

Recognition of General Handwritten Words Using a Hypothesis Generation and Reduction Methodology

John T. Favata and Sargur N. Srihari
Cedar
State University of New York at Buffalo
Buffalo, NY 14260

Abstract

This research develops a recognition algorithm for off-line handwritten word recognition, where a word can be any mixture of discrete characters, cursive components or touching discrete characters. A hypothesis generation and reduction (HGR) paradigm is used as the basis of a solution to the general word recognition problem. Three levels of constraints, character, lexical and spatial, are utilized to hypothesize and rank word interpretations. The first and most general of the constraints are the whole character shapes embedded among the word strokes. An orderly search of the word for embedded characters produces many possible word interpretations. The number of interpretations is reduced by applying the lexical level of constraint, which consists of the character-to-character transition rules of the language. The remaining interpretations are ranked and reduced using inter-character spatial relationships that pertain to the relative character size. Surviving interpretations are compared to the lexicon using string matching so as to finally determine word identity.

1 Introduction

Machine recognition of handwritten words is the problem of transforming a two-dimensional digitized image of a word into a symbolic representation of that word. This research develops a *hypothesis generation and reduction* (HGR) paradigm using character, lexical, and spatial constraints as a very natural and general methodology for the recognition of handwritten words. The HGR model is a bottom-up algorithm which applies several levels of constraints in the recognition process. The most basic of the constraints are the possible character shapes embedded in the strokes of the word. This level by itself produces a large number of plausible interpretations of the word. These interpretations are ranked and reduced in number by the lexical level. This level reduces the number of interpretations by allowing only those interpretations which best satisfy allowable character-to-character transition rules to pass. The last stage reduces the number of interpretations from the lexical stage by ranking the interpretations with respect to the conventional character spatial (height) rules of the language. The best ranked interpretations are output by the algorithm and matched to the lexicon.

There are four very general classes of handwriting that must be handled by any general purpose reading system [Tap86]¹ (see Figure 1): discrete (type 1), purely cursive (type 2),

¹ Tappert defined five basic types of handwriting. The fifth, boxed discrete, is not applicable to this work.

touching discrete (type 3) and mixed (type 4). We decompose the recognition problem into three sub-problems, cursive recognition, discrete recognition, and touching discrete recognition, which are solved by individual HGR algorithms. A general HGR algorithm is developed, which is an extension of the individual algorithms, to provide recognition of any off-line author-independent handwritten word (types 1 to 4) (See Figure 2). The system is able to handle the wide variations in character and word formation that occur in human handwriting. The input to the system is a digitized sample of a handwritten word and a lexicon that contains the unknown word. The output is a ranking of the dictionary according to the analysis of the word. This system is tested using handwritten words from a very large set of authors (postal words) to accurately determine its overall recognition performance.

