

A System to Read Names and Addresses on Tax Forms

Sargur N. Srihari, Yong-Chul Shin, Vemulapati Ramanaprasad and Dar-Shyang Lee

Center of Excellence for Document Analysis and Recognition
and
Department of Computer Science
State University of New York at Buffalo
Buffalo, NY 14260

Abstract

The reading of names and addresses is one of the most complex tasks in automated forms processing. This paper describes an integrated real-time system to read names and addresses on tax forms of the Internal Revenue Service of the United States. The Name and Address Block Reader (NABR) system accepts both machine-printed and hand-printed address block images as input. The application software has two major steps: document analysis (connected component analysis, address block extraction, label detection, hand-print/machine-print discrimination) and document recognition. Document recognition has two non-identical streams for machine-print and hand-print; key steps are: address parsing, character recognition, word recognition and postal database lookup (ZIP+4 and City-State-ZIP files). System output is a packet containing the results of recognition together with database access status file. Real-time throughput (8,500 forms per hour) is achieved by employing a loosely-coupled multiprocessing architecture where successive input images are distributed to available address recognition processors. The functional architecture, software design, system architecture and the hardware implementation are described. Performance evaluation on machine-printed and handwritten addresses are presented.

Index Terms: Address Parsing, Character Recognition, Database, Integrated System, IRS Tax Forms, Multiprocessing, Real-Time Processing.

1 Introduction

The automated processing of forms is of great commercial importance. A large percentage of paper forms in use today require a name and address to be filled-in, in addition to other domain specific data. Names and addresses have wide variability, which makes the reading problem a challenging one. The Internal Revenue Service (IRS) of the United States receives around 200 million tax forms per year. Only ten percent of these returns are filed electronically. Although efforts are under way to increase the use of electronic filing, IRS still has to contend with paper returns for years to come. In an effort to convert the paper returns into electronic form, IRS plans to use a new, intelligent optical character recognition system called *Service Center Recognition/Image processing System (SCRIPS)*. The SCRIPS goal is to replace two current OCR systems at each IRS center to minimize the burden of data entry and hence to speed-up the processing with high accuracy. The current systems in service provide 85% data capture success rate on average 1,700 to 2,000 pages per hour [1][2]. By contrast, SCRIPS is expected to provide a 96% data capture success rate, converting 17,240 documents to images per hour. In essence, SCRIPS consists of five subsystems: image capture, image processing, system control and data management, character recognition and identification, and storage.

In this paper, we describe the design of a subsystem called *Name and Address Block Reader (NABR)*, which is one of the components in the character recognition and identification subsystem of SCRIPS. NABR is expected to match both machine-printed and hand-printed information with cities, states, ZIP Codes and street names stored in its own database. This integrated system was developed at the *Center of Excellence for Document Analysis and Recognition (CEDAR)*, State University of New York at Buffalo. It combines postal address recognition techniques developed at CEDAR over the past decade and current computing power available in the market. NABR is capable of reading both machine-printed and hand-printed address block images extracted from tax forms.

The primary objective of NABR is to recognize each character within the address on a tax form. Recognition is aided by contextual information in the form of various databases, such as city-state-ZIP Code database, street-name database, and personal-name database.

The software and hardware design of NABR posed some special challenges, principal of which were:

- wide range of quality of machine-printed text,
- lack of reliable segmentation and recognition techniques for hand-printed characters,
- demands on processing speed and recognition accuracy,
- size of databases involved (about 1.2 gigabytes compressed), and
- need to provide a cost-effective solution.

Section 2 describes the nature of the images input to the NABR system. Address recognition software modules are described in section 3. The system architecture and hardware design are described in section 4. A performance evaluation of the system is described in section 5.

2 Images input to NABR

The input to the NABR system is a bit-packed binary image in TIFF format. The output is the ASCII form of the address as it is written on the tax form. The address can be either hand-printed or machine-printed. A machine-printed address can be a labeled address (a labeled address is one in which the address is pre-printed on a label and later pasted on top of the tax form) or a typewriter printed address. Some tax forms may have drop-out ink to distinguish between pre-printed text and the actual address. The current NABR system is designed to read 15 different IRS tax forms including 1040EZ, 1040PC, 941 and IRPS forms (1096, 1099-R, 1099-MISC, 5498). It is expected that more forms are added as per future requirement. Examples of hand-printed and machine-printed address blocks from several of the tax forms are shown in Figure 1.

In the label image of Figure 1(b), each of the five lines are handled as follows. The top line, carrier sort line (***** 5-DIGIT 14215) is to be ignored. The next four lines are to be read. Each of the lines is referred to as follows: the first line is the scan line (HZ 123-45-6789 S08 B1), followed by the name line (JOHN J DOE), street address line (123 MAIN STREET), and the city-state-ZIP line (ANYTOWN NY 14260).

3 Address Recognition Software

The overall control flow of the NABR system is shown in Figure 2. It consists of two major modules: document analysis, and document recognition; a structure described in [3]. The detailed design of the modules resembles the design used in postal address reading systems developed at CEDAR, for reading handwritten addresses [4], and for reading either handwritten or machine-printed addresses [5]. In the following sections we describe each of the software modules of NABR in some detail.

3.1 Document Analysis

This involves recognizing the type of address in the input image (hand/machine print, label/no-label) and then segmenting it into individual lines. Figure 3 shows the subtasks in document analysis. If the image is that of a label then the task is to extract that portion of the image corresponding to the label. If the tax form is not printed with drop-out ink then the task is to remove all the pre-printed text from the image.

3.1.1 Connected Component Analysis

Document analysis operations can be speeded-up by performing the operations at the connected component level, rather than at the image level. This insight was previously employed by us in locating address blocks in a real-time address block location system [6]. Thus the first step is to generate connected components of the entire input image.

There are several algorithms to generate connected components. We use the Line Adjacency Graph (LAG) method to generate and represent connected components [7]. This algorithm is quite efficient when implemented in software. The algorithm was further speeded-up by modifying it to generate components directly from the bit-packed image.

Print your name (first, initial, last)

JOHN J. DOE

If a joint return, print spouse's name (first, initial, last)

JANE C. DOE

Home address (number and street). If you have a P.O. box, see page 11. Apt. no.

123 MAIN STREET

City, town or post office, state and ZIP or zip. If you have a foreign address, see page 11.

ANYTOWN, NY 14260

(a)

***** 5-DIGIT 14260

HZ 123-45-6789

S08 B1

JOHN J DOE

123 MAIN STREET

053

ANYTOWN NY

14260

I
R
S

(b)

CENTER FOR DOCUMENT
ANAL.

41194

CEDAR

123456789

321 MAINT ST.

ANYTOWN, NY 14260

(c)

JOHN J. DOE

123 MAIN STREET

ANYTOWN, NY 14260

(d)

Figure 1: Examples of images input to NABR: (a) 1040EZ hand-print with no drop-out ink, (b) 1040 label image with OCR-A Font, (c) 941 hand-print with drop-out ink and (d) 1096 machine-print with drop-out ink.

