

Design and Prototype of a Support System for Archeologists to Decode Scripts on Mokkan

Akihito KITADAI ^a, Kei SAITO ^b, Daisuke HACHIYA ^c, Masaki NAKAGAWA ^d, Hajime BABA ^e,
and Akihiro WATANABE ^f

^{a, d} *Institute of Symbiotic Science and Technology Tokyo University of Agriculture and Technology*
Naka-cho 2-24-16, Koganei
184-8588, Tokyo, JAPAN

^b *Graduate School of Technology Tokyo University of Agriculture and Technology.*
Naka-cho 2-24-16, Koganei
184-8588, Tokyo, JAPAN

^c *Faculty of Technology Tokyo University of Agriculture and Technology*
Naka-cho 2-24-16, Koganei
184-8588, Tokyo, JAPAN

^{e, f} *Independent Administrative Instruction National Research Institute for Cultural Properties, Nara*
Nijyo-cho 2-9-1, Nara
630-8577, Nara, JAPAN

{ak ^a, kei ^b, hachiya ^c}@hands.ei.tuat.ac.jp, nakagawa@cc.tuat.ac.jp ^d, {hajime ^e, akihiro ^f}@nabunken.go.jp

Abstract. This paper describes a design and prototype of a support system for archeologists to read the "mokkan" excavated from Heijyo-kyo, Japan (the ancient palace in the Nara period from AD. 710 to 794). A mokkan is a wooden tablet on which text was written by a brush. Many mokkans were used in the Nara period. Since most of unearthed mokkans have been stained, damaged, and degraded, it is extremely difficult even for archeologists to extract characters from badly blurred or missing ink on a mokkan. We realized binarization functions based on discriminant analysis to extract ink from brownish background. Then, we developed a character recognition function, which outputs candidates for partially missing patterns. The aim of this recognizer is to output candidates even for degraded or partially missing character patterns rather than to read handwritten characters at the maximum speed. We made also a graphical user interface to invoke the above functions, provide experts with suggestions and stimulate their inference. This system supports collaborative work among expert readers by sharing decoding processes and intermediate results of extracted ink images and character recognition.

1. Introduction

Heijyo-kyo is a ruin of an ancient city in Nara, Japan. It was the capital from A.D.710 to 784 and was the center of politics and economy at that time. Since the capital was moved to Kyoto after Nara, Heijyo-kyo was buried under the ground and the area has been used as rice fields, we can now excavate an enormous amount of relics from under the ground. The soil in rice fields has been wet so that even wooden inheritances, which are fragile, oxidized or dried easily, have been kept well under the rice fields.

"Mokkan" is a Japanese generic name to call a wooden tablet on which text was written by a brush in India ink. Since wooden tablets were more accessible than other media to record handwriting and they had enough weatherability, people used them for various usages in the Nara period. Up to now, we have excavated about 350,000 mokkans in Japan and more than 170,000 of them have been excavated from the already excavated part of Heijyo-kyo. The number is increasing as we excavate the larger remaining part in Heijyo-kyo and other areas. By analyzing handwritten contents on excavated mokkans, we can acquire and extend the knowledge on the era. For example, by decoding mokkans used as luggage tags, we are able to know the flow of materials, relations among regions, condition of economy at that period and so on.

Although we find several preceding researches on information technologies for historical documents written on paper (Kim, Cho, Kwag & Kim, 2004; Shi & Govindaraju, 2004; Yan & Leedham, 2004; Gatos, Pratikakis, & Perantonis, 2004), no attempt has been made to use computers to process mokkans, especially such old wooden inheritances.

Since most of mokkans from under the ground have been stained, damaged and degraded, it is difficult even for experts on archeology to extract characters from badly blurred or missing ink on mokkans. Since it is difficult even for human experts, to seek the full automation of extracting and reading characters is impractical. However, image processing and character recognition technologies can be employed and incorporated into an interactive system which provides experts with suggestions and stimulate their inference.

