

# A Hybrid Method for Chinese Character Segmentation

J.F. van Wezel

University of Groningen

September 15, 2017

## Abstract

This work shows a hybrid method between connected components and horizontal vertical projection for Chinese character segmentation with mixed results. There are some problems with the implementation of the system and these problems are located and discussed. Even with these problems the work shows that a hybrid method is able to achieve segmentation with the best of both worlds.

**Character Segmentation - Text recognition - Connected Components - Chinese Characters - Horizontal vertical projection - Computer vision - Machine Learning**

## 1 Introduction

Character recognition is an important part of the fields computer vision and artificial intelligence. For a machine to recognize sentences, words ,and characters in a given image is useful for multiple goals. For example to digitize handwritten texts ([8]). digitized texts have a number of advantages over physical texts. An important advantage is that digital texts are more easily distributed and reproduced. It is also easier to search in digital texts. But probably the greatest advantage is storage. It is easier and takes less room to save and backup a few bytes on a disk than it is to physically backup archive the same texts. As is often the case in the fields computer vision and artificial intelligence, compression is the end goal.

An other practical application of character recognition is text translation ([3]). Text recognition for the use of translation is useful for example tourists who use their smart-phone to take a picture of a sign and have that sign translated in their native language by an translation app. Or even better, have their text translated and inserted into the camera feed directly. This form of augmented reality is still in its infancy but being able to recognize and digitize texts on walls, signs ,and posters for all kinds of purposes is not that hard to imagine.

Usage of character recognition can also be imagined in the retail industry. Here it can be used for label recognition on certain products, currency ,and identification purposes for restricted product groups (like liquor and tobacco).

The automotive industry uses character recognition in their autonomous cars ([6]). They read speed limits and other road signs to form a state of awareness of their surroundings They abstract what the next course of action should be based partly on that gathered state. A normal road sign is not a character but a computer does not know the difference between ‘human’ characters and a specific image. The important part is that the goal is to find meaning in an image by ‘reading’ parts of that image. Sign recognition and character recognition are therefor closely related. They differ in that characters tend to come in a sequence forming sentences. Where road signs are mostly isolated.

Staying on the road character recognition is used for numberplate recognition ([5]). It is used by law enforcers to find stolen or unregistered vehicles but also as an automated form of speed

control. Here a mounted or portable camera is used to take pictures of cars with their number plates and the software recognizes the plates and the characters on them to identify the vehicle.

A character recognition system can be roughly divided into three main components: segmentation, feature extraction, and classification.

Segmentation is where an algorithm tries to find the location of the character or sentences in a given image. Before it can start this task the image often is preprocessed. During preprocessing, an image can be binarized and filtered to reduce small noise and other unwanted large components that might be present in the image. Preprocessing can also include rotating, shearing or other morphological operations.

Feature extraction is where the features of the segmented character images are found and measured. For example if an image is black and white a feature could be the ratio between black and white pixels. An other example of a feature could be the amount of horizontal edges in the image by using an edge detector. The type of features that should be extracted depends heavily on the type of dataset. A western handwritten scroll houses different characters and features than Egyptian hieroglyphs.

Classification is where the, until then unknown, characters get assigned a label. This task can be done by a number of algorithms from the simple K-nearest neighbors, used for classification and regression ([1]), to complicated multilayer convolutional neural networks. An example of such a CNN is LeNet-5 by Lecun ([9]). Which algorithm yields the best results depends on the dataset, the segmentation, and the feature extraction.

This paper will focus on the building of pragmatic preprocessing and segmentation system for a Chinese character dataset. There has been work done in the past on Chinese character segmentation by others. From this work we know there are multiple ways of segmenting the Chinese characters. Examples of methods used in the past are: vertical and horizontal projection based, recognition based, skeleton based, and connected component based [4].

