

Use of Character Recognition and Address Syntax in Address Block Location

Jenchyou Lii, Paul W. Palumbo, and Sargur N. Srihari

Center of Excellence for Document Analysis and Recognition

State University of New York at Buffalo

226 Bell Hall

Buffalo, New York 14260

email: lii@cedar.buffalo.edu

Abstract

Our previous work on address block location has demonstrated that it is possible to achieve a 90% location rate on difficult cases, using low-resolution (100 PPI) images. Such a system does not perform any character recognition. In order to increase the location performance level to, say 95% or higher, it is necessary to incorporate a higher degree of intelligence into the system. One such effort is to incorporate character recognition and address parsing into address block location to supplement the previously pure statistical approach previously developed. The format of an address block offers a mix of syntax as well as context information which, once extracted, will be helpful to do block segmentation refinement and block identification.

1 Introduction

The automated sorting process mainly consists of two parts: the address block location (ABL) and the address interpretation. Although it is assumed that address interpretation is performed after the address block is located, it does not exclude the possibility of using some address recognition to aid the location. For address block location, a complete interpretation of the address need not be performed to determine delivery point encoding, but address recognition can be used in the decision of selecting the desired block from among several candidate blocks.

To read an address block entails recognizing characters and parsing the address. Address syntax itself is a salient feature which can be taken to discriminate a destination address block (DAB) from other text blocks on a mail piece. By recognizing characters and analyzing the address syntax, the system is able to select a more plausible DAB candidate block in a complex environment.

In addition, it is sometimes inevitable that the preliminary DAB blocks, produced by the block segmentation tools, contain irrelevant components and are larger than they should be. To improve DAB segmentation, better DAB boundaries can be

determined with the help of the address syntax. For example, the critical information, defined as the combination of city, state, and ZIP Code, constitutes the last part of an address. The bottom boundary of an address block will become obvious once the critical information is located.

Since the machine-printed mail dominates all letter mail and is easier to apply address interpretation, we focus on improving the ABL performance in this problem domain. In this paper, we illustrate our approach and its advantages in using character recognition and address syntax to locate the destination address block on machine-addressed mail.

2 System Review

For the convenience of readers to appreciate the usage of the address interpretation in address block location, a little background about our earlier system is provided [1][2][3][4].

Figure 1 shows the layout of the previous system. The adaptive thresholding tool (ADTH) first converts a grey-level image into a binary image. Afterwards, two modules, MSEG1 and HSEG1, group neighboring black pixels to produce two types of connected components for machine-printed and handwritten addresses, respectively. The connected component data is then sent to the block segmentation module to cluster into lines and blocks.

In total, there are five component grouping tools to cluster connected components into lines and blocks. One tool (GSEG) takes a top-down approach and the other four tools (MSEG/HSEG/LSEG/PSEG) take a bottom-up approach. GSEG is a top-down segmentation tool applied on acceptable mail pieces. An acceptable mail piece is defined to contain three or fewer blocks of information, i.e. return address, stamp and cancellation mark, and destination address block. The destination address block produced by GSEG will bypass the rest of the system and be considered as the system's final output. The other four bottom-up segmentation tools work on the same mail piece and each of them produces its own set of preliminary blocks. In sequence, four sets of preliminary blocks are sent to the Block Analysis Tools to be refined in recovering imperfect blocks as well as to be analyzed in ranking the DAB candidate blocks. Particularly, MSEG is tuned to process machine-printed addresses, HSEG is tuned to process handwritten addresses, whereas LSEG is designed to form better address text lines for both print methods, and PSEG mainly deals with postcards.

Blocks output by the four bottom-up segmentation tools are separately processed by the Block Analysis Tools (BAT). The BAT performs two tasks. One is to

improve block segmentation and the other is to rank DAB candidate blocks using statistical methods.

On one hand, the BAT consists of a few specialized tools to extract block features and improve the block segmentation, including the Refinement Tool (REFN), the Block Feature Analysis Tool (BFAT), the Block Splitting Tool (BLCS) and the ZIP Code Merging Tool (ZIPM) to adjust the segmentation results. The Refinement Tool improves component segmentation by either merging broken components, or recovering lost components previously thought to be noise. The Block Feature Analysis Tool is responsible for extracting word, line, and block level features from connected component data. Useful results are stored in a global data structure where it can later be accessed by other tools, such as BLCS and ZIPM. The Block Splitting Tool improves block segmentation by removing extra material from a block, e.g. return address, postage, cancellation mark, bar code and so on. The ZIP Code Merging Tool is to find the separate ZIP Code block, due to extra space, and merge it with the rest of the DAB block.