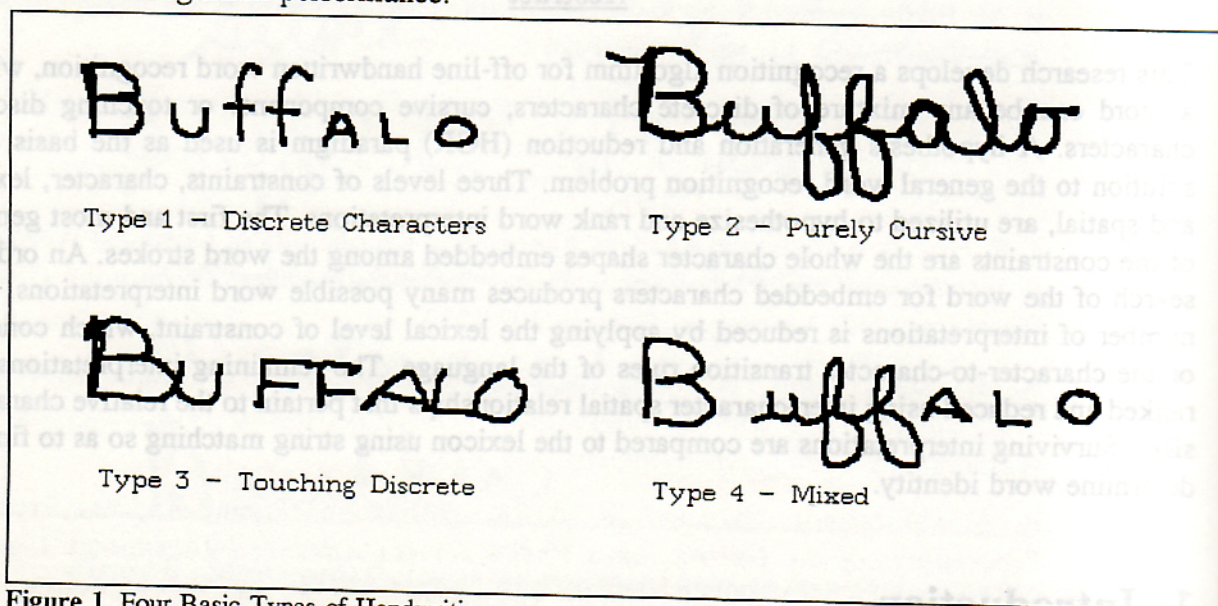


Figure 1. Four Basic Types of Handwriting

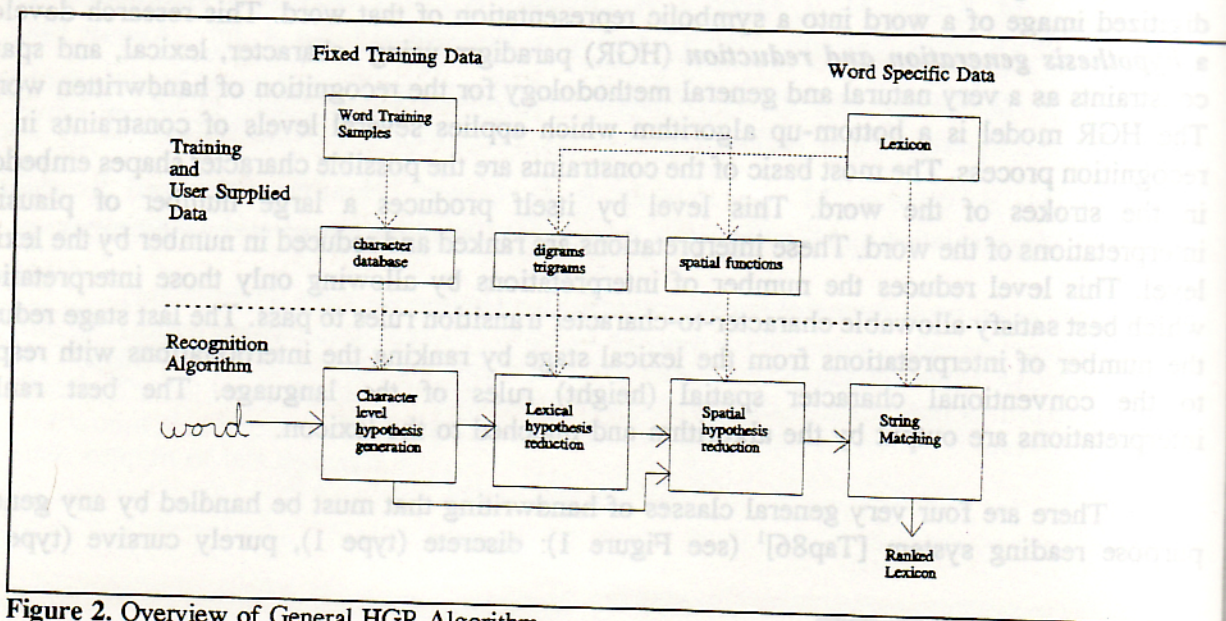


Figure 2. Overview of General HGR Algorithm

2 General Approach to Recognition

The fundamental goal of our recognition system is to segment, isolate and recognize the characters which make up the word. This task is relatively simple for the case of type 1 (discretely printed) words because the segmentation is natural. For types 2, 3 and 4, the problem is much more difficult because finding the correct segmentation is difficult and forced. This research takes the strategy of first producing a (possibly large) number of segmentations (and possible identities) of the word and then progressively refining the segmentations (identities) based on the constraints imposed by our word model assumptions.

We generate a number of plausible word hypotheses by segmenting the word (naturally or forced) and search for characters among pairs of segmentation points. This application of character knowledge to the word is constraint satisfaction based on local constraints (basic stroke shape). These hypotheses will be evaluated and ranked by applying intra-component constraints and global constraints using the lexical and spatial rules. The last step in the recognition process is matching the best hypotheses to a user-supplied lexicon (dictionary) and ranking the lexicon. In certain places in the algorithm, we will incorporate heuristics and certain a priori knowledge about the word to aid in recognition.

3 Basic Recognition Modules

This section will give an outline of the basic recognition modules of the system. These modules handle cursive recognition and touching discrete recognition. In section 4 we will combine these modules along with a discrete character recognizer to achieve general word recognition.

3.1 HGR Algorithm for Cursive Word Recognition

This section gives the details of the HGR algorithm for purely cursive words. The input to the system is a binary image of the word and a lexicon of k words which contains the unknown word. The system analyzes the word and then ranks the lexicon according to the results of this analysis. The system processing has several distinct stages. The first stage is the preprocessing stage which removes noise and normalizes the image for further processing (this stage will be discussed in section 4). It also generates a number of features which are used in later stages. The second step is to generate the segmentations of the word. This is done using a number of different rules and then eliminating redundant points. The third step is the generation of all plausible interpretations of the word using the results of the segmentations and the character recognition system. The final step matches all of the interpretations to the lexicon and computes a ranking of each lexicon word.

3.1.1 Feature Generation

After preprocessing, this part of the system extracts a number of features from the word. These

features are critical for segmentation and detection of certain characters. They are computed from a chaincode description of the upper and lower contours of the word. The features currently used are Across, Ascending, Descending and Tee Strokes. In addition, several other features computed: Upper masses, Lower masses and Holes (See Figure 3).

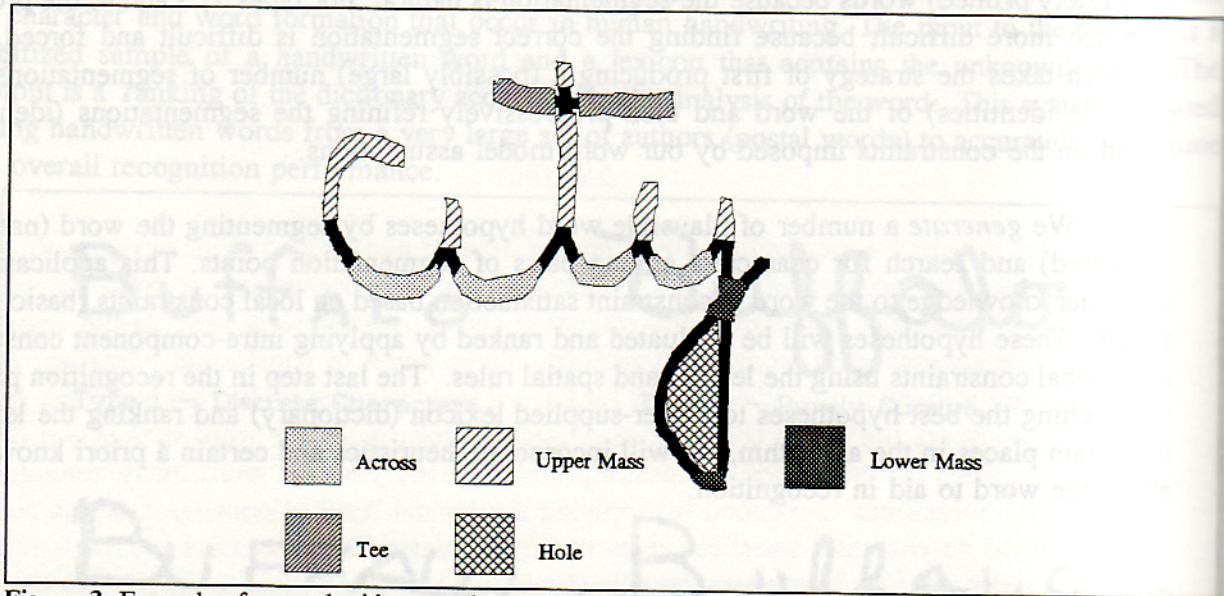


Figure 3. Example of a word with some detected features superimposed.

