p3_verticle = line_verticle(p0p3, p1)
343
344             new_p0 = line_cross_point(p0p3, p0p3_verticle)
345             p1p2 = fit_line([p1[0], p2[0]], [p1[1], p2[1]])
346             p1p2_verticle = line_verticle(p1p2, p3)
347
348             new_p2 = line_cross_point(p1p2, p1p2_verticle)
349             return np.array([new_p0, p1, new_p2, p3], dtype=np.float32)
350
351
352 def sort_rectangle(poly):
353     # sort the four coordinates of the polygon, points in poly should be sorted clockwise
354     # First find the lowest point
355     p_lowest = np.argmax(poly[:, 1])
356     if np.count_nonzero(poly[:, 1] == poly[p_lowest, 1]) == 2:
357         # X, p0
358         p0_index = np.argmin(np.sum(poly, axis=1))
359         p1_index = (p0_index + 1) % 4
360         p2_index = (p0_index + 2) % 4
361         p3_index = (p0_index + 3) % 4

```

```

362         return poly[[p0_index, p1_index, p2_index, p3_index]], 0.
363     else:
364         #
365         p_lowest_right = (p_lowest - 1) % 4
366         p_lowest_left = (p_lowest + 1) % 4
367         angle = np.arctan(-(poly[p_lowest][1] - poly[p_lowest_right][1])/(poly[p_lowest][0] -
            ↪ poly[p_lowest_right][0]))
368         # assert angle > 0
369         if angle <= 0:
370             print(angle, poly[p_lowest], poly[p_lowest_right])
371         if angle/np.pi * 180 > 45:
372             # p2
373             p2_index = p_lowest
374             p1_index = (p2_index - 1) % 4
375             p0_index = (p2_index - 2) % 4
376             p3_index = (p2_index + 1) % 4
377             return poly[[p0_index, p1_index, p2_index, p3_index]], -(np.pi/2 - angle)
378         else:
379             # p3
380             p3_index = p_lowest
381             p0_index = (p3_index + 1) % 4
382             p1_index = (p3_index + 2) % 4
383             p2_index = (p3_index + 3) % 4
384             return poly[[p0_index, p1_index, p2_index, p3_index]], angle
385
386
387 def restore_rectangle_rbox(origin, geometry):
388     d = geometry[:, :4]
389     angle = geometry[:, 4]
390     # for angle > 0
391     origin_0 = origin[angle >= 0]
392     d_0 = d[angle >= 0]
393     angle_0 = angle[angle >= 0]
394     if origin_0.shape[0] > 0:
395         p = np.array([np.zeros(d_0.shape[0]), -d_0[:, 0] - d_0[:, 2],
396             d_0[:, 1] + d_0[:, 3], -d_0[:, 0] - d_0[:, 2],
397             d_0[:, 1] + d_0[:, 3], np.zeros(d_0.shape[0]),
398             np.zeros(d_0.shape[0]), np.zeros(d_0.shape[0]),
399             d_0[:, 3], -d_0[:, 2]])
400         p = p.transpose((1, 0)).reshape((-1, 5, 2)) # N*5*2
401

```

```

402     rotate_matrix_x = np.array([np.cos(angle_0), np.sin(angle_0)]).transpose((1, 0))
403     rotate_matrix_x = np.repeat(rotate_matrix_x, 5, axis=1).reshape(-1, 2,
    ↪ 5).transpose((0, 2, 1)) # N*5*2
404
405     rotate_matrix_y = np.array([-np.sin(angle_0), np.cos(angle_0)]).transpose((1, 0))
406     rotate_matrix_y = np.repeat(rotate_matrix_y, 5, axis=1).reshape(-1, 2,
    ↪ 5).transpose((0, 2, 1))
407
408     p_rotate_x = np.sum(rotate_matrix_x * p, axis=2)[:, :, np.newaxis] # N*5*1
409     p_rotate_y = np.sum(rotate_matrix_y * p, axis=2)[:, :, np.newaxis] # N*5*1
410
411     p_rotate = np.concatenate([p_rotate_x, p_rotate_y], axis=2) # N*5*2
412
413     p3_in_origin = origin_0 - p_rotate[:, 4, :]
414     new_p0 = p_rotate[:, 0, :] + p3_in_origin # N*2
415     new_p1 = p_rotate[:, 1, :] + p3_in_origin
416     new_p2 = p_rotate[:, 2, :] + p3_in_origin
417     new_p3 = p_rotate[:, 3, :] + p3_in_origin
418
419     new_p_0 = np.concatenate([new_p0[:, np.newaxis, :], new_p1[:, np.newaxis, :],
420                               new_p2[:, np.newaxis, :], new_p3[:, np.newaxis, :]],
    ↪ axis=1) # N*4*2
421
422     else:
423         new_p_0 = np.zeros((0, 4, 2))
424         # for angle < 0
425         origin_1 = origin[angle < 0]
426         d_1 = d[angle < 0]
427         angle_1 = angle[angle < 0]
428         if origin_1.shape[0] > 0:
429             p = np.array([-d_1[:, 1] - d_1[:, 3], -d_1[:, 0] - d_1[:, 2],
430                           np.zeros(d_1.shape[0]), -d_1[:, 0] - d_1[:, 2],
431                           np.zeros(d_1.shape[0]), np.zeros(d_1.shape[0]),
432                           -d_1[:, 1] - d_1[:, 3], np.zeros(d_1.shape[0]),
433                           -d_1[:, 1], -d_1[:, 2]])
434             p = p.transpose((1, 0)).reshape((-1, 5, 2)) # N*5*2
435
436             rotate_matrix_x = np.array([np.cos(-angle_1), -np.sin(-angle_1)]).transpose((1, 0))
437             rotate_matrix_x = np.repeat(rotate_matrix_x, 5, axis=1).reshape(-1, 2,
    ↪ 5).transpose((0, 2, 1)) # N*5*2
438
439             rotate_matrix_y = np.array([np.sin(-angle_1), np.cos(-angle_1)]).transpose((1, 0))

```

```

439     rotate_matrix_y = np.repeat(rotate_matrix_y, 5, axis=1).reshape(-1, 2,
    ↪    5).transpose((0, 2, 1))
440
441     p_rotate_x = np.sum(rotate_matrix_x * p, axis=2)[:, :, np.newaxis] # N*5*1
442     p_rotate_y = np.sum(rotate_matrix_y * p, axis=2)[:, :, np.newaxis] # N*5*1
443
444     p_rotate = np.concatenate([p_rotate_x, p_rotate_y], axis=2) # N*5*2
445
446     p3_in_origin = origin_1 - p_rotate[:, 4, :]
447     new_p0 = p_rotate[:, 0, :] + p3_in_origin # N*2
448     new_p1 = p_rotate[:, 1, :] + p3_in_origin
449     new_p2 = p_rotate[:, 2, :] + p3_in_origin
450     new_p3 = p_rotate[:, 3, :] + p3_in_origin
451
452     new_p_1 = np.concatenate([new_p0[:, np.newaxis, :], new_p1[:, np.newaxis, :],
453                               ↪    new_p2[:, np.newaxis, :], new_p3[:, np.newaxis, :]],
    ↪    axis=1) # N*4*2
454
455     else:
456         new_p_1 = np.zeros((0, 4, 2))
457
458     return np.concatenate([new_p_0, new_p_1])
459
460 def restore_rectangle(origin, geometry):
461     return restore_rectangle_rbox(origin, geometry)
462
463 def generate_rbox(im_size, polys, tags):
464     h, w = im_size
465     poly_mask = np.zeros((h, w), dtype=np.uint8)
466     score_map = np.zeros((h, w), dtype=np.uint8)
467     geo_map = np.zeros((h, w, 5), dtype=np.float32)
468     # mask used during training, to ignore some hard areas
469     training_mask = np.ones((h, w), dtype=np.uint8)
470     for poly_idx, poly_tag in enumerate(zip(polys, tags)):
471         poly = poly_tag[0]
472         tag = poly_tag[1]
473
474         r = [None, None, None, None]
475         for i in range(4):
476             r[i] = min(np.linalg.norm(poly[i] - poly[(i + 1) % 4]),
477                       ↪    np.linalg.norm(poly[i] - poly[(i - 1) % 4]))

```

```

478     # score map
479     shrunked_poly = shrink_poly(poly.copy(), r).astype(np.int32)[np.newaxis, :, :]
480     cv2.fillPoly(score_map, shrunked_poly, 1)
481     cv2.fillPoly(poly_mask, shrunked_poly, poly_idx + 1)
482     # if the poly is too small, then ignore it during training
483     poly_h = min(np.linalg.norm(poly[0] - poly[3]), np.linalg.norm(poly[1] - poly[2]))
484     poly_w = min(np.linalg.norm(poly[0] - poly[1]), np.linalg.norm(poly[2] - poly[3]))
485     if min(poly_h, poly_w) < FLAGS.min_text_size:
486         cv2.fillPoly(training_mask, poly.astype(np.int32)[np.newaxis, :, :], 0)
487     if tag:
488         cv2.fillPoly(training_mask, poly.astype(np.int32)[np.newaxis, :, :], 0)
489
490     xy_in_poly = np.argwhere(poly_mask == (poly_idx + 1))
491     # if geometry == 'RBOX':
492     #
493     fitted_parallelograms = []
494     for i in range(4):
495         p0 = poly[i]
496         p1 = poly[(i + 1) % 4]
497         p2 = poly[(i + 2) % 4]
498         p3 = poly[(i + 3) % 4]
499         edge = fit_line([p0[0], p1[0]], [p0[1], p1[1]])
500         backward_edge = fit_line([p0[0], p3[0]], [p0[1], p3[1]])
501         forward_edge = fit_line([p1[0], p2[0]], [p1[1], p2[1]])
502         if point_dist_to_line(p0, p1, p2) > point_dist_to_line(p0, p1, p3):
503             # p2
504             if edge[1] == 0:
505                 edge_opposite = [1, 0, -p2[0]]
506             else:
507                 edge_opposite = [edge[0], -1, p2[1] - edge[0] * p2[0]]
508         else:
509             # p3
510             if edge[1] == 0:
511                 edge_opposite = [1, 0, -p3[0]]
512             else:
513                 edge_opposite = [edge[0], -1, p3[1] - edge[0] * p3[0]]
514         # move forward edge
515         new_p0 = p0
516         new_p1 = p1
517         new_p2 = p2
518         new_p3 = p3

```

