

Cherry Blossom: A System for Japanese Character Recognition

Sargur N. Srihari, Tao Hong and Zhixin Shi

Center of Excellence for Document Analysis and Recognition

State University of New York at Buffalo

Buffalo, New York 14228

email: srihari@cedar.buffalo.edu

Abstract

A general purpose Japanese character recognition system, Cherry Blossom, has been developed at CEDAR in past years. It is designed to recognize Japanese document images in low resolution or with poor print quality. The system includes modules for page skew correction, document segmentation, text segmentation, character recognition and postprocessing. The API code for each module has been developed so that each module can be tested as a stand-alone program or can be called in large application systems such as a document indexing and retrieval system.

In the character recognition module, two classification methods, the nearest-neighbor classifier and the subspace method, have been integrated in an efficient way. The speed of character recognition is about six characters per second from the original one character per second. New techniques, including dynamic feature selection, incremental nearest prototype search and visual similarity analysis, have been developed to speed up character classification while keeping high accuracy. A linguistic postprocessing module improves the accuracy of character recognition by exploiting linguistic context. It was trained on a Japanese text corpus with above 70 million characters.

The JOCR system has also been adapted to recognize Chinese documents. Because the similarity between Chinese and Japanese character sets, the Chinese recognizer can be developed rapidly by simply training on font images and by sharing the information of the well-trained Japanese classifier. A Unicode-based OCR for Far East Languages is being developed based on the JOCR system.

We also conducted experiments on representing Japanese OCR results in HTML files for evaluation, error correction and reading purposes. A real-time GUI program which demonstrates different aspects of the system for Japanese and Chinese document recognition is now available on SunOS.

1 Introduction

The Japanese optical character recognition project at CEDAR is to develop a prototype system for machine-printed Japanese document recognition [33]. Intensive development and research have been conducted to improve accuracy of document recognition on degraded documents such as facsimile pages or images scanned at low resolution.

The Japanese character set is a combination of several scripts: Kanji (Chinese characters), Kana (Japanese consonant and syllabary characters), Roman characters and Arabic numerals. Kanji are Chinese ideographs which were introduced into Japanese. In the Japanese standard, JISx0208-1990 [25], there are 6,879 characters. Among them, 6,355 are Kanji, divided into two distinct sections – JIS Level *I* and Level *II*; where Level *I* spans 2,965 of the most frequently used Kanji characters.

Comparing with the task of developing an English document recognition system, there are more difficulties on developing a Japanese OCR system. Some of those difficulties are:

1. Large variety of page layouts in Japanese documents; For example, the direction of a text line can be horizontal or vertical.
2. Large number of character categories; As described above, there are more than 3,000 character categories which are needed to be recognized.
3. Structural complexity of character patterns; Most Japanese characters are Kanji. The stroke structure of a Kanji usually is complex. Many characters have very similar stroke structure.
4. Large variety of character shapes due to different typeface and image quality in machine-printed documents.

As in systems for English document recognition, a Japanese document recognition system has three

major components: layout analysis, character segmentation/recognition, and postprocessing. Layout analysis techniques for Japanese documents [2, 26] and multi-lingual documents [18, 35] have been developed. Given an image of a document page, layout analysis will locate text blocks and further segment them into text lines which can be either vertically or horizontally oriented. Because word boundaries are usually not visible, character images are directly segmented from the text lines. Connected-component analysis, projection profile analysis and feedback of character recognition result can be used for character segmentation [2, 3, 27].

Because Kanji characters are complex in their stroke structure and many characters share structural similarities, feature descriptors are typically with large dimensions in order to effectively discriminate so many Kanji characters. Stroke features (such as the direction contribution features [2]) and background features (such as surrounding area features [32, 30]) are two types of statistical features which have been proved very useful. Many sophisticated similarity/distance measures have been proposed to achieve more reliable performance of character recognition. Methods such as the multiple similarity analysis [17] and the modified quadratic discriminant functions [20] have been applied to hand-printed and machine-printed character recognition.

Two classifiers have been designed for character recognition in the system [14]. One is the Minimum Error Subspace method; another is the Nearest Neighbor classifier [11]. The features are called "Direction Contributivity" features [2, 28]. The NN classifier has been generalized to use the distance measurement defined in the Minimum Error Subspace method which has been proven to be more accurate than Euclidean distance. Because its confidence calculation is trustworthy, for a character image, the NN classifier provides only the first decision if its confidence score is high; otherwise, the top 5 choice list is output. Classifiers were trained on 200 ppi character images from CEDAR's Japanese character image database [16], and were tested on 200, 300 and 400 ppi images. On 200 ppi character images from CEDAR's test set, the accuracy of the NN classifier is 95.9%. The speed of character recognition varies from 3 to 6 characters per second, depending on image quality.

A linguistic postprocessing module has been added to the system in order to further improve the accuracy of character recognition by exploiting linguistic context. It was trained on a Japanese text corpus with above 70 million characters.

The JOCR system has also been adapted to recognize Chinese documents. Because the similarity between Chinese and Japanese character sets, the

Chinese recognizer can be developed rapidly by simply training on font images and by sharing the information of the well-trained Japanese classifier. A Unicode-based OCR for Far East Languages is being developed based on the JOCR system.

In the system, we also developed methods for visual context analysis, such as image-based character image clustering and inter-character similarity analysis [12, 9]. The results of visual context analysis can be used for OCR error correction [10].

The system is implemented on a SPARC Workstation under SunOS. Given a TIFF image which is a text block or a document page, the system will execute its modules and save the details of analysis and recognition into a file in a special format.

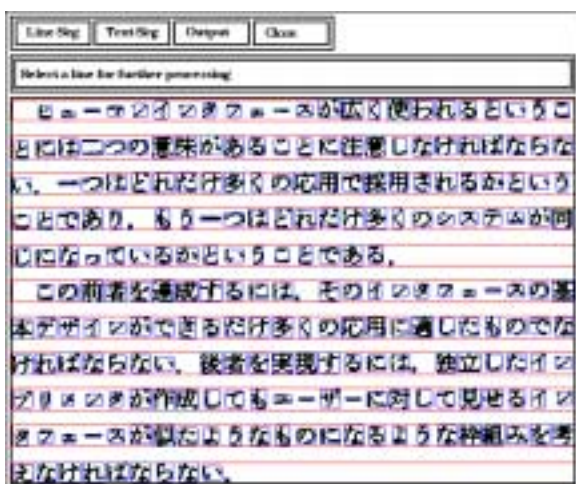
Two X-Window-based GUI programs for system evaluation have been designed using C/C++ and Motif. By using those programs, deskewing, text segmentation, line segmentation and character segmentation can be performed either manually or automatically. For each segmented character image, the top 5 choices are provided by a character classifier. The programs allow a user to use a mouse to manually correct mistakes in segmentation and recognition.

Although these GUI programs are very useful for system evaluation, much work is needed to modify them for the purpose of evaluating different aspects of our OCR system. In order to rapidly prototype visualization and evaluation tools which are more flexible and less platform-dependent, we recently developed a framework to convert the recognition results into a set of HTML files. HTML handles symbols and images uniformly and allows hyperlinks between them. A Web browser such as Netscape Navigator comes with multi-lingual support and can be run on many platforms. Scripts can be easily written to convert recognition results in any level into HTML files so that different aspects of the system can be examined within a local network (Intranet) or on the Internet. Much effort in GUI design, multi-lingual and multi-platform support can therefore be avoided.

