

Control Structure for Interpreting Handwritten Addresses

Edward Cohen, Jonathan J. Hull, and Sargur N. Srihari

Abstract—This correspondence describes the control structure for an intelligent handwritten address interpretation system. The system takes a grey-level address image, segments the address into lines and words, parses the address into meaningful syntactic categories, recognizes words using dynamically generated lexicons, and determines the destination code with the aid of postal directories.

Index Terms—Text processing, handwriting recognition, image processing

I. INTRODUCTION

The Handwritten Address Interpretation System (HWAIS) takes an off-line handwritten postal address image and determines a unique mail delivery point (e.g., a mailbox). In the United States, each

Manuscript received May 20, 1992; revised February 21, 1994. Recommended for acceptance by Associate Editor T. Dean.

E. Cohen is with Accu-Sort Systems Inc., 511 School House Road, Telford, PA 18969 USA; e-mail: cohen@cis.udel.edu.

J.J. Hull, and S.N. Srihari are with the Center for Document Analysis and Recognition (CEDAR), State University of New York at Buffalo, 226 Bell Hall, Buffalo, NY 14260 USA; e-mail: hull@cs.buffalo.edu and srihari@cs.buffalo.edu.

IEEE Log Number 9405251.

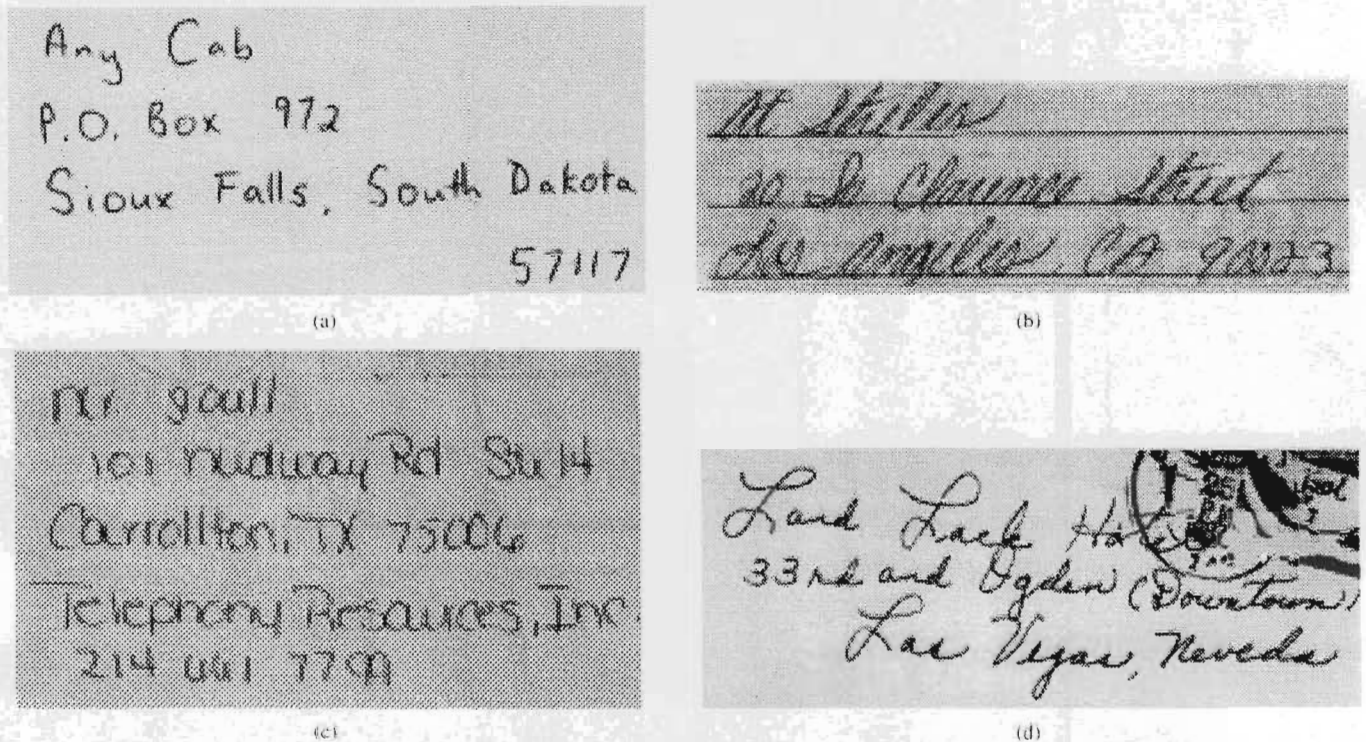


Fig. 1. Examples of handwritten address blocks extracted from live mail.