Recognition based uses a learning algorithm to find the locations of the characters, Guo et al. proposes a self learning license plate method ([7]). Skeleton based uses the skeleton of the image to find the individual characters in a image. A variant of this technique has been demonstrated by Nikolaou et al. ([10]). Component based uses the connected components of an image to find the characters. An example of this method is shown by Chen et al. ([4])

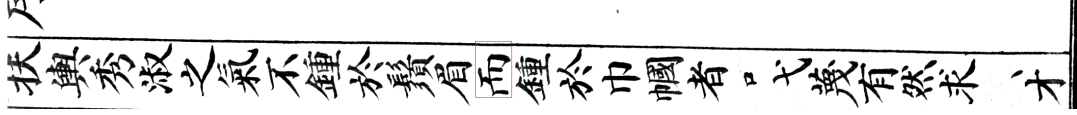
For the created system a combination of various preprocessing techniques, horizontal projection, and connected components was used to segment the dataset. This results in a pragmatic segmentation system that was designed for the given dataset but should also be useful for other Chinese character datasets. The segmented images are used for feature extraction and will be shown in his paper by Lennart van de Guchte and classification with a CNN of the data with the segmented images by Leon Pater.

This work will use labeled data from the given dataset to validate the resulting segmented images. The use of already used methods that were shown to have worked in the shown related work will likely result in a reliable segmentation system.

This paper will continue by explaining a bit more about the dataset. Then the used methods and used parameters will be further explained. After, the type of experiment and how the segmented images were validated will be explained in the Experiment section. From the experiment the results will be shown and discussed. This paper will end with a conclusion on the work done.

## 2 Dataset

This section will shed light on the used Chinese character dataset. It will show some metrics, examples of the noise encountered, and other problems that needed to be solved to achieve the segmented images.



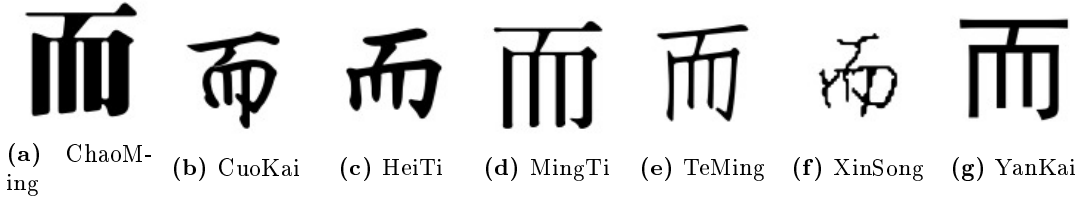
**Figure 1:** Example of a typical image in the dataset

**Table 1:** Dataset metrics

| Size | Labels | Unique labels | Unique characters | Fonts | $\mu_{lw}$ | $\sigma_{lw}$ |
|------|--------|---------------|-------------------|-------|------------|---------------|
| 6000 | 27026  | 589           | 6500              | 7     | 73.5944 px | 8.4463 px     |

The metrics of the used dataset for segmentation

Table 1 shows some metrics from the dataset. The size (6000) is smaller than the number of labels (27026). The size is based on the number of images the dataset holds. The images in the dataset consist of unsegmented lines of Chinese characters. An example of a typical image is shown in figure 1. The red box in the image shows the location of an a labeled character. The locations of these labeled characters are given in xml files along with the line images. The  $\mu_{lw}$  and  $\sigma_{lw}$  signify the mean with of the labels and the standard variation in width of the labels in the dataset.



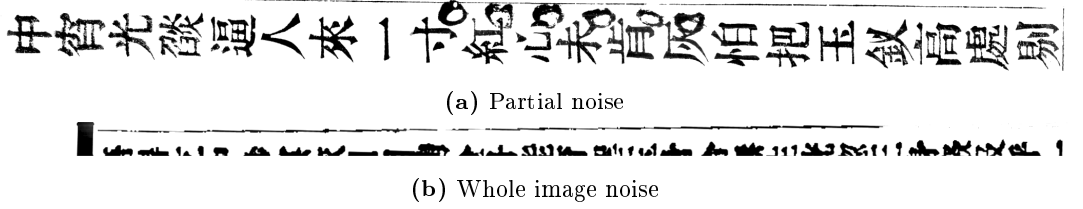
**Figure 2:** The Chinese printed fonts present in the dataset.