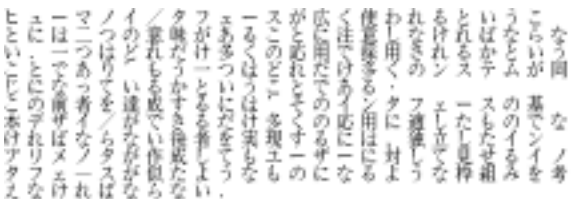
Section 2 gives an overview of the JOCR system design, including its modules for page deskewing, page segmentation, text alignment detection, line segmentation and character segmentation. Section 3 describes the design of two classifiers for character recognitions. Section 4 discusses extending the JOCR system to be multi-lingual recognition system. In section 5, a postprocessing module is described and in section 6, we present an experiment on evaluating Japanese document recognition using HTML. Finally, conclusions and future directions of this project are briefly discussed in section 7.

2 The Overview of the System Design

Cherry Blossom is a prototype system for Japanese document recognition developed at CEDAR in past years. This system is designed to meet several challenges of machine-printed Japanese document recognition: the wide variations in data quality (facsimile, photocopy, newspaper etc.), low scan resolution (e.g., 200ppi), large character set (about 3,000 Kanji characters in the first level JIS code) and the variety of font and point sizes. Given a Japanese document page, it can deskew the image, segment the page into blocks, discriminate text/non-text blocks, determine the alignment of a text block, segment text regions into lines and further into character images, and recognize character images. The system consists of API modules including deskewing, page segmentation, text segmentation, character recognition and postprocessing. In the system, OCR results can be encoded in JIS, Shift-JIS, EUC and Unicode.



(a)



(b)

Figure 1: Character segmentation and recognition. (a) shows line and character segmentation result of a text block. (b) is the recognition result displayed in a *kterm* window.

The design of the modules for deskewing, page segmentation and text segmentation are described as follows. The character recognition module will be presented in the next section.

To detect and correct skew in a scanned document image, a fast directional profile analysis method is used[23]. The range of skew angle handled by the method is about ± 30 degree. After detecting skew angle, a simple affine transformation is applied for skew correction. The skew detection program works well on document images which often contains graphics, tables and figures.

The skew detector was tested with 467 artificially skewed images, with varying amounts of noise, fragmentation, black pixel density and ratio of text to graphics. The document analyzer achieves 97% accuracy within ± 1 degree of the actual skew on good quality document images, and 86% on noisy or complex documents. The average time required for skew detection is 2 cpu seconds on a SPARC-10.

The objective of page segmentation is to decompose a scanned document image into regions which contain homogeneous entities, such as text, graphics and half-tones. A *local-to-global segmentation algorithm* has been designed in *Cherry Blossom* [21]. The page segmentation module was tested with 200 Japanese and English documents. The performance of the document segmentation ranges between 92% completely correct segmentations and 94% partially correct segmentations. On average this module requires 27 cpu seconds on a SPARC-10.

New methods for page segmentation are being investigated. They include smoothed-run-length analysis, background pattern analysis and minimum spinning tree analysis methods. A hybrid method which can achieve better performance on different types of Japanese document pages is being developed.

The text segmentation module is designed with the objective of segmenting a given text block into line images and further into character images [22]. In the system, a five-stage procedure – image quality estimation, character size estimation, text alignment detection, line segmentation and character segmentation, are implemented to fulfill the task of text segmentation.

At the stage of image quality estimation, the system determines aspects of image quality, which will be useful for line and character segmentation. At the stage of character size estimation, character size is estimated globally based on histogram plots of connected components. In clean images, the highest peak in the component histogram corresponds to the width of full pitch characters. In noisy images, this component histogram is smoothed and the mean width of the histogram is evaluated. At the

stage of text alignment detection, methods from two different approaches are implemented to determine whether text is vertically or horizontally aligned. The first approach is based on the analysis of a minimum spanning tree of connected components [18]; and the second one is based on projection profile analysis.

Projection profile analysis is also used at the stage of line segmentation. For character segmentation, character components are initially located by performing a projection profile analysis of black pixels and identifying white spaces as component separators. After local character size estimation, a rule-based approach is then used to partition components into three groups. These groups identify single component characters, partial character components and touching characters (or multi-character components). The sub-character components are merged based on a few simple rules to form complete characters. Touching characters are split based on further analysis and heuristics. At this stage a character recognizer has also been used to validate proposed splitting points.

The text and character segmentation is evaluated with respect to several criteria. Text segmentation performance varies between 98.5% for good quality images and 94% for degraded images. Simple connected component analysis is used for high quality printed text, while more complex analysis is applied for low quality document images.

3 Japanese Character Recognition

Two independent character recognition methods have been developed to recognize more than 3,300 Japanese characters. One classifier uses the minimal error subspace method; another is a fast nearest-neighbor classifier. After implementing and testing many types of feature sets reported in Japanese OCR literature, two feature sets were chosen for further testing because their performance on low quality images was encouraging. These are the Local Stroke Direction (LSD) and Gradient, Structural & Concavity (GSC) feature sets. These feature sets are described below, followed by discussing the design of the minimal error subspace classifier and the nearest neighbor classifier. Results of both classifiers as well as their integration will also be reported.

3.1 Local Stroke Direction Features

Local stroke direction (LSD) features are based on direction contributivity analysis [2]. When given a character image, its LSD feature vector can be computed as follows [28] (see Figure 2 for example): first, for each black pixel, compute its directional run-length for each of four directions in Figure 2 (b) and normalize it as a ratio to the total run-length in all

directions; second, partition the image into $n \times n$ areas and compute the directional run-length of each area as an average of the pixels in the area. Here, we choose that n equals 8. Therefore, the size of LSD feature vector is $4 \times 8 \times 8 = 256$. In Figure 2 (c), the values in the LSD feature vector are scaled to integers at the range of 0 and 255 so that they can be visualized as a grayscale image.

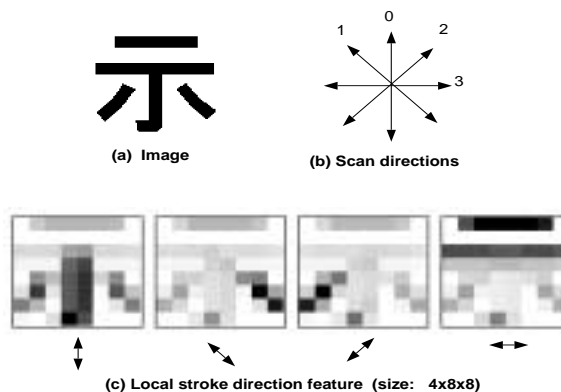


Figure 2: A Kanji character image and its *local stroke direction* (LSD) feature vector

3.2 Gradient, Structural and Concavity Features

The gradient and structural features encode local structure, while the concavity features are global descriptors extracted from binarized images. A gradient map is constructed from the normalized image, by estimating gradient value and direction at each pixel. Histograms of directional features are recorded in each region – indicating the presence (or absence) of a small number of oriented ranges of gradients in the region. Directional histograms from each region are concatenated into a fixed-length *gradient feature* vector. *Structural features* are computed from the gradient map by examining similarities in gradient direction in a local neighborhood of each pixel. These features record curvature in an approximate fashion. Curvature histograms indicating presence of changing orientation of character contours are estimated in each region. These histograms are concatenated from each region in a fixed-length structural feature vector. Gradient and structural features have been described in more detail in [34]. *Concavity features* are coarse global descriptors and are of three kinds: pixel density, large stroke, and true concavity. The large stroke features encode horizontal and vertical strokes in the image. Concave regions enclosed with a character are also recorded.

Concavity features have been described in detail in [5].