```

519     new_p2 = line_cross_point(forward_edge, edge_opposite)
520     if point_dist_to_line(p1, new_p2, p0) > point_dist_to_line(p1, new_p2, p3):
521         # across p0
522         if forward_edge[1] == 0:
523             forward_opposite = [1, 0, -p0[0]]
524         else:
525             forward_opposite = [forward_edge[0], -1, p0[1] - forward_edge[0] * p0[0]]
526     else:
527         # across p3
528         if forward_edge[1] == 0:
529             forward_opposite = [1, 0, -p3[0]]
530         else:
531             forward_opposite = [forward_edge[0], -1, p3[1] - forward_edge[0] * p3[0]]
532     new_p0 = line_cross_point(forward_opposite, edge)
533     new_p3 = line_cross_point(forward_opposite, edge_opposite)
534     fitted_parallelograms.append([new_p0, new_p1, new_p2, new_p3, new_p0])
535     # or move backward edge
536     new_p0 = p0
537     new_p1 = p1
538     new_p2 = p2
539     new_p3 = p3
540     new_p3 = line_cross_point(backward_edge, edge_opposite)
541     if point_dist_to_line(p0, p3, p1) > point_dist_to_line(p0, p3, p2):
542         # across p1
543         if backward_edge[1] == 0:
544             backward_opposite = [1, 0, -p1[0]]
545         else:
546             backward_opposite = [backward_edge[0], -1, p1[1] - backward_edge[0] *
547                                     ↪ p1[0]]
548     else:
549         # across p2
550         if backward_edge[1] == 0:
551             backward_opposite = [1, 0, -p2[0]]
552         else:
553             backward_opposite = [backward_edge[0], -1, p2[1] - backward_edge[0] *
554                                     ↪ p2[0]]
555     new_p1 = line_cross_point(backward_opposite, edge)
556     new_p2 = line_cross_point(backward_opposite, edge_opposite)
557     fitted_parallelograms.append([new_p0, new_p1, new_p2, new_p3, new_p0])
558     areas = [Polygon(t).area for t in fitted_parallelograms]

```



```

557     parallelogram = np.array(fitted_parallelograms[np.argmin(areas)][:-1],
    ↪     dtype=np.float32)
558     # sort the polygon
559     parallelogram_coord_sum = np.sum(parallelogram, axis=1)
560     min_coord_idx = np.argmin(parallelogram_coord_sum)
561     parallelogram = parallelogram[
562         [min_coord_idx, (min_coord_idx + 1) % 4, (min_coord_idx + 2) % 4, (min_coord_idx
    ↪         + 3) % 4]]
563
564     rectangle = rectangle_from_parallelogram(parallelogram)
565     rectangle, rotate_angle = sort_rectangle(rectangle)
566
567     p0_rect, p1_rect, p2_rect, p3_rect = rectangle
568     for y, x in xy_in_poly:
569         point = np.array([x, y], dtype=np.float32)
570         # top
571         geo_map[y, x, 0] = point_dist_to_line(p0_rect, p1_rect, point)
572         # right
573         geo_map[y, x, 1] = point_dist_to_line(p1_rect, p2_rect, point)
574         # down
575         geo_map[y, x, 2] = point_dist_to_line(p2_rect, p3_rect, point)
576         # left
577         geo_map[y, x, 3] = point_dist_to_line(p3_rect, p0_rect, point)
578         # angle
579         geo_map[y, x, 4] = rotate_angle
580     return score_map, geo_map, training_mask
581
582
583     def generator(input_size=512, batch_size=32,
584                   background_ratio=3./8,
585                   random_scale=np.array([0.5, 1, 2.0, 3.0]),
586                   vis=False):
587         image_list = np.array(get_images())
588         print('{} training images in {}'.format(
589             image_list.shape[0], FLAGS.training_data_path))
590         index = np.arange(0, image_list.shape[0])
591         while True:
592             np.random.shuffle(index)
593             images = []
594             image_fns = []
595             score_maps = []

```

```

596     geo_maps = []
597     training_masks = []
598     for i in index:
599         try:
600             im_fn = image_list[i]
601             im = cv2.imread(im_fn)
602             # print im_fn
603             h, w, _ = im.shape
604             txt_fn = im_fn.replace(os.path.basename(im_fn).split('.')[1], 'txt')
605             if not os.path.exists(txt_fn):
606                 print('text file {} does not exists'.format(txt_fn))
607                 continue
608
609             text_polys, text_tags = load_annoataion(txt_fn)
610
611             text_polys, text_tags = check_and_validate_polys(text_polys, text_tags, (h,
        ↪      w))
612             # if text_polys.shape[0] == 0:
613             #     continue
614             # random scale this image
615             rd_scale = np.random.choice(random_scale)
616             im = cv2.resize(im, dsize=None, fx=rd_scale, fy=rd_scale)
617             text_polys *= rd_scale
618             # print rd_scale
619             # random crop a area from image
620             if np.random.rand() < background_ratio:
621                 # crop background
622                 im, text_polys, text_tags = crop_area(im, text_polys, text_tags,
        ↪      crop_background=True)
623                 if text_polys.shape[0] > 0:
624                     # cannot find background
625                     continue
626                 # pad and resize image
627                 new_h, new_w, _ = im.shape
628                 max_h_w_i = np.max([new_h, new_w, input_size])
629                 im_padded = np.zeros((max_h_w_i, max_h_w_i, 3), dtype=np.uint8)
630                 im_padded[:new_h, :new_w, :] = im.copy()
631                 im = cv2.resize(im_padded, dsize=(input_size, input_size))
632                 score_map = np.zeros((input_size, input_size), dtype=np.uint8)
633                 geo_map_channels = 5 if FLAGS.geometry == 'RBOX' else 8

```

```

634         geo_map = np.zeros((input_size, input_size, geo_map_channels),
        ↪ dtype=np.float32)
635         training_mask = np.ones((input_size, input_size), dtype=np.uint8)
636     else:
637         im, text_polys, text_tags = crop_area(im, text_polys, text_tags,
        ↪ crop_background=False)
638         if text_polys.shape[0] == 0:
639             continue
640         h, w, _ = im.shape
641
642         # pad the image to the training input size or the longer side of image
643         new_h, new_w, _ = im.shape
644         max_h_w_i = np.max([new_h, new_w, input_size])
645         im_padded = np.zeros((max_h_w_i, max_h_w_i, 3), dtype=np.uint8)
646         im_padded[:new_h, :new_w, :] = im.copy()
647         im = im_padded
648         # resize the image to input size
649         new_h, new_w, _ = im.shape
650         resize_h = input_size
651         resize_w = input_size
652         im = cv2.resize(im, dsize=(resize_w, resize_h))
653         resize_ratio_3_x = resize_w/float(new_w)
654         resize_ratio_3_y = resize_h/float(new_h)
655         text_polys[:, :, 0] *= resize_ratio_3_x
656         text_polys[:, :, 1] *= resize_ratio_3_y
657         new_h, new_w, _ = im.shape
658         score_map, geo_map, training_mask = generate_rbox((new_h, new_w),
        ↪ text_polys, text_tags)
659
660     if vis:
661         fig, axs = plt.subplots(3, 2, figsize=(20, 30))
662         # axs[0].imshow(im[:, :, :-1])
663         # axs[0].set_xticks([])
664         # axs[0].set_yticks([])
665         # for poly in text_polys:
666         #     poly_h = min(abs(poly[3, 1] - poly[0, 1]), abs(poly[2, 1] - poly[1,
        ↪ 1]))
667         #     poly_w = min(abs(poly[1, 0] - poly[0, 0]), abs(poly[2, 0] - poly[3,
        ↪ 0]))
668         #     axs[0].add_artist(Patches.Polygon(

```

```

669         #         poly * 4, facecolor='none', edgecolor='green', linewidth=2,
        ↪         linestyle='-', fill=True))
670     #     axs[0].text(poly[0, 0] * 4, poly[0, 1] * 4,
        ↪     '{:.0f}--{:.0f}'.format(poly_h * 4, poly_w * 4),
671     #         color='purple')
672     # axs[1].imshow(score_map)
673     # axs[1].set_xticks([])
674     # axs[1].set_yticks([])
675     axs[0, 0].imshow(im[:, :, :-1])
676     axs[0, 0].set_xticks([])
677     axs[0, 0].set_yticks([])
678     for poly in text_polys:
679         poly_h = min(abs(poly[3, 1] - poly[0, 1]), abs(poly[2, 1] - poly[1,
        ↪         1]))
680         poly_w = min(abs(poly[1, 0] - poly[0, 0]), abs(poly[2, 0] - poly[3,
        ↪         0]))
681         axs[0, 0].add_artist(Patches.Polygon(
682             poly, facecolor='none', edgecolor='green', linewidth=2,
        ↪             linestyle='-', fill=True))
683         axs[0, 0].text(poly[0, 0], poly[0, 1], '{:.0f}--{:.0f}'.format(poly_h,
        ↪         poly_w), color='purple')
684     axs[0, 1].imshow(score_map[:, :, :])
685     axs[0, 1].set_xticks([])
686     axs[0, 1].set_yticks([])
687     axs[1, 0].imshow(geo_map[:, :, 0])
688     axs[1, 0].set_xticks([])
689     axs[1, 0].set_yticks([])
690     axs[1, 1].imshow(geo_map[:, :, 1])
691     axs[1, 1].set_xticks([])
692     axs[1, 1].set_yticks([])
693     axs[2, 0].imshow(geo_map[:, :, 2])
694     axs[2, 0].set_xticks([])
695     axs[2, 0].set_yticks([])
696     axs[2, 1].imshow(training_mask[:, :, :])
697     axs[2, 1].set_xticks([])
698     axs[2, 1].set_yticks([])
699     plt.tight_layout()
700     plt.show()
701     plt.close()
702
703     images.append(im[:, :, :-1].astype(np.float32))

```

```

704         image_fns.append(im_fn)
705         score_maps.append(score_map[:, :, 4, np.newaxis].astype(np.float32))
706         geo_maps.append(geo_map[:, :, 4, :].astype(np.float32))
707         training_masks.append(training_mask[:, :, 4, np.newaxis].astype(np.float32))
708
709         if len(images) == batch_size:
710             yield images, image_fns, score_maps, geo_maps, training_masks
711             images = []
712             image_fns = []
713             score_maps = []
714             geo_maps = []
715             training_masks = []
716         except Exception as e:
717             import traceback
718             traceback.print_exc()
719             continue
720
721
722 def get_batch(num_workers, **kwargs):
723     try:
724         enqueueuer = GeneratorEnqueueuer(generator(**kwargs), use_multiprocessing=True)
725         enqueueuer.start(max_queue_size=24, workers=num_workers)
726         generator_output = None
727         while True:
728             while enqueueuer.is_running():
729                 if not enqueueuer.queue.empty():
730                     generator_output = enqueueuer.queue.get()
731                     break
732                 else:
733                     time.sleep(0.01)
734             yield generator_output
735             generator_output = None
736         finally:
737             if enqueueuer is not None:
738                 enqueueuer.stop()
739
740
741
742 if __name__ == '__main__':
743     pass

```

Appendix C - model.py