Average Stroke Thickness - The whole image is examined and a determination of the average stroke width is made. This average width is used to determine the presence of various stroke types (see below). The algorithm scans the image in the horizontal and vertical directions and finds all stroke run-lengths less than a specified thickness. The average of all of the stroke run-lengths is used as the average stroke.

Across Strokes - This feature is the ligatures of the word image. These strokes are defined to be more-or-less horizontal in trajectory and are computed using the upper and lower (smoothed) chaincodes. Three criteria are used to determine if a stroke is a ligature: the lower chaincode direction, the upperchain code direction and the stroke thickness. Ascending and Descending strokes are computed in a similar manner but with different criteria.

Tee Strokes - This feature consists of detecting the horizontal strokes that are used to cross the character t. They are computed using the upper chaincode slope and stroke thickness.

Upper Masses - This feature corresponds to upward pointing strokes that are found in the characters i, u, w, for example. They are computed using the upper chaincode slope and stroke thickness.

Lower Masses - This feature corresponds to downward pointing strokes that are found in the characters n and m for example. They are computed using the lower chaincode slope and stroke thickness.

Holes - This feature is formed by strokes which touch other strokes and enclose a whitespace area. The determination of this feature is done by performing a connected component analysis on whitespace (background). Each hole is a separate component. The boundaries and masses of each hole are determined and saved.

3.1.2 Character Segmentation

The segmentation algorithm used by the system is built using a number of separate modules which generate segmentation points based on features found in the word image. All of the modules are executed and the results of each algorithm are stored in a table. The points are next coalesced and redundant points are removed. The resulting points are the test segmentations for the word. The segmentation algorithm has separate systems for lower-case character segmentation and upper-case character segmentation. Lower-case segmentation points are placed on the bottom contour of the word and cut the stroke immediately above the point. The upper-case segmentation system attempts to trim the leading strokes of the word (assuming it begins with an upper-case character). This separate strategy was found necessary because of the fact that people tend to be very creative with the beginning upper-case character by exaggerating the leading stroke. This can cause recognition problems and it is necessary to carefully trim the leading stroke(s) to normalize the character and thereby improve recognition.

3.1.2.1 Lower-case Character Segmentation

This section describes the algorithms that are used for the generation of segmentation points for lower-case characters. Each algorithm generates segmentation points based on certain features present in the word and are independent of each other.

Ligature Segmentation

A series of segmentation points are generated on the detected ligatures (Across Strokes) in the image. This research departs from other research in that several points are generated per ligature rather than one. This strategy allows the system to be trained with normalized characters, that is, characters which are very carefully trimmed of excess ligature. This strategy enhances character recognition. Several carefully chosen segmentation points allows the system some flexibility (redundancy) for minor mismatching with the normalized training characters (See Figure 4).

Deep Valley Segmentation

This heuristic is used to segment a special case of character formation which occurs when two character are written but no computable ligature exists between the characters. This is caused by the writer making very fast up and down strokes between characters but generating no real horizontal ligature. The only feature that is generated is a deep valley formed between the strokes.

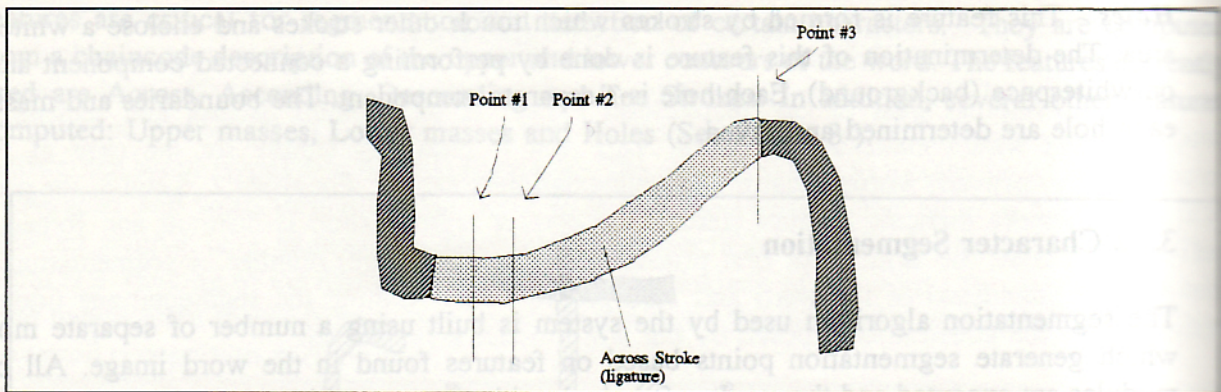


Figure 4. Ligature segmentation points

Double Hole Segmentation

This heuristic is used to segment a special case of touching characters which frequently occur in handwriting. This special case is the double o, which occurs when words with two sequential o characters are written (such as in *wood*). For some authors, the two o's are touching and are carefully segmented for recognition.

Left Hole Segmentation

This heuristic is used to segment another case of character formation in which two characters are written but there is no detectable ligature between the characters. This phenomena was observed in a sufficient number test words to warrant this heuristic. The heuristic places a segmentation point at the left of every valid hole in the word.

Touching Tee Stroke Segmentation

This heuristic is used to segment another case of character formation in which a horizontal t cross stroke touches an adjacent character. This is a very common occurrence. In particular, it is common for most authors to cross a word with two sequential t characters (such as *little*) with one stroke. It is necessary to detect the occurrence of this condition and split the double character.