3.3 Minimum Error Subspace Classifier

Subspace methods have been a major field of study in pattern recognition [29, 6], particularly for classification and feature subset selection. The minimum error subspace classifier is a discriminant function derived from the Karhunen-Loeve expansion. It is given by:

$$g_j(X) = |X - M_j|^2 - \sum_{i=1,k} \{\phi_{i_j}^T(X - M_j)\}^2.$$

Here, $g_j(X)$ defines the discriminant classifier for the j^{th} class, given an input measurement/feature vector X ; M_j is the mean vector and $\phi_{i_j}^T$ is the transpose of the i^{th} eigenvector of the covariance matrix for the j^{th} class; k is chosen as the dimension of the subspace defined by the dominant k eigenvectors. Based on our experiences, k is set to 5 in our experiments. An unknown test feature vector is evaluated with these discriminant functions for each class, that is, $g_j()$, for $j = 1, 2, \dots, C$, where C is the number of classes. The class of X is determined as that of the discriminant function for which the residual error is least, that is: $g_J(X) \leq g_j(X)$, $\forall j \in \{1, 2, \dots, C\}, j \neq J$.

Projections of an unknown test vector on the subspaces defined by the dominant eigenvectors of each class are subtracted from the projections on the entire eigen-space. The subspace of the class which best represents the unknown feature vector, results in the least residual error (and results in the least difference in projections between the whole space and the subspace of the principal eigenvectors). The class corresponding to the subspace with least residual error is chosen as the class-identity of the unknown test vector.

This classifier can also be derived from the modified quadratic discriminant function (MQDF)[19], and can be interpreted as a quadratic discriminant classifier.

The training phase of this classifier requires the computation of the covariance matrices and their eigen-decomposition. For each class, the eigenvectors are then sorted in decreasing order of the corresponding eigenvalues. The mean vector M_j is evaluated for each class. Only the dominant eigenvectors (k , in this case) are used in the subsequent processing. During testing, the unknown or test feature vector is projected onto the subspace defined by these dominant vectors, for each class. The term $|X - M_j|^2$ defines the projection of the test vector onto the whole eigen-space. Hence, the class of an

unknown vector is determined by the best subspace projection, given by the least residual error.

The classifier had been trained and tested on the CEDAR dataset. Using the LSD feature and the GSC feature, results of the ME (Minimum Error) Subspace Classifier are listed in Table 1 (a).

(a) ME Subspace Classifier ($k = 5$)					
LSD			GSC		
Top 1	Top 5	Top 10	Top 1	Top 5	Top 10
93.48%	97.86%	98.45%	93.93%	97.63%	98.19%
(b) NN Classifier (using Euclidean distance)					
LSD			GSC		
Top 1	Top 5	Top 10	Top 1	Top 5	Top 10
92.73%	97.33%	98.01%	93.06%	97.88%	98.60%

Table 1: Character recognition performance of different classifiers on CEDAR’s 200 *ppi* dataset (57,571 images from 3,319 classes).

3.4 A Nearest-Neighbor Classifier

A NN (nearest neighbor) classifier has been designed for Japanese character recognition. New algorithms for prototype reduction, hierarchical prototype organization and fast NN search are implemented in the classifier [11, 14]. In these algorithms, k -nearest/farthest neighbor lists are used to estimate the distribution of samples in the feature space. Given a set of samples, $P = \{p_0, p_1, \dots, p_{n-1}\}$, for each sample p_i , its k -nearest neighbor list, $N_i = \{n_{i0}, n_{i1}, \dots, n_{i,k-1}\}$, and its k -farthest neighbor list, $F_i = \{f_{i0}, f_{i1}, \dots, f_{i,k-1}\}$, have to be pre-computed.

Given a training set, prototype reduction leads to the selection of a subset of samples which can represent the whole training set. Figure 3 illustrates the process of prototype selection for a class. The subset can be generated by deleting redundant samples from the training set. By checking the pre-computed nearest neighbor list of each training sample, the prototype reduction algorithm decides whether a training sample can be deleted without negative effect on the correct classification of itself and other samples. Previous methods, such as Hart’s *condensed nearest neighbor rule* (CNN) and Gates’ *reduced nearest neighbor rule* (RNN), have to call the procedure of classification iteratively to test whether the deletion(or adding) of a prototype affects the correct classification of the other prototypes[8, 7]. The advantage of the method here is to avoid such an iterative process.

The NN classifier using LSD feature and Euclidean distance measure was trained and tested on the data sets from CEDAR’s CDROM. There are

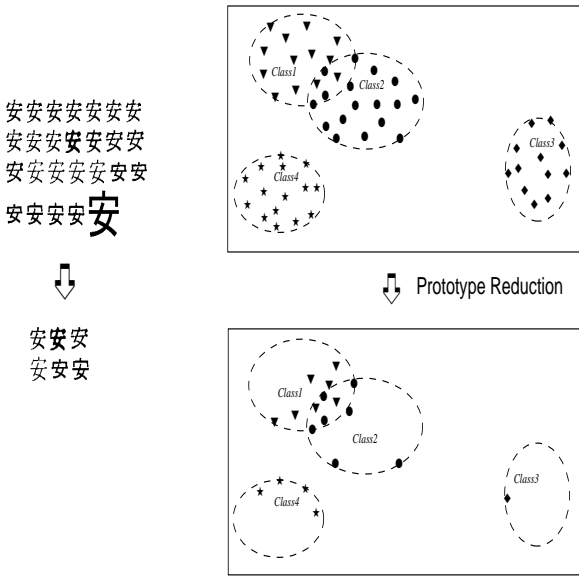


Figure 3: Prototype reduction

118,417 character samples from 3,354 classes in the training set. The number of prototypes derived by prototype reduction is 24,273.

A brute-force NN algorithm suffers from its computational complexity because every prototype has to be matched with the input pattern to see whether it is close to the input pattern. To speed up NN classification, several methods are proposed. An inter-character visual similarity constraints is firstly used. The preliminary experiment shows very promising. Given a text block with 1,721 segmented character images, the system takes less than 17 minutes to recognize it (originally it needs about 29 minutes). In another word, we can improve our classifier to 0.58 second per character from originally 1.0 second per character with the only change.

The idea behind the new method is simple. In a Japanese document page, many characters occur more than one time. An image clustering procedure can be used to group those different image entries of the same character. After image clustering, the NN classifier will recognize characters in a group. After recognizing the first entry of a group by matching with all prototypes, the system will recognize the rest of group members by matching with a small number of prototypes which are in the best match list of the first entry. The size of the reduced prototype set can be less than one tenth of the number of all prototypes. Sometimes the clustering procedure may cause errors so that visually similar characters from different JIS categories are grouped together because in Japanese there are many characters which

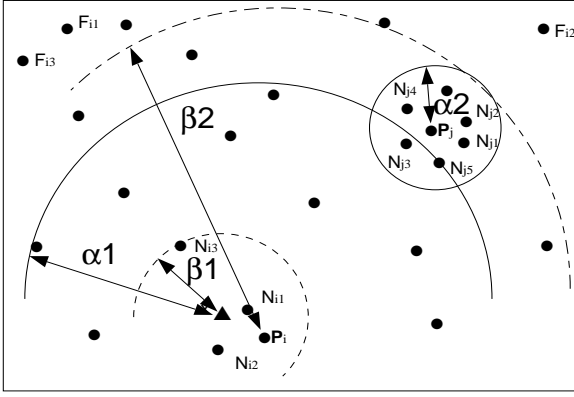
are very similar in shape. The method can tolerate this kind of clustering errors very well. On the example text block with 1,721 segmented characters, 508 clusters were generated. While it takes 0.95 second to recognize the first entry of each group, it needs only 0.25 second to recognize each of the other entries in a group. On average, 0.58 second is needed for each character. The accuracy of classifier is almost the same as the original classifier.