```
1  import tensorflow as tf
2  import numpy as np
3  from tensorflow.contrib import slim
4  tf.app.flags.DEFINE_integer('text_scale', 512, '')
5  from nets import resnet_v1
6  FLAGS = tf.app.flags.FLAGS
7
8
9  def unpool(inputs):
10     return tf.image.resize_bilinear(inputs, size=[tf.shape(inputs)[1]*2,
11         ↪ tf.shape(inputs)[2]*2])
12
13
14 def mean_image_subtraction(images, means=[123.68, 116.78, 103.94]):
15     '''
16     image normalization
17     :param images:
18     :param means:
19     :return:
20     '''
21     num_channels = images.get_shape().as_list()[-1]
22     if len(means) != num_channels:
23         raise ValueError('len(means) must match the number of channels')
24     channels = tf.split(axis=3, num_or_size_splits=num_channels, value=images)
25     for i in range(num_channels):
26         channels[i] -= means[i]
27     return tf.concat(axis=3, values=channels)
28
29 def model(images, weight_decay=1e-5, is_training=True):
30     '''
31     define the model, we use slim's implementation of resnet
32     '''
33     images = mean_image_subtraction(images)
34
35     with slim.arg_scope(resnet_v1.resnet_arg_scope(weight_decay=weight_decay)):
36         logits, end_points = resnet_v1.resnet_v1_50(images, is_training=is_training,
37             ↪ scope='resnet_v1_50')
```

```

37
38     with tf.variable_scope('feature_fusion', values=[end_points.values]):
39         batch_norm_params = {
40             'decay': 0.997,
41             'epsilon': 1e-5,
42             'scale': True,
43             'is_training': is_training
44         }
45         with slim.arg_scope([slim.conv2d],
46                             activation_fn=tf.nn.relu,
47                             normalizer_fn=slim.batch_norm,
48                             normalizer_params=batch_norm_params,
49                             weights_regularizer=slim.l2_regularizer(weight_decay)):
50             f = [end_points['pool5'], end_points['pool4'],
51                  end_points['pool3'], end_points['pool2']]
52             for i in range(4):
53                 print('Shape of f_{} {}'.format(i, f[i].shape))
54             g = [None, None, None, None]
55             h = [None, None, None, None]
56             num_outputs = [None, 128, 64, 32]
57             for i in range(4):
58                 if i == 0:
59                     h[i] = f[i]
60                 else:
61                     c1_1 = slim.conv2d(tf.concat([g[i-1], f[i]], axis=-1), num_outputs[i], 1)
62                     h[i] = slim.conv2d(c1_1, num_outputs[i], 3)
63                 if i <= 2:
64                     g[i] = unpool(h[i])
65                 else:
66                     g[i] = slim.conv2d(h[i], num_outputs[i], 3)
67                 print('Shape of h_{} {}, g_{} {}'.format(i, h[i].shape, i, g[i].shape))
68
69         # here we use a slightly different way for regression part,
70         # we first use a sigmoid to limit the regression range, and also
71         # this is do with the angle map
72         F_score = slim.conv2d(g[3], 1, 1, activation_fn=tf.nn.sigmoid,
73                                ↪ normalizer_fn=None)
74         # 4 channel of axis aligned bbox and 1 channel rotation angle
75         geo_map = slim.conv2d(g[3], 4, 1, activation_fn=tf.nn.sigmoid,
76                                ↪ normalizer_fn=None) * FLAGS.text_scale

```

```

75         angle_map = (slim.conv2d(g[3], 1, 1, activation_fn=tf.nn.sigmoid,
    ↪      normalizer_fn=None) - 0.5) * np.pi/2 # angle is between [-45, 45]
76         F_geometry = tf.concat([geo_map, angle_map], axis=-1)
77
78     return F_score, F_geometry
79
80
81 def dice_coefficient(y_true_cls, y_pred_cls,
82                     training_mask):
83     '''
84     dice loss
85     :param y_true_cls:
86     :param y_pred_cls:
87     :param training_mask:
88     :return:
89     '''
90     eps = 1e-5
91     intersection = tf.reduce_sum(y_true_cls * y_pred_cls * training_mask)
92     union = tf.reduce_sum(y_true_cls * training_mask) + tf.reduce_sum(y_pred_cls *
    ↪      training_mask) + eps
93     loss = 1. - (2 * intersection / union)
94     tf.summary.scalar('classification_dice_loss', loss)
95     return loss
96
97
98 def loss(y_true_cls, y_pred_cls,
99         y_true_geo, y_pred_geo,
100         training_mask):
101     '''
102     define the loss used for training, containing two parts,
103     the first part we use dice loss instead of weighted logloss,
104     the second part is the IoU loss defined in the paper
105     :param y_true_cls: ground truth of text
106     :param y_pred_cls: prediction of text
107     :param y_true_geo: ground truth of geometry
108     :param y_pred_geo: prediction of geometry
109     :param training_mask: mask used in training, to ignore some text annotated by ###
110     :return:
111     '''
112     classification_loss = dice_coefficient(y_true_cls, y_pred_cls, training_mask)
113     # scale classification loss to match the iou loss part

```



```

114     classification_loss *= 0.01
115
116     # d1 -> top, d2->right, d3->bottom, d4->left
117     d1_gt, d2_gt, d3_gt, d4_gt, theta_gt = tf.split(value=y_true_geo, num_or_size_splits=5,
118     ↪     axis=3)
119     d1_pred, d2_pred, d3_pred, d4_pred, theta_pred = tf.split(value=y_pred_geo,
120     ↪     num_or_size_splits=5, axis=3)
121     area_gt = (d1_gt + d3_gt) * (d2_gt + d4_gt)
122     area_pred = (d1_pred + d3_pred) * (d2_pred + d4_pred)
123     w_union = tf.minimum(d2_gt, d2_pred) + tf.minimum(d4_gt, d4_pred)
124     h_union = tf.minimum(d1_gt, d1_pred) + tf.minimum(d3_gt, d3_pred)
125     area_intersect = w_union * h_union
126     area_union = area_gt + area_pred - area_intersect
127     L_AABB = -tf.log((area_intersect + 1.0)/(area_union + 1.0))
128     L_theta = 1 - tf.cos(theta_pred - theta_gt)
129     tf.summary.scalar('geometry_AABB', tf.reduce_mean(L_AABB * y_true_cls * training_mask))
130     tf.summary.scalar('geometry_theta', tf.reduce_mean(L_theta * y_true_cls * training_mask))
131     L_g = L_AABB + 20 * L_theta
132
133     return tf.reduce_mean(L_g * y_true_cls * training_mask) + classification_loss

```

Appendix D - multigpu_train.py

```
1  import time
2  import numpy as np
3  import tensorflow as tf
4  from tensorflow.contrib import slim
5
6  tf.app.flags.DEFINE_integer('input_size', 512, '')
7  tf.app.flags.DEFINE_integer('batch_size_per_gpu', 14, '')
8  tf.app.flags.DEFINE_integer('num_readers', 16, '')
9  tf.app.flags.DEFINE_float('learning_rate', 0.0001, '')
10 tf.app.flags.DEFINE_integer('max_steps', 100000, '')
11 tf.app.flags.DEFINE_float('moving_average_decay', 0.997, '')
12 tf.app.flags.DEFINE_string('gpu_list', '1', '')
13 tf.app.flags.DEFINE_string('checkpoint_path', '/tmp/east_resnet_v1_50_rbox/', '')
14 tf.app.flags.DEFINE_boolean('restore', False, 'whether to restore from checkpoint')
15 tf.app.flags.DEFINE_integer('save_checkpoint_steps', 1000, '')
16 tf.app.flags.DEFINE_integer('save_summary_steps', 100, '')
17 tf.app.flags.DEFINE_string('pretrained_model_path', None, '')
18
19 import model
20 import icdar
21
22 FLAGS = tf.app.flags.FLAGS
23
24 gpus = list(range(len(FLAGS.gpu_list.split(','))))
25
26
27 def tower_loss(images, score_maps, geo_maps, training_masks, reuse_variables=None):
28     # Build inference graph
29     with tf.variable_scope(tf.get_variable_scope(), reuse=reuse_variables):
30         f_score, f_geometry = model.model(images, is_training=True)
31
32         model_loss = model.loss(score_maps, f_score,
33                                geo_maps, f_geometry,
34                                training_masks)
35         total_loss = tf.add_n([model_loss] +
36                                ↪ tf.get_collection(tf.GraphKeys.REGULARIZATION_LOSSES))
37
38     # add summary
```

```

38     if reuse_variables is None:
39         tf.summary.image('input', images)
40         tf.summary.image('score_map', score_maps)
41         tf.summary.image('score_map_pred', f_score * 255)
42         tf.summary.image('geo_map_0', geo_maps[:, :, :, 0:1])
43         tf.summary.image('geo_map_0_pred', f_geometry[:, :, :, 0:1])
44         tf.summary.image('training_masks', training_masks)
45         tf.summary.scalar('model_loss', model_loss)
46         tf.summary.scalar('total_loss', total_loss)
47
48     return total_loss, model_loss
49
50
51 def average_gradients(tower_grads):
52     average_grads = []
53     for grad_and_vars in zip(*tower_grads):
54         grads = []
55         for g, _ in grad_and_vars:
56             expanded_g = tf.expand_dims(g, 0)
57             grads.append(expanded_g)
58
59         grad = tf.concat(grads, 0)
60         grad = tf.reduce_mean(grad, 0)
61
62         v = grad_and_vars[0][1]
63         grad_and_var = (grad, v)
64         average_grads.append(grad_and_var)
65
66     return average_grads
67
68
69 def main(argv=None):
70     import os
71     os.environ['CUDA_VISIBLE_DEVICES'] = FLAGS.gpu_list
72     if not tf.gfile.Exists(FLAGS.checkpoint_path):
73         tf.gfile.MkDir(FLAGS.checkpoint_path)
74     else:
75         if not FLAGS.restore:
76             tf.gfile.DeleteRecursively(FLAGS.checkpoint_path)
77             tf.gfile.MkDir(FLAGS.checkpoint_path)
78

```