Right Lower Mass Segmentation

This segmentation rule is used to segment words that contain descender to ascender character pairs (such as the word *right*). In general, there is a rapid transition from the descender character to the ascender character (a nearly vertical stroke). This rapid transition does not fit the form of a ligature. Therefore, it is necessary to detect the possibility of this condition and place a segmentation point between these characters.

Removal of Redundant points

The segmentation points from all of the segmentation algorithms are concatenated together and sorted according to spatial position. The next step is to remove redundant points in the set. The strategy considers ligature points as having the highest priority and removes all special segmentation points that are within a specified radius of the ligature point. The result of this strategy is that all ligature points plus those special points that are sufficiently far away from the ligature points are retained. This is the segmentation set that is used during the recognition cycle.

3.1.2.2 Upper-case Segmentation

A series of special segmentation points are necessary for the segmentation and recognition of upper-case characters. The rules are used to trim the exaggerated strokes that are sometimes present on these characters. The goal of this part of the system is to normalize the character by trimming these excess strokes. The strategy used for this trimming process is to generate a number of left segmentation points using uniform segmentation at the beginning of the character. Some of these segmentation points will trim the excess ligature and produce (hopefully) a recognizable character. The procedure starts from a right point (RP) on some part of the word and then generates the left uniform segmentation points (LP) at the beginning of the character for a predefined interval. During recognition, the character piece between the current LP and the RP is extracted and sent to character recognition (See Figure 5).

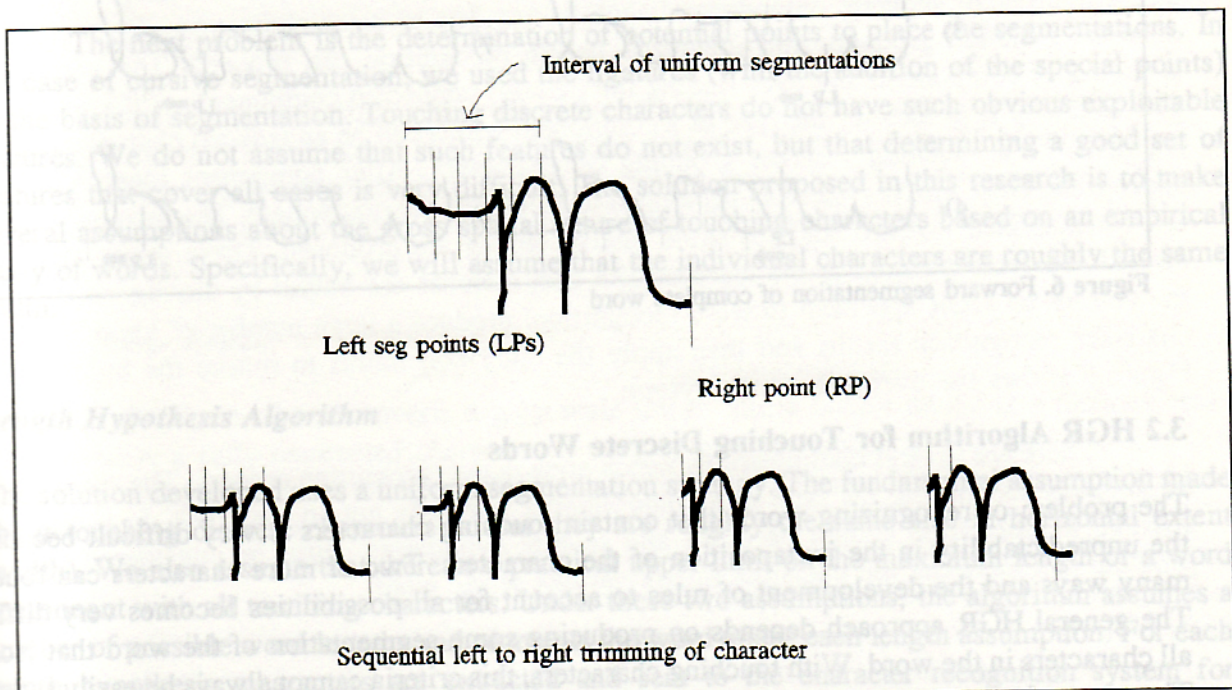


Figure 5. Uppercase trimming

3.1.3 Basic Cursive Character Scan of Word

After the valid segmentation points are determined, the next step is to determine valid characters embedded in the word. This process is accomplished by performing a basic recognition cycle (BRC) at each valid segmentation point in the image (see Figure 6). The basic recognition cycle of the system starts at a left point (LP) in the word. All segmentation points between the LP and the next N ($N = 4$) ligatures, including all special segmentation points, become the right points (RPs). The stroke between the LP and each RP is segmented, extracted and passed to the character recognition system for classification (See Figure 7). The character recognition system returns a confidence for each of the 26 characters. The system keeps track of the results for each segmentation and recognition cycle which includes the confidence and the corresponding segmentation point.