delivery point has been assigned a number, a delivery point code (DPC), which is a ZIP Code with appended digits. The HWAIS must normally locate and read the ZIP Code and street address or P.O. Box number to determine the DPC. City and state names may be used to resolve ambiguities. Examples of address taken from live mail (shown in Fig. 1) illustrate some of the challenges. This correspondence outlines the HWAIS control structure with examples. See [1], [2] for greater details and extensions to other domains.

The United States Postal Service (USPS) funded the research described here (CEDAR) and two similar efforts (ERIM and CGA). While the other groups had comparable performance of ZIP Code location and recognition (66% vs. CEDAR's 76%), our research was the first attempt to determine DPC from handwritten addresses (the other groups have not published figures for comparable systems).

This correspondence is organized as follows. Section II describes algorithms. Section III details how the system processes two handwritten address examples. Section IV describes performance evaluation and experimental results.

II. HWAIS ALGORITHM DESCRIPTION

The control flow in the current HWAIS has ten operational stages.

- A) Preprocessing;
- B) Text line separation based on positional information (separating an address into text lines);
- C) Feature identification;
- D) Text line separation using feature identification (combining B and C);
- E) Text line segmentation (dividing lines into words);
- F) Word classification;
- G) Grading match of word classification to spatial syntax;
- H) Recognizing words that have a high syntax-classification match;
- I) Using semantics to constrain word recognition;
- J) If a reliable semantic evaluation is found, derive an interpretation, otherwise, return to step H until no reliable syntax-classification matches are available.

Preprocessing

Input images are address blocks from live mail scanned at 300 pixels per inch with 256 shades of grey. Two preprocessing steps (thresholding [3] and guide line removal) are performed to produce a binary image (where text pixels are black).

Initial Text Line Separation

Dividing the address into text lines requires two stages. Positional information separates "easy-to-place" components; stage D uses feature identification to refine those results.

Line separation (Fig. 2) divides a binary image into vertical strips. Each strip is converted to a set of blocks, where blocks indicate the presence of text. Overlapping blocks from adjacent strips are connected to form a topological graph representing text lines (Fig. 3), where blocks and block connections are represented as nodes and node connections respectively. Heuristics adjust the graph to correct for typical situations (e.g., ascenders).

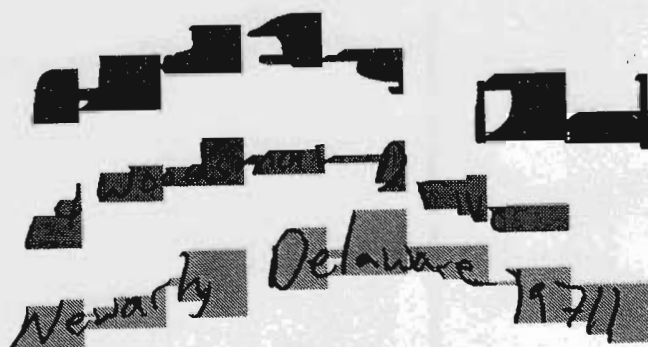
After finalizing the block arrangement, connected components are assigned text lines based on their corresponding blocks. For example, a component located in blocks from a single text line is assigned to that line. A component whose blocks overlap several text lines is assigned to a set of lines. Small components (considered noise) are not assigned to lines. All line assignments (including noise components) may be refined in stage D.

Feature Identification

Feature identification requires shape identification and spatial confirmation. To improve speed, most processing requiring image information (e.g., shape identification) is done first; spatial confirmation requires only symbolic information. Each connected component is initially identified as having (or not having) a particular shape (i.e., dashes, commas, disconnected 5-hats and digits). During this stage,

¹ Handwritten 5's are often written in two strokes. The vertical bar and the bottom curve are written in one stroke and the horizontal bar is written as a separate stroke. When the horizontal stroke is not connected to the other stroke, the horizontal stroke is considered to be a disconnected 5-hat. See Fig. 4(a).