```

79     input_images = tf.placeholder(tf.float32, shape=[None, None, None, 3],
    ↪     name='input_images')
80     input_score_maps = tf.placeholder(tf.float32, shape=[None, None, None, 1],
    ↪     name='input_score_maps')
81     if FLAGS.geometry == 'RBOX':
82         input_geo_maps = tf.placeholder(tf.float32, shape=[None, None, None, 5],
    ↪         name='input_geo_maps')
83     else:
84         input_geo_maps = tf.placeholder(tf.float32, shape=[None, None, None, 8],
    ↪         name='input_geo_maps')
85     input_training_masks = tf.placeholder(tf.float32, shape=[None, None, None, 1],
    ↪     name='input_training_masks')
86
87     global_step = tf.get_variable('global_step', [], initializer=tf.constant_initializer(0),
    ↪     trainable=False)
88     learning_rate = tf.train.exponential_decay(FLAGS.learning_rate, global_step,
    ↪     decay_steps=10000, decay_rate=0.94, staircase=True)
89     # add summary
90     tf.summary.scalar('learning_rate', learning_rate)
91     opt = tf.train.AdamOptimizer(learning_rate)
92     # opt = tf.train.MomentumOptimizer(learning_rate, 0.9)
93
94
95     # split
96     input_images_split = tf.split(input_images, len(gpus))
97     input_score_maps_split = tf.split(input_score_maps, len(gpus))
98     input_geo_maps_split = tf.split(input_geo_maps, len(gpus))
99     input_training_masks_split = tf.split(input_training_masks, len(gpus))
100
101     tower_grads = []
102     reuse_variables = None
103     for i, gpu_id in enumerate(gpus):
104         with tf.device('/gpu:%d' % gpu_id):
105             with tf.name_scope('model_%d' % gpu_id) as scope:
106                 iis = input_images_split[i]
107                 isms = input_score_maps_split[i]
108                 igms = input_geo_maps_split[i]
109                 itms = input_training_masks_split[i]
110                 total_loss, model_loss = tower_loss(iis, isms, igms, itms, reuse_variables)
111                 batch_norm_updates_op = tf.group(*tf.get_collection(tf.GraphKeys.UPDATE_OPS,
    ↪                 scope))

```

```

112         reuse_variables = True
113
114         grads = opt.compute_gradients(total_loss)
115         tower_grads.append(grads)
116
117     grads = average_gradients(tower_grads)
118     apply_gradient_op = opt.apply_gradients(grads, global_step=global_step)
119
120     summary_op = tf.summary.merge_all()
121     # save moving average
122     variable_averages = tf.train.ExponentialMovingAverage(
123         FLAGS.moving_average_decay, global_step)
124     variables_averages_op = variable_averages.apply(tf.trainable_variables())
125     # batch norm updates
126     with tf.control_dependencies([variables_averages_op, apply_gradient_op,
127         ↪ batch_norm_updates_op]):
128         train_op = tf.no_op(name='train_op')
129
130     saver = tf.train.Saver(tf.global_variables())
131     summary_writer = tf.summary.FileWriter(FLAGS.checkpoint_path, tf.get_default_graph())
132
133     init = tf.global_variables_initializer()
134
135     if FLAGS.pretrained_model_path is not None:
136         variable_restore_op = slim.assign_from_checkpoint_fn(FLAGS.pretrained_model_path,
137             ↪ slim.get_trainable_variables(),
138                                     ignore_missing_vars=True)
139     with tf.Session(config=tf.ConfigProto(allow_soft_placement=True)) as sess:
140         if FLAGS.restore:
141             print('continue training from previous checkpoint')
142             ckpt = tf.train.latest_checkpoint(FLAGS.checkpoint_path)
143             saver.restore(sess, ckpt)
144         else:
145             sess.run(init)
146             if FLAGS.pretrained_model_path is not None:
147                 variable_restore_op(sess)
148
149     data_generator = icdar.get_batch(num_workers=FLAGS.num_readers,
150                                     input_size=FLAGS.input_size,
151                                     batch_size=FLAGS.batch_size_per_gpu * len(gpus))

```

```

151     start = time.time()
152     for step in range(FLAGS.max_steps):
153         data = next(data_generator)
154         ml, tl, _ = sess.run([model_loss, total_loss, train_op], feed_dict={input_images:
↪     data[0],
155
↪     input_score_maps:
↪     data[2],
156
↪     input_geo_maps:
↪     data[3],
157
↪     input_training_masks:
↪     data[4]})
158     if np.isnan(tl):
159         print('Loss diverged, stop training')
160         break
161
162     if step % 10 == 0:
163         avg_time_per_step = (time.time() - start)/10
164         avg_examples_per_second = (10 * FLAGS.batch_size_per_gpu *
↪     len(gpus))/(time.time() - start)
165         start = time.time()
166         print('Step {:06d}, model loss {:.4f}, total loss {:.4f}, {:.2f}
↪     seconds/step, {:.2f} examples/second'.format(
167             step, ml, tl, avg_time_per_step, avg_examples_per_second))
168
169     if step % FLAGS.save_checkpoint_steps == 0:
170         saver.save(sess, FLAGS.checkpoint_path + 'model.ckpt',
↪     global_step=global_step)
171
172     if step % FLAGS.save_summary_steps == 0:
173         _, tl, summary_str = sess.run([train_op, total_loss, summary_op],
↪     feed_dict={input_images: data[0],
174
↪     input_score_maps:
↪     data[2],
↪     input_geo_maps:
↪     data[3],

```

176

↪ input_training_masks

↪ data[4])

177 summary_writer.add_summary(summary_str, global_step=step)

178

179 if __name__ == '__main__':

180 tf.app.run()

Appendix E - eval.py

```
1  # encode=utf-8
2  # text detection
3  import time
4  import os
5  import cv2
6  import numpy as np
7  import tensorflow as tf
8  import model
9  from preproc import get_images, resize_image
10 from detect import detect
11 from recognize import recognize
12 import box_modifier
13 import enchant
14 import locality_aware_nms as nms_locality
15
16 tf.app.flags.DEFINE_string('test_data_path', './images/', '')
17 tf.app.flags.DEFINE_string('model_path', './model_CRNN/', '')
18 tf.app.flags.DEFINE_string('gpu_list', '0', '')
19 tf.app.flags.DEFINE_string('checkpoint_path', './east_icdar2015_resnet_v1_50_rbox/', '')
20 tf.app.flags.DEFINE_string('output_dir', './images/', '')
21
22 FLAGS = tf.app.flags.FLAGS
23 dictEN = enchant.Dict("en_US")
24 dictFR = enchant.Dict("fr_FR")
25
26
27 def dict_check(sim_pred):
28     if dictEN.check(sim_pred) is True:
29         return sim_pred + '_EN'
30     elif dictFR.check(sim_pred) is True:
31         return sim_pred + '_FR'
32
33
34 def main(argv=None):
35     os.environ['CUDA_VISIBLE_DEVICES'] = FLAGS.gpu_list
36     try:
37         os.makedirs(FLAGS.output_dir)
38     except OSError as e:
```



```

39         if e.errno != 17:
40             raise
41
42     with tf.get_default_graph().as_default():
43         input_images = tf.placeholder(tf.float32, shape=[None, None, None, 3],
44             ↪ name='input_images')
45
46         global_step = tf.get_variable('global_step', [],
47             ↪ initializer=tf.constant_initializer(0), trainable=False)
48
49         f_score, f_geometry = model.model(input_images, is_training=False)
50
51         variable_averages = tf.train.ExponentialMovingAverage(0.997, global_step)
52         saver = tf.train.Saver(variable_averages.variables_to_restore())
53
54     with tf.Session(config=tf.ConfigProto(allow_soft_placement=True)) as sess:
55         ckpt_state = tf.train.get_checkpoint_state(FLAGS.checkpoint_path)
56         model_path = os.path.join(FLAGS.checkpoint_path,
57             ↪ os.path.basename(ckpt_state.model_checkpoint_path))
58         print('Restore from {}'.format(model_path))
59         saver.restore(sess, model_path)
60
61     im_fn_list = get_images()
62     for im_fn in im_fn_list:
63         im = cv2.imread(im_fn)[: , :, :-1]
64         start_time = time.time()
65         im_resized, (ratio_h, ratio_w) = resize_image(im)
66
67         timer = {'net': 0, 'restore': 0, 'nms': 0}
68         start = time.time()
69         score, geometry = sess.run([f_score, f_geometry], feed_dict={input_images:
70             ↪ [im_resized]})
71         timer['net'] = time.time() - start
72
73         boxes, timer = detect(score_map=score, geo_map=geometry, timer=timer)
74         print('{} : net {:.0f}ms, restore {:.0f}ms, nms {:.0f}ms'.format(
75             im_fn, timer['net'] * 1000, timer['restore'] * 1000, timer['nms'] *
76             ↪ 1000))
77
78     if boxes is not None:
79         boxes = boxes[:, :8].reshape((-1, 4, 2))
80         boxes[:, :, 0] /= ratio_w

```

```

75         boxes[:, :, 1] /= ratio_h
76
77         duration = time.time() - start_time
78         print('[timing] {}'.format(duration))
79
80         text = []
81         box_coor = []
82         for box in boxes:
83             # to avoid submitting errors
84             box = box_modifier.sort_poly(box.astype(np.int32))
85             if np.linalg.norm(box[0] - box[1]) < 5 or np.linalg.norm(box[3] - box[0])
↳             < 5:
86                 continue
87             # cropping the boxes
88             for _ in box:
89                 _box = im[(np.amin(box, 0)[1]):(np.amax(box, 0)[1]),
90                          (np.amin(box, 0)[0]):(np.amax(box, 0)[0])]
91                 _box_coor = [(np.amin(box, 0)[0]), (np.amax(box, 0)[0]),
↳                 (np.amin(box, 0)[1]),
92                          (np.amax(box, 0)[1])]
93                 if recognize(_box) is not None:
94                     text.append(recognize(_box))
95                     box_coor.append(_box_coor)
96
97         matches = box_modifier.match_check(box_coor)
98
99         for i in range(len(matches)):
100             text[matches[i][0]] = text[matches[i][0]] + ' ' + text[matches[i][1]]
101
102         box_modifier.comb_boxes(box_coor, matches)
103         box_location = box_modifier.get_box_location(box_coor, im)
104
105         return text, box_location
106
107
108 if __name__ == '__main__':
109     tf.app.run()

```

Appendix F - preproc.py

```
1  import cv2
2  import os
3  import tensorflow as tf
4
5  FLAGS = tf.app.flags.FLAGS
6
7  def get_images():
8      '''
9      find image files in test data path
10     :return: list of files found
11     '''
12     files = []
13     exts = ['.jpg', '.png', '.jpeg', '.JPG']
14     for parent, dirnames, filenames in os.walk(FLAGS.test_data_path):
15         for filename in filenames:
16             for ext in exts:
17                 if filename.endswith(ext):
18                     files.append(os.path.join(parent, filename))
19             break
20     print('Found {} images'.format(len(files)))
21     return files
22
23
24 def resize_image(im, max_side_len=2400):
25     '''
26     resize image to a size multiple of 32 which is required by the network
27     :param im: the resized image
28     :param max_side_len: limit of max image size to avoid out of memory in gpu
29     :return: the resized image and the resize ratio
30     '''
31     h, w, _ = im.shape
32
33     resize_w = w
34     resize_h = h
35
36     # limit the max side
37     if max(resize_h, resize_w) > max_side_len:
```