On the other hand, the Block Analysis Tools is in charge of ranking the candidate destination address blocks from each segmentation tool based on a statistical approach. Many relatively specialized tools are employed to compute geometrical features, such as block size, block layout, block location, block overlapping, block consistency verification, and so on. Each of them subsequently generates a corresponding piece of evidence. The evidence combination tool (EVHP), a scheme based on the Dempster-Shafer theory, is then used to combine these pieces of evidence and provides a single confidence value which is used to rank blocks based on their likelihood of being the actual DAB.

The structure of each segmentation tool being followed by a BAT is called 4-stream BAT, since it has four computational flows. Each stream independently generates its own ranking of candidate blocks. Another critical process, the ranking Combination Tool (COMB) is designed to draw a consensus from these four independent ranking results. The consensus is based on (1) confidence values and the ranks of candidate blocks in each ranking and (2) determining if two blocks in different rankings actually correspond to the same block. A combined ranking is thus produced and the top ranked block is taken as the most likely destination address block. Finally, the Postprocessing Tool (POST), placed at the end of the system, attempts to correct for minor segmentation problems in the top-ranked candidate DAB. It removes extra material such as the words "To" and "Pay to the order of".

3 Use Character Recognition and Address Syntax in the Block Analysis Tools

It is not surprising to see a destination address block embedded in a complex environment on a machine-printed mail piece. Figure 3 is a typical example. In this case, several text blocks are present on a mail piece and they pose a challenge for a system to select the correct DAB. A human identifies the address block by means of mixed knowledge. Best of all is the postal address syntax. A text block with address syntax distinguishes itself from non-address blocks. Therefore, it is attractive to bring character recognition and address syntax analysis into an address block location system to help DAB identification.

Our previous work on address block location did not perform any character recognition. To complement the statistical method used in our earlier system, we have incorporated the syntactic approach to improve the system performance. As soon as the system reads all characters in a text block, word segmentation and word string matching are executed to assign a syntactic tag to each word. Important postal syntactic features, such as the critical information and the street syntax, are extracted by parsing these syntactic tags. Several tools, including the ZIP Code Merging Tool, the Syntax Tool, and the Postprocessing Tool, use the syntactic features to either improve the block segmentation or better select the destination address block. The ZIP Code Merging Tool assumes the responsibility of recovering the lost ZIP Code and has been made more powerful due to the availability of character recognition and the syntactic features. The Syntax Tool has been added to the system to parse the candidate address blocks. Ranking priority will, therefore, be given to those blocks which demonstrate the address syntax. The Postprocessing Tool improves the block segmentation mainly relying on the results of address syntax analysis.

In the following subsections, we elaborate on the methodologies of performing character recognition, computing syntactic features, and using address syntax in individual tools in the system. Subsection 3.1 describes how the system achieves character segmentation and recognizes characters. Subsection 3.2 discusses word level processing, i.e. word segmentation and word string matching. The detection of the critical information is represented in subsection 3.3. The last three subsections are devoted to three applications of using the syntactic features derived from character recognition. Subsection 3.4 narrates how to use character recognition in the ZIP Code Merging tool. Subsection 3.5 details the way of using character recognition in the Syntax tool. Using address syntax to improve DAB segmentation is introduced in subsection 3.6.

3.1 Character Segmentation and Character Recognition

Character segmentation is primarily accomplished by MSEG1 and HSEG1, where each connected component presumably corresponds to one character. The Refinement Tool further operates on these connected components to fix possibly broken characters. Subsequently, the Block Feature Analysis Tool (BFAT) calls the Character Recognition Unit (CRU) to recognize characters of every line in a block. The relationship between these tools is drawn in Figure 2.

Character recognition is performed in ABL at 100 PPI resolution images, instead of 300 PPI resolution images. There are two reasons for doing this. The first reason is 100 PPI is the resolution used in other portions of the ABL system to segment the candidate blocks for locating the destination address. It is consistent to have a character recognition algorithm work on the same resolution as well. The second reason is character recognition performed at 300 PPI requires about nine times the computation time than it does at 100 PPI. To facilitate the implementation of a real-time system, 100 PPI character recognition is preferred.

The current methods to perform character recognition at 100 PPI is the gradient based Bayesian classification approach. Since this algorithm appears to work well at 300 PPI in the Machine-Print Character Recognition project at CEDAR, the CEDAR ABL group explored the use of this algorithm at lower resolution. It is possible to switch to another algorithm, e.g. neural network based classification, if the new method is proven better.

Unlike the binary images, on the gray level images it is possible to take into account the intensity variations. The gradient-based approach utilizes the grey-level change in the intensity of various pixels to perform machine-printed character recognition. For each character image, a gradient map of the entire image is calculated. The gradient map depicts the gradient magnitude and direction of each pixel which in turn depicts the magnitude and direction of maximum change in gray level intensity for each image point. The Sobel operator is used to calculate the gradient map of the image. Once the gradient map is obtained, the "character" image is partitioned along both horizontal and vertical directions and is thus divided into a fixed number of rectangular cells. The gradient directions are also quantized into a fixed number, k , of gradient angle bins and each bin covers a range of $360/k$ gradient angles.