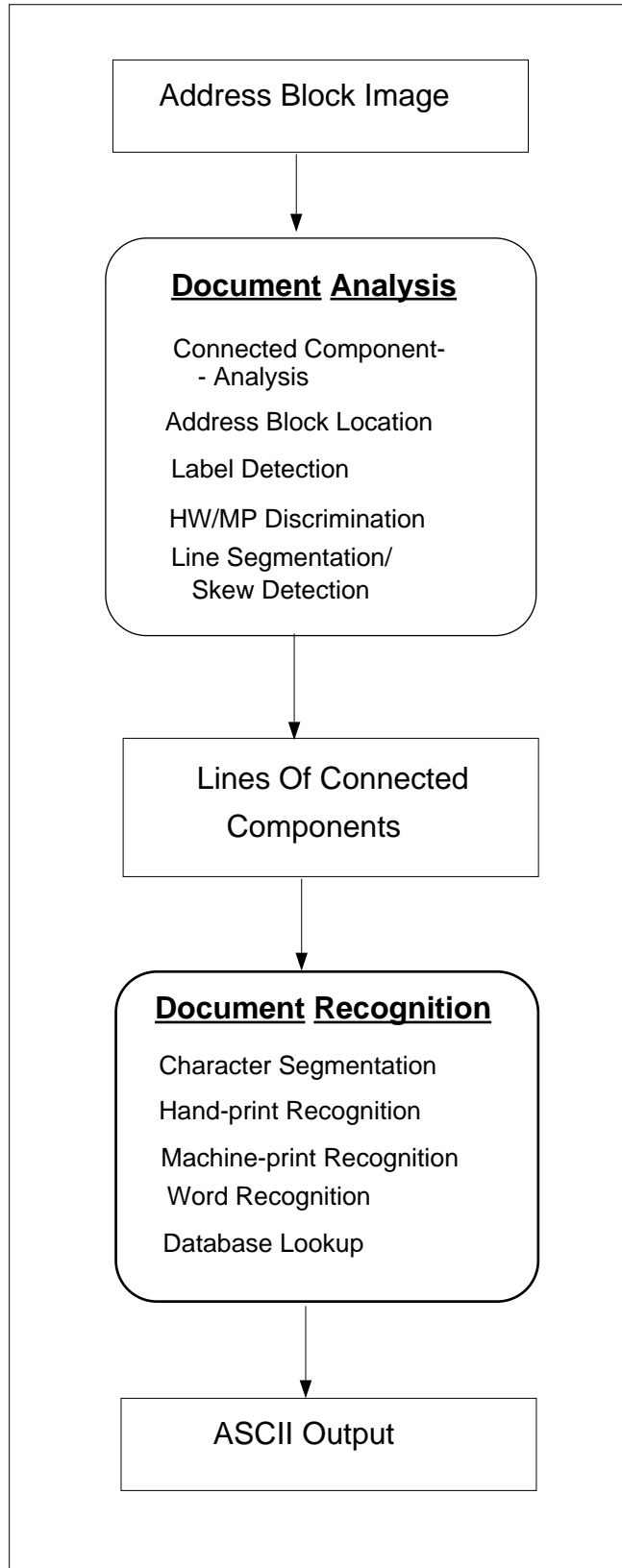


Figure 2: Top-level Modular Organization

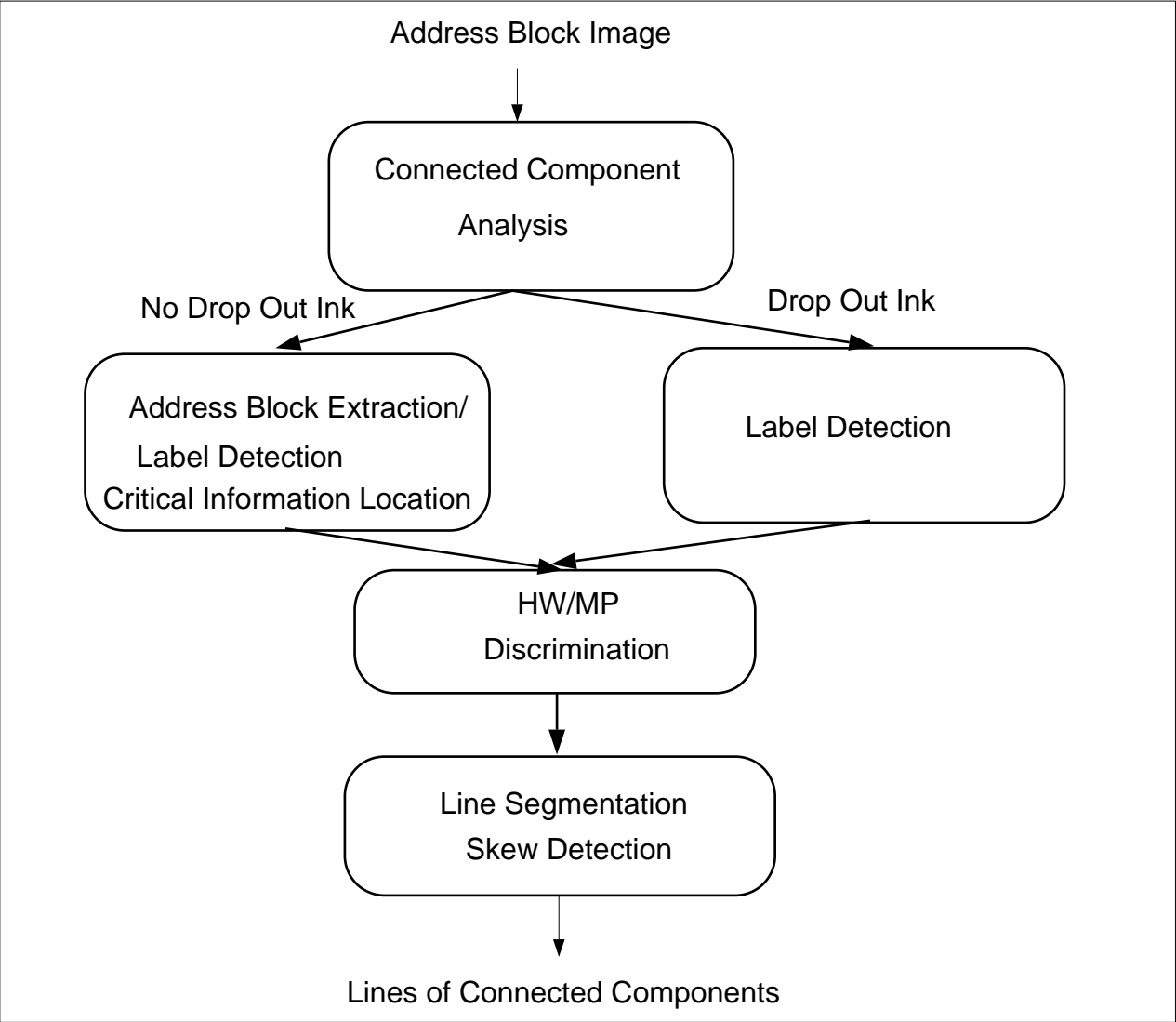


Figure 3: Document Analysis.

3.1.2 Address Block Extraction/Label Detection

If a tax form is not pre-printed in drop-out ink (this depends on the type of form being processed), then certain pre-determined registration marks are used to delete the pre-printed text. These registration marks correspond to areas of high pixel density in the address block. By examining regions around these registration marks, we determine whether there is pre-printed text to be deleted or whether a label has been pasted on top of the tax form.

In case of tax forms that have drop-out ink, the spacing between lines is used to determine the presence of a label. In case of a labeled image the spacing is much smaller than in the case of a non-labeled image. The detection of the presence of a label facilitates parsing the lines into name line, street line and ZIP Code line because of the fact that non-labeled images have fixed parsing.

3.1.3 Hand-print/Machine-print Discrimination

It is necessary to determine whether the address is hand-printed or machine-printed, since different character segmentation and recognition algorithms are used for each script type; such a strategy is also used in real-time postal mail address reading [5]. The hand/machine print discrimination is performed using six symbolic features extracted from the connected components:

1. standard deviation of connected component widths,
2. standard deviation of connected component heights,
3. average component density,
4. aspect ratio,
5. distinct different heights, and
6. distinct different widths.

The classifier employed is a Fisher's linear discriminant [8]. During training, the equation of the hyperplane that separates the features into two distinct clusters is found. For recognition, the unknown sample is classified by determining the side of the hyperplane its feature vector lies in. The result is a fast, high-performance discriminator: on a test set of 800 postal addresses, the correct discrimination rate was 95% [9].

3.1.4 Line Segmentation/Skew Detection

One approach to line segmentation is to use a vertical histogram. This method can be very slow as it is pixel-based. Also histogram method will fail badly if the input image is skewed. To overcome these two problems, a clustering algorithm that works directly on the connected components of the image, is used. This decreases the processing time for line segmentation. Also vertical overlap between adjacent components is used to cluster the components into lines. The skew is computed as the clusters are determined. The clustering criterion (vertical overlap) makes the algorithm more robust to skew in the input image.

3.2 Document Recognition

After preprocessing, we have lines of connected components. We also have an indication whether the image is hand-printed (HP) or machine-printed (MP). Depending on the type of the image (HP or MP), different segmentation algorithms are used. The process flow for character segmentation and character recognition is shown in Figure 4.