A dynamic feature selection for feature matching has been developed. In our implementation of the LSD feature, there are 256 integer elements. The size of the feature vector is fairly large. The NN classifier runs slowly because it has to do comparison on each of those elements between the test pattern and a prototype. It will be much faster if only a small subset of feature elements can be selected for matching. Given a character image, after extracting its LSD feature vector, we can dynamically select such a subset of element according to the pivotal directions in the feature extraction algorithm.

we have also developed a fast NN search algorithm, *FNN*. The basic idea is to avoid unnecessary comparisons. Suppose there is a test sample X to be classified (see the triangle in Figure 4). It has to compare with prototypes from the prototype set. After the comparison between the input pattern X and a prototype P_i , we know that the distance between them is very small. For those prototypes in P_i 's k -farthest neighbor list, such as F_{i1}, F_{i2} and F_{i3} shown in Figure 4, without going further to compare each of them with the input sample X , we know the distance will be very large and therefore they can not be in X 's k -NN list. Similarly, after the comparison between X and a prototype P_j , we know that the distance between them is very large. For those prototypes in P_j 's k -nearest neighbor list, such as N_{j1}, N_{j2}, N_{j3} and N_{j4} shown in Figure 4, without going further to compare each of them with the input sample X , we know the distance will be very large and therefore they can not be in X 's k -NN list.

An incremental feature match algorithm was also designed which firstly dynamically select a small number of feature elements (about one sixth of the feature size) and do NN search, then use a larger number of feature elements to refine the NN search result on a small number of prototypes with small match distance.

In order to further speed up the search process, The prototype set can be organized into a hierarchical representation. The prototype set is divided into two levels. The first level is the so-called "*centroid set*." For a centroid, there may be a set of prototypes which are stored in the second level. The set for a centroid is called as its "*package*." A centroid can approximately represent those prototypes in its



Note: ● denotes prototype; ▲ is an input pattern x.

Figure 4: Speed up nearest neighbor search by avoiding redundant comparisons

package. The process to create such a hierarchy is to pack prototypes for selected centroids.

After packing prototypes, the k -nearest neighbor list and k -farthest neighbor list of each prototype in the centroid set C will be computed. The NN classifier can be described as a two-step procedure. Given a test sample X , its approximate k -nearest neighbors \overline{N}_X in the centroid set C can be calculated using the search algorithm described above. Then the prototypes in the packages of those centroids in \overline{N}_X will be compared with X to generate the final k -nearest neighbors N_X for the sample.

In the NN classifier, distance between a test pattern and a prototype can be measured in many different ways, from simple metrics to sophisticated ones. We first tried the Euclidean distance, then changed to use both the Euclidean distance and the distance defined by the *ME Subspace* method.

Using the Euclidean distance, the performance of the *FNN* classifier on the test set is shown in Table 1 (b). The prototype set used in our NN classifier has 24,273 prototypes. The number of centroids is 1,919 if we define the package size as 50. By comparing with 2,195 prototypes on average the system can achieve 92.7% top 1 accuracy. Performance of the classifier using the GSC feature set is also reported in the table.

By using the Euclidean distance, the accuracy of the NN classifier is lower than that of the subspace classifier although there are multiple prototypes per class in the NN classifier while there is only one prototype for each class in the subspace method. It seems to us that the distance measurement defined in the subspace method is much more accurate than the Euclidean distance. Therefore, we designed a new NN classifier to use the subspace distance. Because the subspace method is time-consuming, we

proposed to use both the Euclidean distance and the subspace distance in the classifier. Given a test pattern P and n prototypes, steps to recognize P are:

- Step 1 – find first n_1 closest prototypes of P among n prototypes based on Euclidean distance;
- Step 2 – refine the distance scores of the pattern P with n_1 prototypes by using the subspace distance;
- Step 3 – re-rank n_1 prototypes by their refined distance scores and output the top 10 choices.

In our experiments, n is 24,273 and n_1 is set to 100. In the design, instead of computing the subspace distance between the test pattern and each of 3,300 classes, only n_1 subspace distances have to be computed. The Euclidean-distance-based NN algorithm for *Step1* can be a brute-force NN algorithm or the fast NN search algorithm described above.

Originally each class in *ME Subspace* method has a set of eigen vectors and a mean vector. In the modified NN classifier, we also generalized the subspace method to use multiple prototypes for each class. In the new classifier, each class has a set of eigen vectors and a set of prototypes. This extension gives the capability of using the subspace method in NN classifier so that reliable decisions can be made based on majority voting.

The new NN classifier was designed in order to integrate the NN search method with the *ME Subspace* method. It has also been trained and tested on 200 ppi datasets from CEDAR’s character image database. Table 2 lists results of the classifier in which *Step1* was implemented as a brute-force NN classifier and the FNN algorithm described above. By using the FNN algorithm, each test image was compared with about 6,203 prototypes out of 23,449 prototypes. The LSD feature is used in the classifier here.

To generate reliable confidence information for OCR decision, we proposed a method to compute the confidence score for each choice from the NN classifier [14]. In the method, confidence scores are calculated based on both majority voting result and distance information.

Given a test pattern t , a NN classifier can provide two types of candidate lists: one is the nearest prototype list

$$P = p_1, p_2, \dots, p_i, \dots, p_n$$

where p_i is a prototype. The identity of p_i is $JIS(p_i)$ and the distance between p_i and t is $DIST(p_i)$. In

	Number of comparisons	Accuracy		
		top1	top5	top10
Brute-force NN search	23,449	95.92%	98.09%	98.57%
		(55,221/57,566)	(56,471/57,566)	(56,748/57,566)
Fast NN search	6,203	95.83%	97.97%	98.43%
		(55,166/57,566)	(56,401/57,566)	(56,667/57,566)

Table 2: Performance of NN classifiers using Euclidean distance and Subspace method. The only difference between two classifiers is that in *Brute-force NN* search the test pattern has to compare with every prototype while only a portion of prototypes are compared in the *Fast NN* search algorithm.

P , prototypes may have the same identity. Another is the choice list

$$C = c_1, c_2, \dots, c_j, \dots, c_m$$

where c_j is a JIS code. In C , choices have to be distinct. For each c_j , $DIST(c_j)$ is defined as $DIST(p_i)$, where p_i is the first prototype in P so that $JIS(p_i) = c_j$.

For each c_j , we can compute its confidence as

$$CONF(c_1) = k/3.0$$

where k satisfies

$$JIS(p_1) = \dots = JIS(p_k) \neq JIS(p_{k+1})$$

and otherwise

$$CONF(c_j) = DIST(c_1)/DIST(c_j)$$

Based on above experiments, a new system configuration is implemented for the character recognition. Significant improvement on classification speed has been achieved. The improvement is the result of the integration of a flexible control structure, image clustering algorithm, prototype packing method, dynamic sub-feature selection method, incremental nearest-neighbor searching algorithm and the modified quadratic discriminant function. This progress makes our system a step closer to a real-time system. Figure 6 is the flowchart of the character recognition module in the current system. As illustrated in the figure, although there are seven steps for character recognition, character images do not have to go through all of them. For a character passes all seven steps, it takes about 0.4 second. Many character images in large clusters only have to pass Steps 1, 2 and 6 and only have to compare their feature vectors with 500 prototypes. This takes less than 0.1 second.