This paper describes the design and implementation of the support system that helps archeologists to read mokkans. Section 2 presents the basic idea of applying the information processing for them to read mokkans. Section 3 shows the architecture of the support system. Section 4 describes each component of the support system. Section 5 draws conclusion.

2. Basic idea of a support system using information processing technologies

The experts on reading and decoding handwriting extract ink parts from a mokkan or its picture first. However, Very often, ink has been blurred, damaged or missing because:

- Color of ink parts has been faded out or decolorized.
- Color boundaries between ink parts and the background (skin of wood with grain) have been vague since the surfaces of wooden tablets have been turned darkish and stained.
- Some regions of a mokkan containing ink parts have been lost or broken.

For these reasons, the experts have to make conjectures or hypotheses on the lacked ink parts.

We consider that image processing, handwritten character pattern recognition (HCPR) and some technologies of information processing may assist the experts' work. At least, the results of the information processing can stimulate the inference of the experts.

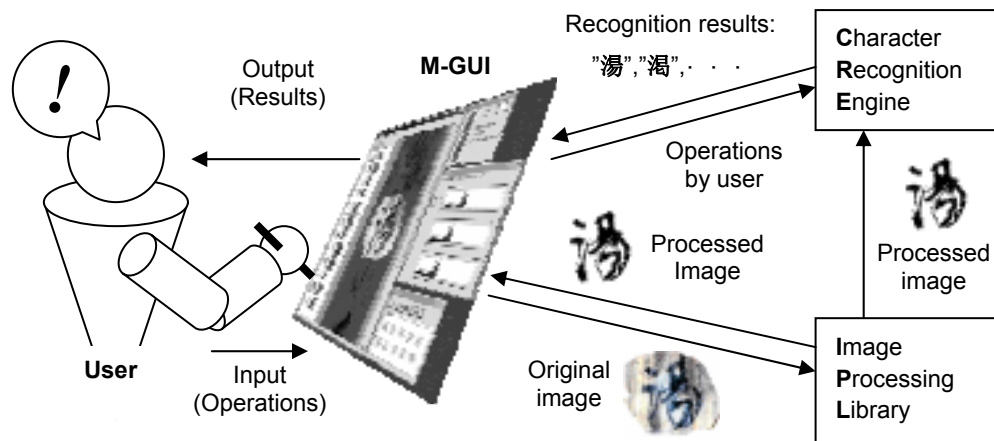


Figure 1. Architecture of the support system.

3. Architecture of a support system

We show the architecture of the support system in figure 1. The support system consists of three components. "Image processing library" (IPL) supplies the functions of fundamental image processing. "Character recognition engine" (CRE) provides HCPR for old Japanese characters used in Nara. Mokkan-GUI or M-GUI in short is the graphical user interface and it enables users to use IPL and CRE interactively. We describe the details of each component below.

4. Components of the support system

4.1 Image processing library

IPL has several functions to extract ink parts from the digital color image of a mokkan. As the results of the extraction, IPL outputs digital binary images in which ink parts are expressed as black pixels. The binary image is not only necessary for CRE but also helpful for expert users to read and decode handwriting on a mokkan. Also, IPL provides simple image processing methods to transform, enhance or shrink its image.

Since discriminant analysis (DA) can be used to binarize digital color and gray scale images, we have employed DA as the basic algorithm of ink extraction from morkkans with darkish grain, stained surface, and ink parts faded out.

(1) DA for gray scale images

In this method, DA is applied to a gray scale image generated from a color image of a mokkan. This is the simplest method of the ink extraction provided by IPL, but mostly effective for images without heavy stain and darkish grain.