The fonts used in the dataset are shown in figure 2. There are seven fonts in total, which show a similarity with the characters in the images in the dataset. It seems assumable from this observation that the characters were printed instead of handwritten.

The images in the dataset are rows of characters but as shown in figure 1 contain sometimes an extra row with one or more characters in it. Sometimes there is whitespace on one or all sides of the image and sometimes the whole image is noise. The character rows are more or less given but the exact location of the the character row needs still to be found in the image.

The Chinese language is generally written from top to bottom. The orientation of the images is however in the western orientation of the language, from left to right. Because the whole dataset was in the western orientation the dataset will be handled as such. The horizontal line present in the image is to separate the character rows. This line is not present in all images. The image also shows that the rotation of the image is not straight, it is rotated by a few degrees. The characters them self however do not seem to show a curve. This is also not true for all images, some images do show a slight curvature. The rotation is probably a result of the scanning process. This is mostly done by hand and prone to error. The line could be abused to find the correct rotation of the image. But as stated the line is not present in all images and a different approach would be needed for the images without the line.

Figure 1 also shows a ‘fat’ vertical line on the far right side of the image. This line is also present in most images but not all. Sometimes this line shows up on the far left side of the image or on both sides. This line is not a character and needs to be handled as noise, but it can also be used to find the starting and end point of a character row.



**Figure 3:** Examples of noise found in the dataset

The image shown in figure 3a illustrates noise on the center characters often found in this dataset. This noise might be the result of water damage or perforation of the sheet where the characters were originally on. This noise could be resolved with morphological operations and using a dataset specific filter. The dataset also contains some images that are by them self noise. Figure 3b shows an image that is complete noise. This specific image seems to be located to low during the scanning process. This resulted in only half of the original characters showing in the dataset.

### 3 Method

This section will show what methods where used to segment the dataset. It will do so in a chronological order based on the example image shown in section 2 in figure 1.

#### 3.1 Segmentation

##### 3.1.1 Rotation and binarization

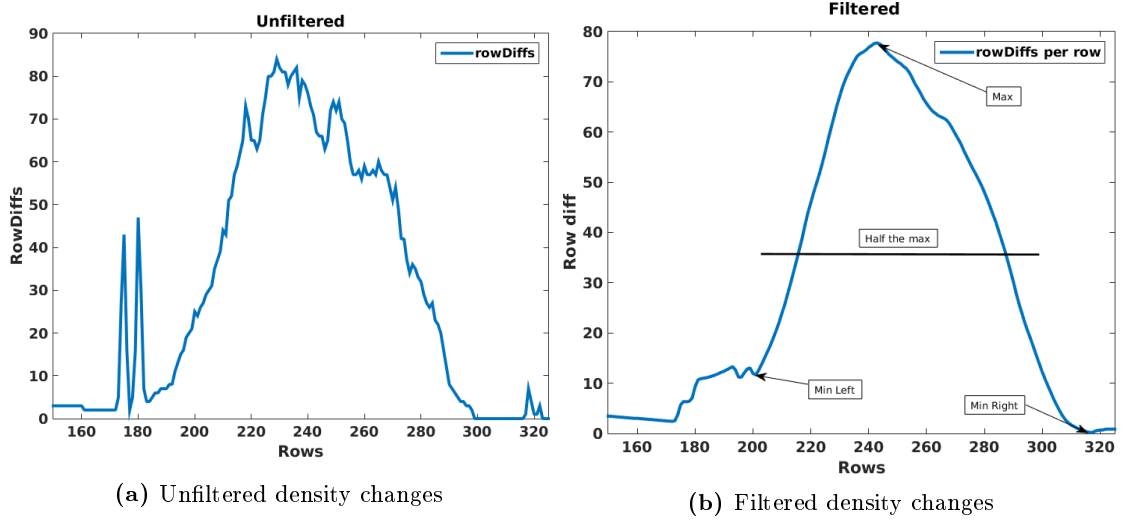
As discussed in the introduction (1) the Chinese characters were supplied in greyscale strokes of characters in the western orientation (from left to right).

