

Identification of forgeries in handwritten petitions for ballot propositions

Sargur Srihari, Veshnu Ramakrishnan, Manavendar Malgireddy, and Gregory R. Ball

Department of Computer Science and Engineering
Center of Excellence for Document Analysis and Recognition (CEDAR)
UB Commons
520 Lee Entrance, Suite 202
Amherst, NY 14228

ABSTRACT

Many governments have some form of “direct democracy” legislation procedure whereby individual citizens can propose various measures creating or altering laws. Generally, such a process is started with the gathering of a large number of signatures. There is interest in whether or not there are fraudulent signatures present in such a petition, and if so what percentage of the signatures are indeed fraudulent. However, due to the large number of signatures (tens of thousands), it is not feasible to have a document examiner verify the signatures directly. Instead, there is interest in creating a subset of signatures where there is a high probability of fraud that can be verified. We present a method by which a pairwise comparison of signatures can be performed and subsequent sorting can generate such subsets.

Keywords: Signature verification, ballot propositions

1. INTRODUCTION

Many local and state governments in the United States have form of a system for creating ballot propositions. These propositions, a form of “direct democracy,” can have the effect of directly creating or affecting existing laws based upon individual citizens’ votes. More than half of all American states and cities have initiative and referendum processes, and more than seventy percent of the population lives in such a state or city.¹ Furthermore, ten European nations, six former Soviet states, the European Union, and Taiwan all allow initiatives, referendums, or both. Initiatives and referendums are generally implemented as methods by which citizens can propose new laws or vote on existing laws, respectively, by collecting a number of signatures. Subsequently, a popular vote can take place on the issue in question. The topics of ballot measures vary widely. Examples include obtaining popular approval for bond measures, granting research money, creating or repealing taxes, determining issues on marriage, services, and minimum wage law.

2. BALLOT PETITION DATASETS

Petition required for the initiative or referendum process must be individually signed by each person on the list, and are generally accompanied by handwritten mailing addresses. Legally, the address block often can be written by a party other than the individual who signs, although a the address block may be used to help narrow down the list of forgery candidates. The number of individuals involved can be near one hundred thousand or more, and this poses a challenge for efficient analysis of handwriting. An individual signing in as a impostor is termed as a forgery (i.e., some agent signing another individual’s name falsely). Computationally speaking, a pairwise comparison of all signatures can be an intractable problem. Manual verification of such huge numbers of pairs is infeasable for human evaluators (100,000 signatures creates billions of pairwise comparisons). Example segments from actual ballot petitions are found in Figure 1. We approach this problem by attempting to create a pairwise list of candidate forgery pairs for a human examiner to verify. Some forgeries may be missed, and some

Further author information:

E-mail: srihari@cedar.buffalo.edu, Telephone: (716) 645-6164

non-forgeries may be included in the list. However, the intent of the list is to generate a more manageable list for a human examiner (perhaps from billions to thousands), allowing the examiner to locate an elevated percentage of any fraud present by analysing a fraction of the pairs.

15	Michele Perry	76 Valmar St		
16	Michele Perry			
11	Walter Shetter	37 Pleasant St.		
12	Oscar Loren	16 Lincoln St		
11	W. Winslow	41 Silvermine Rd		

Figure 1. Example images extracted from actual ballot measure petitions

An efficient solution to this problem has widescale implications in the real world. Verification systems using handwriting have been well researched and give good results when there is sufficient handwritten information present (half page/full page).^{2,3} However in this specific problem, the amount of information is restricted to single line handwritten address and signature. Further the presence of margins, lines, and logos in the documents interferes with analysis. We discuss below our proposed solution in four main steps. We first extract the handwritten address block. Second, we utilize transcript mapping to extract appropriate discriminatory features. Next, we use parallel processing to generate a list of overly-similar pairs. Last, we cluster the pairs to form clusters of handwritten elements that are all similar to one another.

3. METHODOLOGY

3.1 Preprocessing

A scanned document typically consists of several (approximately 20) rows of handwritten addresses and signatures of individuals. The first task in our method is to extract each handwritten address automatically. This involves the identification of the cell where the address is present. This task is a special case of line segmentation that uses the ruled line information present in the image. Horizontal and vertical projection profiles are used to find the candidate rule line positions. The average width of each line as well as the density information are useful features that are computed for this purpose. The result of preprocessing are a image snippets of handwritten addresses (several from each scanned image). Processed results are shown in Figure 2.