Every gradient angle bin in each cell corresponds to one feature. The total number of features is equal to the product of the number of vertical divisions, the number of horizontal divisions and the number of bins. A feature is said to be present if its corresponding bin contains pixels which have a gradient magnitude greater than a threshold value. Otherwise, the feature is considered absent. It is the presence or absence of these features which is used as a feature vector for character recognition.

Once the features are extracted, a probability value of the test "character" is obtained by multiplying the feature vector with the weight matrix which is calculated based upon the training set. The probability value will then be used by Bayesian classification in determining the class of the test "character" among sixty-two classes, including fifty-two upper and lower case characters plus ten digits.

3.2 Word Segmentation and Word String Matching

Based on the gap between connected components, word segmentation is carried out by BFAT. The String Matching Tool (STRM), invoked by BFAT (see Figure 2), will process each input word and assign it a proper tag (e.g. state tag, numeral tag, street suffix tag, and so on). If the majority of the characters in an input word are digits, a numeral tag will be attached to that word. Otherwise, the String Matching Tool proceeds to match the input word against the keywords stored in a database and labels it with an appropriate tag when a match is found. For short words with 2, 3, or 4 characters, an exact match of the character recognition results and a keyword is required. For longer words, the longest common subsequence algorithm is used to compare the input word against the keywords. Tolerance is allowed for one or two characters in the long word matching.

The way the keywords are stored in the database is not purely in lexicographical order. Instead they are first sorted according to character length. The words of the same length constitute a partition of the database. And then words inside each partition are arranged in lexicographical order. Also, an array is used to record the starting positions of every partition. The advantage of making such an arrangement is that not every keyword needs to be compared with a given input word. Only keywords whose length approximates the length of the input word deserve matching consideration. As we know, it does not make too much sense to compare a two-character word with a ten-character word. Therefore, the search domain in the database is reduced to a few continuous partitions, compared to the whole database otherwise. To reduce the computation time further, a binary search is used to look up a particular partition of the database whenever exact matching is required.

If more than one keyword is acceptable, a ranking mechanism is activated to determine the best matched word. However, if no keyword is close enough to be compatible with the input word, an unknown tag will be assigned. So far, we have used the state tag, the numeral tag, and tags denoted for P.O. Box, street suffix, and title in the keyword database as seen in Table 1. There are about 500 entries in the keyword database. More keywords and syntactic tags will be added into the database as the system is being developed.

3.3 Detect the Critical Information

Through analyzing those syntactic tags, a text line is parsed to detect the partial/complete critical information. As a result, it falls into one of the four categories as shown in Table 2. In the first case, a numeral tag with 5 or 9 characters is found to follow a state tag. The second case indicates a state tag is detected, but no numeral tag. The ZIP Code may stand alone in one line as seen in the third case. Otherwise, it belongs to the fourth case. Notice that the city name is optional and may contain up to three words if it exists.

Our earlier approach to detect the critical information was built on the pseudo word recognition, where the connected component counts of "words" were used to locate the possible state abbreviation and the ZIP Code. Using the character recognition results tremendously enforced our address syntax analysis. At present, not only the state abbreviation and the ZIP Code can be strongly verified, but also the full state name can be identified.

3.4 Use Character Recognition in the ZIP Code Merging Tool

Machine-printed address blocks often have extra space between the city/state names and the state/ZIP Code. During the connected component grouping, this excessive space results in text being excluded from the main body of the address. The ZIP Code Merging Tool recovers from this inevitable exclusion of the state and/or the ZIP Code by checking a set of conditions. Among the blocks produced by each segmentation tool, a list of DAB candidate blocks and a list of one-line text blocks are chosen by the ZIPM tool. These one-line text blocks are suspicious of holding the hypothesized missing information. Hereafter, they are called ZIP Code candidate blocks though they possibly comprise the city and the state as well. Generally speaking, three conditions need checking in order to pair match DAB and ZIP Code candidate blocks.

First, both DAB and ZIP Code candidate blocks have the same print method, either machine-printed or handwritten. There is a tool available in the system, called the handwriting machine-printing discrimination tool (HWMP). The HWMP tool classifies a segmented block and labeled it as being printed either by hand or by machine. By inspecting the handwriting/machine-printing tags of both DAB and ZIP Code candidate blocks, the ZIP Code Merging Tool will not merge two blocks with different print methods.