In order to segment the characters from an image a combination of methods were used. First the image is binarized. Binarization is done by applying a Gaussian filter on the image to remove some of the pixel noise. Then Otsu’s [11] method of finding a threshold is used. The image is then binarized using this threshold.

After binarization the image is rotated. This is done because most images in the dataset are titled slightly. As discussed in the dataset section (2) This is probably due to some variations in the orientation of the paper during the scanning process. Rotation is done by resizing the image in the vertical direction to half its size. The idea behind this is that the characters will be squeezed together and form a horizontal line, this line can then be used to find the optimal rotation. This is done by taking the horizontal densities. The maximum densities are saved for the different rotations (between  $+1^\circ$  and  $-1^\circ$ ). The rotation with the largest maximum density is where the line is likely to be the most horizontal. This is taken as the optimal rotation. The image is rotated with this optimal angle.

##### 3.1.2 Vertical location

After the characters are rotated the location of the characters is determined on the vertical axes. The used method to determine this location is horizontal projection on the changes in pixels from black to white. On the densities an averaging filter is first applied. Each row is averaged with ten of its neighbor rows. To illustrate this the changes in densities were taken on the example image



**Figure 4:** Method of finding the vertical location of the characters

shown in the dataset section (1). The effect of the averaging is clearly visible in figure 4. In figure 4a the unfiltered image is shown. The plot has many peaks and also the top black line is clearly visible in the densities on the left side of the image. In figure 4b the filtered image is shown. Only the global peaks stay and the noise peak from the black line has almost disappeared.

$$S = \begin{bmatrix} -1 & 0 & 1 \end{bmatrix}$$

**Figure 5:** Sobel operator

In order to find the location of the characters, the global maximum of changes in densities is taken. This is always the top of the highest peak of the characters. To find the bases on both sides of the peak a Sobel operator is used (shown in figure 5). Sobel's operators can be used to find the gradient in a signal, demonstrated in a paper by Richard O. Duda and Peter E. Hart ([2]) and later by Sobel self ([13]). To find these minima the half of the max is taken and from this half the first occurrence where the gradient crosses zero in the left and right direction is taken as a base of the mountain. The half of the max is used because on the top of the plot sometimes goes up and down even with the average filtering. The half max makes sure the whole height of the characters is taken. After determining the vertical location of the characters the image is cropped accordingly. Whitespace on the left and right side of the image is also removed.

扶輿秀淑之氣不鍾於鬚眉而鍾於巾幗者。弋戩有然求 才

**Figure 6:** Rotated and cropped image

### 3.1.3 Connected components

We are left with the row of characters. An example is shown in figure 6. This image is rotated and cropped from figure 1. In order to find the characters in this row the connected components are taken. These are taken with the method described by Sedgewick ([12]). Sedgewick's method will give us a list of all the connected components in the image. If multiple components are above

each other we assume that these components are part of the same character. These components are merged. Some components are not directly on top of each-other. It is determinant how much the components overlap vertically. If the overlapping width is larger or the same as 0.4 times the width of the smaller of the two components, they are merged. This is done to make sure the components are on top of each other and not some small part of one component is on top of another component.