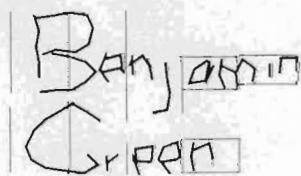
Robert James Dussal
 14 Woodsman Drive
 Newark Delaware 19711



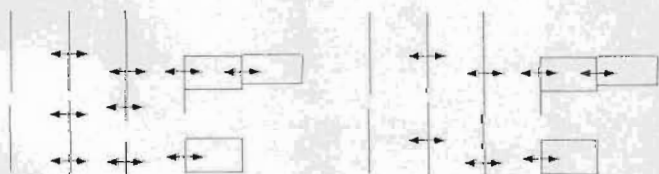
(a)

(b)

Fig. 2. Line separation shading algorithm. (a) Original binary image. (b) Blocks formed from shading. For each block, the right side shows the block height and the left side shows a text component used to determine the height. Each text line in the address has a slightly different color shading, indicating that the shading algorithm determined the line structure correctly.



(a)



(b)

(c)



(d)

Fig. 3. Line separation heuristics stages: (a) original image and the blocks created from it, (b) topological graph created from original block set (double arrows indicate connections between nodes), (c) topological graph after heuristic modifications, (d) final block set which allows each word to be assigned to a separate line.

components may have multiple identifications based on shapes. For instance, a short horizontal bar may be identified as having the shape of a dash and a disconnected 5-hat (as in Fig. 4).

Spatial information refines decisions. For example, a component with a disconnected 5-hat shape must have a component (C) to its immediate left in the same text line. C must either be identified as the digit 5 or as an unrecognized digit. C must not be a comma, dash or disconnected 5-hat (see Fig. 4).

Final Line Separation

Component position information (used in stage B) can be insufficient to separate lines. Components from one text line may be positioned closer to another (e.g., commas are often closer to the text line beneath their own). Heuristics are used to make new line assignments for such components.

Other components (typically characters from adjacent text lines that touch) must be forcibly segmented (by segmentation lines) and placed in more than one text line. This segmentation line is determined by



(a)

(b)

Fig. 4. Examples of feature identification: (a) horizontal bar has shape of disconnected 5-hat or dash, but location next to a 5 shaped digit identifies bar as a disconnected 5-hat. (b) horizontal bar has shape of disconnected 55-hat or dash, but location eliminates disconnected 5-hat possibility and suggests a dash.

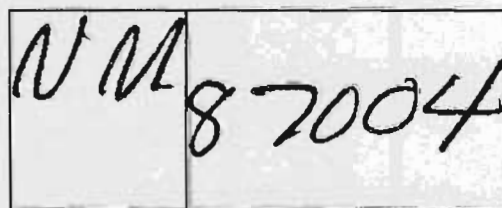
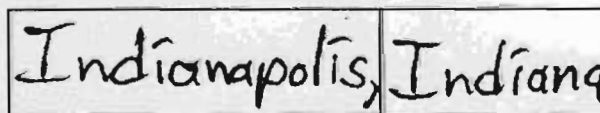
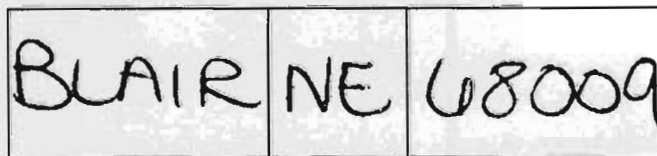


Fig. 5. Different types of word breaks.

minimizing the mean square error in a first-order polynomial whose data points are corners of bounding boxes of components already assigned to a single line. Each segment of a split component is assigned to its nearest text line.