3.2.1 Machine-Print Segmentation

It was found that most of the machine-printed addresses on the tax forms have fixed-pitch font. This information is very useful in breaking touching characters as it makes the segmentation task trivial. The fixed-pitch font determination is done by examining the spacing between the connected components that have similar width and height as that of normal characters and also by looking at the word gaps. The first step in character segmentation is to group connected components that satisfy some pre-determined criteria (like horizontal overlap, distance between components, etc). The objective of the grouping is to cluster broken characters into one group. Each of the groups can contain one or more components. By comparing the width and height of a group of components to the rest of the character like components, the number of characters in that group are estimated. If the group is considered to contain a complete character, a character image will be generated. If the group contains more than one character, then it is segmented into two or more characters. If the image has fixed-pitch font then the segmentation is a straight-forward task. Otherwise the vertical profile of the group of components is used to locate the break points. If the vertical profile analysis fails, split-and-classify technique is used to find the break points. In the last technique, the image is segmented into different sequences of regions based on probable local minima and estimated number of characters in the image. Each region is recognized and assigned a confidence measure by a classifier. The sequence of regions with highest confidence measure is chosen as the segmentation decision.

3.2.2 Hand-Print Segmentation

The hand-print segmentation algorithm is based on an iterative row-partitioning of the image [10], where each partition consists of a set of contiguous rows. The number of rows in a partition decreases for each iteration. At each iteration, overlapping partitions are grouped into a blob. Each of these blobs is tested with a character recognizer. If the recognition yields a high confidence character, then that blob is removed in the next iteration. After the last iteration (when the number of rows per partition is one) each blob corresponds to a connected component. At this stage each of these blob can contain one or more character. An estimation is made about the approximate number of characters (n) in each of the blobs remaining. Character splitting techniques are used to split the blob into either $n - 1$, n or $n + 1$ characters. The decision to choose either $n - 1$, n or $n + 1$ characters is made based on the combined character confidences.

Four types of splitting techniques are used to split blobs into individual characters. In *histogram* method vertical density projections are calculated, and the column with minimal value is chosen for vertical split. In the *upper-lower contour* method the split is made by connecting the valley in the upper profile and a peak in the lower profile. *Upper contour* method uses valleys in the upper contour for split and *lower contour* method uses peaks in the lower contour.

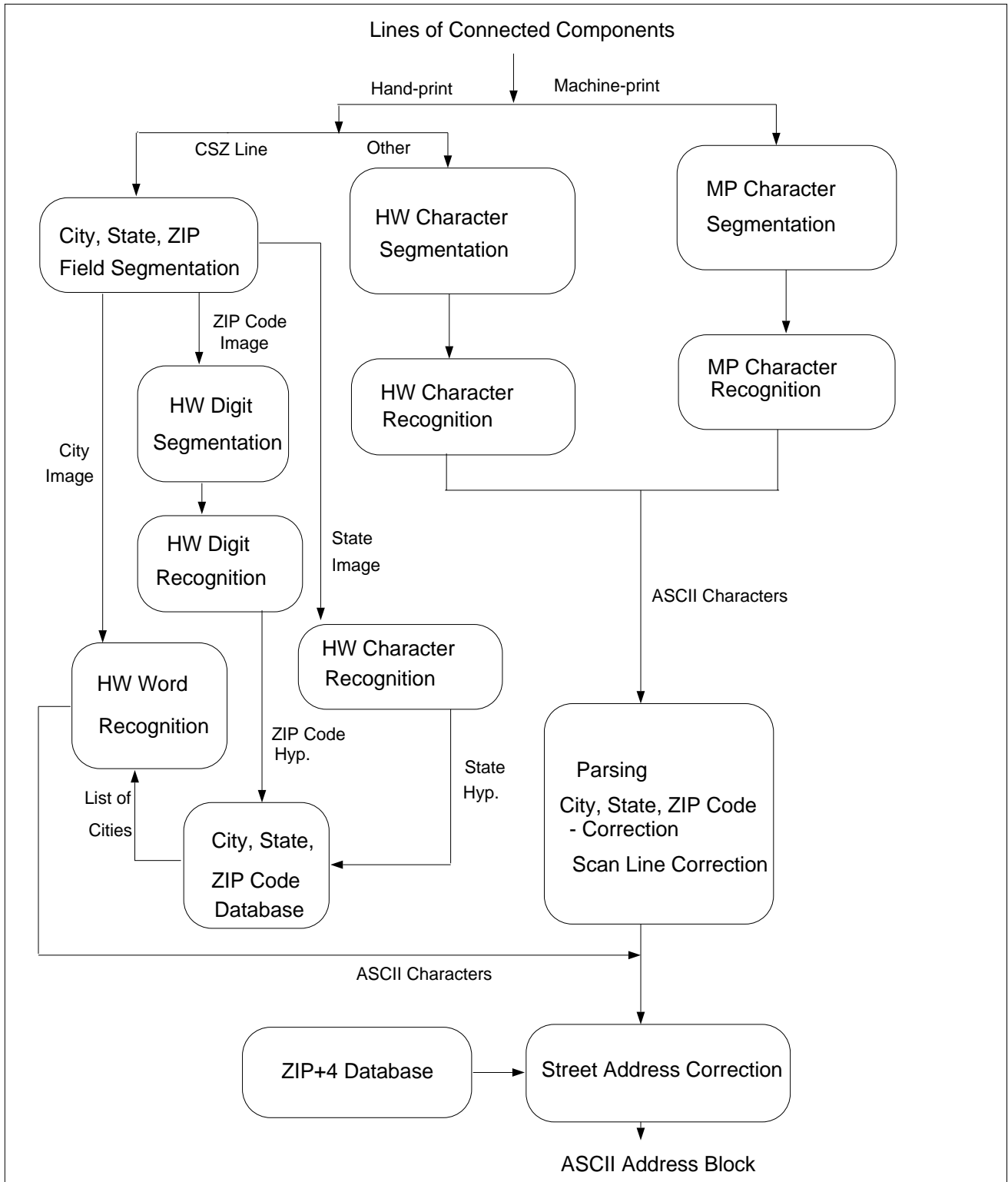


Figure 4: Document Recognition

3.2.3 Character Recognition

Three basic OCR engines are used for character recognition, each OCR engine specifically trained for a different purpose (machine-print, hand-printed alphabets, and hand-printed digits) [11]. All three engines use an artificial neural network classifier consisting of one hidden layer, trained using backpropagation. A special OCR-A engine is based on a two-layer neural network. The OCR-A and handprint digit modules consist of a combination of different algorithms.

The *training dataset* for machine-print (98,308 images) was from tax forms and postal envelopes. The data sets for handprinted alphabets (113,694 images) were from various sources, including NIST data sets. The OCR-A engine was trained on 3,620 images. The data set for hand-printed digits was from a database created at CEDAR from postal envelopes.

OCR Modules. The basic features are contour (or gradient) features [12]. The contour direction for each edge pixel is calculated by indexing to a lookup table with its 3x3 neighborhood. The basic algorithm calculates the chaincode direction for each pixel by convolving a kernel over the image; the actual chaincode sequence is not constructed. This process can be implemented efficiently using register shifts and bit masking operations, eliminating repetitive gradient calculations at every pixel position. This implementation (on a SUN) constitutes significant saving in feature extraction time. Since opposite gradient directions signify the same stroke direction, these directions are folded from 0 to 180 degrees. The directions are then quantized into 4 general directions. These general directions of all edge pixels within each of the 4x4 partitioned region in the image are histogrammed and normalized with respect to image size. The feature vector is of length 64. Architectures of the network vary in different modules: with 64 input units, 150 hidden units and 71 output units for machine-print, and 64-100-71 for hand-print.

The version with floating weights timed more than 200 characters per second on a Sparc10 for a 71 class problem. For the implementation running inside NABR system, network weights are converted to integer values to improve execution speed on the address recognition processors. This algorithm, supplemented by various contextual postprocessing schemes, provides suitable recognition accuracy for machine printed characters. For example, for the name field in a 1040PC, the number of classes is 31 (26 upper-case characters and punctuation marks), and so the other character classes are filtered out of the output choices. Some commonly confused characters and several punctuation marks are recognized by simple special purpose classifiers. For example, U and W have a special discriminator for OCR-A fonts when the character strokes are thick. It computes the depth of the middle region of the character. Other special discriminators are used for (-,.) and (K, <).