Secondly, certain constraints are imposed on the relative positions of DAB and ZIP Code candidate blocks. An acceptable region is defined to be around the south-

eastern outskirts of a DAB candidate block. A ZIP Code candidate block must be present in this region in order to be accepted.

Thirdly, the address syntax between DAB and ZIP Code candidate blocks should be consistent. The detection of the critical information in the DAB and ZIP Code candidate blocks are performed as we mentioned in the previous subsection and pieces of their critical information need to match. For example, if a DAB candidate is hypothesized to be missing the ZIP Code and a ZIP Code block demonstrates containing a 5 or 9 digit number, then they are said to be consistent in the syntax.

Without using address syntax, the system has a hard time doing ZIP Code merging in the following cases. When a block resembling a ZIP Code is present around the DAB candidate block, how does the system decide whether the DAB candidate block has lost its ZIP Code? If more than one small machine-printed block is hanging around a DAB candidate block, which one should be selected as the ZIP Code block for merging? Using character recognition and address syntax analysis offers a good solution. By reading the text blocks, the system is able to perform better in recovering the lost ZIP Code block and at the same time avoids merging irrelevant material.

3.5 Use Character Recognition in the Syntax Tool

If all machine-printed mail pieces simply contains three types of objects, i.e. return address, destination address, and stamp/cancellation mark, analyzing the address syntax seems unnecessary for serving the purpose of locating the destination address. On the contrary, since quite a percentage of machine-printed mail pieces contain other unpredictable text and nontext blocks, it becomes uncertain as to selecting the correct destination address block **without** analyzing the address syntax. For this reason, the Syntax Tool is incorporated into the system (see Figure 2) to parse the entire contents of a block based on the address syntax. When a block shows good postal syntactic features with its critical information being detected, it will be favored as the desired block over other competitors. Internally, the Syntax Tool consists of three modules, i.e. the Syntax Feature Analysis Tool, the Syntax Ranking Tool, and the Syntax Combination Tool. The first tool does the actual address syntax analysis. The second tool tries to select from each stream of the system a good DAB candidate whose address syntax can be detected. The last tool is a decision tool which takes into account of evidence from the statistical portion of the system for address syntax. More details of their functionality are described as follows.

Because the critical information is the most important part of an address and its unique structure makes it easy to recognize, the Syntax Feature Analysis Tool examines a block and decides if the block contains the critical information. Recall

that every word in a block has been attached a tag by the Sting Matching Tool, the fact that a state tag and a numeral tag are labeled consecutively either on the same line or on neighboring lines strongly suggests the existence of the critical information. A block can also be considered as a plausible address block if it has been created by the ZIPM Tool, and the critical information has been recovered. Therefore, a block stands out among others by showing the evidence of possessing the state and ZIP Code.

As we mentioned earlier, there are four computational streams corresponding to four segmentation tools in our system. At the end of each stream, a ranked list of DAB candidate blocks is computed based on statistical combination of diverse pieces of evidence. To compensate for the possibility of misranking a non-address block as the top ranked block, the Syntax Ranking Tool plays the role of rearranging the order of candidates if another more promising block is detected by the Syntax Feature Analysis Tool. In fact, the Syntax Ranking Tool picks up the first block which shows the critical information from the ranking list of each stream. Though it may happen that no block is selected due to the unavailability of the address syntax, up to four DAB candidates, one from each stream, may be collected at this stage. These favored blocks are further passed to the Syntax Combination Tool where the final choice is made.

If there is more than one block provided by the Syntax Ranking Tool, the Syntax Combination Tool first classifies blocks into groups based on block overlapping and selects the group close to the lower right corner on the mail piece. The best looking block from that group will then become the final choice. If there is only one block provided by the Syntax Ranking Tool and its location is reasonable, the block will be taken as the destination address block. Otherwise, the Syntax Combination Tool will output the block generated from the Combination Tool as it was in the previous system.

3.6 Use Address Syntax to Improve DAB Segmentation

Though an address block which contains a small amount of extra text or non-text components may not be absolutely unacceptable, it is still worth while finding a way to improve the block segmentation. After the DAB is selected by the Syntax Combination Tool, the Postprocessing Tool uses the address syntax to improve the DAB segmentation. The objective is to locate correct the DAB boundaries and consequently eliminate undesirable components from the DAB to reduce the number of minor oversegmentation problems produced by our system.

It is advantageous to begin with the critical information, if it can be located. The critical information defines the last portion of an address block and, therefore,

anything below the state/city/ZIP Code can be ignored. Once the state and the ZIP Code are identified, the left boundary will be determined with the help of the locations of the possible city words preceding the state name. As the city name may consist of up to three words, the first three words located before the state name suggest three vertical cutting edges. The POST tool needs to examine the positions of these potential cutting lines and uses other information to judge the left boundary of a DAB block. The fact that most of machine-printed address blocks are left aligned is a good clue to follow. If a potential cutting line intersects words in a text block, it is not considered as a valid left boundary. Whereas, if a potential cutting line just touches the head of a word in every line, it is favored as a left boundary. Also, if it happens that P.O. Box can be recognized, their starting position will be directly taken as the left boundary.