Figure 5 shows character recognition result of the NN classifier on images of different quality. In the figure, only the first choice is provided if its confidence is high; otherwise, the top five choice list is output and further postprocessing is needed to make a decision.

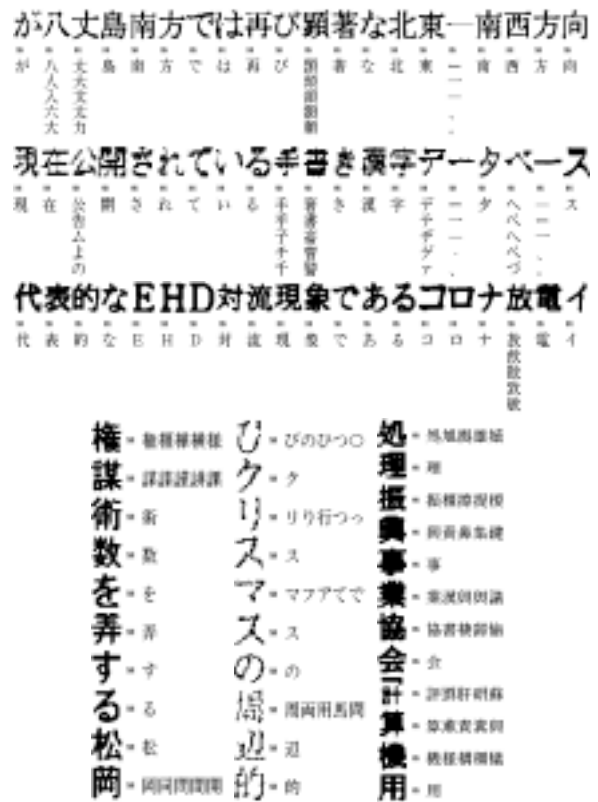


Figure 5: Character recognition. The NN classifier is used. For each character image, the classifier gives a candidate list. Because the confidence information from the NN classifier is reliable, the classifier provides only top first candidate if its confidence is high.

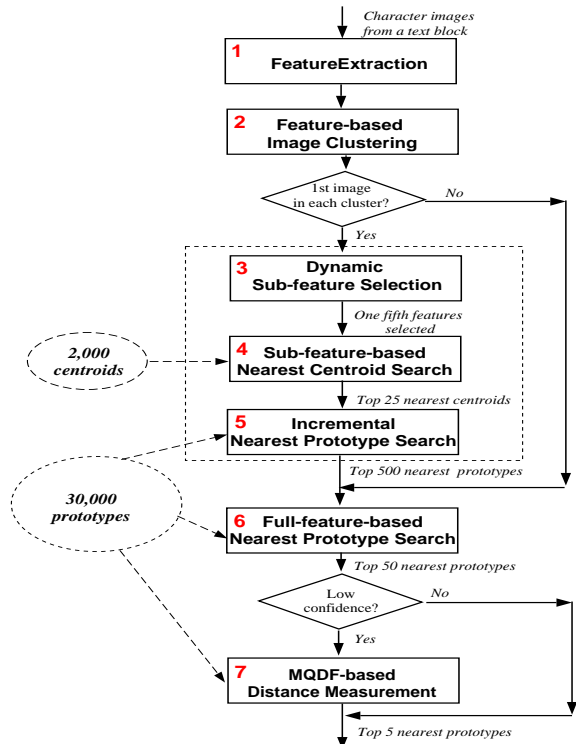


Figure 6: The framework of character recognition in the JOCR system

4 Multi-Lingual Document Recognition

Based on our developed Japanese document recognition system, we also investigated recognition for documents in other languages such as Chinese and Korean.

A new method is proposed for rapidly prototyping a recognition system for documents written in a new language (such as Chinese) based on the well-trained JOCR system. The new recognition system uses only the font samples and does not need to be trained on real world image samples. As a natural result, a Unicode-based OCR for Oriental languages (such as Chinese, Japanese and Korean) has been developed based on our JOCR system.

It usually takes long time and high cost to develop a high accuracy OCR for documents in a different language. Much effort is usually on training the OCR on many scanned document images. Although font images are easier to be collected than real world document samples, OCR using only font images usually can not achieve high accuracy because the number of font samples is very limited. An accurate distance measurement such as MQDF needs more than

100 character samples for each class.

Can we adapt our JOCR for other Oriental languages such as Chinese by using only the available font images? We did some experiments to answer the question. We collected 4 different Chinese fonts from Internet. They are Song, FangSong, Hei and Kai. By using them as prototypes, the accuracy of the NN classifier is not high. Based on the fact that the Kanji characters in Japanese are very similar to the Hanzi characters in Chinese, we proposed a method to use the information from our well-trained Japanese OCR to improve the performance of the Chinese OCR. Firstly, we conducted visual similarity analysis between Chinese and Japanese font images. A mapping table has been derived. For those Chinese characters which have visual equivalents in Japanese, their corresponding Japanese prototypes can be used to recognize Chinese characters.

The method above has been implemented so that the system can also achieve high recognition accuracy on document in Chinese as same as in Japanese. A multi-lingual document recognition system using the international standard, Unicode, is being also developed.

5 Post Processing

Character segmentation and character classification are error-prone. To improve recognition result, post-processing methods can be applied. During the post-processing stage, various contextual constraints can be used to detect and correct OCR errors. In our JOCR system, we are developing methods for two specific tasks: recognition error correction and segmentation error correction.

In the task of recognition error correction, a top first choice has to be generated for each character which is suspected to be classified incorrectly. The confidence score provided by the classifier can be used to decide whether postprocessing is needed for the character. Two types of constraints can be considered: one is the linguistic constraints; the other is the visual similarity constraints.

One of the simplest types of linguistic constraints is the character bigram. A character bigram can be defined as a tuple, (c_1, c_2, f) , where (c_1, c_2) is a character pair, and f is the frequency of c_1 and c_2 being together. Bigram data can be trained from a large text corpus. For a character image with several alternatives generated by a classifier, bigram information can be very useful to find a choice from the alternative list which best fit its local context (e.g., the characters before or after the character).

In past years, we have conducted experiments on using bigram information for English recognition error correction. The methods have been implemented in our JOCR system. For Japanese, a bi-

gram database with 242,682 bigram entries has been collected from a large Japanese news corpus.

The post processing is done only when it is needed. For a text segment in the original image, if there is no character with a high enough confidence can be found in the fonts, then we say that the “truth value” of this text segment is not yet determined and that the post processing needs to be performed on this text segment. Because the top five choices chosen as possible candidates for any text segment are always available, so the post processing can be done as the following. First of all, we take the first choice of the five candidates and compare it with each of the five choices of its left or right neighbor text segment respectively, if there is such a text segment. Then we can get five “distance values”. These distance values are defined according to the frequency of the two characters appearing together in documents such as books and newspapers of the language. Once we have these distance values, we can use a strategy to find a most possible case that the candidate characters should appear in the text of our special document.

Visual similarity constraints can also be useful for recognition error correction. The intuition behind is that: visually equivalent character images should be recognized as the same characters. By examining OCR result, we found that this constraint sometimes is violated. Different alternative lists can be provided by a character classifier for several visually almost equivalent character images. Previously, we developed methods to cluster character images based on their visual similarities [12]. The clustering algorithm has also implemented as a module in the JOCR system. After clustering, a procedure of consistency checking will be applied to those character images in the same cluster to correct potential errors.

Another types of errors can be corrected by post-processing is character segmentation errors. For those character images with low recognition confidence and with very large or small width/height aspect ratio, segmentation errors may happen. Re-segmentation can be conducted by splitting or merging, then character classification can be applied to get new recognition results. Progress has been made to design a postprocessing procedure to correct segmentation errors.