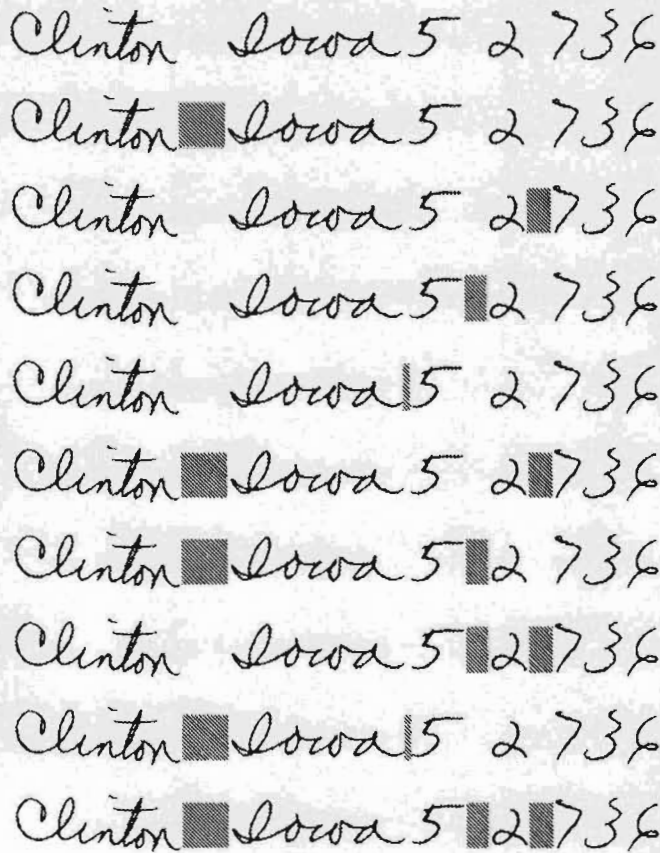


Fig. 6. Different ways of separating a text line into words. Two largest gaps give sixth grouping from top. Correct choice (first and fourth largest gaps) is second from bottom.

Text Line Segmentation

Components of each text line are clustered into words to create different hypotheses. Words are formed by ordering a text line's components from left to right and selecting word breaks between components. Word-break locations are hypothesized based on horizontal spaces between components, a shift in the mean vertical position of components, and the location of commas (see Fig. 5).

Word-break locations are ordered based on their distances (a comma-gap distance is computed as the sum of distances on either side of the comma; a shift in mean vertical position is always a word break). Different word-break groupings (Fig. 6) are chosen to separate a line into 1-5 words (e.g., for three words, use the top two widest word breaks for one set, use first and third widest word breaks for another set).

Word Classification

Each word hypothesized by text line segmentation is classified (into one or more categories with assigned confidence values) based on features and spatial locations of its components. Word categories are *characters* (including a number-of-characters estimate), *city*, *state-abbreviation*, *digits* (including a number-of-digits estimate), *digits-followed-by-dash*, *ZIP+4 Codes*, *barcodes*, *P.O.*, *Box*, and *noise*.

Grading Word Categories to Syntax

Our syntax is two-dimensional, since words are horizontally ordered along a text line and text lines are vertically ordered. One common handwritten address syntax has the bottom-most line containing three words, city, state, and ZIP Code (CSZ). For this syntax to be present, the bottom-most line must be divisible into three words and each must be matched to its expected form. The first word must

TABLE I
RESULTS OF READING ZIP CODE WRITTEN ON HANDWRITTEN ADDRESSES

Correctly read images found in the USPS directory	746	74.0%
Correctly read images not found in the USPS directory	18	1.8%
Total Correct	764	75.80%
Rejected images that contained ZIP Codes	221	21.9%
Rejected images that contained no ZIP Code	11	2.2%
Total Reject	232	23.0%
Errors in images that contained ZIP Codes	11	1.1%
Errors in images that did not contain ZIP Codes	1	0.1%
Total Error	12	1.2%
Total images	1008	100.0%

have characteristics of a "city" word, meaning it must be classified as a "city" or "text", but the "city" classification would have a higher confidence match. State name and ZIP Code are matched similarly.

Matches between image text and syntax are graded by ranking high-confidence syntax matches with a highly probable syntax over low-confidence syntax matches with less probable syntax. The HWAIS system matches against five different syntaxes: *P.O. + Box + number* (all on one line), *street-number + street-name* (all on one line), *ZIP Code* (on 1-2 lines, 5-digit ZIP Code and 4-digit add-on may be on separate lines), *city name + state name + ZIP Code* (on 1-4 lines), and *state name + ZIP Code* (on 1-3 lines). The HWAIS tries to find the most complete match possible (e.g., a *city name + state name + ZIP Code* match will be chosen over a *ZIP Code* match).