Furthermore, using the critical information as the base line, the POST tool checks the distance between text lines from the bottom up. Big gaps usually indicates a possible top boundary with the exception of the first title line. In certain situation, the top boundary of an address block becomes obvious, if title (e.g. Mr., Mrs.) can be identified. In addition, the information regarding street number, P.O. Box and street suffix reveals the layout of an address block. The syntax of the street information can be obtained by observing either P.O. Box immediately followed by a number or a number followed by a street suffix with a few words apart. Whenever the street syntax is confirmed, the POST tool will be careful to include other text into the top of the street line.

4 Experiments and Results

The algorithm used for character recognition was tested with slight changes in the parameters (i.e. number of image partitions and number of gradient angle bin) for 100 PPI and the results were then compared with that for 300 PPI. The training set consisted of 47565 images and the testing set of 9500 images. Each image was partitioned into 5 vertical and 5 horizontal divisions and 18 gradient angle bins. The total number of features is $5 * 5 * 18 = 450$. The recognition success rate for 100 PPI and 300 PPI was 76.5% and 87.37% using 62 recognition classes. It is projected that once folding of similar classes is performed (e.g. O, 0, and o), the recognition performance should increase by about 7%.

Based on experiments with 365 machine-printed images, we estimate the incorporation of character recognition and its usage in ABL can increase the system performance by more than 3%. We use an example to illustrate how to use character recognition to help locate the destination address block. The test mail piece and its block segmentation results from MSEG are shown in Figure 3. Table 3 shows the

character recognition and the string matching results for block #1, where each word is labeled with either a syntactic tag or an unknown tag. Since the critical information can be detected by observing a state tag followed by a numeral tag in the last line, block #1 is favored as the destination address block. The final system result is shown in Figure 4.

5 Conclusions

We have shown many advantages of using postal syntactic features in address block location. Performing character recognition at low resolution (100 PPI) is desirable for a real-time system. Efficient string matching algorithm is designed to match words against a dictionary of about 500 keywords and assign a syntactic tag to each word in a text block. The critical information and the street syntax are extracted by parsing the syntactic tags. The usage of the postal syntactic features has been introduced in many aspects such as recovering lost information, selecting a more likely destination address block, and improving DAB block segmentation. With character recognition is available in the system, the utilization of other syntax and context information on the mail pieces becomes feasible. The initial testing results of the incorporation of the syntactic approach was positive. We are looking forward to doing more testing and maturing the system.

References

- [1] P.W. Palumbo, Jung Soh, S.N. Srihari, R. Sridhar, V. Demjanenko, Postal Address Block Location in Real Time, *IEEE Computer*, 34-42, July 1992.
- [2] P.W. Palumbo, Jung Soh, S.N. Srihari, V. Demjanenko, R. Sridhar, Real-Time Address Block Location using Pipelining and Multiprocessing, *Proceedings of USPS Advanced Technology Conference*, Vol. 1, 73-87, 1990.
- [3] S.N. Srihari, C.-H. Wang, P.W. Palumbo, J.J. Hull, Recognizing address blocks on mail pieces: specialized tools and problem-solving architecture, *AI Magazine*, 8(4), 25-40 (Winter 1987).
- [4] C.-H. Wang, S.N. Srihari, A framework for object recognition in a visually complex environment and its application to locating address blocks on mail pieces, *International Journal of Computer Vision*, 2(2), 125-152 (1988).