(2) DA for each elemental channel of primary colors

IPL provides a function of DA for each channel of primary colors used to extract ink parts from the image of a mokkan with heavy stain or darkish grain. First, IPL generates three images by extracting each element of color channels: red (R), green (G) and blue (B). Each image is converted to a binary image by DA. By overlapping three binary images, IPL generates an eight-colored image consists of R, G, B, R+G, G+B, B+R, R+G+B (white), and black (black pixels in every binary image) domains. Users can remove any color domains among the eight colors to obtain the image of ink parts. Figure 3 shows an example: since the principle color element of skin of wood containing darkish grains is brown, we can obtain a fine image of ink parts by removing R and R+G domains.

We also employ CMY or CMYK color channels as well as RGB (C: cyan, M: magenta, Y: yellow and K: black). Fig. 4 shows a successful case of ink extraction. Since the skin of the wood contains yellow strongly, we can generate a good ink part image by adapting DA to the image of the yellow channel element extracted from the color image of the mokkan.

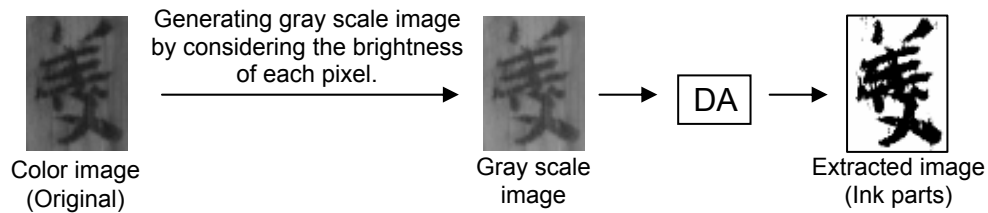


Figure 2. DA for gray scale image.

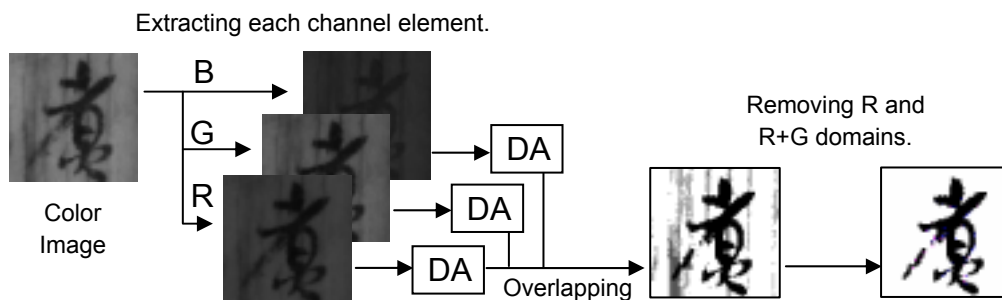


Figure 3. Ink parts extraction by using DA for each color element.

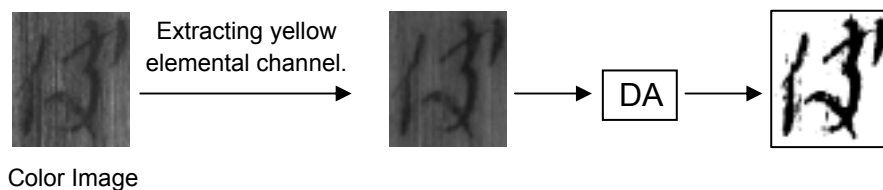


Figure 4. Ink extraction by using DA for the yellow channel element in CMYK.

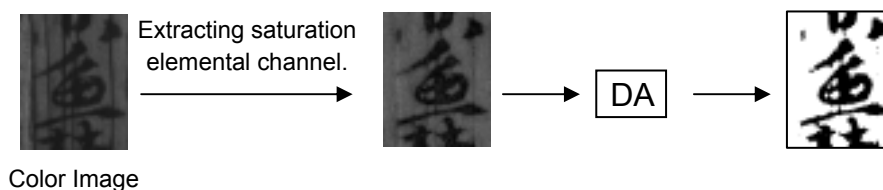


Figure 5. Ink parts extraction by using DA for saturation channel element in HSV.