```

38         ratio = float(max_side_len) / resize_h if resize_h > resize_w else
           ↪ float(max_side_len) / resize_w
39     else:
40         ratio = 1.
41     resize_h = int(resize_h * ratio)
42     resize_w = int(resize_w * ratio)
43
44     resize_h = resize_h if resize_h % 32 == 0 else (resize_h // 32 - 1) * 32
45     resize_w = resize_w if resize_w % 32 == 0 else (resize_w // 32 - 1) * 32
46     im = cv2.resize(im, (int(resize_w), int(resize_h)))
47
48     ratio_h = resize_h / float(h)
49     ratio_w = resize_w / float(w)
50
51     return im, (ratio_h, ratio_w)

```

Appendix G - box_modifier.py

```
1  import itertools
2  import numpy as np
3
4
5  def sort_poly(p):
6      min_axis = np.argmin(np.sum(p, axis=1))
7      p = p[[min_axis, (min_axis + 1) % 4, (min_axis + 2) % 4, (min_axis + 3) % 4]]
8      if abs(p[0, 0] - p[1, 0]) > abs(p[0, 1] - p[1, 1]):
9          return p
10     else:
11         return p[[0, 3, 2, 1]]
12
13
14     # look for matches itertools permutations
15     def match_check(box_coor):
16         comb_box = []
17         for i, j in itertools.permutations(box_coor, 2):
18             if (j[0] <= i[1] and j[2] <= i[3]) and (i[0] <= j[1] and i[2] <= j[3]):
19                 comb_box.append((box_coor.index(i), box_coor.index(j)))
20         matches = list(set(tuple(sorted(l)) for l in comb_box))
21         for k in range(len(matches)):
22             if box_coor[matches[k][0]][2] <= box_coor[matches[k][1]][3]:
23                 matches[k] = (matches[k][1], matches[k][0])
24             elif box_coor[matches[k][0]][0] <= box_coor[matches[k][1]][1]:
25                 matches[k] = (matches[k][1], matches[k][0])
26             else:
27                 pass
28         return matches
29
30
31     def comb_boxes(box_coor, matches):
32         for i in matches:
33             box_coor[i[0]] = [(np.minimum(box_coor[i[0]][0], box_coor[i[1]][0])),
34                               (np.maximum(box_coor[i[0]][1], box_coor[i[1]][1])),
35                               (np.minimum(box_coor[i[0]][2], box_coor[i[1]][2])),
36                               (np.maximum(box_coor[i[0]][3], box_coor[i[1]][3]))]
37
38         return box_coor
```

```

39
40
41 def get_box_location(box_coor, im):
42     im_size = np.shape(im)
43     box_location = []
44     for i in range(len(box_coor)):
45         box_x = (box_coor[i][0] + box_coor[i][1]) / 2
46         box_y = (box_coor[i][2] + box_coor[i][3]) / 2
47         print(float(box_x/im_size[0]), float(box_y/im_size[1]))
48         box_location.append([box_y / im_size[0], box_x / im_size[1]])
49
50     return box_location

```

Appendix H - crnn_main.py

```
1  from __future__ import print_function
2  import argparse
3  import random
4  import torch
5  import torch.backends.cudnn as cudnn
6  import torch.optim as optim
7  import torch.utils.data
8  from torch.autograd import Variable
9  import numpy as np
10 from warpctc_pytorch import CTCLoss
11 import os
12 import utils
13 import dataset
14
15 import models.crnn as crnn
16
17 parser = argparse.ArgumentParser()
18 parser.add_argument('--trainroot', required=True, help='path to dataset')
19 parser.add_argument('--valroot', required=True, help='path to dataset')
20 parser.add_argument('--workers', type=int, help='number of data loading workers', default=2)
21 parser.add_argument('--batchSize', type=int, default=64, help='input batch size')
22 parser.add_argument('--imgH', type=int, default=32, help='the height of the input image to
    ↳ network')
23 parser.add_argument('--imgW', type=int, default=100, help='the width of the input image to
    ↳ network')
24 parser.add_argument('--nh', type=int, default=256, help='size of the lstm hidden state')
25 parser.add_argument('--niter', type=int, default=25, help='number of epochs to train for')
26 parser.add_argument('--lr', type=float, default=0.01, help='learning rate for Critic,
    ↳ default=0.00005')
27 parser.add_argument('--beta1', type=float, default=0.5, help='beta1 for adam. default=0.5')
28 parser.add_argument('--cuda', action='store_true', help='enables cuda')
29 parser.add_argument('--ngpu', type=int, default=1, help='number of GPUs to use')
30 parser.add_argument('--crnn', default='', help="path to crnn (to continue training)")
31 parser.add_argument('--alphabet', type=str, default='0123456789abcdefghijklmnopqrstuvwxyz')
32 parser.add_argument('--experiment', default=None, help='Where to store samples and models')
33 parser.add_argument('--displayInterval', type=int, default=500, help='Interval to be
    ↳ displayed')
```

```

34 parser.add_argument('--n_test_disp', type=int, default=10, help='Number of samples to display
    ↳ when test')
35 parser.add_argument('--valInterval', type=int, default=500, help='Interval to be displayed')
36 parser.add_argument('--saveInterval', type=int, default=500, help='Interval to be displayed')
37 parser.add_argument('--adam', action='store_true', help='Whether to use adam (default is
    ↳ rmsprop)')
38 parser.add_argument('--adadelat', action='store_true', help='Whether to use adadelat (default
    ↳ is rmsprop)')
39 parser.add_argument('--keep_ratio', action='store_true', help='whether to keep ratio for
    ↳ image resize')
40 parser.add_argument('--random_sample', action='store_true', help='whether to sample the
    ↳ dataset with random sampler')
41 opt = parser.parse_args()
42 print(opt)
43
44 if opt.experiment is None:
45     opt.experiment = 'expr'
46 os.system('mkdir {0}'.format(opt.experiment))
47
48 opt.manualSeed = random.randint(1, 10000) # fix seed
49 print("Random Seed: ", opt.manualSeed)
50 random.seed(opt.manualSeed)
51 np.random.seed(opt.manualSeed)
52 torch.manual_seed(opt.manualSeed)
53
54 cudnn.benchmark = True
55
56 if torch.cuda.is_available() and not opt.cuda:
57     print("WARNING: You have a CUDA device, so you should probably run with --cuda")
58
59 train_dataset = dataset.lmdbDataset(root=opt.trainroot)
60 assert train_dataset
61 if not opt.random_sample:
62     sampler = dataset.randomSequentialSampler(train_dataset, opt.batchSize)
63 else:
64     sampler = None
65 train_loader = torch.utils.data.DataLoader(
66     train_dataset, batch_size=opt.batchSize,
67     shuffle=True, sampler=sampler,
68     num_workers=int(opt.workers),
69     collate_fn=dataset.alignCollate(imgH=opt.imgH, imgW=opt.imgW, keep_ratio=opt.keep_ratio))

```



```

70  test_dataset = dataset.lmdbDataset(
71      root=opt.valroot, transform=dataset.resizeNormalize((100, 32)))
72
73  nclass = len(opt.alphabet) + 1
74  nc = 1
75
76  converter = utils.strLabelConverter(opt.alphabet)
77  criterion = CTCLoss()
78
79
80  # custom weights initialization called on crnn
81  def weights_init(m):
82      classname = m.__class__.__name__
83      if classname.find('Conv') != -1:
84          m.weight.data.normal_(0.0, 0.02)
85      elif classname.find('BatchNorm') != -1:
86          m.weight.data.normal_(1.0, 0.02)
87          m.bias.data.fill_(0)
88
89
90  crnn = crnn.CRNN(opt.imgH, nc, nclass, opt.nh)
91  crnn.apply(weights_init)
92  if opt.crnn != '':
93      print('loading pretrained model from %s' % opt.crnn)
94      crnn.load_state_dict(torch.load(opt.crnn))
95  print(crnn)
96
97  image = torch.FloatTensor(opt.batchSize, 3, opt.imgH, opt.imgH)
98  text = torch.IntTensor(opt.batchSize * 5)
99  length = torch.IntTensor(opt.batchSize)
100
101  if opt.cuda:
102      crnn.cuda()
103      crnn = torch.nn.DataParallel(crnn, device_ids=range(opt.ngpu))
104      image = image.cuda()
105      criterion = criterion.cuda()
106
107  image = Variable(image)
108  text = Variable(text)
109  length = Variable(length)
110

```

```

111  # loss averager
112  loss_avg = utils.averager()
113
114  # setup optimizer
115  if opt.adam:
116      optimizer = optim.Adam(crnn.parameters(), lr=opt.lr,
117                              betas=(opt.betas, 0.999))
118  elif opt.adadelta:
119      optimizer = optim.Adadelta(crnn.parameters(), lr=opt.lr)
120  else:
121      optimizer = optim.RMSprop(crnn.parameters(), lr=opt.lr)
122
123
124  def val(net, dataset, criterion, max_iter=100):
125      print('Start val')
126
127      for p in crnn.parameters():
128          p.requires_grad = False
129
130      net.eval()
131      data_loader = torch.utils.data.DataLoader(
132          dataset, shuffle=True, batch_size=opt.batchSize, num_workers=int(opt.workers))
133      val_iter = iter(data_loader)
134
135      i = 0
136      n_correct = 0
137      loss_avg = utils.averager()
138
139      max_iter = min(max_iter, len(data_loader))
140      for i in range(max_iter):
141          data = val_iter.next()
142          i += 1
143          cpu_images, cpu_texts = data
144          batch_size = cpu_images.size(0)
145          utils.loadData(image, cpu_images)
146          t, l = converter.encode(cpu_texts)
147          utils.loadData(text, t)
148          utils.loadData(length, l)
149
150          preds = crnn(image)
151          preds_size = Variable(torch.IntTensor([preds.size(0)] * batch_size))

```