Hand-printed Digit Recognition. Since recognition of handprinted addresses relies heavily on accurately reading the ZIP Codes, a different algorithm is used for handwritten digit recognition. A combination of two algorithms are used (Figure 5). The first one is based on chaincode features. Directions and curvatures at each point on the digit contour is first calculated. Then a histogram similar to the one described above is constructed for each image partition to form a feature vector of 296 dimension. A multilayer perceptron is used for classification. The second recognizer uses the polynomials of pixels in a normalized image as features. The polynomial recognizer uses 1241 experimentally determined coefficients [13]. We have also found that combining results of multiple recognition of the same image based on different preprocessing improves accuracy. The raw image is normalized to a fixed size using two different algorithms: a proportional size normalization and a deskewing moment normalization. The confidence values of these two recognition processes are summed before combined with those of the chaincode based method. An alternative for handwritten

digit module, which is based on combination of features extracted at multiple levels of resolution [14], is also being evaluated.

3.2.4 Parsing and Contextual Post-Processing

At this stage, we have lines of words/characters with OCR results. The objective of the postprocessing module is to correct OCR results using contextual information. Parsing involves classifying each line into City, State, ZIP line or Address Line or Name Line and also identifying the most frequently used words in a postal address (street, south etc.). Contextual constraints imposed by postal directories, such as the City-State-ZIP Code database and the ZIP+4 database are useful to correct OCR results [15].

3.2.5 Parsing/Word Category Assignment

Before contextual postprocessing can be performed, the OCR results have to be parsed [16] [17]. There are two levels of parsing. One is *vertical parsing*. Vertical parsing involves classifying each line as name line, street line, CSZ (city, state, ZIP Code) line, etc. Vertical parsing is trivial for non labeled images because the parsing information can be determined from the location of each line with respect to the tax form. In case of labeled image, vertical parsing is done by using spatial constraints and the location of each line with respect to the others. The second type of parsing is called *horizontal parsing*. It involves labeling each word in a line with appropriate tags. Tags are used to denote special types of words that occur in postal addresses (e.g., street number, street name, ZIP Code, street suffix, etc.). Horizontal parsing is accomplished by comparing each word in a line with a set of keywords that occur very frequently in postal addresses (e.g., South, North, Mr., Street, etc.). Each keyword has a set of tags associated with it (e.g., North can be a personal name or a street prefix). When a word in a line matches with a keyword then that word is labeled with the tags associated with the matched keyword. From this list of tags, an appropriate tag is chosen based on the vertical parsing information (i.e., NORTH in name line will be a family name and in address line it will be labeled as street prefix).

Once an unambiguous horizontal parsing has been done, the databases are used to correct OCR results. The ZIP Code database is used to correct city names and the city name database to correct ZIP Codes and state names. We also have a database of all street names in the entire country. This database is indexed by ZIP Code and a street number and it is used to correct street names, assuming that the ZIP Code and street number are recognized correctly [18].

3.2.6 City, State and ZIP Code Postprocessing

City-State-ZIP database is used to correct the OCR results of the CSZ line. This database is a comparatively small database (about 5 MB) and resides entirely in main memory. It supports three types of queries.

- Given a ZIP Code, return all (City, State) pairs. This is the fastest of the three queries.
- Given a city name substring, return all city names which contain *exactly* those characters in the substring. This query is slightly slower than the first one.

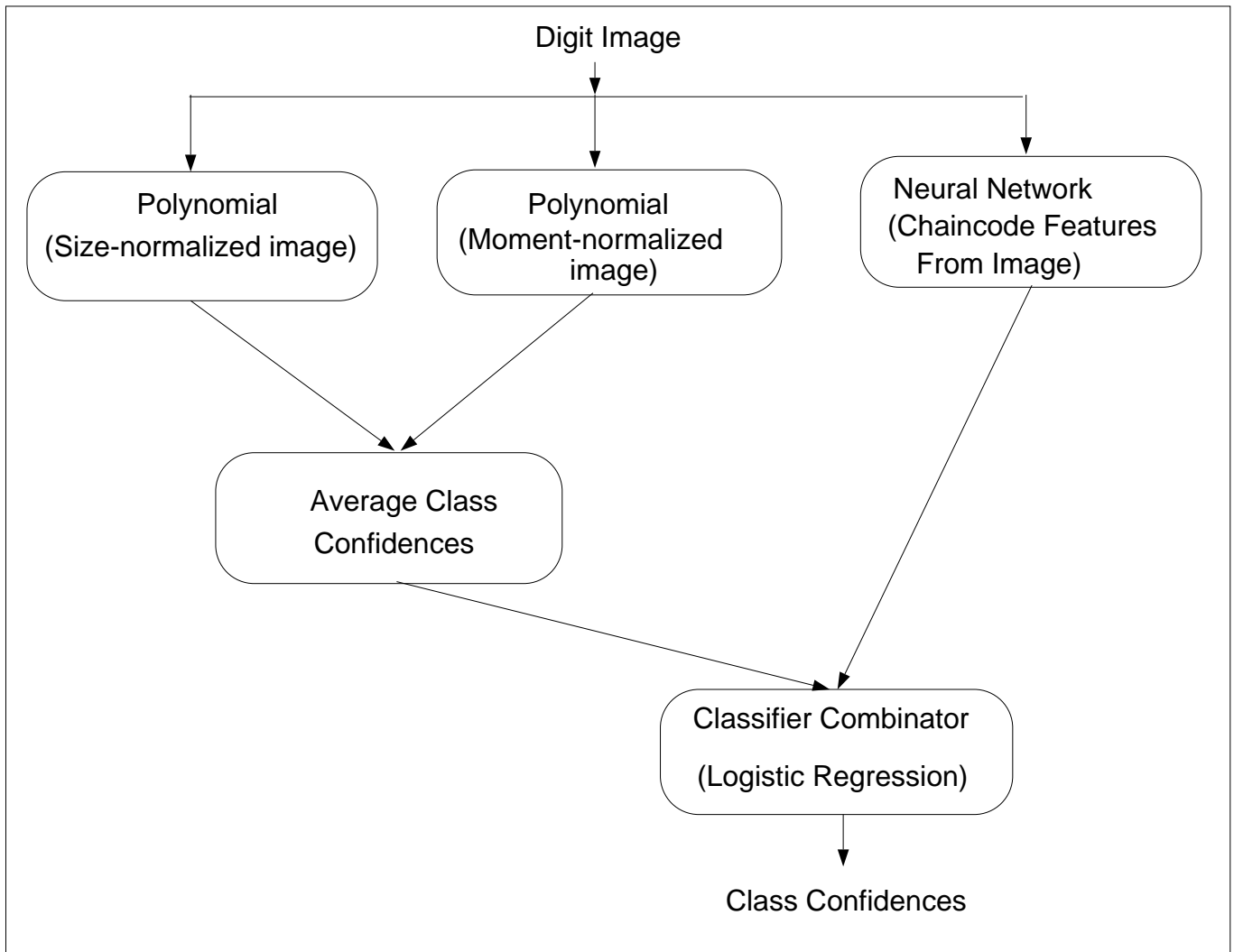


Figure 5: Hand-printed digit recognizer

- Given a city name substring, return all city names which *contain* the substring. This is the slowest of all queries.

Since digits are detected more accurately than alpha, first the ZIP Code is used to get all the city names from the database. Each of these city names from the database are then compared with the OCR results of the recognized city name. If there is no good match, then all the characters in the city name are used to find the corresponding ZIP Codes. These ZIP Codes are compared with the OCR results of the OCR results of the ZIP Code. If a good match is still not found, then a city name substring match is used.

3.2.7 Street Address Postprocessing

The street address database, also called ZIP+4 database, contains street name information of the entire country. ZIP Code and street number are used as indices to get a list of street names from this database. Each of these street names is compared against the OCR results of the recognized street names. If there is a good match then the OCR results are appropriately corrected. The uncompressed ZIP+4 database occupies 4GB of memory space. In order to reduce the storage requirements, the database is compressed, at the same time keeping the database access reasonably fast.