3.2 Feature extraction

We compute a set of discriminating features to help distinguish two individual's handwriting. We use a similar set of features as in earlier research² which can be categorized as macro and micro features. Macro features are global features of the document and micro features are character level features. Typical macro features include entropy (gray scale), connectedness, average slant, average word gap etc. Extracting micro features is more involved in that it requires the recognition of characters. Once the characters in each address snippet are recognized, the characters can be compared against other address snippets. For example, all the images of the character 'a' in one sample can be compared with those in another. To recognize such characters, we utilize a "transcript mapping" technique. A transcript mapper⁴ maps the given text of the address to the image snippet. For the particular problem at hand, the transcripts of each address is made available either by manual truthing or by using existing registration information from the voters list. A transcript mapper works by oversegmenting a given snippet into several segments. Combination of segments form character images. Stroke ligatures are used for segmentation. After segmentation, a dynamic programming algorithm is used to optimally combine the segments to best match the given transcript.

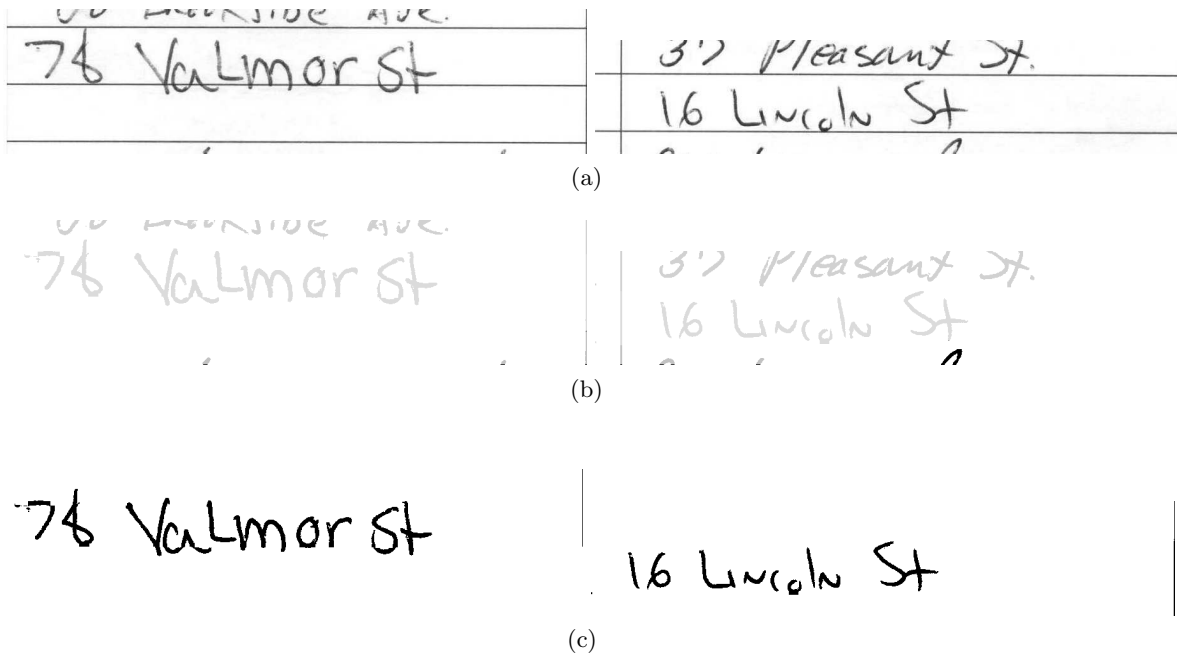


Figure 2. Automatically extracted address blocks (a) Identified address blocks, (b) Ruled line removal, (c) Final address image extraction

3.2.1 Transcript mapping

Before running the transcript mapping procedure, preprocessing including image binarization, line separation, and automatic word segmentation are performed. Again, errors may be produced in each of these steps. The result is a set of auto-segmented word images W , as shown in Figure 3, defined as follows:

$$W = \langle w_1, w_2, \dots, w_i, \dots, w_n \rangle \quad (1)$$

where n is the total number of word images, and w_i represents one word image segment. Normally we will have three types of situations for w_i . (i) It may contain just one word, such as w_1 and w_5 in the example. In this case word segmentation is correct. (ii) Or it may contain more than one word (i.e. undersegmentation error). For example, w_4 grouped “of” and “the” as one word segment, or (iii) it may only contain a part of a word (i.e. oversegmentation error), such as w_2 and w_3 in the example, which should be combined together and mapped to only one word, “cosponsor”, in the transcript.