(3) DA for each elemental channel of HSV

IPL provides a function that adapts DA to each channel of hue, saturation, and value (HSV) composing the color image of a mokkan with heavy stain or darkish grain. Fig. 5 shows an example that DA for the saturation elemental channel can extract a fine image of ink parts from the color image of a mokkan containing extremely grayed grain. Such conversions of color channels worked well for historical paper documents (Shi & Govindaraju, 2004).

(4) Other image processing functions

IPL has control functions of brightness, contrast, and color balance of digital images used to generate eye-friendly images for experts.

4.2 Character recognition engine

We can find some effective method to recognize complete character pattern (Liu, Liu, & Dai, 1997; Okamoto & Yamamoto, 2004; Velek & Nakagawa, 2002). However, the character recognition engine in this system should output candidates even for degraded or partially missing character patterns and stimulate imagination of experts. The most important problem is to output candidates even for character patterns with missing ink parts. A realistic solution is for experts to augment missing ink parts roughly although precise augmentation cannot be expected. Therefore, we extended our HCPR method to accept a ternary image in which ink parts are expressed in black, background in white and augmented parts in gray.

Fig. 6 shows an extreme example of augmentation. Since the right half of the mokkan has been lost, we can conjecture that all character patterns in the mokkan have lost their right-hand sides. In this case, the users can only augment the right-hand side by drawing a rectangle in gray by pen or mouse. In other cases, the users can trace blurred ink.

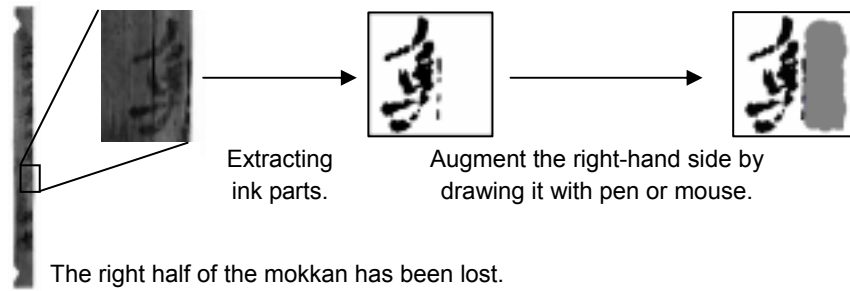


Figure 6. Augmenting information to the character pattern with lacked parts.

Our CRE first apply nonlinear normalization to a ternary image. Second, it divides the normalized image into an array of cells. Third, it extracts 8-directional features from each pixel as shown in Fig. 8. For each pixel P_i , the directional feature of P_i in each direction $F(P_i, d)$ is defined by eq. (1) where d is from 0 to 7, P_x denotes P_i or one of its 8-neighbor pixels P_d and $C(P_x)$ takes the value 0 (white), 1 (gray), or 2 (black). Fourth, within each cell, the directional feature $F(P_i, d)$ of each pixel P_i is summed up for every direction.

$$F(P_i, d) = \left\{ 1 - \frac{|C(P_i) - C(P_d)|}{2} \right\} \times \{C(P_i) + C(P_d)\} \quad (1)$$

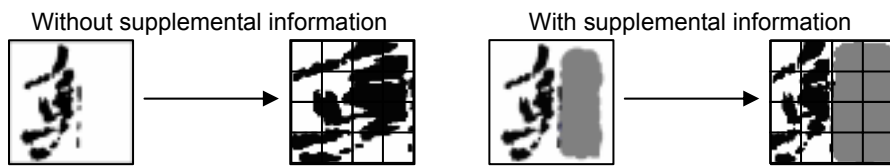


Figure 7. Nonlinear normalization with/without supplemental information.