Recognizing Words

Numbers (ZIP Code, P.O. Box number or street number) are segmented into digits and recognized. For state-abbreviations, a segment-and-classify algorithm designed for characters is used; and for other word types, a dictionary is used to assist in recognition. See [1] for details.

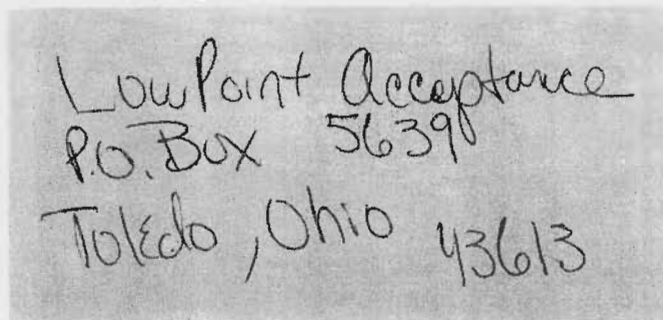
Using Semantic Information to Constrain Word Recognition

Once all relevant words have passed through the initial recognition phase, semantic information is used to improve recognition performance. ZIP Codes are checked against a USPS directory. Only 44% of possible 5-digit ZIP Codes are assigned; so roughly 56% of errors are detected.

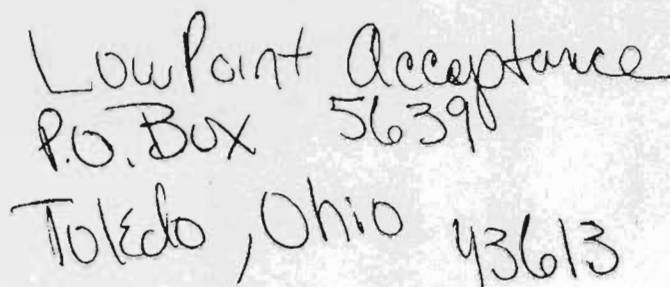
If the state name is recognized, it is compared to the first two ZIP Code digits and several actions can be taken. When state name and ZIP Code agree, the ZIP Code is accepted as correct even if digits are recognized with low confidence. If the state name is recognized with high confidence and the first and/or second digit of the ZIP Code do not have high-confidence recognition, the ZIP Code will be corrected to correspond to the state name if possible. If the state name and ZIP Code are recognized with high confidence and they conflict, the ZIP Code is accepted as correct since its recognition is more reliable.

The HWAIS uses information from street numbers and P.O. Box numbers. If a 9-digit ZIP Code is recognized, its 4-digit add-on can be confirmed by recognizing the P.O. Box number. If a 5-digit ZIP Code is recognized, its 4-digit add-on can be determined by recognizing the P.O. Box number. The HWAIS can locate and recognize P.O. Box numbers in 57% of the images which contain them.

If a ZIP Code and street number of an address are determined, a dictionary of possible street names is extracted from a USPS directory. When a street name is chosen, the HWAIS uses the recognized ZIP Code and street address (street number and name) to assign a 4-digit add-on. Adding the last two digits of the street number to the ZIP+4 Code creates the DPC. On addresses containing street addresses whose ZIP Codes were automatically located and recognized, the HWAIS locates and recognizes 41% of the street addresses.



(a)



(b)

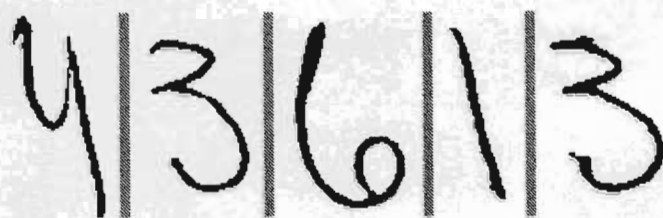
Toledo, Ohio 43613

city 1.00
7 chars 1.00

PO 1.00
Box 1.00
3 chars 1.00
4 digits 0.66

5 digits 1.00
Box 1.00
5 chars 0.80

(c)