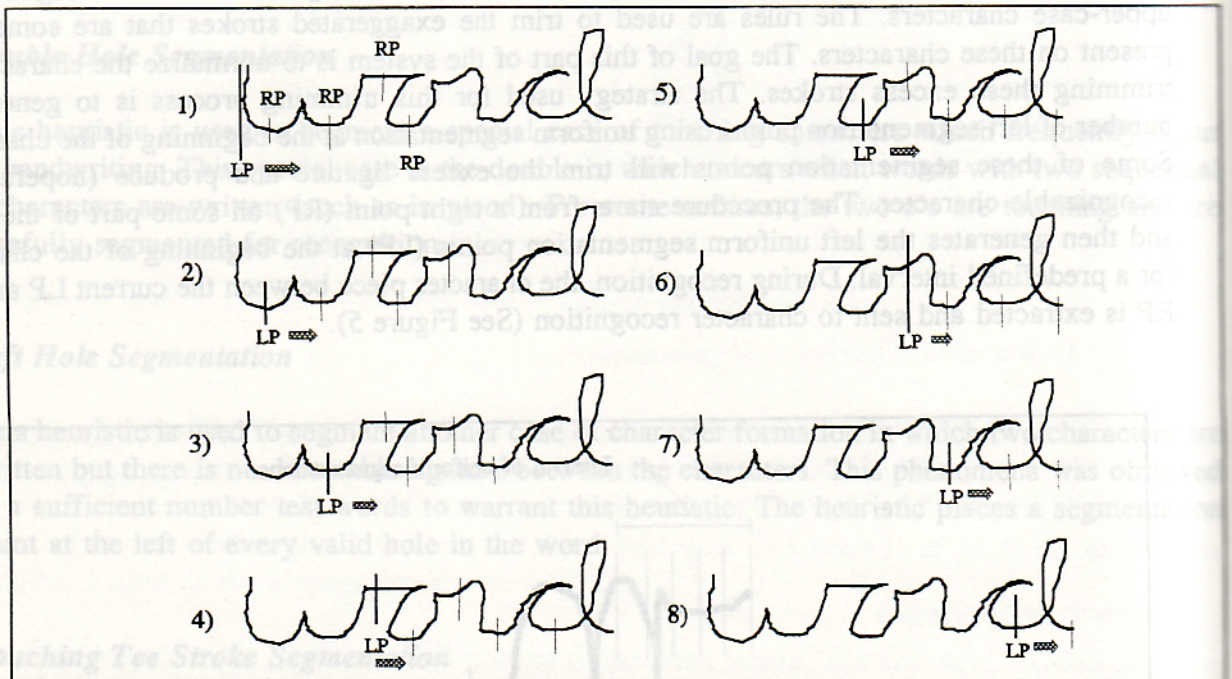


Figure 6. Forward segmentation of complete word

3.2 HGR Algorithm for Touching Discrete Words

The problem of recognizing words that contain touching characters is very difficult because of the unpredictability in the juxtaposition of the characters. Two or more characters can touch in many ways and the development of rules to account for all possibilities becomes very difficult. The general HGR approach depends on producing some segmentation of the word that isolates all characters in the word. With touching characters, this criteria cannot always be easily satisfied. The types of segmentations required to segment a word may require segmentations at an oblique angle, or even worst, a trajectory that is not a straight line. Fortunately, most handwriting will not require exotic segmentations to isolate characters.

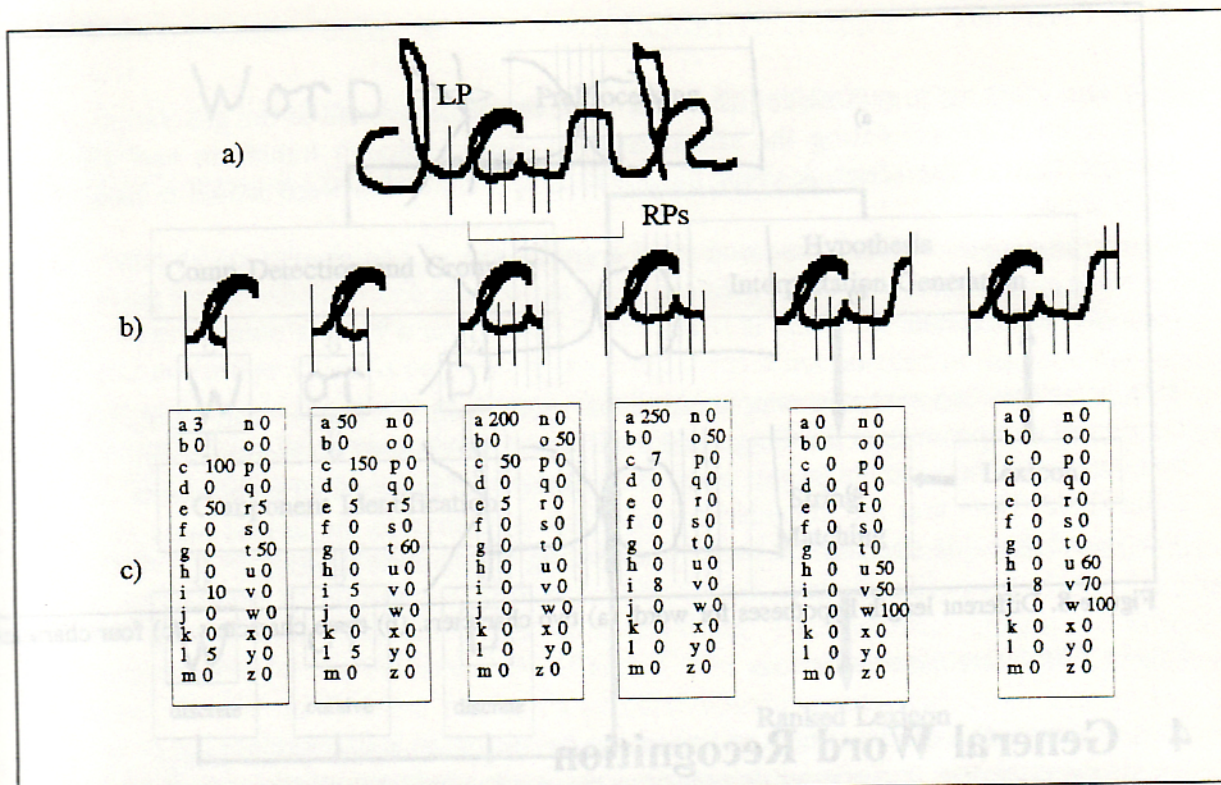


Figure 7. Example of basic cycle. (a) word with region to be recognized. (b) sequential strokes extracted from word. (c) character recognition results of (b).

The next problem is the determination of potential points to place the segmentations. In the case of cursive segmentation, we used the ligatures (with the addition of the special points) as the basis of segmentation. Touching discrete characters do not have such obvious exploitable features. We do not assume that such features do not exist, but that determining a good set of features that cover all cases is very difficult. The solution proposed in this research is to make several assumptions about the gross spatial nature of touching characters based on an empirical study of words. Specifically, we will assume that the individual characters are roughly the same width.