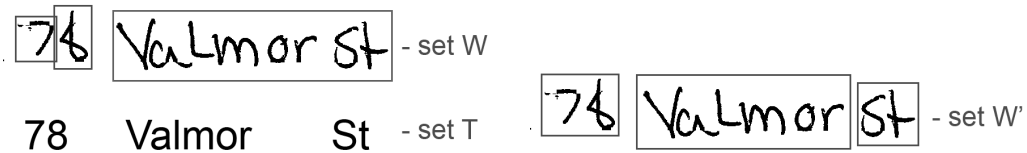


Figure 3. Examples of different sets in the transcript matching problem definition.

On the other hand, for a corresponding document image we have a truth transcript T , as shown in Figure 3, that is an ordered list of textual words, expressed as:

$$T = \langle t_1, t_2, \dots, t_j, \dots, t_m \rangle \quad (2)$$

where m is the total number of textual words in the transcript and t_j is one text word.

Since the proposed algorithm does not only align the transcript with the document image, but also tries to improve the word segmentation, we also define the improved set of word images W' , as shown in the in Figure 3, which is supposed to be more closer to the “ideal” segmentation result, as follows:

$$W' = \langle w'_1, w'_2, \dots, w'_j, \dots, w'_m \rangle \quad (3)$$

Similar to the situations for W , (i) each w'_j may be the same as some w_i , which indicates a correct auto-segmentation; or (ii) it may be part of w_i , which indicates a fix of undersegmentation; or (iii) it may be a combination of w_i and w_{i+1} , which indicates a fix of oversegmentation.

In addition, we set the size of set W' also to be m , which is the same as the set T . That is because the goal of our alignment algorithm is to assign one improved word segment to each of the textual words in T , and in this way, we can find an optimal mapping between a sequence of truth words, T (i.e. the transcript), and another sequence of better segmented word images W' .

A brief overview of the transcript mapping algorithm is shown in Figure 4.

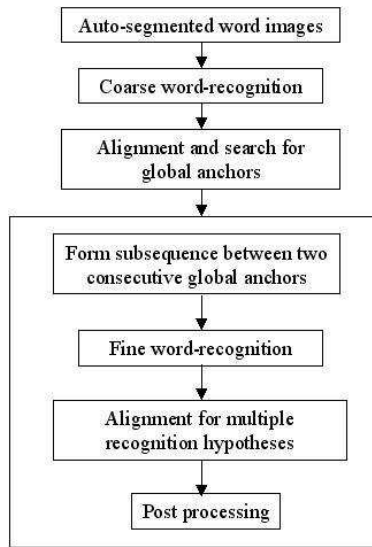


Figure 4. Diagram of the transcript mapping algorithm

In our current system, a connected component based clustering algorithm is used to perform the line separation followed by a neural network based word segmentation algorithm processing each segmented line image. The procedure is described more fully by Kim et al.⁵ A small modification made at the word segmentation step is to add a constraint on the total number of auto-segmented word images; from the transcript we already know that the total number of textual words is m and use this by requiring that the total number of segmented word images n should be bounded by

$$m(1 - 0.15) \leq n \leq m(1 + 0.15) \quad (4)$$

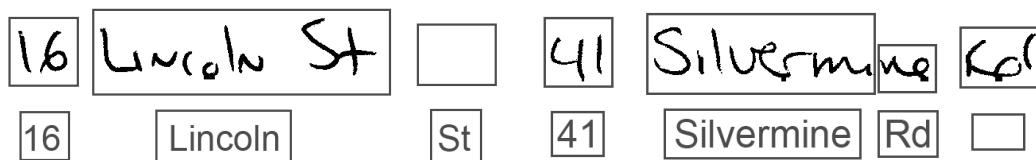
Coarse word recognition for all the word images is performed followed by coarse alignment on the entire document. The goal of this step is to seek for a set of global anchors. That is, for each auto-segmented word image w_j , we generate a lexicon based on its position information (i.e. the index number “j”), and then do word recognition for this word image. Here the size of the lexicon is chosen to be less than or equal to 20 (depending on the total number of words), which is a tradeoff between lexicon coverage and the recognition accuracy. The word recognizer will return an ordered list of lexicon words ranked in descending order of the dissimilarity value (i.e. distance measure). In this step, only the top 1 choice of the returned word list will be considered as the truth of the word image.