```

152         cost = criterion(preds, text, preds_size, length) / batch_size
153         loss_avg.add(cost)
154
155         _, preds = preds.max(2)
156         preds = preds.squeeze(2)
157         preds = preds.transpose(1, 0).contiguous().view(-1)
158         sim_preds = converter.decode(preds.data, preds_size.data, raw=False)
159         for pred, target in zip(sim_preds, cpu_texts):
160             if pred == target.lower():
161                 n_correct += 1
162
163         raw_preds = converter.decode(preds.data, preds_size.data, raw=True)[:opt.n_test_disp]
164         for raw_pred, pred, gt in zip(raw_preds, sim_preds, cpu_texts):
165             print('%-20s => %-20s, gt: %-20s' % (raw_pred, pred, gt))
166
167         accuracy = n_correct / float(max_iter * opt.batchSize)
168         print('Test loss: %f, accuray: %f' % (loss_avg.val(), accuracy))
169
170
171     def trainBatch(net, criterion, optimizer):
172         data = train_iter.next()
173         cpu_images, cpu_texts = data
174         batch_size = cpu_images.size(0)
175         utils.loadData(image, cpu_images)
176         t, l = converter.encode(cpu_texts)
177         utils.loadData(text, t)
178         utils.loadData(length, l)
179
180         preds = crnn(image)
181         preds_size = Variable(torch.IntTensor([preds.size(0)] * batch_size))
182         cost = criterion(preds, text, preds_size, length) / batch_size
183         crnn.zero_grad()
184         cost.backward()
185         optimizer.step()
186         return cost
187
188
189     for epoch in range(opt.niter):
190         train_iter = iter(train_loader)
191         i = 0
192         while i < len(train_loader):

```

```

193     for p in crnn.parameters():
194         p.requires_grad = True
195     crnn.train()
196
197     cost = trainBatch(crnn, criterion, optimizer)
198     loss_avg.add(cost)
199     i += 1
200
201     if i % opt.displayInterval == 0:
202         print('%d/%d [%d/%d] Loss: %f' %
203               (epoch, opt.niter, i, len(train_loader), loss_avg.val()))
204         loss_avg.reset()
205
206     if i % opt.valInterval == 0:
207         val(crnn, test_dataset, criterion)
208
209     # do checkpointing
210     if i % opt.saveInterval == 0:
211         torch.save(
212             crnn.state_dict(), '{0}/netCRNN_{1}_{2}.pth'.format(opt.experiment, epoch,
213                               ↵ i))

```

Appendix I - adaptor.cpp

```
1  #include "pybind11/pybind11.h"
2  #include "pybind11/numpy.h"
3  #include "pybind11/stl.h"
4  #include "pybind11/stl_bind.h"
5
6  #include "lanms.h"
7
8  namespace py = pybind11;
9
10
11 namespace lanms_adaptor {
12
13     std::vector<std::vector<float>> polys2floats(const std::vector<lanms::Polygon>
14 ↪      &polys) {
15         std::vector<std::vector<float>> ret;
16         for (size_t i = 0; i < polys.size(); i++) {
17             auto &p = polys[i];
18             auto &poly = p.poly;
19             ret.emplace_back(std::vector<float>{
20                 float(poly[0].X), float(poly[0].Y),
21                 float(poly[1].X), float(poly[1].Y),
22                 float(poly[2].X), float(poly[2].Y),
23                 float(poly[3].X), float(poly[3].Y),
24                 float(p.score),
25             });
26         }
27         return ret;
28     }
29
30
31     /**
32     *
33     * \param quad_n9 an n-by-9 numpy array, where first 8 numbers denote the
34     *               quadrangle, and the last one is the score
35     * \param iou_threshold two quadrangles with iou score above this threshold
36     *               will be merged
37     *
```

```

38     * \return an n-by-9 numpy array, the merged quadrangles
39     */
40     std::vector<std::vector<float>> merge_quadrangle_n9(
41         py::array_t<float, py::array::c_style | py::array::forcecast>
42         ↪ quad_n9,
43         float iou_threshold) {
44         auto pbuf = quad_n9.request();
45         if (pbuf.ndim != 2 || pbuf.shape[1] != 9)
46             throw std::runtime_error("quadrangles must have a shape of (n, 9)");
47         auto n = pbuf.shape[0];
48         auto ptr = static_cast<float*>(pbuf.ptr);
49         return polys2floats(lanms::merge_quadrangle_n9(ptr, n, iou_threshold));
50     }
51 }
52
53 PYBIND11_PLUGIN(adaptor) {
54     py::module m("adaptor", "NMS");
55
56     m.def("merge_quadrangle_n9", &lanms_adaptor::merge_quadrangle_n9,
57         "merge quadrangles");
58
59     return m.ptr();
60 }

```

Appendix J - recognize.py

```
1  # coding=utf-8
2
3  import models.crrn as crrn
4  import torch
5  from torch.autograd import Variable
6  import utils
7  import dataset
8  import enchant
9  import unicodedata
10 from PIL import Image
11
12 model_path = './crrn.pth'
13 alphabet = '0123456789ABCDEFGHIJKLMNOPQRSTUVWXYZabcdefghijklmnopqrstuvwxyz'
14 dictEN = enchant.Dict("en_US")
15 dictFR = enchant.Dict("fr_FR")
16
17
18 #
19     ↪ https://stackoverflow.com/questions/517923/what-is-the-best-way-to-remove-accent-in-a-python-unicode-string
20 def remove_accents(input_str):
21     nfkd_form = unicodedata.normalize('NFKD', input_str)
22     only_ascii = nfkd_form.encode('ASCII', 'ignore')
23     return only_ascii.decode('utf-8')
24
25 def recognize(box_coord):
26     '''
27
28     '''
29     model = crrn.CRNN(32, 1, 37, 256)
30     if torch.cuda.is_available():
31         model = model.cuda()
32     # print('loading pretrained model from %s' % model_path)
33     model.load_state_dict(torch.load(model_path))
34
35     converter = utils.strLabelConverter(alphabet)
36     transformer = dataset.resizeNormalize((100, 32))
37
```

```

38     image = Image.fromarray(box_coor).convert('L')
39     image = transformer(image)
40     if torch.cuda.is_available():
41         image = image.cuda()
42     image = image.view(1, *image.size())
43     image = Variable(image)
44     model.eval()
45     preds = model(image)
46     _, preds = preds.max(2)
47     # preds = preds.squeeze(2)
48     preds = preds.transpose(1, 0).contiguous().view(-1)
49     preds_size = Variable(torch.IntTensor([preds.size(0)]))
50     raw_pred = converter.decode(preds.data, preds_size.data, raw=True)
51     sim_pred = converter.decode(preds.data, preds_size.data, raw=False)
52     if len(sim_pred) > int(1):
53         if dictEN.check(sim_pred) is True:
54             return sim_pred
55         elif dictFR.check(sim_pred) is True:
56             return sim_pred
57         else:
58             sugg = dictFR.suggest(sim_pred)
59             for i in range(0, len(sugg)):
60                 sugg[i] = remove_accents(sugg[i])
61             if sim_pred in sugg:
62                 return sugg[0]
63             else:
64                 pass
65     else:
66         pass

```


Appendix K - crnn.py

```
1  import torch.nn as nn
2
3  class BidirectionalLSTM(nn.Module):
4
5      def __init__(self, nIn, nHidden, nOut):
6          super(BidirectionalLSTM, self).__init__()
7
8          self.rnn = nn.LSTM(nIn, nHidden, bidirectional=True)
9          self.embedding = nn.Linear(nHidden * 2, nOut)
10
11      def forward(self, input):
12          recurrent, _ = self.rnn(input)
13          T, b, h = recurrent.size()
14          t_rec = recurrent.view(T * b, h)
15
16          output = self.embedding(t_rec)  # [T * b, nOut]
17          output = output.view(T, b, -1)
18
19          return output
20
21
22  class CRNN(nn.Module):
23
24      def __init__(self, imgH, nc, nclass, nh, n_rnn=2, leakyRelu=False):
25          super(CRNN, self).__init__()
26          assert imgH % 16 == 0, 'imgH has to be a multiple of 16'
27
28          ks = [3, 3, 3, 3, 3, 3, 2]
29          ps = [1, 1, 1, 1, 1, 1, 0]
30          ss = [1, 1, 1, 1, 1, 1, 1]
31          nm = [64, 128, 256, 256, 512, 512, 512]
32
33          cnn = nn.Sequential()
34
35          def convRelu(i, batchNormalization=False):
36              nIn = nc if i == 0 else nm[i - 1]
37              nOut = nm[i]
38              cnn.add_module('conv{0}'.format(i),
```

```

39         nn.Conv2d(nIn, nOut, ks[i], ss[i], ps[i]))
40     if batchNormalization:
41         cnn.add_module('batchnorm{0}'.format(i), nn.BatchNorm2d(nOut))
42     if leakyRelu:
43         cnn.add_module('relu{0}'.format(i),
44                         nn.LeakyReLU(0.2, inplace=True))
45     else:
46         cnn.add_module('relu{0}'.format(i), nn.ReLU(True))
47
48     convRelu(0)
49     cnn.add_module('pooling{0}'.format(0), nn.MaxPool2d(2, 2)) # 64x16x64
50     convRelu(1)
51     cnn.add_module('pooling{0}'.format(1), nn.MaxPool2d(2, 2)) # 128x8x32
52     convRelu(2, True)
53     convRelu(3)
54     cnn.add_module('pooling{0}'.format(2),
55                     nn.MaxPool2d((2, 2), (2, 1), (0, 1))) # 256x4x16
56     convRelu(4, True)
57     convRelu(5)
58     cnn.add_module('pooling{0}'.format(3),
59                     nn.MaxPool2d((2, 2), (2, 1), (0, 1))) # 512x2x16
60     convRelu(6, True) # 512x1x16
61
62     self.cnn = cnn
63     self.rnn = nn.Sequential(
64         BidirectionalLSTM(512, nh, nh),
65         BidirectionalLSTM(nh, nh, nclass))
66
67
68     def forward(self, input):
69         # conv features
70         conv = self.cnn(input)
71         b, c, h, w = conv.size()
72         assert h == 1, "the height of conv must be 1"
73         conv = conv.squeeze(2)
74         conv = conv.permute(2, 0, 1) # [w, b, c]
75
76         # rnn features
77         output = self.rnn(conv)
78
79     return output

```

Appendix L - utils.py

```
1  #!/usr/bin/python
2  # encoding: utf-8
3
4  import torch
5  import torch.nn as nn
6  from torch.autograd import Variable
7  import collections
8
9
10 class strLabelConverter(object):
11     """Convert between str and label.
12
13     NOTE:
14
15     Insert `blank` to the alphabet for CTC.
16
17     Args:
18
19     alphabet (str): set of the possible characters.
20     ignore_case (bool, default=True): whether or not to ignore all of the case.
21     """
22
23     def __init__(self, alphabet, ignore_case=True):
24         self._ignore_case = ignore_case
25         if self._ignore_case:
26             alphabet = alphabet.lower()
27         self.alphabet = alphabet + '-' # for `-1` index
28
29         self.dict = {}
30         for i, char in enumerate(alphabet):
31             # NOTE: 0 is reserved for 'blank' required by wrap_ctc
32             self.dict[char] = i + 1
33
34     def encode(self, text):
35         """Support batch or single str.
36
37     Args:
38
39     text (str or list of str): texts to convert.
40
41     Returns:
```