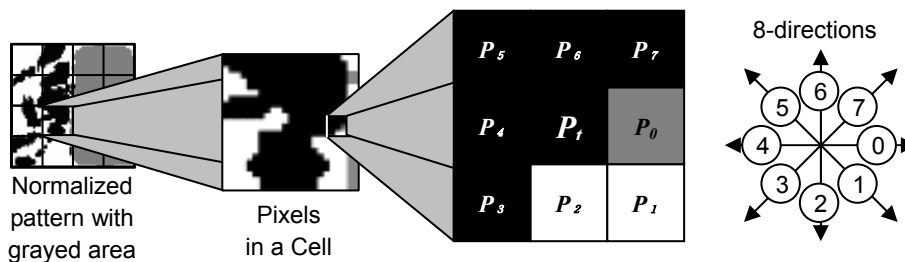


Figure 8. Directions and pixels.

Augmented gray ink has two roles. One is for non-linear normalization to normalize the original ink parts properly without expanding it to the whole character size as shown in Fig. 7. As a result, each directional feature will be extracted from the cell at correct place. The other role is to provide white noise to missing directional features. This is better than null features to guess original character patterns.

4.3 M-GUI

M-GUI provides the functions of IPL and CRE to users of the support system. Moreover, it manages the process of reading and decoding handwriting on mokkans. Since M-GUI can save, manage, and load the workspaces consisting of HCPR results, edited images, et al., the users of the support system can suspend and resume their work at any time and can share the workspaces for their collaborative work. Moreover, M-GUI can export the HCPR results and edited images to word processors, image editors and other application software. Fig. 9 shows a screen shot of M-GUI.



Figure 9. M-GUI.

5. Conclusion

We described the support system that helps experts to read mokkans. Ink extraction by image processing and candidate suggestion by character recognition are the main results. Since we are now constructing a pattern database of old characters, the evaluation of the support system and each component are our future work.

Acknowledgment

This work is supported by Grant-in-Aid for Scientific Research under the contract number S:15102001

References

- Gatos, B., & Pratikakis, I., & Perantonis, S. J. (2004). An Adaptive Binarization Technique for Low Quality Historical Documents. *Proc. of the 6th Workshop on Document Analysis Systems -DAS 2004*, Florence, Italy, September 8-10, pp. 102-113.
- Kim, M., S., Cho, K., T., Kwag, H., K., & Kim, J., H., (2004). Segmentation of Handwritten Characters for Digitalizing Korean Historical Documents. *Proc. of the 6th Workshop on Document Analysis Systems -DAS 2004*, Florence, Italy, September 8-10, pp. 114-124.
- Liu, C., L., Liu, Y., J., & Dai, R., W. (1997). Preprocessing and Statistical/Structural Feature Extraction for Handwritten Numeral Recognition. *Progress of Handwriting Recognition*, A.C. Downton and S. Impedovo eds., World Scientific, pp. 161-168.
- Okamoto, M., & Yamamoto, K. (1999). On-line Handwritten Character Representation using Directional Features and Direction-Change Features. *Journal of IEE Japan*, 119(3), pp. 358-366. (in Japanese).
- Shi, Z., & Govindaraju, V. (2004). Historical Document Image Enhancement Using Background Light Intensity Normalization. *Proc. of the 17th Int. Conf. on Pattern Recognition -ICPR2004*, Cambridge, UK, August 23-26, (2aP.Mo-ii).
- Velek, O., & Nakagawa, M. (2002). The Impact of Large Training Sets on the Recognition Rate of Off-Line Japanese Kanji Character Classifiers. *Proc. of the 5th Workshop on Document Analysis Systems -DAS'02*, Princeton, US, August 19-21, pp. 106-109.
- Yan, C., & Leedham, G. (2004). Decompose-Threshold Approach to Handwriting Extraction in Degraded Historical Document Images. *Proc. of the 9th Int. Workshop on Frontiers in Handwriting and Recognition -IWFHR-9*, Tokyo, Japan, October 26-29, pp. 239-244.