After the coarse word recognition, each word image has been assigned a text word. Then a coarse alignment will be performed on the entire document in order to find a set of global anchors. Here the alignment problem is solved by a dynamic programming algorithm. A word image will be chosen as a global anchor if the following conditions are both satisfied: (i) this word is in the longest common subsequence that is generated from the dynamic programming algorithm, (ii) its associated distance value from the recognition result should be small (in our current system, it should be less than a trained threshold value).

After we get a set of global anchors, the next step is to segment the entire document into several subsequences. For every two consecutive global anchors, w_i and w_j , the set of all the word images in between is just a subsequence of the entire document, and we call it W_{ij} . Then each of such subsequence can be treated as a shorter document.

For each such subsequence, we perform a fine word recognition for each word image with a smaller lexicon (size of the lexicon is about 10). But this time, we may not only keep the top 1 word from the returned list, but may also keep the top 2 or top 3 as the possible candidates. The criteria for whether keeping more candidates are as follows. We choose top 2 only if the distance value of top-1 is large and the difference between the distance value of top-1 and top-2 is small (both criteria use some threshold values that were estimated using the training data). Therefore multiple hypotheses for the recognition string could be generated by choosing different choice, if they are available. After that, the same dynamic programming algorithm is used to find the optimal alignment between the truth string and each of the recognition string. Finally the hypothesis with the highest alignment score is selected.

Post-processing is performed after each subsequence gets a best alignment from all the hypotheses. It looks through the actual alignments and try to fix any mismatching alignment caused by segmentation errors. Normally when undersegmentation happens, there will be an empty word in the result sequence as illustrated by Figure 5(a). In this case, usually it should be another mismatch in its neighbor (either before it or after it), so we can try to combine the missed word lexicon with the mismatched one to form a new word and add it into the current lexicon, and do the word recognition again using the new lexicon, if the new word is recognized as the top 1 choice, then this word image will be segmented into characters and then grouped as two words based on the truth words, and then a cut will be made at the boundary of these two word images; two text words are assigned to them respectively.



(a) Example of undersegmentation error

(b) Example of oversegmentation error

In the case of oversegmentation, usually there will be a space in the truth sequence, as illustrated by Figure 5(b). In this case, there are two possibilities. One is that the corresponding image piece is not a word image and is even not a part of word image. It could be a punctuation, such as “?” and “!”. The other is that it is a part of word. We can distinguish these two cases by doing the following. We attach it to its neighbors (to the end of the word before it, or to the front of the word after it), then we do word recognition for the new word image using the same lexicon again. If the returned text word of the top-1 choice is matched or its distance value is lower than the previous one, then we keep the new combined image (we may just fix a oversegment). If the returned text is not matched or its distance value is greater than the previous one, then we leave that piece of image out.

String alignment is also a common problem in many research areas, such as in bioinformatics.⁶ Here, given two textual sequences P and Q : one is from the transcription (the truth) and the other is formed by the results from word recognizer, we have a similar alignment problem. The difference is that the elements of the string are no longer characters but words and the “string” here is really a sentence (i.e. a sequence of words). We design our alignment algorithm as follows.

Similar to the algorithm for string matching, we need to introduce a special word “-”, which represent the insertion of an empty word (or a gap). Given two strings P and Q , with $|P| = n$ and $|Q| = m$, in order to compute an optimal alignment of P and Q , we first define a *score function* as follows.

If p_i and q_j are each a single word or an empty word, then $\sigma(p_i, q_j)$ denotes the score of aligning p_i and q_j . In our problem, we define $\sigma(p_i, q_j)$ as

$$\sigma(p_i, q_j) = \begin{cases} 1 + c_j & \text{if } p_i = q_j \\ -1 & \text{otherwise} \end{cases} \quad (5)$$

where c_j is a confidence value for q_j and $0 \leq c_j \leq 1$. This confidence value is computed from the distance measure associated with each choice returned by the word recognizer.

Then we define $V(i, j)$ as the value of an optimal alignment of the string $\langle p_1, \dots, p_i \rangle$ and $\langle q_1, \dots, q_j \rangle$. The value of an optimal alignment of P and Q is then $V(n, m)$. And it's solved by using the following recurrence formula:

$$V(i, j) = \max \left\{ \begin{array}{ll} V(i-1, j-1) & + \sigma(p_i, q_j), \\ V(i-1, j) & + \sigma(p_i, -), \\ V(i, j-1) & + \sigma(-, q_j) \end{array} \right\} \quad (6)$$

After we get the score of an optimal alignment, the actual alignments can be recovered by retracing the dynamic programming steps back from the $V(n, m)$ entry.