After merging vertically we are left with a list of components that might be outliers in terms of width in comparison with the mean width of the characters in the supplied labeled dataset. In order to recognize outliers the mean and standard deviation was calculated from the characters in the labeled dataset, also shown in section 2 and table 1. A component is deemed a small outlier if:

$$C_w < \mu_{lw} + 2\sigma_{lw} \quad (1)$$

and a big outlier if:

$$C_w > \mu_{lw} + 2\sigma_{lw} \quad (2)$$

Here  $C_w$  is the component's width,  $\mu_{lw}$  is the mean width, and  $\sigma_{lw}$  is the standard deviation in width. Both were taken from the labels of the dataset

Small component can be noise or a character. A small component is more likely to be a character if it has a big amount of whitespace left and right from it. Except when it is the first or the last component in the list. To check for this the amount of whitespace on both sides of the small component is measured and added to the width of that component. If the component is the first or the last component in the list the white space is only measured on one side and is added twice to the width. With the added width the component is reevaluated for being an outlier.

The list with components still holds noise components. There are different types of noise present in the images. This algorithm tries to find unwanted 'blobs' and too small components. The unwanted 'blobs' are detected by taking the square of the component and filling it with white pixels. The component itself is then printed on this square in black pixels. The component is determined to be noise if the rectangle consists of more than 40 percent black pixels. This number is based on the labels of the dataset where the amount of black pixels in the rectangle was always lower than 40 percent of the pixels. This method removes components like the black bar shown in figure 6 on the outer right side of the image. An image is deemed noise if the width of the character is smaller than 15 pixels.

The small components that are left in the list are now merged with the smallest of its two neighbor components. If a neighboring component is deemed big a small component will not be merged with it. If a component has a big neighbor to the left and a big neighbor to the right, This small component will not be merged. After each merger the components widths will be evaluated. If two merged components are still small, the algorithm will try to merge them again with one of the (new) smallest neighbors.

After merging all small components are dealt with. The component list still contains big outliers. Some might even have been created during merging. These components are split by projecting the variations in white pixels to black pixels (or the amount of fluctuations in white pixels). Where there are less variations there is more white space. The characters are split on where the fluctuations of the whitespace is minimal. After splitting the components are reevaluated on width. If a component is determined big it is split again. This is done until the component is not big anymore.

### 3.1.4 Finding the character box

With the split big components the segmentation part is completed. In order to validate the segmented images the location of the characters need to be determined in the original image. Because the image is rotated before segmentation the locations of the segmented images need to be rotated to point to the correct location in the original image. The location of the segmented

characters is given by taking the  $x, y, width$  and  $height$  variables. This creates a rectangle on the original image where a character is determined to be. But because the characters outputted by the segmentation system are (merged) connected components there might be parts of other characters present in the rectangle that are not present in the character outputted by the segmentation system.

The approximation of the rectangle is made by calculating all corner coordinates of the rectangle by taking the max and min row and column of the characters in the rotated image. Then these coordinates are rotated with the following formulas:

$$x' = y \times \sin(\alpha) + x \times \cos(\alpha) \quad (3)$$

$$y' = y \times \cos(\alpha) - x \times \sin(\alpha) \quad (4)$$

After rotating the minimum value for  $x'$  and  $y'$  is taken to determine the upper left corner of the rectangle. Then the width is determined by:

$$width = abs(min(x_1, x_2, x_3, x_4) - max(x_1, x_2, x_3, x_4)) \quad (5)$$

The height is calculated in a similar fashion:

$$height = abs(min(y_1, y_2, y_3, y_4) - max(y_1, y_2, y_3, y_4)) \quad (6)$$

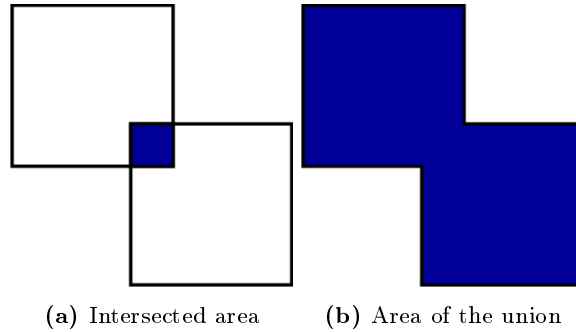
These equations give the  $x, y, width$  and  $height$  variables needed to reconstruct the rectangle. After approximating the rectangle the characters are saved for use in the feature extractor and later the classifier of the character segmentation system.