```

39         torch.IntTensor [length_0 + length_1 + ... length_{n - 1}]: encoded texts.
40         torch.IntTensor [n]: length of each text.
41     """
42     if isinstance(text, str):
43         text = [
44             self.dict[char.lower() if self._ignore_case else char]
45             for char in text
46         ]
47         length = [len(text)]
48     elif isinstance(text, collections.Iterable):
49         length = [len(s) for s in text]
50         text = ''.join(text)
51         text, _ = self.encode(text)
52     return (torch.IntTensor(text), torch.IntTensor(length))
53
54     def decode(self, t, length, raw=False):
55         """Decode encoded texts back into strs.
56
57         Args:
58             torch.IntTensor [length_0 + length_1 + ... length_{n - 1}]: encoded texts.
59             torch.IntTensor [n]: length of each text.
60
61         Raises:
62             AssertionError: when the texts and its length does not match.
63
64         Returns:
65             text (str or list of str): texts to convert.
66         """
67         if length.numel() == 1:
68             length = length[0]
69             assert t.numel() == length, "text with length: {} does not match declared length:
70                 ↳ {}".format(t.numel(), length)
71             if raw:
72                 return ''.join([self.alphabet[i - 1] for i in t])
73             else:
74                 char_list = []
75                 for i in range(length):
76                     if t[i] != 0 and (not (i > 0 and t[i - 1] == t[i])):
77                         char_list.append(self.alphabet[t[i] - 1])
78                 return ''.join(char_list)
79         else:

```

```

79         # batch mode
80         assert t.numel() == length.sum(), "texts with length: {} does not match declared
           ↪ length: {}".format(t.numel(), length.sum())
81         texts = []
82         index = 0
83         for i in range(length.numel()):
84             l = length[i]
85             texts.append(
86                 self.decode(
87                     t[index:index + l], torch.IntTensor([l]), raw=raw))
88             index += l
89         return texts
90
91
92     class averager(object):
93         """Compute average for `torch.Variable` and `torch.Tensor`. """
94
95         def __init__(self):
96             self.reset()
97
98         def add(self, v):
99             if isinstance(v, Variable):
100                 count = v.data.numel()
101                 v = v.data.sum()
102             elif isinstance(v, torch.Tensor):
103                 count = v.numel()
104                 v = v.sum()
105
106             self.n_count += count
107             self.sum += v
108
109         def reset(self):
110             self.n_count = 0
111             self.sum = 0
112
113         def val(self):
114             res = 0
115             if self.n_count != 0:
116                 res = self.sum / float(self.n_count)
117             return res
118

```

```

119
120 def oneHot(v, v_length, nc):
121     batchSize = v_length.size(0)
122     maxLength = v_length.max()
123     v_onehot = torch.FloatTensor(batchSize, maxLength, nc).fill_(0)
124     acc = 0
125     for i in range(batchSize):
126         length = v_length[i]
127         label = v[acc:acc + length].view(-1, 1).long()
128         v_onehot[i, :length].scatter_(1, label, 1.0)
129         acc += length
130     return v_onehot
131
132
133 def loadData(v, data):
134     v.data.resize_(data.size()).copy_(data)
135
136
137 def prettyPrint(v):
138     print('Size {0}, Type: {1}'.format(str(v.size()), v.data.type()))
139     print('| Max: %f | Min: %f | Mean: %f' % (v.max().data[0], v.min().data[0],
140                                             v.mean().data[0]))
141
142
143 def assureRatio(img):
144     """Ensure imgH <= imgW."""
145     b, c, h, w = img.size()
146     if h > w:
147         main = nn.UpsamplingBilinear2d(size=(h, h), scale_factor=None)
148         img = main(img)
149     return img

```

Appendix M - dataset.py

```
1  #!/usr/bin/python
2  # encoding: utf-8
3
4  import random
5  import torch
6  from torch.utils.data import Dataset
7  from torch.utils.data import sampler
8  import torchvision.transforms as transforms
9  import lmdb
10 import six
11 import sys
12 from PIL import Image
13 import numpy as np
14
15
16 class lmdbDataset(Dataset):
17
18     def __init__(self, root=None, transform=None, target_transform=None):
19         self.env = lmdb.open(
20             root,
21             max_readers=1,
22             readonly=True,
23             lock=False,
24             readahead=False,
25             meminit=False)
26
27         if not self.env:
28             print('cannot creat lmdb from %s' % (root))
29             sys.exit(0)
30
31         with self.env.begin(write=False) as txn:
32             nSamples = int(txn.get('num-samples'))
33             self.nSamples = nSamples
34
35         self.transform = transform
36         self.target_transform = target_transform
37
38     def __len__(self):
```

```

39         return self.nSamples
40
41     def __getitem__(self, index):
42         assert index <= len(self), 'index range error'
43         index += 1
44         with self.env.begin(write=False) as txn:
45             img_key = 'image-%09d' % index
46             imgbuf = txn.get(img_key)
47
48             buf = six.BytesIO()
49             buf.write(imgbuf)
50             buf.seek(0)
51             try:
52                 img = Image.open(buf).convert('L')
53             except IOError:
54                 print('Corrupted image for %d' % index)
55                 return self[index + 1]
56
57             if self.transform is not None:
58                 img = self.transform(img)
59
60             label_key = 'label-%09d' % index
61             label = str(txn.get(label_key))
62
63             if self.target_transform is not None:
64                 label = self.target_transform(label)
65
66         return (img, label)
67
68
69     class resizeNormalize(object):
70
71     def __init__(self, size, interpolation=Image.BILINEAR):
72         self.size = size
73         self.interpolation = interpolation
74         self.toTensor = transforms.ToTensor()
75
76     def __call__(self, img):
77         img = img.resize(self.size, self.interpolation)
78         img = self.toTensor(img)
79         img.sub_(0.5).div_(0.5)

```



```

80         return img
81
82
83     class randomSequentialSampler(sampler.Sampler):
84
85         def __init__(self, data_source, batch_size):
86             self.num_samples = len(data_source)
87             self.batch_size = batch_size
88
89         def __iter__(self):
90             n_batch = len(self) // self.batch_size
91             tail = len(self) % self.batch_size
92             index = torch.LongTensor(len(self)).fill_(0)
93             for i in range(n_batch):
94                 random_start = random.randint(0, len(self) - self.batch_size)
95                 batch_index = random_start + torch.range(0, self.batch_size - 1)
96                 index[i * self.batch_size:(i + 1) * self.batch_size] = batch_index
97             # deal with tail
98             if tail:
99                 random_start = random.randint(0, len(self) - self.batch_size)
100                 tail_index = random_start + torch.range(0, tail - 1)
101                 index[(i + 1) * self.batch_size:] = tail_index
102
103             return iter(index)
104
105         def __len__(self):
106             return self.num_samples
107
108
109     class alignCollate(object):
110
111         def __init__(self, imgH=32, imgW=100, keep_ratio=False, min_ratio=1):
112             self.imgH = imgH
113             self.imgW = imgW
114             self.keep_ratio = keep_ratio
115             self.min_ratio = min_ratio
116
117         def __call__(self, batch):
118             images, labels = zip(*batch)
119
120             imgH = self.imgH

```