6 Evaluation of OCR Results Using HTML

For rapid prototyping of evaluation tools which are more flexible and less platform-dependent, a framework for integrating the JOCR system with the Internet/Intranet has been outlined (see Figure 8). Under the framework, the output of the JOCR sys-

tem is converted to HTML files, which can be viewed with a Web browser such as Netscape Navigator 2.0/3.0. Simple Perl scripts can be written to convert the OCR results in a detailed format (in the format shown in Figure 7) into HTML files.

Compared with the approach of designing special GUI programs for performance evaluation, the advantages of using HTML and a Web browser are obvious: scripts for converting OCR results into HTML are much easier to write, test and maintain than X-Window-based GUI programs; HTML files for OCR results can be accessed on different platforms; and the Japanese language is automatically supported in Netscape Navigator 2.0/3.0.

For a Japanese document page, the recognition system can generate OCR results for each text block. In current experiments, we focused on the task of converting OCR results for a text block into HTML files. Aspects of system performance to be evaluated are: character segmentation, character recognition and character clustering. Given OCR results of a text block, a Perl script was written to generate six HTML files:

1. **Block_Image**: the HTML file for the input image; The TIFF image is converted to GIF format for use in HTML.
2. **Character_Images**: the HTML file for character segmentation; This file shows character segmentation results. A sequence of character images is displayed. Here, images are organized in their original order.
3. **OCR_Result**: the HTML file for recognition results; This file shows character recognition results. Corresponding to each character image, its top choice is displayed in JIS code.
4. **Symbol_or_Image**: the HTML file for the hybrid representation of recognition results; For each character position, either a JIS code or a GIF image is placed. For a character, a JIS code is given if the OCR confidence is high; otherwise, the original character image is placed there.
5. **Image_Clustering**: the HTML file for image-based clustering; Results of character image clustering are saved in the file. In a cluster, visually equivalent character images are grouped.
6. **OCR_Result_Clustering**: the HTML file for OCR-result-based clustering. Result of JIS-code-based clustering is saved in the file. Character images recognized as the same JIS code are grouped into clusters.

In each of those HTML files except the first one, for every character position (either a JIS code or an

```

Image quality/type:N
Global Char Size(height, width): 18 18
Text Alignment:H
Line (srow, scol, nrow, ncol): 4 0 28 1088
Line (srow, scol, nrow, ncol): 41 0 28 1088
Line (srow, scol, nrow, ncol): 79 0 27 1088
Line (srow, scol, nrow, ncol): 118 0 27 1088
Line (srow, scol, nrow, ncol): 154 0 26 1088
Line (srow, scol, nrow, ncol): 189 0 27 1088
Line (srow, scol, nrow, ncol): 225 0 27 1088
Char (srow, scol, nrow, ncol, type): 4 21 21 19 C
次 CONF: 2.00
Char (srow, scol, nrow, ncol, type): 6 45 17 17 C
に CONF: 16.67
Char (srow, scol, nrow, ncol, type): 19 66 6 3 C
、 CONF: 4.33
Char (srow, scol, nrow, ncol, type): 6 90 17 19 C
八 入 人 兵 大 CONF: 0.67 0.14 0.14 0.10 0.09
Char (srow, scol, nrow, ncol, type): 4 111 21 21 C
丈 大 丈 太 支 CONF: 0.33 0.08 0.08 0.06 0.06
Char (srow, scol, nrow, ncol, type): 4 136 21 19 C
島 CONF: 3.00
Char (srow, scol, nrow, ncol, type): 7 157 19 20 C
以 CONF: 1.67
Char (srow, scol, nrow, ncol, type): 5 180 22 20 C
南 CONF: 2.33
Char (srow, scol, nrow, ncol, type): 9 204 15 17 C
の CONF: 16.67
Char (srow, scol, nrow, ncol, type): 6 225 21 21 C
伊 信 研 倍 供 CONF: 0.33 0.13 0.13 0.12 0.12
Char (srow, scol, nrow, ncol, type): 7 248 19 20 C
豆 CONF: 1.67
Char (srow, scol, nrow, ncol, type): 6 271 21 20 C
海 CONF: 7.33
Char (srow, scol, nrow, ncol, type): 6 293 22 21 C
積 CONF: 2.00
Char (srow, scol, nrow, ncol, type): 7 317 21 20 C
海 CONF: 7.33
Char (srow, scol, nrow, ncol, type): 7 339 20 21 C
底 庇 庭 直 庵 CONF: 0.33 0.09 0.08 0.07 0.07
Char (srow, scol, nrow, ncol, type): 9 364 17 18 C
か CONF: 12.00
Char (srow, scol, nrow, ncol, type): 9 389 18 14 C
6 CONF: 12.67
Char (srow, scol, nrow, ncol, type): 8 407 20 21 C
探 探 挾 扶 係 CONF: 0.67 0.22 0.15 0.14 0.14
Char (srow, scol, nrow, ncol, type): 9 430 20 20 C
--More-- (7%)

```

Figure 7: The details of text recognition results.

image), a hyperlink to its OCR file is provided so that its OCR result can be examined by clicking a mouse on it.

The OCR file of a character is also in HTML. In the file, the character image is extracted from the block image and saved in GIF format; the classification result is also saved in a HTML file.

The Netscape browser provides a mechanism to have several frames simultaneously inside its window. Each frame can display an individual HTML file. Using the browser, a homepage for JOCR results was designed with three frames: the TOC frame, the OCR frame and the MAIN frame (see Figure 9). The top one is the TOC frame which has a table in which each entry is a hyperlink to one of six HTML files described above. The left bottom one is the OCR frame which is used to display the recognition result for a character. The right bottom one is the MAIN frame which is designed for displaying one of the six HTML files defined above. By clicking on a hyperlink in the TOC frame, the linked file will be displayed in the MAIN frame; by clicking on a hyperlink in the MAIN frame, the linked OCR file will be shown in the OCR frame [13].¹

Activated by selecting the entry `Symbol_or_Image` in the TOC frame, a hybrid HTML file for OCR re-

sults can be displayed in the MAIN frame (see Figure 9). Here, each character is either a JIS code or an image. For this text block, most characters are in JIS because their OCR confidence scores are high. For those characters with low OCR confidence, their original images are displayed. Obviously, those positions where images are placed usually have segmentation problems or recognition errors. Although there are many segmentation and recognition errors in the OCR results which make the JIS code text difficult to read, the content in the hybrid representation is quite readable and not much information is lost because of recognition errors. If the user wants to correct those errors, he can load the hybrid HTML file into an HTML editor (such as HotMetal from SoftQuad, Inc or the WYSIWYG editor which is embedded in the Netscape's Navigator Gold for Windows 95 and NT) and replace those images with their equivalent JIS code.

Experiments on visual contextual analysis have been conducted in our system [12]. In the simplest case, character image clustering is performed. After image clustering, images which are visually equivalent will be grouped as clusters. Image clustering results can be converted into HTML format, as shown in Figure 10 (a). By browsing this page, accuracy and efficiency of the clustering method can be evaluated. Because each image in a cluster is also a hyperlink to its OCR file, by checking the OCR result of every character inside a cluster, we can examine whether the recognizer runs consistently on visually

¹The HTML pages presented in the paper can be found at <http://www.cedar.buffalo.edu/JOCR/HTML-demo/w95.html2.frm.html>. More demonstrations for various Japanese and Chinese document pages can be located at <http://www.cedar.buffalo.edu/JOCR/HTML-demo/index.html>.