### 3.2 Feature Extraction

### 3.3 Classification

## 4 Experiment

This section will show by example how the segmented dataset is validated.

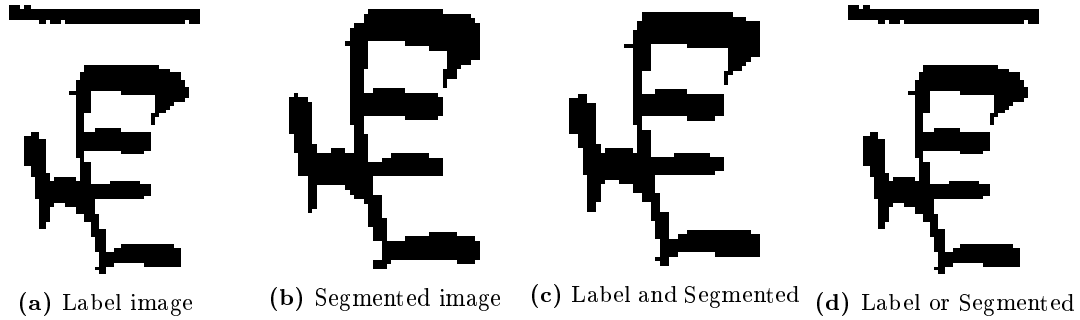


**Figure 7:** Method of validating the segmented data (intersection over union)

In order to validate the segmented images obtained by the segmentation the intersection over union is used. As shown in figure 7a the intersection between two boxes is the area that overlaps between these boxes. The union over two boxes give the area both the boxes occupy (figure 7b). The fraction between the intersection and the union tells us how well the boxes fit on top of each other. With this fraction the performance of the segmentation can be measured.

Because the characters consist mainly of black pixels, and the location of the box is arbitrary, the intersection over union is done only over the black pixels in both boxes. This is done by binerizing the original image and taking its complement. The complement is taken so black pixels

will be represented by ones and white pixels by zeros. Now the labeled box and a segmented box are both printed on a white background the size of the original image. This results in two binary images on which an And is the same as the intersection, an Or is the same as a union. Summing the result from these operations give us the needed scalars to perform the intersection over union fraction.



**Figure 8:** Example of intersection over union on a segmented character in the dataset,  $IoU \approx 0.73$ . The character used is the twelfth character in the image shown in figure 1.

In order to be able to say anything about the performance of the segmentation results will be thresholded from an intersection over union rate from 0.1 to 10. if the intersection over union is lower than 0.1 the boxes are not touching each other enough.



## 5 Results

This section will show the direct results of the intersection over union from the segmentation system. It will also show some cases where the segmentation seemed to have failed and where it succeeded.

**Table 2:** IoU Results

| 0.1   | 0.2   | 0.3   | 0.4   | 0.5   | 0.6   | 0.7   | 0.8   | 0.9   | 1     |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 0.988 | 0.973 | 0.948 | 0.905 | 0.841 | 0.769 | 0.704 | 0.631 | 0.532 | 0.117 |

The results from the segmented characters with the labeled data from the dataset.

## 6 Discussion

This section will discuss the results given in the results section 5, both positive and negative results will be handled to give an inclusive image of the achieved segmentation.

In figure ?? two components are shown that have been successfully merged together by the merging algorithm discussed in section 3. The left most component in figure ?? seems to be a small components merged with the seemingly normal sized component to the right of it forming the shown character.

Figure ?? shows a small segmented character. In figure ?? the location in the original dataset is shown. It stands out that the component is a character on its own and is handled properly by the segmentation system.

The previous two characters where examples of where the segmentation method seemed to have handled correctly. Now we will discuss some examples where the system seems to make mistakes.

The first example is shown in figure ?. Where two components seem to have been merged by the system where they seem to be two individual characters. This is probably due to the combined width of these two characters not being long enough for the system to see them as such.