Length Hypothesis Algorithm

The solution developed uses a uniform segmentation strategy. The fundamental assumption made about touching discrete characters is that they are roughly the same size in horizontal extent (width). We also assume that there is a practical upper limit on the maximum length of a word component with all touching characters. Under these two assumptions, the algorithm assumes a number of possible word lengths and segments the word under each length assumption. For each length hypothesis, characters are extracted and sent to the character recognition system for analysis. If all characters of the segmented word under the current length hypothesis meet acceptance criteria, the segmentation hypothesis is accepted and the character recognition results are retained. All of the length hypotheses that produce valid characters are output at the end of the analysis. This process is performed for length hypotheses of 2, 3, 4, 5 (see Figure 8).

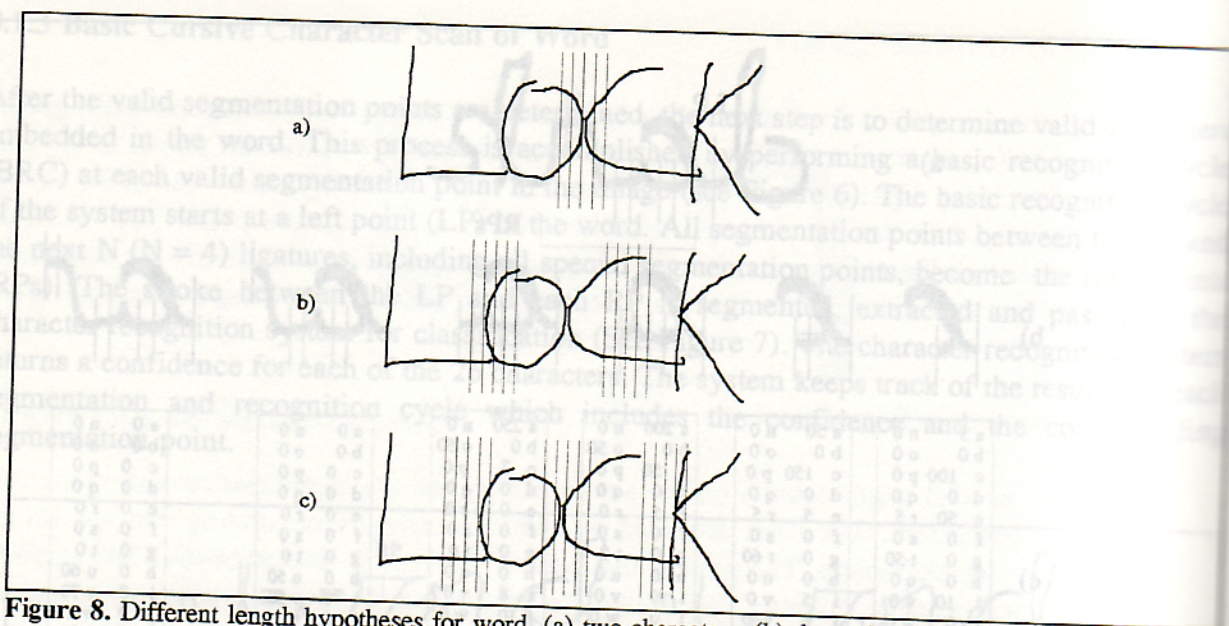


Figure 8. Different length hypotheses for word. (a) two characters. (b) three characters. (c) four characters.

4 General Word Recognition

The problem of recognizing unconstrained general handwritten words is difficult because the combinatorial number of style variations of the word. A handwritten word image can be split into a number of separate components² each having one or more characters. Each component must be isolated, grouped, and identified before recognition can proceed. The isolation and grouping of components involves finding the (nearly) connected components of the word image. The next step is to identify the type of a component which can be an isolated character, a cursive segment, or a touching discrete grouping of characters. After component identification, the appropriate recognition system is dispatched to recognize the component. The recognition results are concatenated together and plausible interpretations are produced. These interpretations are ranked and then matched to the lexicon. The best matches are output by the algorithm. This is the bottom-up HGR approach to word recognition: produce a large number of hypotheses using character recognition results and then apply the other constraints to reduce the hypothesis set. Figure 9 shows the overview of the system.

There are several complications which arise in the recognition process. The main problem of correctly identifying the type of a component is generally not an easy task. In some cases a positive identification cannot be made and several hypotheses must be considered. The next problem is generating the plausible interpretations of the word.

² A component is defined as a stroke or sequence of strokes that are physically connected - that is, there exists a path between any pair of (black) pixels in the component through only black pixels.