Figure 8: Integrating the JOCR system with the Internet/Intranet.

similar images.

Similarly, clustering based on OCR results is also very useful for fast error correction [1]. The MAIN frame of Figure 10 (b) shows the clustering result for the example block. All images which are recognized as the same character are grouped together.

It is easy to modify the Perl script when the appearance of an HTML page needs to be changed. In our experiments, performance of character segmentation, character recognition and image-based clustering were evaluated. New scripts for converting OCR results to HTML files can be designed rapidly when some other aspect of OCR performance needs to be studied.

OCR transforms a given document from its image representation into its symbolic representation. After this process, we obtain a text document which is electronically searchable, indexable and reusable. However, the transformation is error-prone. Therefore, if we just save the pure text document generated by the OCR process, the text may have many errors. Without referring to the original image representation from which they were derived, some of those errors may not be recoverable. As suggested by our experiment shown in Figure 9, a hybrid document which combines a symbolic representation and an image representation may relieve the problem. Such a hybrid document in HTML is searchable and also error-recoverable. If we represent an OCR'd document in such a hybrid format, OCR errors will not have much negative effect on the human reading process. Therefore, the hybrid HTML can be recommended as an ideal format for OCR'd documents in an online digital library [24].

Our Japanese OCR system is designed to recognize the first level of Kanji in the JIS code. When characters from the second level of the JIS code are encountered, OCR results for those characters will always be incorrect. It is best to keep those characters as original images so that the text will still be readable with a browser [15, 13].

Document images can be edited after segmentation. Usually this needs a specially designed GUI program [4]. An HTML file browser or editor can work for this purpose, as we have shown in our experiments. For a text block image, after character segmentation, a sequence of character images arranged in their original linear order can be rendered as readable as the original image, and differ-

ent markup instructions can be applied. This kind of image-based document reconstruction can be very useful. We conducted another experiment to reconstruct documents more efficiently. For a Japanese document, usually there are many visually equivalent character images (as we have shown in our image-based clustering experiment). After image clustering, those visually equivalent character images can be viewed as the same image at different positions. Therefore, when the document image is reconstructed, we do not have to paste every image to its position. Instead, for those equivalent character images, a prototype can be generated to represent all of them at different positions [15, 13].

7 Conclusions and Future Directions

The Japanese OCR system, *Cherry Blossom*, has been designed to recognize Japanese documents digitized at low resolution such as 200 ppi. A wide variety of Japanese document images extracted from newspapers, journals, facsimiles and other sources have been used in the development and testing of this system. Results of character classifiers using CEDAR's Japanese character database have been presented. Experiments on Internet-based system evaluation is also reported. we focused on using HTML to evaluate performance at the block level. Experiments can be generalized to use HTML for physical/logical layout analysis. In this framework, we use hyperlinks and image-embedding functions provided by HTML to integrate components of the OCR results with components of the text image to overcome the difficulties of OCR performance evaluation and error correction. In the general paradigm of multivalent documents [31], we will extend the approach to tackle more sophisticated problems in visual document recognition.

More efforts will be directed towards further improving each module of the system and benchmarking of the performance of the system.

Acknowledgment

Present and past members of the *Cherry Blossom* team include: Geetha Srikantan, Yi L. Chiang, Chi Fang, Brian Grom, Zhou Hong, Jonathan J. Hull, Stephen W. Lam, Dar-Shyang Lee, Qunfeng Liao, Minsio Kariya, Seiichiro Kashimori, Chong O. Kim, Daniel R. Mechanic, Xiaomin Ren, Weh-Chih Wang,

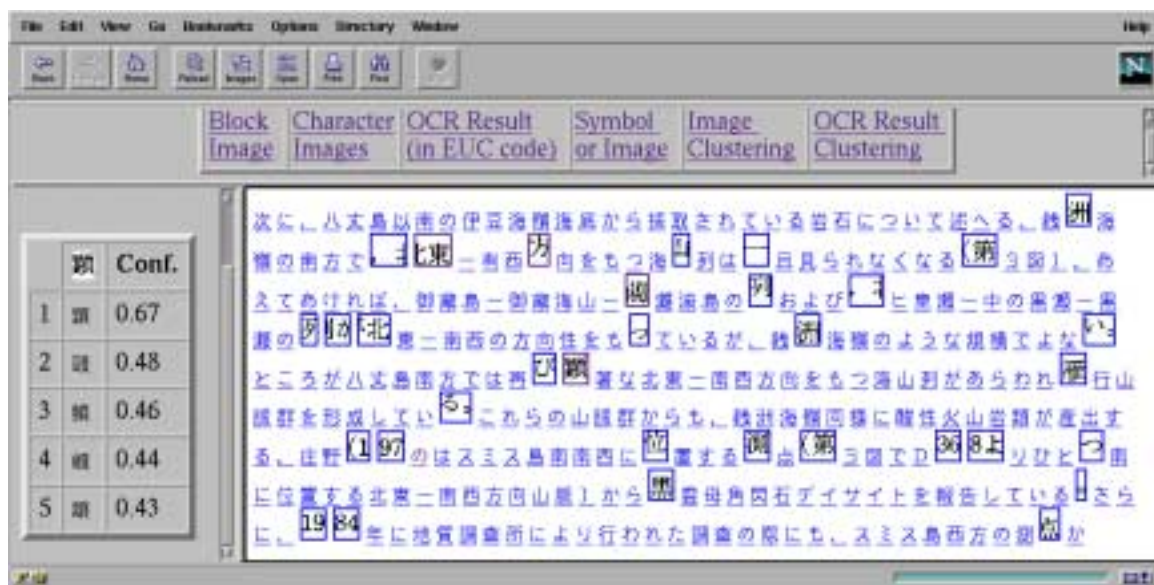


Figure 9: Multi-frame design for presenting OCR results.

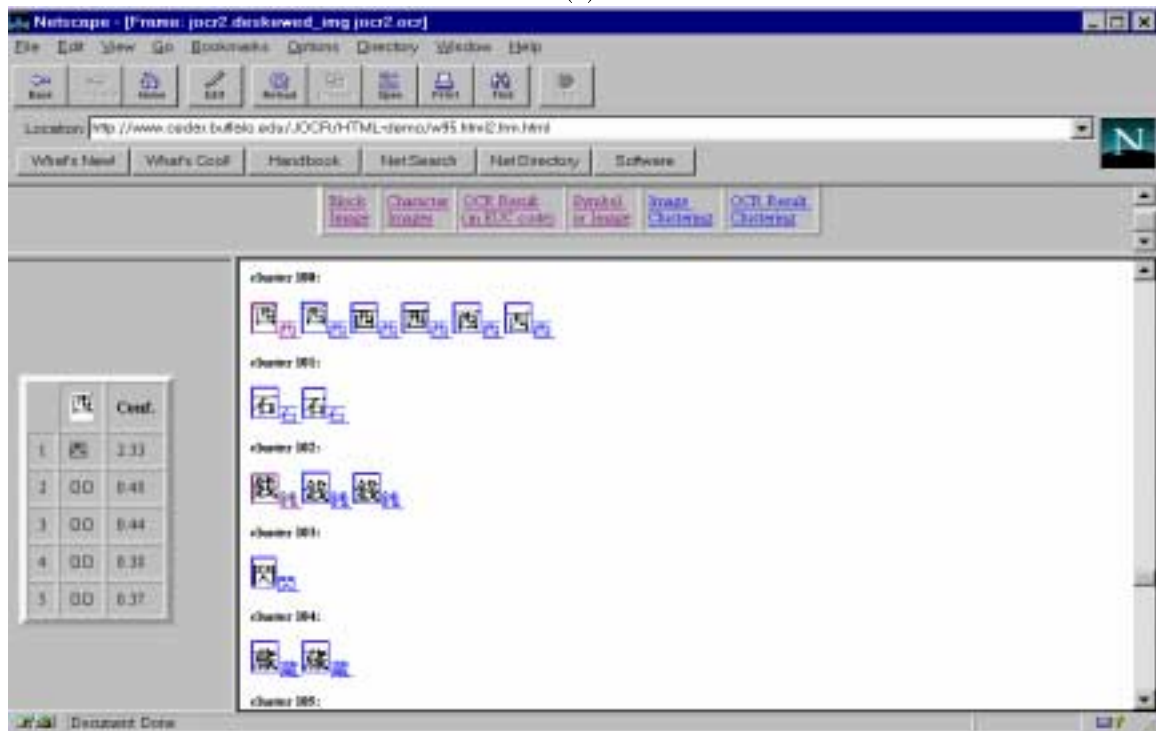
Shu-Fang Wu, Ying Yao, Xiaoming Yang, Victor C. Zandy and Jason Zhu.