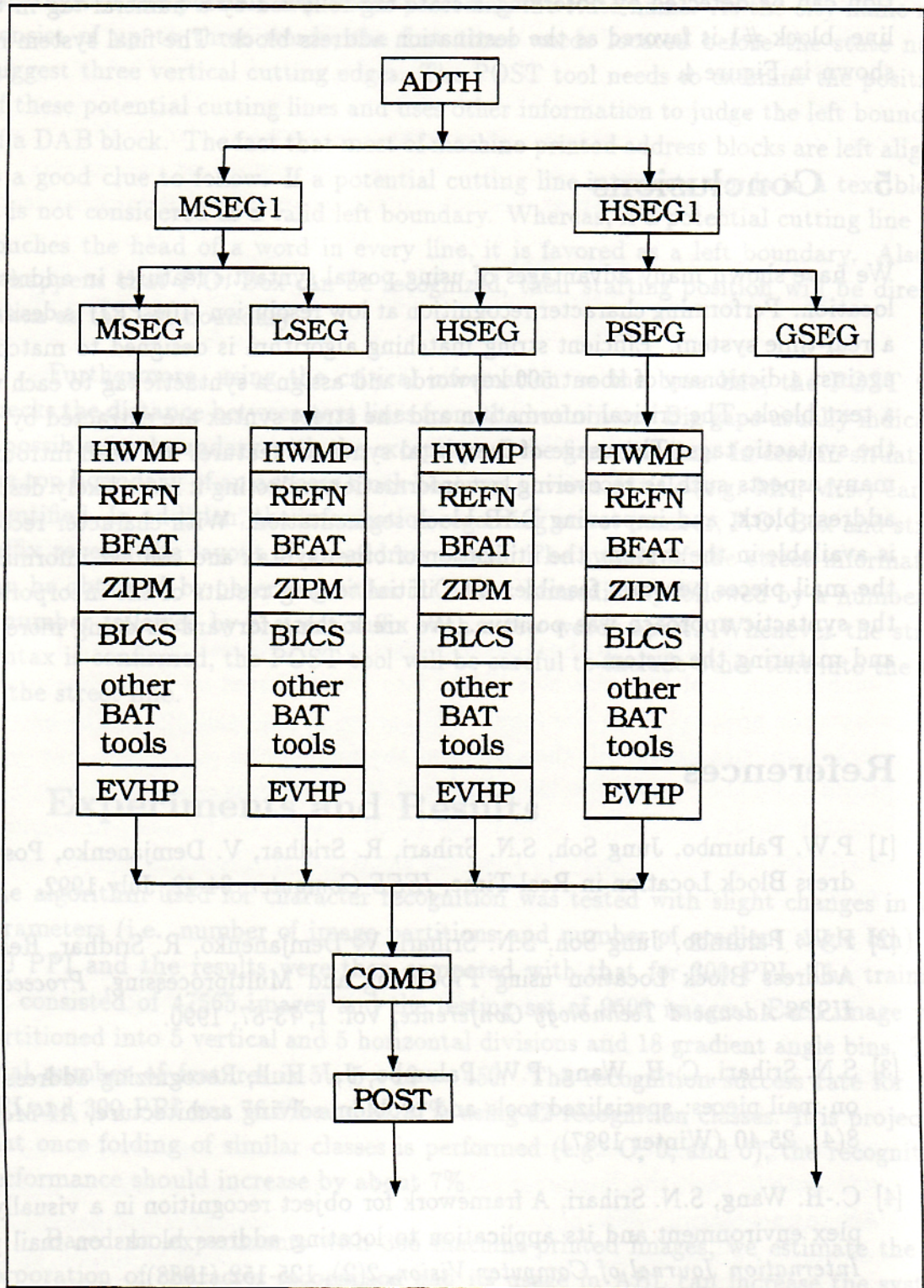


Figure 1: System Layout

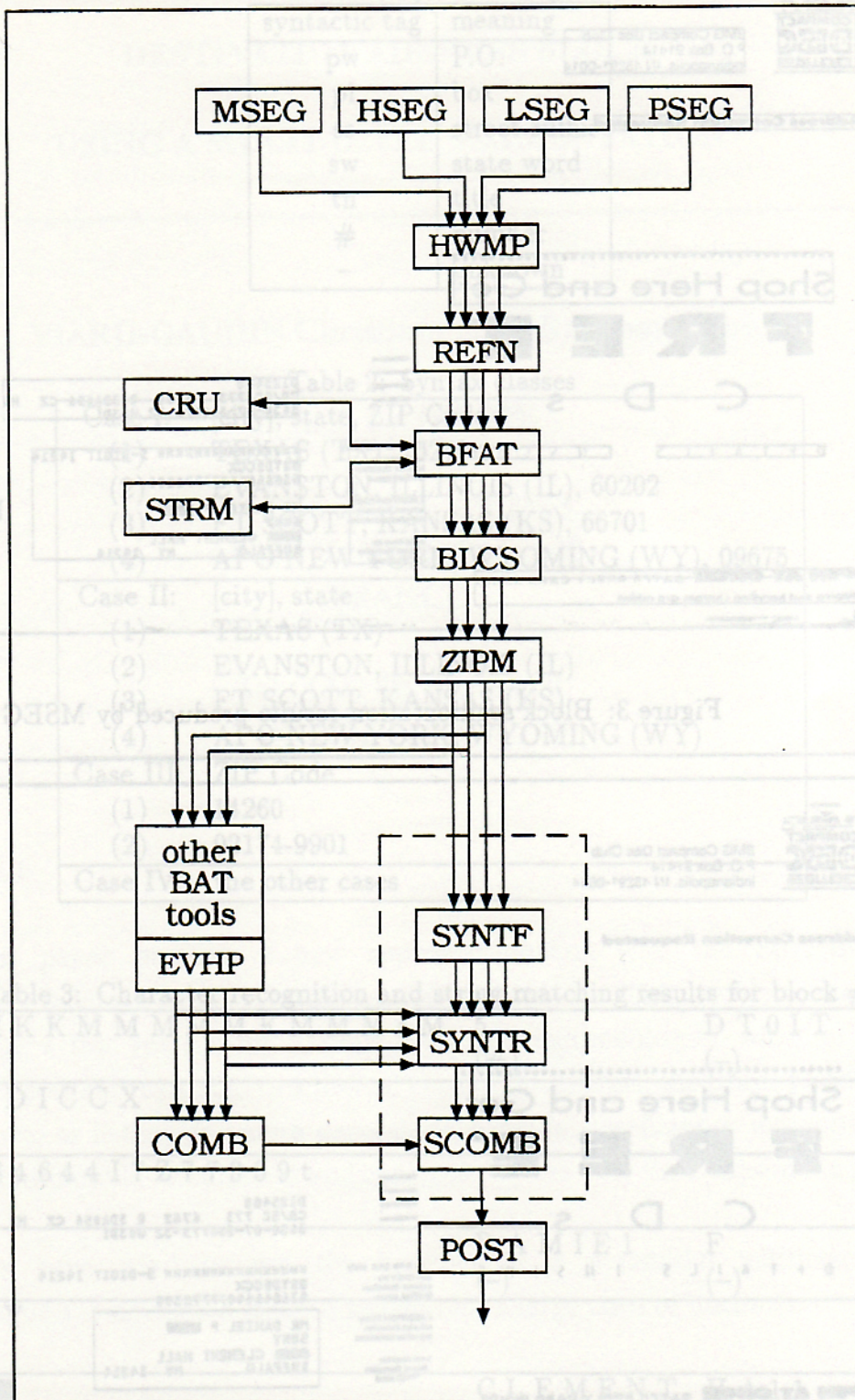


Figure 2: New Configuration of BAT Tools

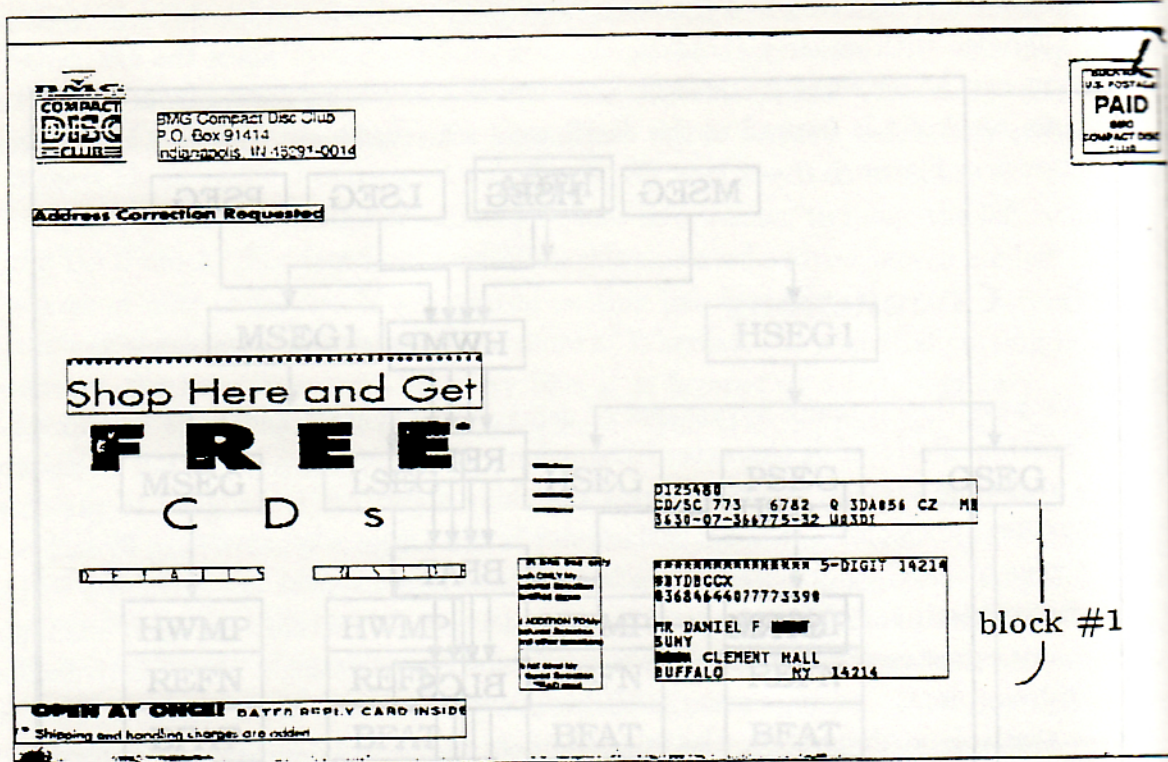


Figure 3: Block segmentation results produced by MSEG

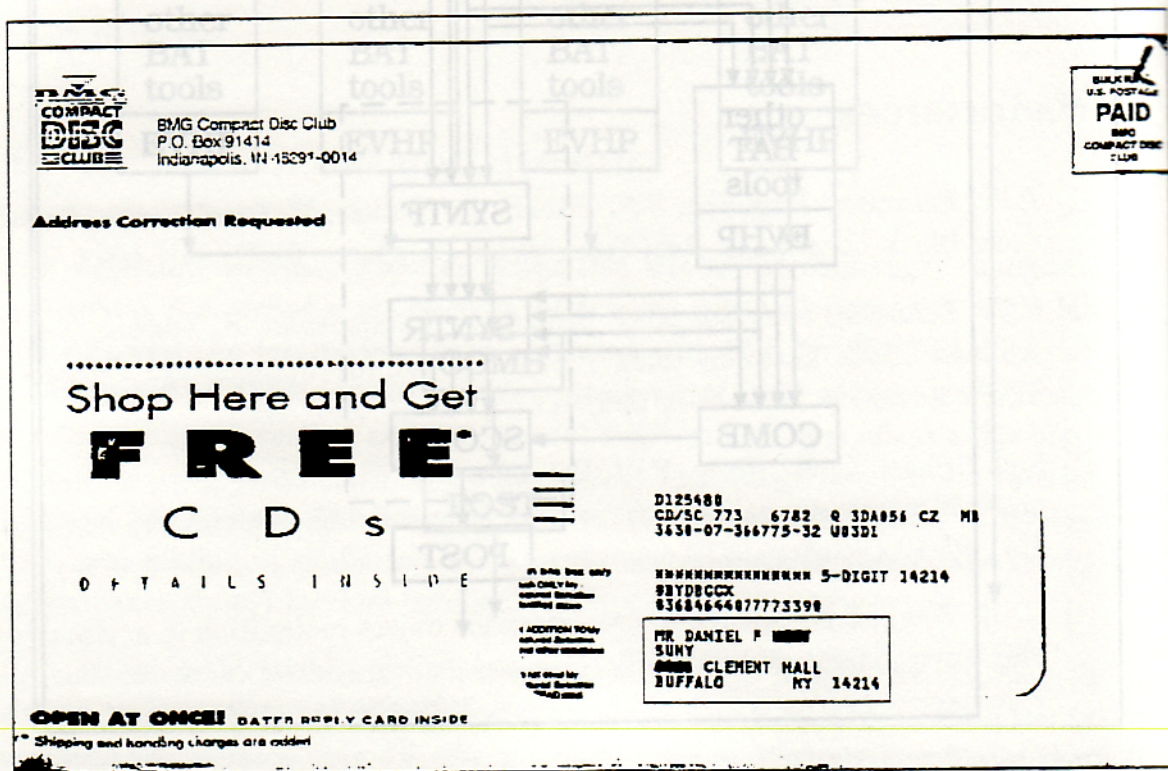


Figure 4: Top-ranked destination address block

Table 1: Syntactic tags

syntactic tag	meaning
