

# Vector Model Based Indexing and Retrieval of Handwritten Medical Forms

Huaigu Cao, and Venu Govindaraju  
Center for Unified Biometrics and Sensors (CUBS)  
Dept. of Computer Science and Engineering  
University at Buffalo, Amherst, NY  
{hcao3, govind}@cubs.buffalo.edu

## Abstract

*A vector model based information retrieval of handwritten medical forms is presented in this paper. In order to improve the IR performance on the erroneous output of handwriting recognition (HR) systems, a variation of the vector model is made to estimate the number of occurrences of terms from word segmentation and recognition probabilities. IR Tests show that our approach outperforms the retrieval of ordinary HR text in terms of mean average precision (MAP), R-Precision, and interpolated 11-point precisions.*

## 1 Introduction

Our work is motivated by the objective of building a system for indexing and retrieving New York State Pre-hospital Care Reports (PCR forms). In New York State all patients who enter the Emergency Medical System (EMS) are tracked through their pre-hospital care to the emergency room using the PCR. The PCR is used to gather vital patient information and has manifold purpose: (i) it is used by the emergency room staff in providing appropriate care, (ii) it is to facilitate record keeping about the patient and the care provided, and (iii) to afford large regional EMS (such as WREMS, Western Regional Emergency Medical Systems,) to gain trend knowledge through macro analysis. Our goal is to automate the collection of data from the PCR and enable efficient maintenance, and dissemination of information. The task proposed is quite challenging for several reasons: (i) handwritten responses are very loosely constrained in terms of writing style, format of response, and choice of text, (ii) medical lexicons of words are very large (about 50,000 entries), (iii) forms are often very noisy due to irrepressible emergency situations. This leads to difficulties in the automatic transcription of forms. Our dataset consists of over 6000 scanned PCR forms containing medicine and healthcare information. Each PCR contains about 100

handwritten words on average so the content is very short. The word recognition rate of the forms using Word Model Recognizer (WMR) [7] is as low as 20% as reported in [10].

Several works have been done to improve the IR performance of OCR'ed text. Researchers [2, 3] have shown that the performance of OCR text retrieval is badly affected when dealing with short or low quality documents. In [6, 11, 12] different approaches modeling typical recognition errors were proposed. In [11] a probabilistic model for misrecognition was proposed and this model was used to design the term-weighting scheme of information retrieval. The approach that generates candidate terms for each "true" search term and adds the retrieval results of candidate terms into the final result was studied in [12]. In [6], a language model that took common recognition errors into account was built. This language model can then be used to approximate an "uncorrupted" version of a particular document, and it can be used for retrieval in a language modeling approach.

The problem of indexing and retrieving handwritten documents has recently been addressed by researchers. Due to low recognition accuracies, It is difficult to use probabilistic modeling of OCR'ed text for indexing and retrieving handwritten documents. A new trend in this research area is to index every word image with word recognition probabilities of all term candidates. Rath and Manmatha[14] proposed an OCR-free approach to historical manuscript retrieval that learned the joint probability of the query word and features of the word image. However probability estimate cannot be accurate when the dimensionality of feature vector increases. A few works have attempted to solve this problem using OCR ranks [4, 5, 8]. Lee *et al.*[8] implemented retrieval on text composed of top- $k$  candidates of character recognition results of Hangul document images. Howe *et al* proposed a probabilistic model that uses top- $k$  word recognition results. Instead of considering top- $k$  candidates equally, the probability of the  $k$ -th candidate is assumed to be inversely proportional to  $k$ . In [4], a vector model that uses top- $k$  OCR output is defined and the posterior proba-