(d)

Original	Modified
4 0.299	4 0.999
3 0.988	3 0.988
6 0.997	6 0.997
1 0.994	1 0.994
3 0.994	3 0.994

(e)

PO BOX 5639

2 digits 1.00	Box 1.00	4 digits 1.00
Box 1.00	state abbrev 1.00	Box 1.00
PO 0.80	2 chars 1.00	4 chars 0.80
state abbrev 0.80	2 digits 0.34	
2 chars 0.80		

(f)

5 0.996
6 0.996
3 0.999
9 0.996

(g)

43613-0639

(h)

Fig. 7. Example of HWAIS processing an image. (a) Grey-level image. (b) Thresholded image. (c) City, State, and ZIP Code Candidates. (d) Segmented ZIP Code. (e) Digits recognition results for ZIP Code (digit result and correct confidence) and adjusted digit recognition confidence results. (f) P.O., Box, and Box-number Candidates. (g) Digit recognition results of P.O Box number. (h) 9-digit ZIP+4 Code (which is a DPC).

III. EXAMPLES OF HWAIS ON AN ADDRESS IMAGES

Figs. 7 and 8 show some intermediate steps of images processed by the HWAIS. The first example shows how a poorly recognized ZIP Code has its confidence increased by the state name recognition results and how the P.O. Box number is located and recognized to determine a DPC. The second example demonstrates how the HWAIS can correctly locate a ZIP Code after initially failing to do so, and how a street address can determine a DPC.

The grey level for the first example (Fig. 7) is image shown in (a), and the thresholded image in (b). No guide line removal is required for this image; and, the address is separated into text lines and words.

Each text line is segmented into several different word groupings. Word classifications for one three-word grouping of the line are shown in (c). The first word is classified as *city* (because it has

text followed by a comma) and *7-chars* with high confidence (1.00). Similarly, the second word is classified as *PO* (a short word with a component recognized as a zero at the end), *Box* (a short word), and *3-chars* (the *i* is considered to be noise) with high confidence. The word is also considered to be *4-digits* with moderate confidence (0.66) because the first and last connected components are recognized as zeros.

This word segmentation results in a bottom-line syntax of *city-3-chars-5-digits* (c) corresponding to *city-state-ZIP*. Based on this match, the system selects the *5-digits* word as the ZIP Code candidate. The ZIP Code candidate is segmented (d) and isolated digits are recognized. Digit recognition results and their confidences are shown in (e). The recognized ZIP Code *43613* is checked against the USPS ZIP Code directory and found to be valid, but it would ordinarily be rejected because the first digit (which looks like a *4* or a *9*) is recognized as a *4* with low confidence (0.299). However, the state