Compression scheme. We take advantage of the vertical redundancy (similarity of information in consecutive records) to achieve desired compression. Attribute coding [19] is one such compression scheme, which takes advantage of the relationships between records in a structured database by storing a record as a set of differences from the previous record. These differences are specified by functions and their arguments. With the aid of statistical analysis of the database, a set of functions is devised for each attribute. These functions use information in the previous records as well as information that has been already computed in the current record, as arguments to calculate the next attribute in the current record. To keep database access reasonably fast, it is not feasible to have every record depend on all the records before it. Then to recover record n , every record before n need to be accessed, making the database access slow. Therefore, the records are compressed in groups of N . The value of N can be adjusted for a space/time tradeoff. Some of the functions used in attribute coding are more likely to appear than others. To take advantage of this non-uniform distribution, the function pointers are encoded using Huffman codes [20]. One advantage of the Huffman code is that decompression can be achieved in constant time using a lookup table. The size of compressed ZIP+4 database is around 1.1 GB. The techniques used in compressing the ZIP+4 database are described in detail in [21].

The ZIP+4 database organization is such that all information corresponding to a ZIP Code is stored in contiguous disk blocks. A global ZIP Code table contains pointers to the street information for each of the ZIP Codes. Street information is divided into blocks with each block containing information for 20 streets. Each block is compressed using attribute and Huffman coding as discussed earlier. Only one type of query is provided for database access: given a street number and a ZIP Code, return all the corresponding street names. The query takes place as follows:

1. obtain pointer to global information from ZIP Code table (basically a disk block offset and length).
2. fetch global information from the disk, and use street number to calculate the block which contains the specific information about the street number.

3. fetch the block containing the street number data, and decompress the block to find street name information corresponding to the specified street number.

3.2.8 Word Recognition

The process flow for machine-print is inadequate for hand-print; when the same flow is used, hand-print address recognition falls short of the desired performance goal. The primary reason is the assumption that addresses contain mostly discrete characters and a few touching characters. Hand-printed addresses have mostly touching characters with few discrete ones. Segmentation of such images into characters is still a research area. Currently we use a lexicon-based word recognizer, where possible. This word recognizer generates a set of segmentation points for each handwritten word such that a character contains at most four segments and no segment contains more than one character. The recognizer uses the lexicon and a character recognizer to group the individual segments into characters. This lexicon-based approach is possible in case of CSZ and street line. In case of the CSZ line, the ZIP Code is recognized first. Using this ZIP Code, a set of possible cities are extracted from the database. This set of cities is passed to the word recognizer as a lexicon along with the image corresponding to the city (Figure 4). Lexicon-based recognizers have been developed at CEDAR for over a decade, with each succeeding method being an improvement over the previous. Using a recent word-model word recognizer [22] we obtained an average of 55% of field correct rate on the CSZ line as shown in Table 3. Similarly one can use ZIP Code and street number to get a set of possible street names. These names can then be used as a lexicon to recognize street names. However, extracting the street name image from the street line is non-trivial, as this involves accurate recognition of street suffix, street prefix, etc.

4 Multiprocessing System Architecture

The speed requirement of SCRIPS is 17,200 forms per hour. Each SCRIPS system employs two NABR units. In order to meet the required processing performance and speed goals (about 8,500 forms per hour), a multiprocessing scheme was chosen as the NABR system architecture. Such an architecture has been previously employed at CEDAR for real-time address block location [6] and for real-time handwritten address interpretation [23]. In this architecture, each of several processors executes identical copies of the address recognition software. A system controller (SC) distributes each successive input image to an available processor. The architecture is illustrated in Figure 6.

4.1 System Controller

Overall system activities are monitored and coordinated by the SC. In order to monitor the system activities, a shared memory provided in the SC's memory space is used. Other components of the system notify their status to the SC by updating the memory contents in the shared memory. Also, the SC can probe each subsystem to detect any failure. The SC can reboot any subsystem selectively based on the probing status. The SC is responsible for distribution of input images to the address recognition processors. If no available processor is found when a new input image is ready, the SC should be able to find a candidate processor to kill a stale job.

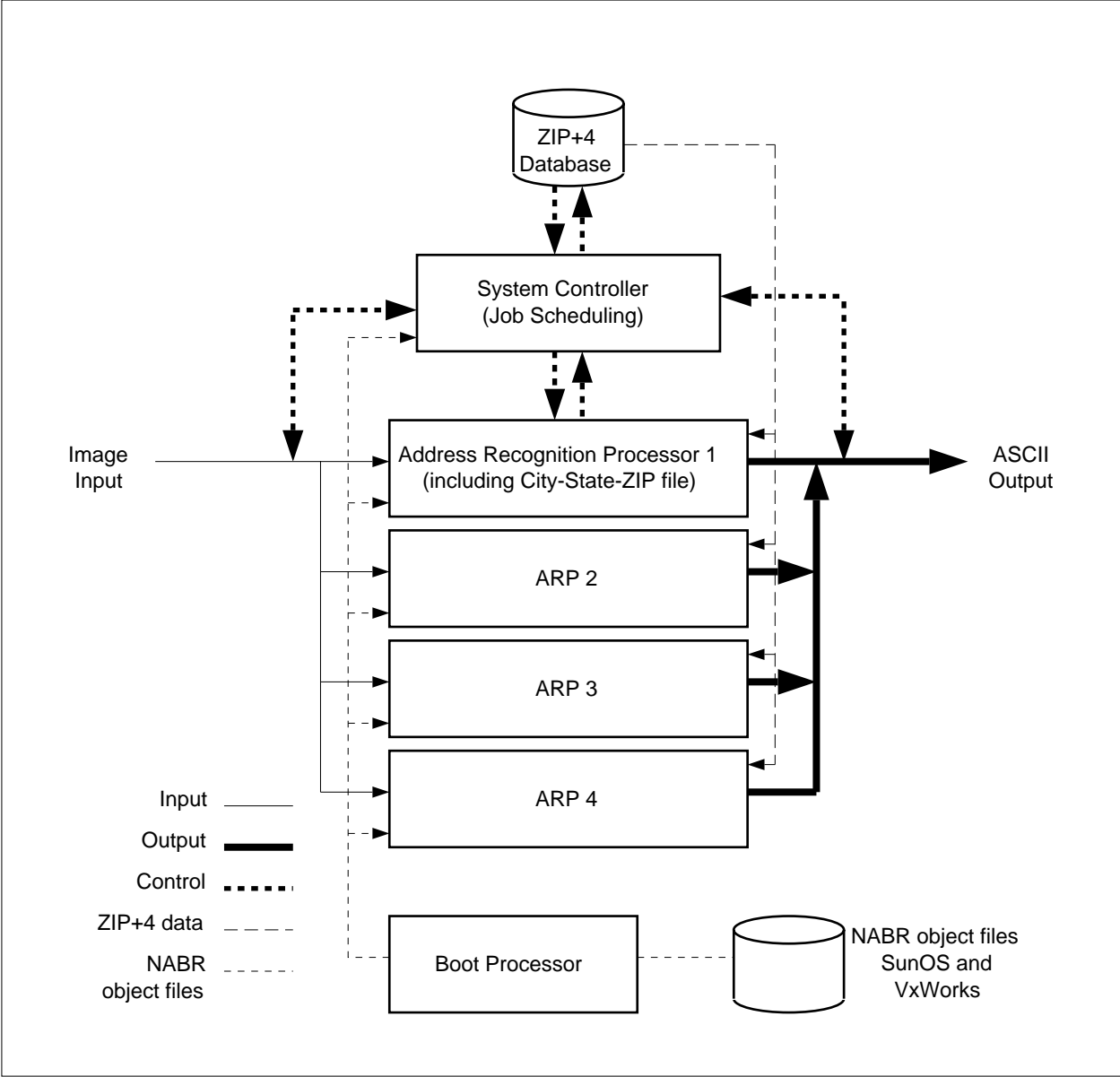


Figure 6: The architecture of NABR system

4.2 Address Recognition Processors

The main processing power of the system is based on four address recognition processors (ARPs) which form a loosely-coupled multiprocessor [24] system. There are two levels of tasks in each ARP. The upper level task is activated when there is no input image to process, and hence the processor is in an idle state, waiting for a start command from the system controller. Once an input image is transferred and stored in a shared memory space on an ARP, the SC issues a start command to the ARP. The start command is recognized by the upper level task and a lower level task is spawned accordingly. The lower level task performs the reading of name and address. Once the lower level task is spawned successfully, the upper level task is put in a sleep mode, monitoring minimal status flags. The upper level task is activated again after the lower level task has been completed.