References

- [1] Illuminator user's manual, release 0.9. Technical report, RAF Technology, Inc., May 31 1995.
- [2] T. Akiyama and N. Hagita. Automated entry system for printed documents. *Pattern Recognition*, 23:1141–1154, 1990.
- [3] S. Ariyoshi. A character segmentation method for Japanese printed documents coping with touching character problems. In *Proceedings of the International Conference on Pattern Recognition*, pages 313–316, 1992.
- [4] S.C. Bagley and G.E. Kopec. Editing images of text. *CACM*, 37(12):63–72, 1994.
- [5] J.T. Favata, G. Srikantan, and S. N. Srihari. Handprinted character/digit recognition using a multiple feature/resolution philosophy. In *International Workshop on the Frontiers of Handwriting Recognition, IWFHR - 4*, 1994.
- [6] K. Fukunaga. *Introduction to Statistical Pattern Recognition*. Academic Press, 1990.
- [7] G.W. Gates. The reduced nearest neighbor rule. *IEEE Transactions on Information Theory*, IT-18(3):431–433, May 1972.
- [8] P.E. Hart. The condensed nearest neighbor rule. *IEEE Transactions on Information Theory*, IT-14(3):515–516, May 1967.
- [9] T. Hong. *Degraded Text Recognition Using Visual and Linguistic Context*. PhD thesis, Computer Science Department of SUNY at Buffalo, 1995.
- [10] T. Hong and J. J. Hull. Improving OCR performance with word image equivalence. In *Symposium on Document Analysis and Information Retrieval*, pages 177–190, April 24–26 1995.
- [11] T. Hong, S.W. Lam, J.J. Hull, and S.N. Srihari. The design of a nearest-neighbor classifier and its use for Japanese character recognition. In *Proceedings of the Third International Conference on Document Analysis and Recognition ICDAR '95*, 1995.
- [12] T. Hong, S.W. Lam, J.J. Hull, and S.N. Srihari. Visual similarity analysis of Chinese characters and its uses in Japanese ocr. In *Proceedings of the Conference on Document Recognition, 1995 SPIE Symposium (SPIE95)*, February 1995.
- [13] T. Hong and S. N. Srihari. Representing ocred documents in html. In *Proceedings of International Conference on Computational Intelligence and Multimedia Applications*, pages 10–12, 1997.



(a)



(b)

Figure 10: (a) The result of image-based clustering is shown in the MAIN frame; (b) The result of JIS-code-based clustering is shown in the MAIN frame.

- [14] T. Hong, G. Srikantan, V.C. Zandy, C. Fang, and S. N. Srihari. Character recognition in a Japanese text recognition system. In *Proceedings of the Conference on Document Recognition, 1996 SPIE Symposium(SPIE96)*, January 1996.
- [15] T. Hong, S. Wu, and S.N. Srihari. Evaluating japanese document recognition in the internet/intranet environment. In *Proceedings of Workshop on Document Analysis Systems*, pages 14–16, 1996.
- [16] <http://www.cedar.buffalo.edu/Database/JOOCR/>. CEDAR's Japanese character image database.
- [17] T. Ijima, H. Genchi, and K. Mori. A theory of character recognition by pattern matching method. In *Proc. First IJ CPR*, pages 587–594, 1973.
- [18] D.J. Ittner and H.S. Baird. Language-free layout analysis. In *Proceedings of the International Conference on Document Analysis and Recognition*, pages 336–340, 1993.
- [19] F. Kimura, M. Shridhar, and Y. Miyake. Relationship among quadratic discriminant functions for pattern recognition. In *International Workshop on the Frontiers of Handwriting Recognition, IWFHR - 4*, 1994.
- [20] F. Kimura, K. Takashina, S. Tsuruoka, and Y. Miyake. Modified quadratic discriminant functions and the application to Chinese character recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 9:149–153, 1987.
- [21] S. W. Lam. A local-to-global approach to complex document layout analysis. In *Proceedings of Machine Vision Applications*, 1994.
- [22] S. W. Lam, Q. F. Liao, and S. N. Srihari. A divide-and-conquer approach to Japanese text segmentation. In *Proceedings of the Conference on Document Recognition, 1995 SPIE Symposium(SPIE95)*, February 1995.
- [23] S. W. Lam and V. C. Zandy. Skew detection using directional profile analysis. In *Proceedings of Machine Vision Applications*, 1994.
- [24] M. Lesk. Substituting images for books: The economics for libraries. In *Symposium on Document Analysis and Information Retrieval*, pages 1–16, April 15-17 1996.
- [25] K. Lunde. *Understanding Japanese Information Processing*. O'Reilly and Associates, Inc., Sebastopol, CA, USA, 1993.
- [26] Q. Luo, T. Watanabe, and N. Sugie. A structure recognition method for Japanese newspapers. In *SDAIR*, pages 217 – 234, 1992.
- [27] Y. Maeda, F. Yoda, K. Matsuura, and H. Nambu. Character segmentation in Japanese hand-written document images. In *Proceedings of the International Conference on Pattern Recognition*, pages 769–772, 1986.
- [28] S. Mori, K. Yamamoto, and M. Yasuda. Research on machine recognition of handprinted characters. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6:386–405, 1984.
- [29] E. Oja. *Subspace Methods in Pattern Recognition*. John Wiley & Sons, 1983.
- [30] R. Oka. Handwritten Chinese character recognition by using cellular feature. In *Proc. 6th Int. Joint Conf. Pattern Recognition*, pages 783–785, 1982.
- [31] T. Phelps and R. Wilensky. Multivalent documents: Inducing structure and behaviors in online digital documents. In *Proceedings of the 29th Hawaii International Conference on System Sciences*, January 3-6 1996.
- [32] K. Sakai, S. Hirai, T. Kawada, S. Amano, and K. Mori. An optical Chinese character reader. In *Proc. 3rd Proceedings of the International Conference on Pattern Recognition*, pages 122–126, 1976.
- [33] S. N. Srihari and Geetha Srikantan. Japanese optical character recognition. In *Proceedings of 1995 Symposium on Document Image Understanding Technology*, pages 236–246, 1995.
- [34] G. Srikantan, S.W. Lam, and S. N. Srihari. Gradient-based contour encoding for character recognition. *Pattern Recognition Journal*, 29(7):1147–1160, July 1996.
- [35] S. Tsujimoto and H. Asada. Major components of a complete text reading system. *Proceedings of the IEEE*, 80:1133–1149, 1992.