pw	P.O.
pi	box
ss	street suffix
sw	state word
tn	title
#	number
-	unknown

Table 2: Syntax classes

Case I:	[city], state, ZIP Code
(1)	TEXAS (TX), 75240
(2)	EVANSTON, ILLINOIS (IL), 60202
(3)	FT SCOTT, KANSAS (KS), 66701
(4)	APO NEW YORK, WYOMING (WY), 09675
Case II:	[city], state
(1)	TEXAS (TX)
(2)	EVANSTON, ILLINOIS (IL)
(3)	FT SCOTT, KANSAS (KS)
(4)	APO NEW YORK, WYOMING (WY)
Case III:	ZIP Code
(1)	14260
(2)	02174-9901
Case IV:	the other cases

Table 3: Character recognition and string matching results for block #1

r M M K K M M M M M K M M M a M 5	D T O I T	I 4 2 1 4
(-)	(#)	(#)
8 B Y D I C C X		
(-)		
t 3 6 8 4 6 4 4 I 7 Z 7 7 3 3 9 t		
(#)		
M R	D A M I E 1	F
(tn)	(-)	(-)
5 u N Y		
(-)		
[REDACTED]	C L E M E N T	H A 1 1
(#)	(-)	(-)
B u F F A 1 0	N Y	I 4 2 1 4
(-)	(sw)	(#)