Word-model Word Recognition (WMR)⁷ takes as inputs a word image and a lexicon, and computes a dissimilarity score (distance) between each lexicon word and the word image. The lexicon words are ranked and returned with the top 1 choice as the best match between the lexicon and the word image.

To match the word image against a lexicon, WMR involves three major phases, i.e., segmentation, feature extraction and matching. The segmentation phase separates a word image into smaller pieces called segments. Each segment represents a character or a sub-character (i.e. a part of a character). During the phase of feature extraction, 74 chain code based features are extracted from all possible combinations of 1-4 consecutive segments (called as super-segments). A super-segment corresponds to a single character in a word of the lexicon.

Given a lexicon word, the matching phase uses a dynamic programming algorithm to match features of the super-segments with the “ideal” features (obtained in the training procedure of WMR) of characters from the lexicon word, and computes a distance as the matching score. The matching phase is repeated for all lexicon words. As a result, the output of WMR consists of a list of lexicon entries ranked in descending order of their matching score values.

While the matching phase determines the segmentation points between segments that correspond to characters in a lexicon word, WMR can also be used to segment character images from a word image if a single true lexicon word is presented. Further more, it will be used in the post processing to fix some of the undersegmentation errors, i.e., segment an undersegmented word image into two or more word images.

3.3 Verification model

The third phase is the verification phase. In this phase we compare (features) pairs of address snippets to obtain a similarity measure between them. A higher similarity measure indicates that the two snippets are potential candidates for a forgery. The final goal is to identify all the snippets that were written by the same person. In essence, this is a clustering problem where each cluster represents address snippets written by the same individual. If N snippets are present, there exists $\binom{N}{2}$ possible comparisons. This results in a $N \times N$ symmetric similarity matrix S . The clustering problem now is that of identifying cliques in a graph with vertices as the snippets and edge weights the values in the S matrix. For very large value of N , the $\binom{N}{2}$ verifications as well as the clique problem are computationally expensive. Hence, it is infeasible to first wait for the verification task to complete. In our approach, we combine the verification problem with the clique problem to obtain intermediate solutions as the algorithm runs. The verification model is described below.

The model used for verification is Naive Bayes. Features between two snippets are compared and thereby obtaining a distance. A large training set of ensemble of pairs of documents were used to learn the parameters. The ensemble of pairs consists of two sets: (i) pairs of documents written by the same individual and (ii) pairs written by different individuals. The distance obtained by comparison of features between each pair in the training set are modeled parametrically using Gaussian or Gamma densities. This represents the class conditionals in the Naive Bayes model. The prior probability of a pair to belong to the same or different writer is considered equal. For a given new pair of document (in our case address snippets), the class (same or different) posterior probabilities are evaluated using Bayes rule. The features are considered independent and this enables ease of use of multiple features in the Naive Bayes model. The result of the model is a log likelihood ratio (LLR) that represents the similarity between two pairs of snippets. A positive LLR indicates the pairs were written by the same writer and a negative LLR indicates the pairs were written by different writers.

Due to the large number of comparison to be made, parallel processing is done. A server is responsible for choosing nodes which should be compared next and assigns those to free clients. Any computer at any time can bind to the server as a client and process a job assigned by the server.

The snippets are nodes in a graph. The pairs of snippets are compared and when a matching pair is found (a significantly high positive LLR), we start growing the graph by breadth first approach with these two nodes, until we exhaust links from these nodes to other nodes. The threshold for matching was obtained using a validation set. These groups of nodes form a primary cluster. This cluster is not a maximal clique in the sense that it is not necessary that every pair in a clique satisfy the threshold LLR similarity. This strict definition of the maximal clique is relaxed. Instead it suffices that only a majority of the nodes are connected with high LLR value to a given node. This clustering goal is achieved using infinite mixture of Gaussians. This is explained further below.

3.4 Clustering

In the fourth phase we cluster these nodes using an infinite mixture of Gaussians. This allows us to model data multimodally. Our data is potentially multimodal because a forgery involves the same individual writing several snippets. For performing this clustering, every node is represented using a feature vector. A binary vector is used in our approach. We indicate by 0's absence of connection between two nodes and by ones the connections. The final cluster that we get out of this process is then weighted proportional to the similarity measure that each node makes within that cluster. The highly weighted clusters are considered as possible forgeries.