The next example is a noise reduction problem. In figure ?? the segmented character is shown and in figure ?? the cropped original image is shown. It is clearly visible that there are two stripes missing. In this case the system classified the other stripes as noise and just removed them from the character list.

In figure ?? and figure ?? two segmented characters are shown. When taking a look at figure ?? we can see that the two segmented characters were segmented wrongly to contain 1.5 character and 0.5 character. This is probably due to an error in the merging process. Possibly because the merging process does not take the whitespace between the components into account.

In the three figures ?? to ?? faulty segmentation due to noise is shown. In figure ?? the black line above the characters merges the components together due to the vertical merging process. The splitting algorithm does not seem to split the image in this case and needs to be investigated.

Figure ?? shows a black dot and is probably noise. However it is difficult to classify this kind of ball as such because its size is the same as a small character.

The next example of noise is however better to deal with and is shown in figure ?. This is the black bar discussed in section 3. The system tries to handle these cases as described in section 3 but the parameter used needs further tweaking.

## 7 Conclusion

Based on the example shown in the results section and discussed in the discussion section we must conclude that the segmentation achieved seems good but could be better when more work is put in to it. Most examples handled in the discussion concluded in the same thing. Namely that some part of the system was faulty or that some parameter still needed to be tweaked more. Almost all problems shown could be handled by the methods implemented in the system. The connected components implementation together with horizontal and vertical projection shows promise and should (if worked out more) be able to achieve better results. This work however shows that a hybrid form can achieve segmentation with the best of both worlds.

## References

- [1] N. S. Altman. An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician*, 46(3):175–185, 1992.
- [2] Abraham Bookstein. Pattern classification and scene analysis. richard o. duda , peter e. hart. *The Library Quarterly*, 44(3):258–259, 1974.
- [3] Chongmu Chen, Da-Han Wang, and Hanzi Wang. *Scene Character and Text Recognition: The State-of-the-Art*, pages 310–320. Springer International Publishing, Cham, 2015.

- [4] Jiun-Lin Chen, Chi-Hong Wu, and Hsi-Jian Lee. *Chinese Handwritten Character Segmentation in Form Documents*, pages 348–362. Springer Berlin Heidelberg, Berlin, Heidelberg, 1999.
- [5] Shan Du, Mahmoud Ibrahim, Mohamed S. Shehata, and Wael M. Badawy. Automatic license plate recognition (alpr): A state-of-the-art review. *IEEE Trans. Circuits Syst. Video Techn.*, 23(2):311–325, 2013.
- [6] M. Y. Fu and Y. S. Huang. A survey of traffic sign recognition. In *2010 International Conference on Wavelet Analysis and Pattern Recognition*, pages 119–124, July 2010.
- [7] J. M. Guo and Y. F. Liu. License plate localization and character segmentation with feedback self-learning and hybrid binarization techniques. *IEEE Transactions on Vehicular Technology*, 57(3):1417–1424, May 2008.
- [8] S. Jaeger, C.-L. Liu, and M. Nakagawa. The state of the art in japanese online handwriting recognition compared to techniques in western handwriting recognition. *Document Analysis and Recognition*, 6(2):75–88, Oct 2003.
- [9] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. In *Intelligent Signal Processing*, pages 306–351. IEEE Press, 2001.
- [10] Nikos Nikolaou, Michael Makridis, Basilis Gatos, Nikolaos Stamatopoulos, and Nikos Papamarkos. Segmentation of historical machine-printed documents using adaptive run length smoothing and skeleton segmentation paths. *Image and Vision Computing*, 28(4):590 – 604, 2010.
- [11] N. Otsu. A threshold selection method from gray-level histograms. *IEEE Transactions on Systems, Man, and Cybernetics*, 9(1):62–66, Jan 1979.
- [12] Robert Sedgewick. *Algorithms in C*. Addison-Wesley, 3rd edition, 1998.
- [13] Irwin Sobel. History and definition of the so-called "sobel operator". 2014.