While the lower level task is activated, a task-delete command can be issued by the SC based on scheduling. The lower level task is deleted upon receiving a task-delete command and the processor becomes free for the next input image.

Each ARP maintains a heart beat, incrementing a specific memory location in shared memory space, so that the SC can identify the status of a processor. The status of the activity is also indicated by a set of LEDs on the front panel.

4.3 Hardware Implementation

The NABR system was implemented using commercial off-the-shelf boards on a VME backplane. The system consists of six SPARC engine boards, three Transputers, one SCSI host adapter, two database disks, one system disk and one tape drive. Figure 7 shows the major components of the NABR system. A view of the actual system is shown in Figure 8.

Among the six SPARC engine boards, four are microSPARC based CPU boards with 47.3 MIPS of computing power [25] for *ARPs*, two are SPARC2 based CPU boards with 28.5 MIPS [26] – one is for *SC*, and another for *boot processor*. The SC and the ARPs are running real-time operating systems VxWorks 5.0.2 and 5.1, respectively. The SPARC2 based SC can map its entire memory to VME memory space, whereas the microSPARC based ARPs can only map maximum one megabyte to VME. The one megabyte mapping window is used for input and output. SC maps each of the one megabyte windows (corresponding to each ARP) into its own address space but at different addresses so that it can differentiate each ARP by looking at the corresponding memory space. The boot processor is running Solaris 1.1 (equivalent to SunOS 4.1.3). The ARPs run same version of the embedded recognition software. For program and local database memory, each board has 64 Mbytes of RAM.

Three Transputers present in the SCSI interface are used to provide a versatile interface between the boot processor and the outside SCSI interface for input and output. One is for the host system interface, the next one is for handling input and output buffer, and the last one for SCSI interface. Differential SCSI interface is employed for reliability in noisy environment.

The ZIP+4 database is stored in two disks and accessed through a SCSI host adapter. The tape drive is used to install a new database, to upgrade recognition software, and to backup system.

When a new input image is fetched and stored in a buffer, this is notified by the SC and a job scheduling is started. Once an available ARP is found as a target, the address of the shared memory space of a target processor is given to the Transputers and the image transfer is initiated. Image transfer operation is performed through the boot processor over the VME interface.

The SC is also responsible for arranging the database access. A database access request generated by an ARP is queued on a first come first service basis. The database access is initiated by the SC through the SCSI host adapter.

Every activity monitored by the SC is saved in a log file for future reference. The file access is performed through Ethernet communication between the SC and the boot processor to minimize any degradation in VME bandwidth that is used for input and output transfer. The Ethernet communication is also used for the SC and ARPs to fetch their own operating system kernels when the system is booted.

4.4 Address Databases

The compressed database is stored in one disk (Barracuda disk with 2.572 Mbytes unformatted capacity). Another disk is provided to hold backup data. The database is accessed on the fly through a SCSI host adapter. Efficient database query is made to enable a fast access to the database. The queries requested by the ARPs are fetched by the SC to provide a centralized access to a single database. The access initiated by the SC is fetched by the SCSI host adapter and the corresponding disk blocks are read and the retrieved information is directly transferred to the shared memory space of the corresponding ARP.

4.5 Boot Processor

The boot processor is based on Solaris 1.1 (SunOS 4.1.3) so as to facilitate system maintenance and to provide versatile communication channels with other system components. Users can login into the SC or the ARPs through the boot processor to find out details of the status. The boot processor has a system disk that contains VxWorks kernel for the system controller and the ARP as well as the Solaris operating systems for itself.

4.6 Input and Output Packets

An input packet consists of a name and address block, and other pertinent information, such as form type and field locations, constructed externally to the NABR. Depending on tax forms, maximum three name and address blocks can be put in an input packet. The input image is a bit-packed binary image in TIFF captured at 200 ppi.

An output packet contains processed information including character recognition results, database access status, and other corrected information associated with a set of confidence values.

4.7 I/O Interface

The differential SCSI interface is used for data transfer between the NABR and the SCRIPS' system control and data management subsystem. The differential SCSI signals received by the NABR are converted to the single ended protocol, then read in by a transputer module with a single ended SCSI interface which runs a code necessary for handling SCSI target transactions, so that the SCSI transputer module looks like a tape drive unit to the outside world for all intents and purposes [27]. Another transputer module [28] interacts with the SCSI transputer module so that the data packet is transferred to and received from the SCSI transputer module. The transputer