This being a real world problem, the evaluation of the method's performance on real data needs to be done by visual inspection of the results. However, for empirical verification, a mock election petition dataset was obtained.

4. DATASET AND RESULTS

Table 1. Zero Classification Border Results

True Positive	89 / 186	47.85%
False Positive	355 / 4185	8.45%
True Negative	3830 / 4185	91.52%
False Negative	97 / 186	52.15%
Precision	89 / 444	20.05%
Recall		47.85%

Table 2. Maximum Likelihood Classification Border Results (threshold 7.02)

True Positive	28 / 186	15.05%
False Positive	34 / 4185	0.81%
True Negative	4151 / 4185	99.19%
False Negative	158 / 186	84.95%
Precision	28 / 62	45.16%
Recall		15.05%

Table 3. Precision/recall Percentage By Number Of Possible Character Comparisons

	0 - 5	5 - 10	10 - 15	> 15
Zero border precision	11.86	26.62	30.43	43.48
Zero border recall	38.89	57.81	56	40
ML border precision	15.79	57.14	58.33	66.6
ML border recall	4.17	25	28	8

The dataset consisted of 95 handwritten address written by 19 authors, with each author contributing 5 samples. Every address in the dataset was different, and authors were told to write in their own natural handwriting. The results with two classifiers—the zero classification border and the maximum likelihood classification border—are summarized in Tables 1 and 2 respectively. The results grouped by the number of possible character comparisons between the two samples are shown in Table 3.

The percentage of correct answers/precision increased as the number of letter matches in common increases. When there was an overlap of 15 characters (e.g. some three characters in one word in common with some five in the other), the precision was about 43% and 67% for the zero border and ML border classifiers respectively.

Recall, on the other hand, was not as straightforwardly linked to the letter match count. The pattern for both classifiers had 0-5 character and > 15 character matches getting the worst results, and the 5-10 and 10-15 character matches getting the best. The zero border did a better job with recall, implying that there were more false positives but it caught more of the true positives.

5. CONCLUSIONS AND FUTURE DIRECTIONS

Precision was about 20%/45% overall for the zero border and ML border classifiers respectively. Although this number is somewhat far from the ideal 100%, this is still an interesting result because the generated candidate match pairs would allow a human examiner to perhaps catch about half of the fraudulent pairs while looking at only 4% of the matches, saving significant time. In order to combat being potentially overwhelmed with non-fraudulent pair “noise,” some post processing into cliques (even small or incomplete ones) can be performed. It is possible that some forgery pairs are missed by the process, but unless there is a very high percentage of forgeries, manual examination of pairwise signatures is unlikely to find a significant number of forgeries due to the very high number of manual pairwise comparisons required. Thus, the process described allows a manual examiner to elevate the number of forgeries detected while still having a manageable list. Furthermore, since the process works breadth-first, partial results will be available for examination while the candidate lists are still being generated.

REFERENCES

- [1] Matsusaka, J. G., “Direct democracy works,” *Journal of Economic Perspectives* **19**, 185–206 (Spring 2005).
- [2] Huang, C. and Srihari, S. N., “Automatic writer verification with handwriting of twins and non-twins,” in [*Proc. International Graphonomics Society Conference*], 184–187 (November 2007).
- [3] Srihari, S. N., Huang, C., and Srinivasan, H., “On the discriminability of the handwriting of twins,” *Journal of Forensic Sciences* **53**, 430–446 (March 2008).
- [4] Huang, C. and Srihari, S. N., “Mapping transcripts to handwritten text,” in [*Proc. International Workshop on Frontiers in Handwriting Recognition (IWFHR-10)*], 15–20 (2006).
- [5] Kim, G., Govindaraju, V., and Srihari, S. N., “An architecture for handwritten text recognition systems,” *International Journal on Document Analysis and Recognition* **2**(1), 37–44 (1999).
- [6] Durbin, R., Eddy, S., Krogh, A., and Mitchison, G., [*Biological Sequence Analysis: Probabilistic Models of Proteins and Nucleic Acids*], Cambridge University Press (2001).
- [7] Kim, G. and Govindaraju, V., “A lexicon driven approach to handwritten word recognition for real-time applications,” *IEEE Transactions on Pattern Analysis and Machine Intelligence* **19**, 366–379 (April 1997).