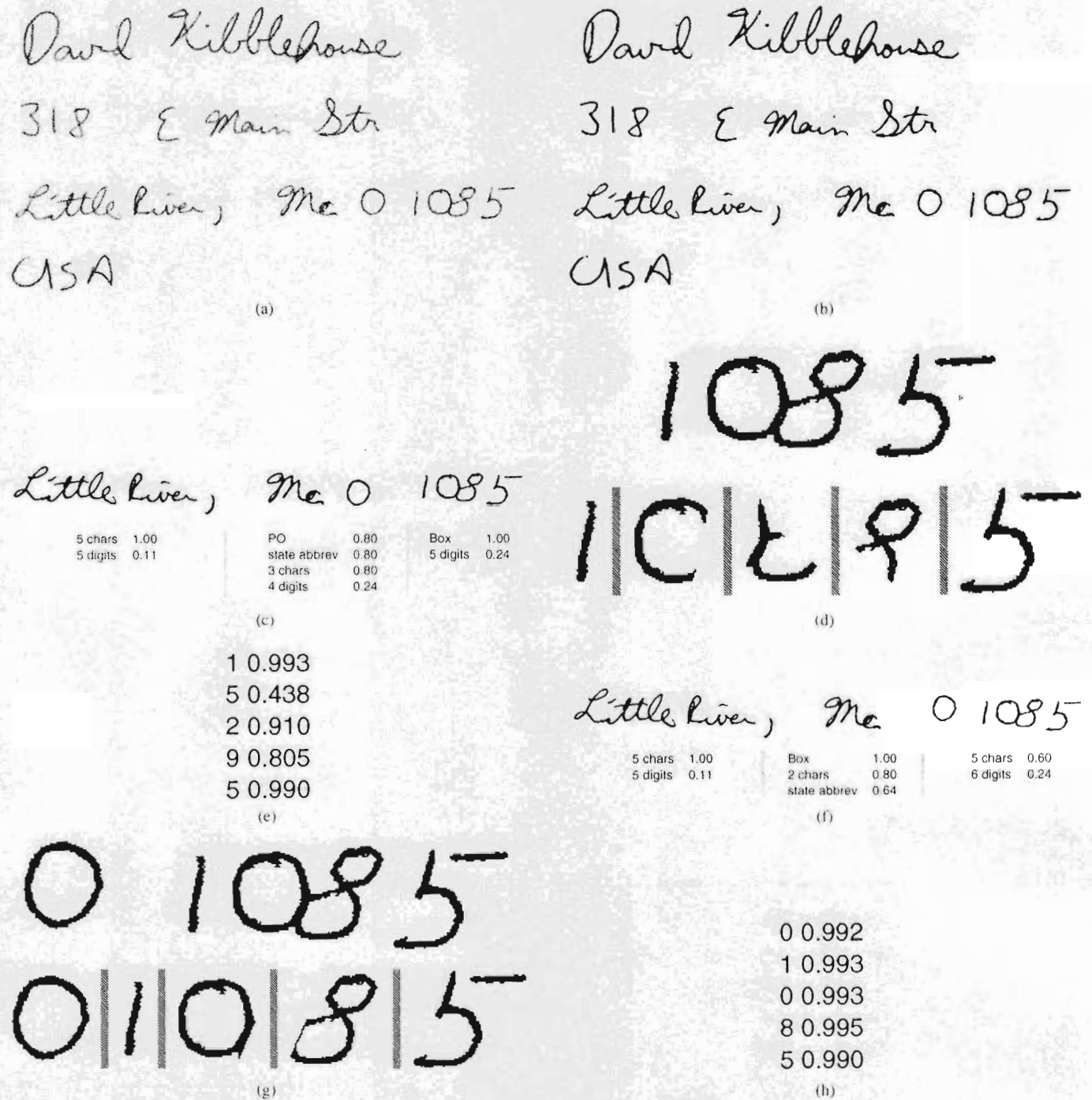


Fig. 8. Example of HWAIS processing an image. (a) Grey-level image. (b) Thresholded image. (c) First hypothesis for city, state, and ZIP Code candidates. (d) First segmented ZIP Code. (e) First digits recognition results for ZIP Code (digit result and correct confidence). (f) Second hypothesis for city, state, and ZIP Code candidates. (g) Second segmented ZIP Code. (h) Second digits recognition results for ZIP Code (digit result and correct confidence).

name is recognized as *Ohio* and since all Ohio ZIP Codes begin with the digit 4, the confidence of the first digit is boosted, and the ZIP Code is accepted.

The HWAIS must also interpret the P.O. Box line to determine a DPC. The middle line of the address is parsed to *P.O.-Box-4-digits* (f) which matches *P.O.-Box-Box-number*. This line could also match to *2-digits-2-chars-4-chars* for a street address line (*street-number-street-name*), but the P.O. Box line match has higher confidence.

The Box-number is recognized as 5639 (g) with high confidence, and the 4-digit add-on 0639 is determined using a lookup in a USPS directory. In this instance, the ZIP Code and the 4-digit add-on form the 9-digit DPC (h).

The second example (Fig. 8) shows the grey-level image (a) and the thresholded image (b). Each text line is separated and word hypotheses within each line are created.

After classification, the top parsing hypothesis for the *city-state-ZIP* for this address image is found in the second line from the bottom. One three-word hypothesis (c) for this line locates *5-chars-state-abbreviation-5-digits* as a syntax match. This match is incorrect, since the first digit of the ZIP Code is included in the state name. This mistake occurs because the ZIP Code contains a connected component (consisting of two large touching digits, 0 and 8) that is estimated to contain three digits (making the last word a 5-digits word). In addition, the middle word has two medium-sized connected



Fig. 8. (continued). Example of HWAIS processing an image. (i) Street number and name candidates. (j) Street number digit recognition results. (k) Original street name dictionary and 4-digit add-ons. (l) Ranked street name dictionary and 4-digit add-ons. (m) 9-digit ZIP+4 Code. (n) Final DPC.

components which provides a good match to the *state abbreviation* classification.