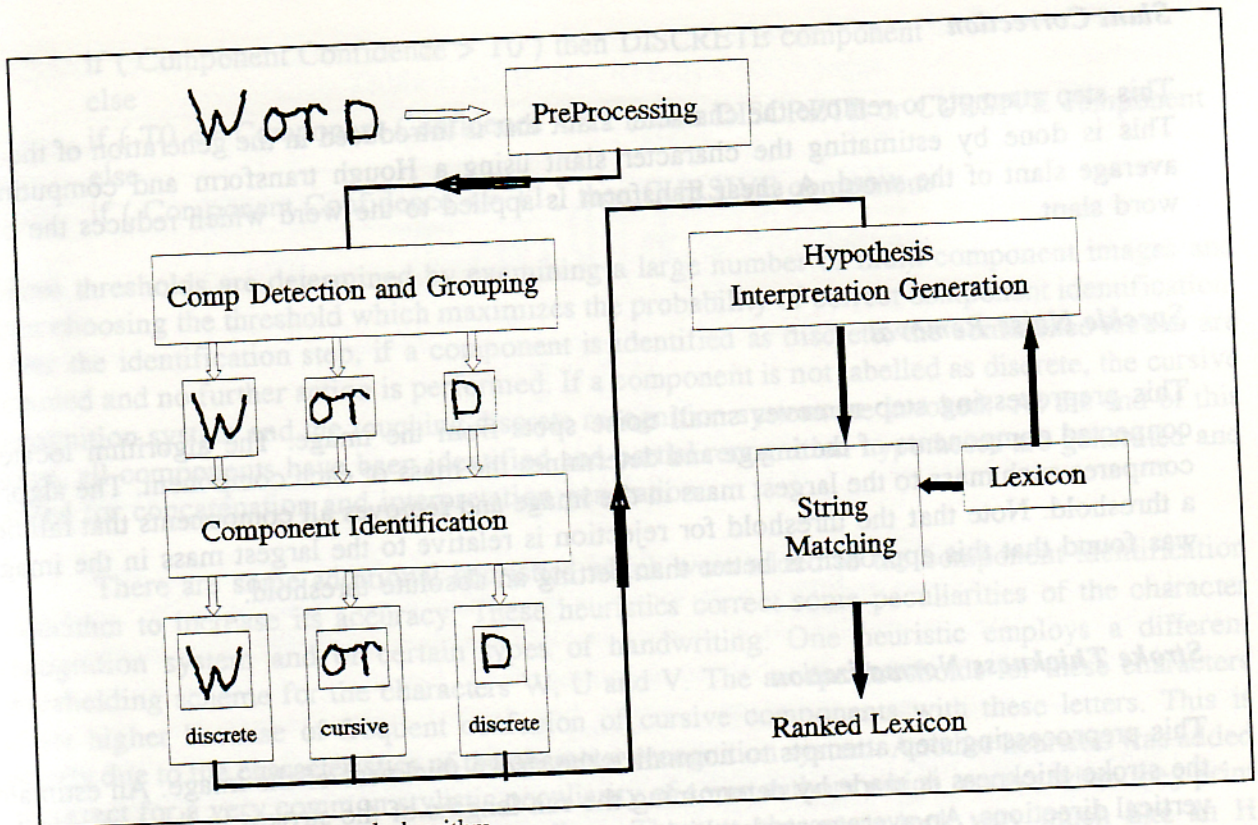


Figure 9. Overview of general algorithm

4.1 Preprocessing

Several general preprocessing steps are done on the composite raw image before word recognition starts. These steps attempt to normalize the image as much as possible. The preprocessing steps are performed once over the whole image before any recognition. When a recognition module is executed to handle a component, the preprocessing steps are not repeated.

Broken Image Repair

Finding a threshold which converts a grey scale image into a binary image is a difficult task. Occasionally, a fully connected character or word is unintentionally broken into two or more parts because of a failure in the thresholding algorithm. This partition of a character or word into disconnected components can cause recognition errors if each subcomponent is treated as a character. Empirical evidence has indicated that when a thresholding error occurs, the resultant breakpoint between the components is very small (1-3 pixels wide). This suggests an algorithm which identifies components and checks the proximity to other components. If there is a small gap between components, the gap is bridged and the components are reconnected.

Slant Correction

This step attempts to remove the character slant that is introduced in the generation of the word. This is done by estimating the character slant using a Hough transform and computing the average slant of the word. A shear transform is applied to the word which reduces the overall word slant.

Speckle Noise Removal

This preprocessing step removes small noise spots from the image. The algorithm locates all connected components of the image and determines the mass of each component. The algorithm compares each mass to the largest mass in the image and removes all components that fall below a threshold. Note that the threshold for rejection is relative to the largest mass in the image. It was found that this approach is better than setting an absolute threshold.

Stroke Thickness Normalization

This preprocessing step attempts to normalize the stroke thickness of the image. An estimate of the stroke thickness is made by determining the run-lengths of the strokes in the horizontal and vertical directions. An average stroke thickness is computed and if it is above a defined threshold, the image strokes are reduced. Reduction is done by finding and marking the outer contour of the image and peeling these contour pixels away. This is done as many times as necessary to reduce the stroke thickness to the defined threshold.

4.2 Component Isolation and Identification

After the preprocessing stage, the logically connected components in the raw image are identified and extracted. The algorithm uses connected component analysis and some spatial heuristics to group components [Coh92] into logical units. It is necessary to use additional heuristics to determine groupings because certain cases of disconnected but logically related components can appear in word images. These cases usually involve disconnected "hats" as in the characters J and T, or disconnected strokes such as in printed "f". In general, logical grouping can be a very difficult problem and misgrouped strokes can cause recognition errors.

Each isolated logically grouped component of the image is passed to a character recognition system that has been trained on discretely printed characters. This recognition system uses the multi-feature template matching approach discussed in [Fav92] and trained with discretely printed characters, cursive upper-case characters, and some cursive lower-case characters. The confidence results of this component recognition are used to determine the status of the component. Two simple thresholds T0 and T1 are employed to determine the type of a component using the following algorithm:

A component is defined as a similar sequence of strokes that are physically connected - that is, there is a path between any pair of (black) pixels in the component through only black pixels.


```

if ( Component Confidence > T0 ) then DISCRETE component
else
if ( T0 <= Component Confidence < T1 ) then DISCRETE or CURSIVE component
else
if ( Component Confidence <= T1 ) then CURSIVE component

```

These thresholds are determined by examining a large number of multi-component images and then choosing the threshold which maximizes the probability of correct component identification. After the identification step, if a component is identified as discrete, the confidence results are retained and no further action is performed. If a component is not labelled as discrete, the cursive recognition system and the touching discrete recognition system are invoked. At the end of this stage, all components have been identified and partial recognition hypotheses are generated and stored for concatenation and interpretation generation.

There are some additional heuristics which were added to the component identification algorithm to increase its accuracy. These heuristics correct some peculiarities of the character recognition system and of certain types of handwriting. One heuristic employs a different thresholding scheme for the characters W, U and V. The accept thresholds for these characters is set higher because of frequent confusion of cursive components with these letters. This is strictly due to the characteristics of the character recognition system. Another heuristic was added to correct for a very common stylistic peculiarity of some authors which occurs when they print words that have a double t in their spelling. These tt characters look very much like an H character and are recognized as such. It is necessary to add an additional test for the occurrence of an H in a component (other than the first component) and force the invocation of the cursive and touching discrete systems for the component.