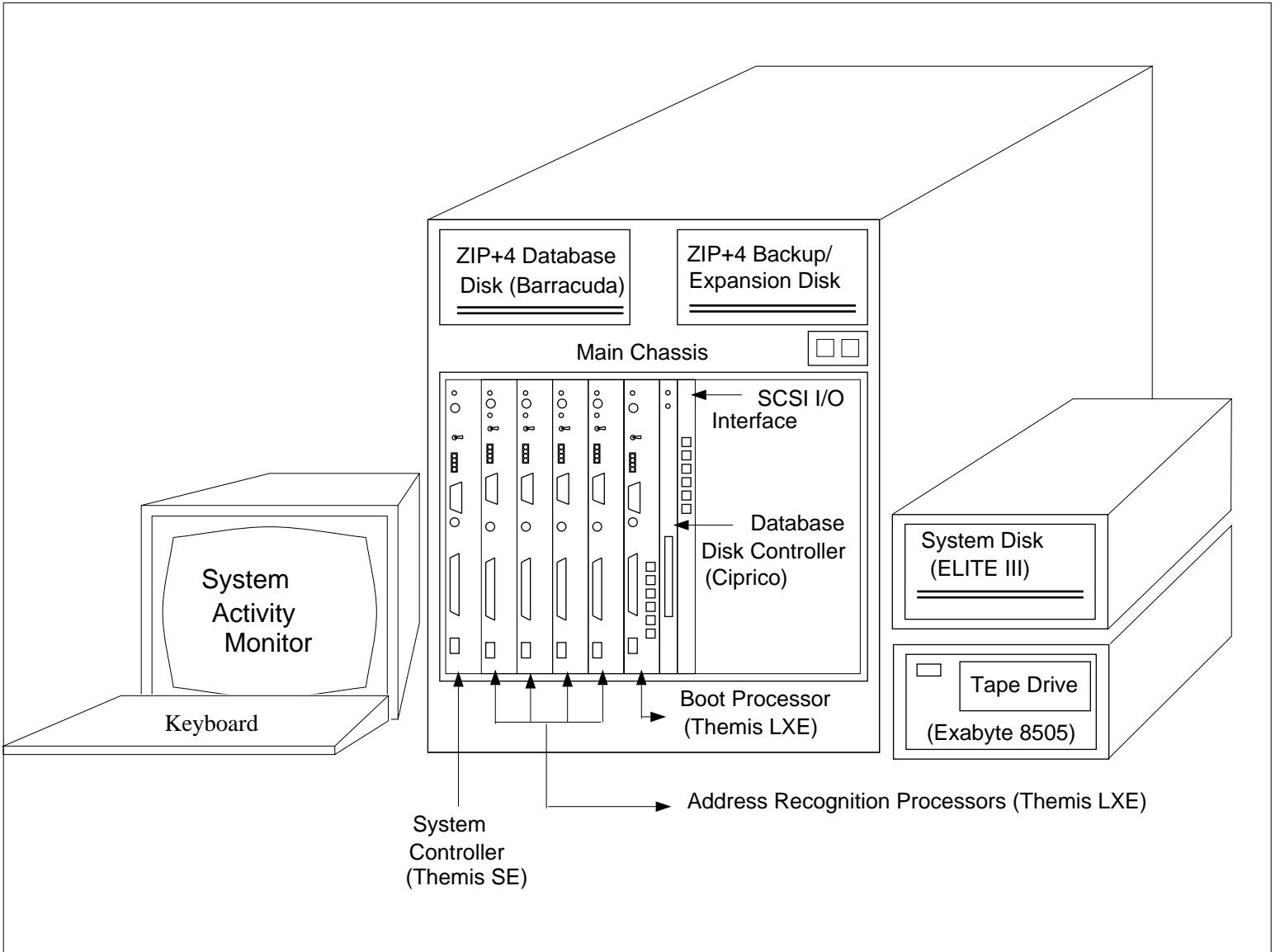


Figure 7: The NABR system

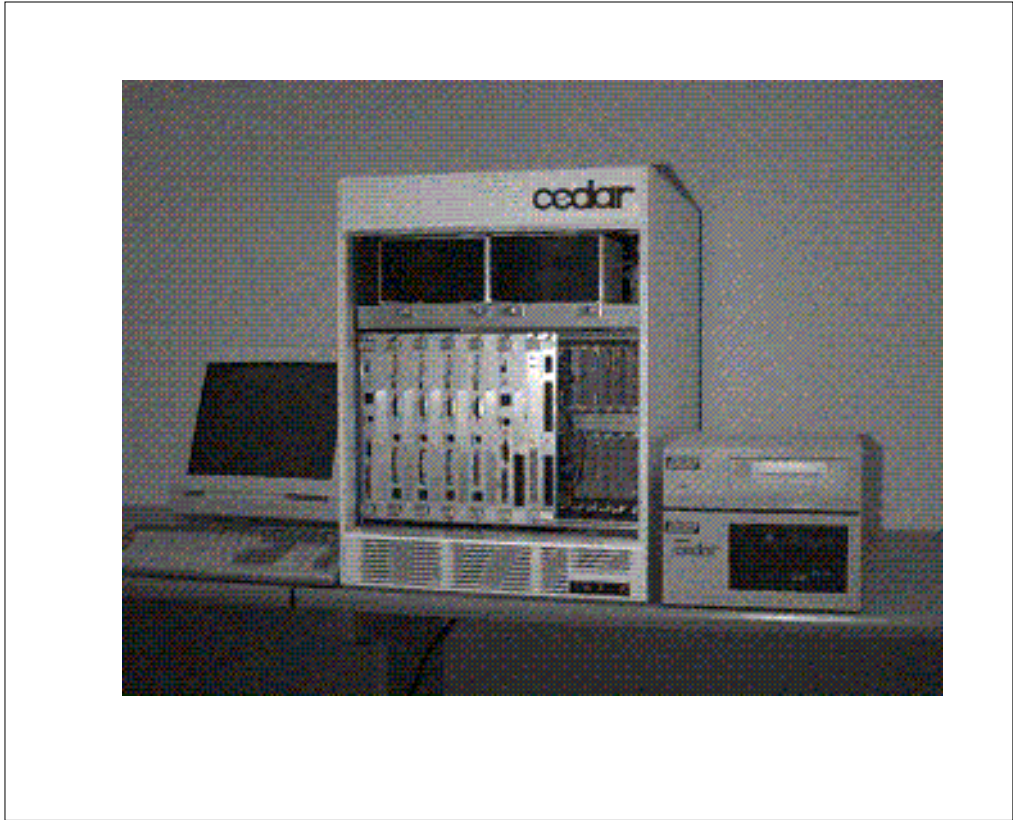


Figure 8: A view of actual system

module also provides extra links so that external transputers can access the NABR system. The transputer module is responsible for queuing images and returning result data. When an address recognition processor is available in the system the next available image in the transputer module is transferred through a HSI/Sbus [29] over the VME bus to the appropriate ARP. The data transfer between the HSI/Sbus and a ARP is performed via programmed I/O.

4.8 Communication

In addition to the data I/O, the NABR also provides an Ethernet LAN to SCRIPS' system control and data management subsystem, so that the overall activities of the NABR can be monitored remotely.

4.9 Scalable Architecture

The entire system is designed to be scalable for any application. Additional computing power can be employed to meet demanding application by simply adding more ARPs in the chassis, without rewriting the control structure.

4.10 Software Complexity

Table 1 gives an idea about the complexity of NABR software as measured by the number of lines of C code for each of modules and sub-modules of the system.

5 Performance

NABR is currently being deployed at several IRS processing centers. For the purpose of reporting performance for this paper, a set of test images was collected for measuring the performance of the current NABR system. The test set consists of 72 hand-printed images and 628 machine-printed images. Table 2 shows the number of images collected per form.

5.1 Performance Measure

Each image has a set of filling areas called fields. For the hand-printed images in the test set, there are five fields; name line (NL), street address line (AD), city (CT), state (ST) and ZIP Code (ZP). Some images have multiple NL fields. For the machine-printed images, there are additional fields such as scan line (SL), quarter date (QT), employer identification number (EIN), and social security number (SS). Field correctness is defined such that all the recognized character should match to the truth of the field. Any mis-recognized, missing or additional character in the output constitutes an incorrect field. Performance on each field has been measured and summarized in Table 3. Performance on those fields (CT, ST,ZP), where contextual information is used, is much higher than on those fields where contextual information is absent. The performance on street address line is not as high as City, State and ZIP fields because for the street address line to be recognized accurately both the ZIP Code and the street number have to be recognized accurately. Also the same street address line can be written in many different variations (like abbreviating (or expanding) the suffix (or prefix), dropping of suffix (or prefix) etc.

Table 1: Software Complexity of the NABR system

Module	Number of Lines of C code
Input/Output Interface	5918
Document Analysis	
Connected Component Analysis	649
Address Block Extraction	966
Line Segmentation	1839
Word Segmentation	1357
HP/MP discrimination	152
Other (Noise Removal, Control Structure etc.)	1005
	Sub Total 5968
Document Recognition	
Machine-print Segmentation	1954
Hand-print Segmentation	3964
Parsing	2884
City, State, Zip- -Postprocessing	1397
Street Line Postprocessing	967
Scan Line Postprocessing	733
Matching Algorithms	663
Control Structure etc	1392
	Subtotal 13415
Character Recognition	
Chaincode Feature Extraction	996
Polynomial Feature Extraction	851
Gradient Feature Extraction	316
Classifier and Other	1261
	Sub Total 3424
Hand-print Word Recognition	5241
CSZ database Interface	3390
ZIP+4 database Interface	1922
Hardware Interface	
Database Disk Controller	1681
ARP Scheduling	3367
System Controller	4112
SCSI I/O Interface	980
Database Loading/Creation	1249
	Sub Total 11389
Total	50667

Table 2: Test Image Tax Forms

Image Class	1040EZ	IRPS	941	1040PC
Hand Printed	15	57	0	0
Machine Printed	125	418	10	75

Character read rate (CRR) is a performance measure obtained using an edit distance between a recognized field and a truth record [30]. Let $R = (R_1, \dots, R_i, \dots, R_n)$ be a list of n recognized fields, and $T = (T_1, \dots, T_i, \dots, T_n)$ be a list of fields in the truth record, where $R_i = \text{null string}$ for some i and $T_i \neq \text{null string}$ for all i . Let W be a list of field pairs, $W_i = (R_i, T_i)$, where $1 \leq i \leq n$, and E_i be an edit distance of W_i . To avoid the character read rate being a negative number, a modified distance, D_i for W_i , is defined by

$$D_i = \begin{cases} E_i & \text{if } \text{strlen}(T_i) \geq \text{strlen}(R_i) \\ \text{strlen}(T_i) - \text{strlen}(R_i) + E_i & \text{otherwise} \end{cases}$$

Then, the character read rate is defined by

$$\text{character read rate} = 1 - \frac{\sum_{i=1}^n D_i}{\sum_{i=1}^n \text{strlen}(T_i)}.$$

The character read rate can be used for obtaining an estimation of the system performance with respect to any human effort to correct the mis-recognized output. However, the character read rate may not reflect the true character recognition performance of the built-in recognizers, since it might include errors from other tools in the system. In other words, the character read rate can be affected and reduced due to segmentation or parsing errors. For example, a 5-character word may be segmented into a 6-character word and hence the resulting recognition cannot be a success. If the word were segmented correctly, it would get the right recognition. Therefore, another performance measure called *OCR correct rate* is defined by considering only the recognition results with same string length in the truth records. Let P be a list of field pairs, $P_j = (R_j, T_j)$, where $\text{strlen}(R_j)$ is equal to $\text{strlen}(T_j)$ for all j , and $1 \leq j \leq p$. Now the OCR correct rate is defined by

$$\text{OCR correct rate} = 1 - \frac{\sum_{j=1}^p E_j}{\sum_{j=1}^p \text{strlen}(T_j)}.$$

Table 4 shows overall performance, including the character read rate and the OCR correct rate.