The ZIP Code for this parse is segmented (d) and the digit results are shown in (e). Digit recognition results are not all high confidence, and the state name is not recognized; so no digit corrections are made and the ZIP Code candidate is rejected.

The second parsing hypothesis chosen for the *city-state-ZIP* is the correct one. This parse of the same line into three words gives a syntax match of *5-chars-state-abbreviation-6-digits*. The *6-digits* match to the ZIP Code is allowed because digit estimation was required, and the HWAIS knows estimation errors may occur. The ZIP Code is segmented (g) and all digits are recognized with high confidence (h). The state name is still not recognized, so it does not affect the digit confidences.

The street address line (i) is correctly located by matching the second line, from the top, to *3-digits 5-chars*. The street number is correctly recognized as 318 with high confidence (j). Using a USPS directory, the street number, and the 5-digit ZIP Code, the HWAIS creates a dictionary (k). This dictionary contains all valid street names in the 01805 ZIP Code area whose street number is 318. Their associated 4-digit add-ons (k) are also shown. The abbreviations in the street names are expanded (e.g., E MAIN ST becomes EAST MAIN ST, EAST MAIN STR, EAST MAIN STREET, E MAIN ST, E MAIN STR, etc.) to create an expanded dictionary. This dictionary and the street name image are passed to the "word recognition with dictionary" algorithm producing a ranked dictionary (l). Confidence levels of the ranked dictionary indicate that *E. Main St.* is the correct street name. A ZIP+4 Code is created by adding the 4-digit add-on to the 5-digit ZIP Code (m). The last two digits of the street number (18) are then added to the ZIP+4 Code to create the final DPC (n).

IV. PERFORMANCE EVALUATION

The performance of the HWAIS is described in two phases, reading the ZIP Code present on the mail piece and assigning a DPC.

The system performance for reading ZIP Codes on a set of 1008 images of live mail addresses is 75.8% correct with a 1.2% error rate

(Table I). Since all DPC assignment by our HWAIS relies on ZIP Code determination, the 75.8% represents an upper bound on DPC assignment.

The system performance on DPC was tested on another set of 973 images (all addressed to locations inside Buffalo). These addresses were chosen because only Buffalo-area DPC directories were available during testing. ZIP Code recognition performance for these images is 67.4% correct with 1.11% errors. The drop in ZIP Code recognition performance was due to the nature of the ZIP Codes (i.e., ZIP Codes in this set contained many 4's and 2's, 14222 was most common; and our digit recognition algorithm performs worst on 4's and 2's).

DPC performance varied from 33.1% correct with 19.6% errors to 24.4% correct with 4.3% errors, depending upon the acceptance threshold level. This performance is the current state of the art for fully automatic systems that determine DPC's from handwritten United States' addresses. While insufficient for practical applications, continuing efforts are underway to improve performance.

ACKNOWLEDGMENT

We would like to thank M.-Y. Chen, A. Commike, P. Cullen, J. Favata, R. Fenrich, E. Kleinberg, C.-C. Kuan, A. Kundu, D.-S. Lee, C. O'Connor, G. Seni, J. Tan and J. Zhou for contributing ideas and code towards the HWAIS. This work was supported by the USPS under Task Order 104230-86M-3990.

REFERENCES

- [1] E. Cohen, "Interpreting handwritten text in a constrained domain," Tech. Rep. 92-06, State Univ. of New York at Buffalo, Feb. 1992.
- [2] E. Cohen, J. J. Hull, and S. N. Srihari, "Understanding handwritten text in a structured environment: Determining zip codes from addresses," in *Character and Handwriting Recognition: Expanding Frontiers*, P. S. P. Wang, Ed. Singapore: World Scientific, 1991, pp. 221-264.
- [3] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Trans. Syst., Man Cybern.*, vol. SMC-9, no. 1, pp. 63-66, Jan. 1979.