4.3 Interpretation Generation

The interpretation generation takes the recognition results from each component and concatenates them. The process of producing interpretations done in two stages. The first stage generates a main path which is the concatenation of a particular set of discrete recognition results and cursive segmentations (paths). A discrete component generates only one path (which may have a substitution set of k characters). A cursive component may generate a number of paths which correspond to each plausible segmentation of the component. Associated with each path is a substitution set for each character (between pairs of segmentation points). The total number of paths generated for a word is the product of the number of paths for each component. After a specific main path is generated, the next step is to generate plausible interpretations of the word. This is done by following subpaths through the substitution sets of the main path. The number of interpretations produced is the product of each substitution set size in the currently examined main path. The algorithm generates all main paths and ranks all interpretations (subpaths). The ranking mechanism uses digram and trigram statistics gathered from the lexicon and employs certain spatial rules of (cursive) characters. These spatial rules reflect the spatial context of pairs of cursive characters, such as "e" is taller than "l" with respect to surrounding characters. A set of compatibility functions is used to rank cursive components with respect to character recognition results. The top N (N=50) interpretations are retained and passed to the matching

algorithm for final matching with the lexicon. Figure 10 shows an example of interpretation generation.

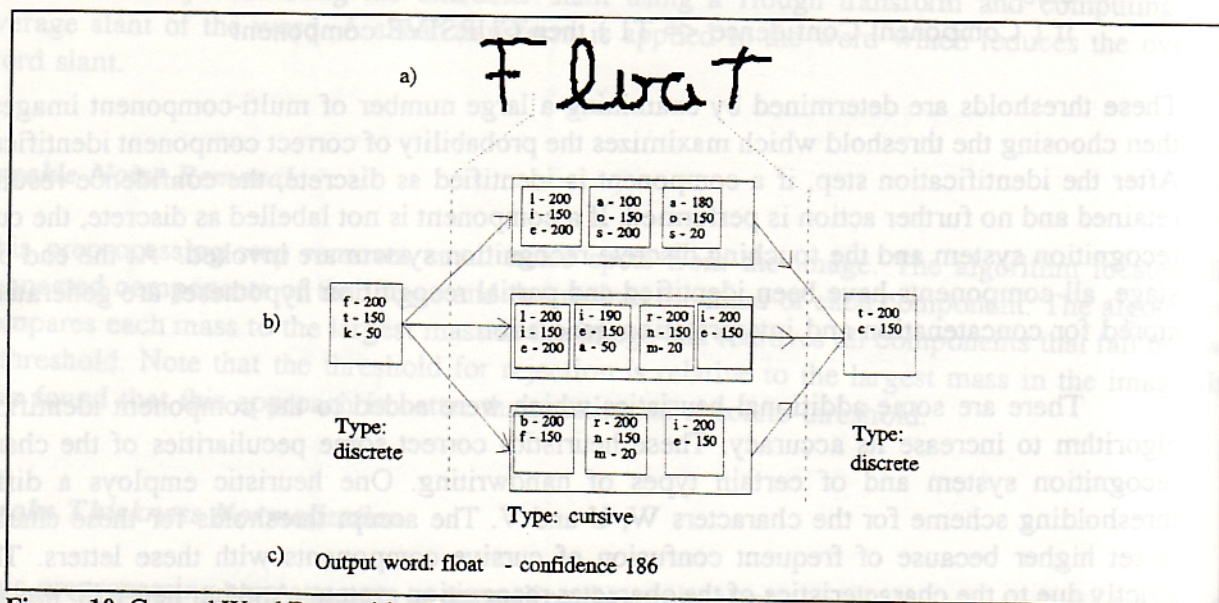


Figure 10. General Word Recognition Example. (a) word image. (b) component type determination and character recognition results. (c) one plausible word from a path through (b).

5 Lexicon Matching

The matching of the interpretations generated by the HGR algorithm is done using string matching techniques. A modified Levenshtein metric [Oku76] is used to match each of the interpretations to each of the lexicon words. The best matching interpretation-lexicon pair is output by the system.

6 Experiments

An implementation of this system was tested on a database of handwritten city words taken from live mail. Lexicons of 10, 100, and 1000 words were prepared and tested on 600 words from the database. The word type and frequency reflect the underlying statistics of the mail stream. Table 1 shows the results of this experiment.

Lexicon	Top 1	Top 2	Top 5	Top 10	Top 20
10	89.5%	95.7%	99.3	100.0%	
100	77.6%	82.8%	89.7%	92.4%	94.9%
1000	60.4%	67.3%	74.7%	79.6%	84.8%

7 Summary

This paper outlines the HGR approach to general word recognition. The method tries to identify the type of each word component and analyzes it under that type assumption. After all components are individually analyzed, the results of each analysis are concatenated and the valid paths in the word are generated. Each path is expanded and a number of interpretations are generated and ranked using lexical and spatial constraints. The matching of the best interpretations with the lexicon produces the final identity of the word.

References

- [Boz89] R. Bozinovic and S. Srihari, "Off-Line Cursive Script Word Recognition", *IEEE PAMI*, Vol. 11, No 1, Jan 1989
- [Coh92] E. Cohen, "Interpreting Handwritten Text in a Constrained Domain", Ph.D. Dissertation, SUNY Buffalo, Dept of Computer Science, TR-92-06, Feb 1992
- [Ehr75] R.W. Ehrich and K.J. "Koehler, Experiments in the Contextual Recognition of Cursive Script", *IEEE Trans. of Computers*, Vol C-24, No. 2, Feb 1975
- [Fav92] J. Favata and S. Srihari, "Offline Recognition of Handwritten Cursive Words", SPIE/IS&T Symposium On Electronic Imaging Technology, San Jose, Ca., Feb. 1992
- [Oku76] T. Okuda et al, "A Method for the Correction of Garbled Words based on the Levenshtein Metric", *IEEE Trans. on Computers*, Vol. c-25, No. 2, Feb. 1976
- [Tap86] C.C Tappert, "An Adaptive System for Handwriting Recognition", in *Graphonomics: Contemporary Research in Handwriting*, H.S.R. Kao (ed.), Elsevier Science Publishers, 1986