6 Conclusion

We have described an integrated system for document processing, specifically for reading text contained in address blocks of tax forms. The system combines a number of software modules

Table 3: Field Correct Rates

Filed Name	Hand Printed		Machine Printed	
	Correct Rate	Total Fields	Correct Rate	Total Fields
SL	n/a	n/a	80.59%	170
QT	n/a	n/a	0.00%	5
EIN	n/a	n/a	100.00%	5
SS	n/a	n/a	97.53%	81
NL	8.64%	81	63.43%	689
AD	15.28%	72	71.96%	624
CT	47.22%	72	93.79%	628
ST	59.72%	72	98.41%	627
ZP	58.33%	72	96.82%	628

Table 4: Overall Performance of NABR system

	Hand-Printed	Machine-Printed
Average field correct rate (# of fields)	37.13% (369)	84.50% (3457)
All-field correct rate (# of images)	2.78% (72)	44.43% (628)
Character read rate (dist/chars)	62.68% (1208/3237)	95.44% (1482/32497)
OCR correct rate (dist/chars)	89.53% (147/1404)	97.86% (643/30022)

that CEDAR has developed over the past decade, as well as new techniques developed specifically for tax form processing. The NABR system achieves a throughput of 8,500 forms per hour. It has been implemented using a cost effective real-time processing architecture. This architecture, derived from other real-time architectures developed at CEDAR for address-block location and handwritten address interpretation, was the first system that CEDAR deployed in the field. The design can be easily migrated into other similar document analysis applications. Both hardware and software design details have been presented and the performance of the system was discussed. More research on hand-printed image analysis and recognition is necessary for the development of a higher performing system.

Acknowledgement

The NABR system has evolved from over a decade of address reading research done with the sponsorship of the United States Postal Service. During this period, over two hundred individuals at CEDAR have contributed ideas and code to the system. Several of them are individuals who are authors of the references in this paper. The implementation of the NABR system and its deployment at several IRS sites was made possible by a subcontract from Grumman Data Systems.

References

- [1] T. Minahan, "IRS System Turns Taxpayers' Scrawls into Neat Images," *Government Computer News*, p. 51, Jan. 10, 1994.
- [2] "Document Image Capture Cuts Tax Gobbledygook," *Signals*, vol. 48, no. 10, pp. 26-28, June 1994.
- [3] S. N. Srihari, "Character Recognition," in *Encyclopedia of Artificial Intelligence (Second Edition)* (S. Shapiro, ed.), pp. 138-150, John Wiley, 1992.
- [4] S. N. Srihari, E. Cohen, J. J. Hull, and L. Kuan, "A System to Locate and Recognize ZIP Codes in Handwritten Addresses," *International Journal of Research and Engineering: Postal Applications*, pp. 37-65, 1989.
- [5] S. N. Srihari, "High Performance Reading Machines," *Proceedings of the IEEE*, 80(7), pp. 1120-1132, 1992.
- [6] P. W. Palumbo, S. N. Srihari, J. Soh, R. Sridhar, and V. Demjanenko, "Postal address block location in real time," *IEEE Computer*, pp. 34-42, 1992.
- [7] T. Pavlidis, *Algorithms For Graphics and Image Processing*. Computer Science Press, Rockville, MD, 1982.
- [8] R. O. Duda and P. E. Hart, *Pattern Classification and Scene Analysis*. Wiley-Interscience, 1973.
- [9] J. F. Giattino, "Fisher's Pairwise Linear Discriminant Method for Font Type Classification," M. S. Project Report, Department of Computer Science, SUNY at Buffalo, April 1993.

- [10] R. Fenrich and S. N. Srihari, "System for Recognizing Handwritten Character Strings Containing Overlapping and/or Broken Characters," *United States Patent No. 5,321,768*, June 14, 1994.
- [11] S. N. Srihari, "Recognition of Handwritten and Machine-Printed Text for Postal Address Interpretation," *Pattern Recognition Letters*, vol. 14, pp. 291–302, 1993.
- [12] G. Srikantan, S. W. Lam, and S. N. Srihari, "Gradient-based contour encoding for grayscale character recognition," in *Proceedings, SPIE Symposium on Electronic Imaging*, (San Jose, California), Jan-Feb 1993.
- [13] U. Srinivasan, "Polynomial Discriminant Method of Handwritten ZIP Code Recognition," M. S. Project Report, Department of Computer Science, SUNY at Buffalo, November 1989.
- [14] J. T. Favata, G. Srikantan, and S. N. Srihari, "Hand-printed digit/character recognition using a multiple feature/resolution philosophy," in *Proceedings, International Workshop on Frontiers in Handwriting Recognition*, (Taipei, Taiwan), pp. 57–66, 1994.
- [15] S. N. Srihari, S. W. Lam, P. B. Cullen, and T. K. Ho, "Document Image Analysis and Recognition," in *In Artificial Intelligence Methods and Applications* (N. G. Bourbakis, ed.), pp. 590–617, World Scientific, 1992.
- [16] P. W. Palumbo and S. N. Srihari, "Text parsing using spatial information for recognizing addresses in mail pieces," in *Proceedings, International Conference on Pattern Recognition*, (Paris, France), pp. 1068–1070, 1986.
- [17] M. Prussak and J. J. Hull, "A multi-level pattern matching method for text image parsing," in *Proceedings, IEEE-CS Conference in Applications of Artificial Intelligence*, (Miami Beach, FL), pp. 183–189, 1991.
- [18] J. J. Hull, V. Demjanenko, R. Sridhar, and S. N. Srihari, "Architectural design for address interpretation in a Laboratory Address Recognition Unit (LARU)," in *Proceedings, USPS Advanced Technology Conference*, pp. 1313–1323, November 1992.
- [19] Technology Resource Center, Arthur D. Little, Inc. , "Specification of directory modifications," Technical Report Task Order 90-08, Contract No. 104230-83-D-1902, United States Postal Service, December 1990.
- [20] D. Huffman, "A method for the construction of minimum-redundancy codes," in *Proceeding of IEEE*, pp. 1098–1101, 1952.
- [21] P. Filipiski, J. Hull, and S. N. Srihari, "Compression for Fast Read-Only Access of a Large Database," in *Proceedings, the United States Postal Service Advanced Technology Conference*, (Washington, D.C.), pp. 1113–1127, Nov-Dec 1992.
- [22] G. Kim and V. Govindaraju, "Handwritten Word Recognition for real-time applications," in *Proceedings, Third International Conference on Document Analysis and Recognition*, (Montreal, Canada), August 1995.
- [23] V. Govindaraju, S. N. Srihari, and A. Shekhawat, "Interpretation of Handwritten Addresses in U.S. Mail Stream," in *Proceedings, Second International Conference on Document Analysis and Recognition*, (Tsukuba Science City, Japan), pp. 291–294, October 1993.

- [24] K. Hwang and F. A. Briggs, *Computer Architecture and Parallel Processing*. Mc-Graw Hill, New York, NY, 1984.
- [25] Themis Computer, *SPARC LXE User's Manual*, 1993.
- [26] Themis Computer, *SPARC 2SE/PLUS Hardware Reference Manual*, 1993.
- [27] Alta Technology, *SCSITRAM – SCSI 2 TRAM with Optional SCSI Differential Driver Board – Installation Guide and User Reference Manual*, 1994.
- [28] Alta Technology, *CTRAM2 Data Sheet*, 1994.
- [29] Alta Technology, *HSI/SBus Installation and User Reference*, 1994.
- [30] R. A. Wagner and M. J. Fischer, “The String-to-String Correction Problem,” *Journal of the ACM*, vol. 21, pp. 168–173, Jan 1974.

Legends

Figure 1 Examples of images input to NABR. (a) 1040EZ Hand-print with no drop-out ink, (b) 1040 Label image with OCRA Font, (c) 941 Hand-print with drop-out ink and (d) 1096 Machine-print with drop-out ink.

Figure 2. Top-level Modular Organization.

Figure 3. Document Analysis.

Figure 4. Document Recognition.

Figure 5. Hand-printed digit recognizer

Figure 6. The architecture of NABR System

Figure 7. The NABR System

Figure 8. A view of actual system