```

121     imgW = self.imgW
122     if self.keep_ratio:
123         ratios = []
124         for image in images:
125             w, h = image.size
126             ratios.append(w / float(h))
127         ratios.sort()
128         max_ratio = ratios[-1]
129         imgW = int(np.floor(max_ratio * imgH))
130         imgW = max(imgH * self.min_ratio, imgW) # assure imgH >= imgW
131
132     transform = resizeNormalize((imgW, imgH))
133     images = [transform(image) for image in images]
134     images = torch.cat([t.unsqueeze(0) for t in images], 0)
135
136     return images, labels

```

References

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Gregory S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian J. Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Józefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Gordon Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul A. Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda B. Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *CoRR*, abs/1603.04467, 2016.
- [2] Ouais Alsharif and Joelle Pineau. End-to-end text recognition with hybrid HMM maxout models. *CoRR*, abs/1310.1811, 2013.
- [3] B. Bai, F. Yin, and C. L. Liu. A fast stroke-based method for text detection in video. In *2012 10th IAPR International Workshop on Document Analysis Systems*, pages 69–73, March 2012.
- [4] A. Bissacco, M. Cummins, Y. Netzer, and H. Neven. Photoocr: Reading text in uncontrolled conditions. In *2013 IEEE International Conference on Computer Vision*, pages 785–792, Dec 2013.
- [5] Jeffrey R. Blum, Mathieu Bouchard, and Jeremy R. Cooperstock. Whats around me? spatialized audio augmented reality for blind users with a smartphone.
- [6] J. Canny. A computational approach to edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-8(6):679–698, Nov 1986.
- [7] H. Chen, S. S. Tsai, G. Schroth, D. M. Chen, R. Grzeszczuk, and B. Girod. Robust text detection in natural images with edge-enhanced maximally stable extremal regions. In *2011 18th IEEE International Conference on Image Processing*, pages 2609–2612, Sept 2011.
- [8] Xiangrong Chen and A. L. Yuille. 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.*, volume 2, pages II–366–II–373 Vol.2, June 2004.

- [9] Adam Coates, Blake Carpenter, Sanjeev Satheesh, Bipin Suresh, Tao Wang, David J. Wu, and Andrew Y. Ng. Text detection and character recognition in scene images with unsupervised feature learning.
- [10] R. Collobert, K. Kavukcuoglu, and C. Farabet. Torch7: A matlab-like environment for machine learning. In *BigLearn, NIPS Workshop*, 2011.
- [11] M. Donoser, H. Bischof, and S. Wagner. Using web search engines to improve text recognition. In *2008 19th International Conference on Pattern Recognition*, pages 1–4, Dec 2008.
- [12] Tefilo Emdio de Campos, Bodla Rakesh Babu, and Manik Varma. Character recognition in natural images. 2:273–280, 01 2009.
- [13] B. Epshtein, E. Ofek, and Y. Wexler. Detecting text in natural scenes with stroke width transform. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2963–2970, June 2010.
- [14] Alex Graves, Santiago Fernandez, and Faustino Gomez. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *In Proceedings of the International Conference on Machine Learning, ICML 2006*, pages 369–376, 2006.
- [15] S. M. Hanif and L. Prevost. Text detection and localization in complex scene images using constrained adaboost algorithm. In *2009 10th International Conference on Document Analysis and Recognition*, pages 1–5, July 2009.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [17] Pan He, Weilin Huang, Yu Qiao, Chen Change Loy, and Xiaoou Tang. Reading scene text in deep convolutional sequences. 06 2015.
- [18] Sepp Hochreiter and Jrgen Schmidhuber. Long short-term memory, 1997.
- [19] Sanghoon Hong, Byung-Seok Roh, Kye-Hyeon Kim, Yeongjae Cheon, and Minje Park. Pvanet: Lightweight deep neural networks for real-time object detection. *CoRR*, abs/1611.08588, 2016.
- [20] Jan Hendrik Hosang, Rodrigo Benenson, and Bernt Schiele. Learning non-maximum suppression. *CoRR*, abs/1705.02950, 2017.
- [21] Lichao Huang, Yi Yang, Yafeng Deng, and Yinan Yu. Densebox: Unifying landmark localization with end to end object detection. *CoRR*, abs/1509.04874, 2015.

- [22] R. Huang, P. Shivakumara, and S. Uchida. Scene character detection by an edge-ray filter. In *2013 12th International Conference on Document Analysis and Recognition*, pages 462–466, Aug 2013.
- [23] M. R. Islam, C. Mondal, M. K. Azam, and A. S. M. J. Islam. Text detection and recognition using enhanced ms-er detection and a novel ocr technique. In *2016 5th International Conference on Informatics, Electronics and Vision (ICIEV)*, pages 15–20, May 2016.
- [24] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Synthetic data and artificial neural networks for natural scene text recognition. In *Workshop on Deep Learning, NIPS*, 2014.
- [25] Vedaldi A. Jaderberg M. and Zisserman A. Deep feature for text spotting. pages 512–528, 2014.
- [26] Yingying Jiang, Xiangyu Zhu, Xiaobing Wang, Shuli Yang, Wei Li, Hua Wang, Pei Fu, and Zhenbo Luo. R2CNN: rotational region CNN for orientation robust scene text detection. *CoRR*, abs/1706.09579, 2017.
- [27] D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. Ghosh, A. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu, F. Shafait, S. Uchida, and E. Valveny. Icdar 2015 competition on robust reading. In *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*, pages 1156–1160, Aug 2015.
- [28] Kwang In Kim, Keechul Jung, and Jin Hyung Kim. Texture-based approach for text detection in images using support vector machines and continuously adaptive mean shift algorithm. volume 25, pages 1631–1639, Dec 2003.
- [29] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- [30] Tao Kong, Anbang Yao, Yurong Chen, and Fuchun Sun. Hypernet: Towards accurate region proposal generation and joint object detection. *CoRR*, abs/1604.00600, 2016.
- [31] J. J. Lee, P. H. Lee, S. W. Lee, A. Yuille, and C. Koch. Adaboost for text detection in natural scene. In *2011 International Conference on Document Analysis and Recognition*, pages 429–434, Sept 2011.
- [32] Minhua Li and Chunheng Wang. An adaptive text detection approach in images and video frames. In *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*, pages 72–77, June 2008.

- [33] Jian Liang, David S. Doermann, and Huiping Li. Camera-based analysis of text and documents: a survey. *International Journal of Document Analysis and Recognition (IJDAR)*, 7:84–104, 2004.
- [34] M. Liao, B. Shi, and X. Bai. TextBoxes++: A Single-Shot Oriented Scene Text Detector. *ArXiv e-prints*, January 2018.
- [35] Minghui Liao, Baoguang Shi, Xiang Bai, Xinggang Wang, and Wenyu Liu. Textboxes: A fast text detector with a single deep neural network. *CoRR*, abs/1611.06779, 2016.
- [36] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott E. Reed, Cheng-Yang Fu, and Alexander C. Berg. SSD: single shot multibox detector. *CoRR*, abs/1512.02325, 2015.
- [37] X. Liu, D. Liang, S. Yan, D. Chen, Y. Qiao, and J. Yan. FOTS: Fast Oriented Text Spotting with a Unified Network. *ArXiv e-prints*, January 2018.
- [38] X. Liu and W. Wang. Robustly extracting captions in videos based on stroke-like edges and spatio-temporal analysis. *IEEE Transactions on Multimedia*, 14(2):482–489, April 2012.
- [39] J. Matas, O. Chum, M. Urban, and T. Pajdla. Robust wide baseline stereo from maximally stable extremal regions. In *In Proc. BMVC*, pages 384–393, 2002.
- [40] K. Mikolajczyk, T. Tuytelaars, C. Schmid, A. Zisserman, J. Matas, F. Schaffalitzky, T. Kadir, and L. Van Gool. A comparison of affine region detectors. *International Journal of Computer Vision*, 65(1):43–72, Nov 2005.
- [41] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. *CoRR*, abs/1606.04797, 2016.
- [42] A. Mishra, K. Alahari, and C. V. Jawahar. Scene text recognition using higher order language priors. In *BMVC*, 2012.
- [43] Anand Mishra, Karteek Alahari, and C V. Jawahar. Scene text recognition using higher order language priors. 09 2012.
- [44] L. Neumann and J. Matas. Real-time scene text localization and recognition. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3538–3545, June 2012.
- [45] Tatiana Novikova, Olga Barinova, Pushmeet Kohli, and Victor Lempitsky. Large-lexicon attribute-consistent text recognition in natural images.

- [46] Y. F. Pan, X. Hou, and C. L. Liu. A hybrid approach to detect and localize texts in natural scene images. *IEEE Transactions on Image Processing*, 20(3):800–813, March 2011.
- [47] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [48] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *CoRR*, abs/1506.01497, 2015.
- [49] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. *CoRR*, abs/1505.04597, 2015.
- [50] Zohra Saidane and Christophe Garcia. Automatic scene text recognition using a convolutional neural network. In *In Proceedings of the Second International Workshop on Camera-Based Document Analysis and Recognition (CBDAR, 2007*.
- [51] Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *CoRR*, abs/1507.05717, 2015.
- [52] Baoguang Shi, Xinggang Wang, Pengyuan Lv, Cong Yao, and Xiang Bai. Robust scene text recognition with automatic rectification. *CoRR*, abs/1603.03915, 2016.
- [53] Cunzhao Shi, Chunheng Wang, Baihua Xiao, Yang Zhang, Song Gao, and Zhong Zhang. Scene text recognition using part-based tree-structured character detection. pages 2961–2968, 06 2013.
- [54] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014.
- [55] A. Tabassum and S. A. Dhondse. Text detection using mser and stroke width transform. In *2015 Fifth International Conference on Communication Systems and Network Technologies*, pages 568–571, April 2015.
- [56] Zhi Tian, Weilin Huang, He Tong, Pan He, and Yu Qiao. Detecting text in natural image with connectionist text proposal network. volume 9912, pages 56–72, 10 2016.
- [57] Andreas Veit, Tomas Matera, Lukas Neumann, Jiri Matas, and Serge J. Belongie. Coco-text: Dataset and benchmark for text detection and recognition in natural images. *CoRR*, abs/1601.07140, 2016.

- [58] Kai Wang, B. Babenko, and S. Belongie. End-to-end scene text recognition. In *2011 International Conference on Computer Vision*, pages 1457–1464, Nov 2011.
- [59] T. Wang, D. J. Wu, A. Coates, and A. Y. Ng. End-to-end text recognition with convolutional neural networks. In *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, pages 3304–3308, Nov 2012.
- [60] Tao Wang, David J Wu, Adam Coates, and Andrew Y Ng. End-to-end text recognition with convolutional neural networks. 02 2018.
- [61] W. Wu, D. Chen, and J. Yang. Integrating co-training and recognition for text detection. In *2005 IEEE International Conference on Multimedia and Expo*, pages 4 pp.–, July 2005.
- [62] C. Yao, X. Bai, W. Liu, Y. Ma, and Z. Tu. Detecting texts of arbitrary orientations in natural images. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1083–1090, June 2012.
- [63] Cong Yao, Xiang Bai, Nong Sang, Xinyu Zhou, Shuchang Zhou, and Zhimin Cao. Scene text detection via holistic, multi-channel prediction. *CoRR*, abs/1606.09002, 2016.
- [64] Q. Ye and D. Doermann. Text detection and recognition in imagery: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(7):1480–1500, 2015.
- [65] C. Yi and Y. Tian. Localizing text in scene images by boundary clustering, stroke segmentation, and string fragment classification. *IEEE Transactions on Image Processing*, 21(9):4256–4268, Sept 2012.
- [66] X. C. Yin, Z. Y. Zuo, S. Tian, and C. L. Liu. Text detection, tracking and recognition in video: A comprehensive survey. *IEEE Transactions on Image Processing*, 25(6):2752–2773, June 2016.
- [67] Matthew D. Zeiler. ADADELTA: an adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.
- [68] Zheng Zhang, Zhang Chengquan, Wei Shen, Cong Yao, Wenyu Liu, and Xiang Bai. Multi-oriented text detection with fully convolutional networks. pages 4159–4167, 06 2016.
- [69] X. Zhou, C. Yao, H. Wen, Y. Wang, S. Zhou, W. He, and J. Liang. East: An efficient and accurate scene text detector. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2642–2651, July 2017.

- [70] Karel Zimmermann. A new class of learnable detectors: Cser- class specific extremal regions.

List of Abbreviations

- IEEE:** Institute of Electrical and Electronics Engineers, Inc.
- WHO:** World Health Organization
- CVPR:** The Conference on Computer Vision and Pattern Recognition
- ICCV:** International Conference on Computer Vision
- ICDAR:** International Conference on Document Analysis and Recognition
- OCR:** Optical Character Recognition
- CC:** Connected Component
- MSER:** Maximally Stable Extremal Regions
- SWT:** Stroke Width Transform
- CRF:** Conditional Random Field
- NN:** Neural Network
- DNN:** Deep Neural Network
- CNN:** Convolutional Neural Network
- DCNN:** Deep Convolutional Neural Network
- FCN:** Fully Convolutional Network
- RNN:** Recurrent Neural Network
- EAST:** Easy Accurate Scene Text
- RCNN & CRNN:** Convolutional Recurrent Neural Network
- NMS:** Non Maximum Suppression
- LSTM:** Long Short Term Memory
- CTC:** Connectionist Temporal Classification
- GPU:** Graphics Processing Unit
- CPU:** Central Processing Unit

GHz: Gigahertz, 1 billion hertz

SDK: Software Development Kit

JSON: JavaScript Object Notation

GPS: Global Positioning System