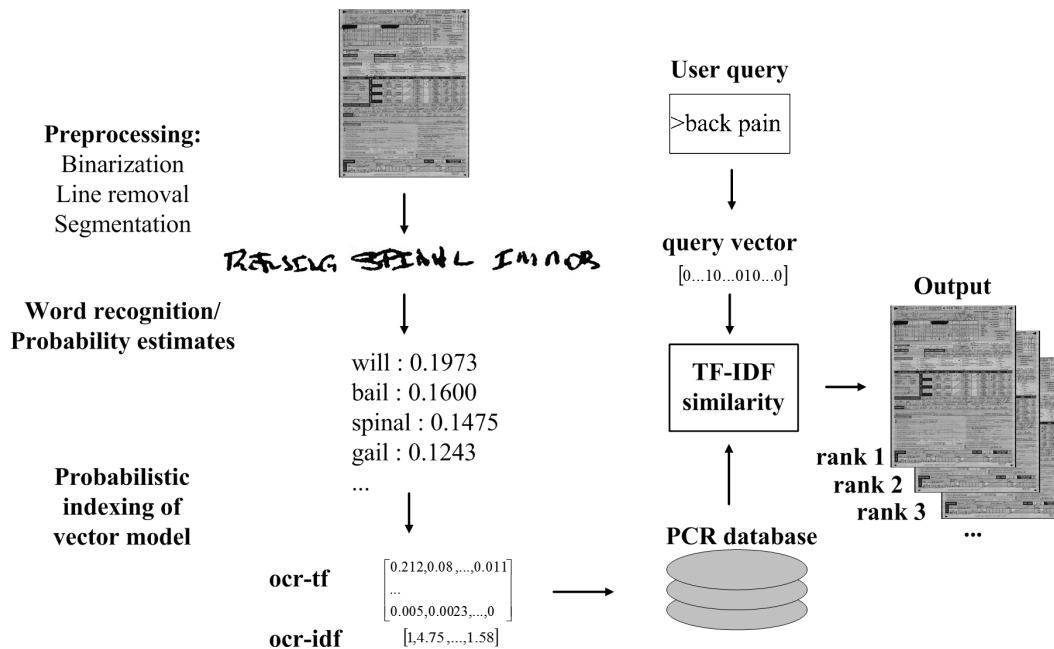


Figure 1. Flowchart of the search engine.

bility of word recognition at rank  $k$  is learnt from a training set rather than assigned empirically.

Besides new approaches to handwritten document retrieval, keyword spotting as an alternative approach of indexing and retrieving handwritten documents has been proposed in [9, 13, 16]. The idea is to search the document for a certain keyword by feature matching instead of recognition.

This paper describes a PCR handwritten form retrieval algorithm based on a vector model using multiple word recognition choices of top ranks. The algorithm in this paper is different from previous work [4, 5, 8, 14] in that word segmentation errors are modeled by segmentation probabilities and that the posterior probability of word recognition is modeled as a gaussian function of the edit distance produced by word recognizer, so it is not simply a function of rank regardless of the similarity between the word image and the entry in the lexicon.

The organization of this paper is as follows. The PCR retrieval system using OCR-rank based vector model is described in Section 2. Experiment results are presented in section 3. Finally section 4 concludes our work.

## 2 The Vector Model based Search Engine for PCR Forms

### 2.1 System overview

The flowchart of the search engine for the handwritten PCR forms is shown in figure 1. When indexing the

database, every PCR form is binarized, the form layout is analyzed, borderlines are removed, and lines of text are segmented into word images. Then a word recognizer produces ranked word recognition results with estimated poster probabilities. In the final step of indexing, OCR-related term frequency (TF) and inverted document frequency (IDF) are calculated and stored. When searching the database for relevant documents, the user input query is converted to a query vector and the cosine similarity of vector model is calculated for each document. Documents are ranked in the decreasing order of similarity and top documents are returned. Although the form segmentation and recognition methods are also very important, this paper focuses on the information retrieval method, i.e., the vector model which is discussed in the following section.

Owing to the fact that a PCR form has a fixed layout, form layout analysis is done by template matching. Line separation is done by a connected component based clustering algorithm. To ensure most word images are retrieved, the clustering algorithm is performed twice, each time with a different threshold of maximum distance between components allowed, and the produced word image hypotheses are combined. Therefore it is possible that multiple hypotheses overlap. The word recognition algorithm we use is based on over-segmentation and matching to minimize Euclidian distance between test image and cluster centers of the training set[7].

## 2.2 The Vector Model for Low-Accuracy OCR Results

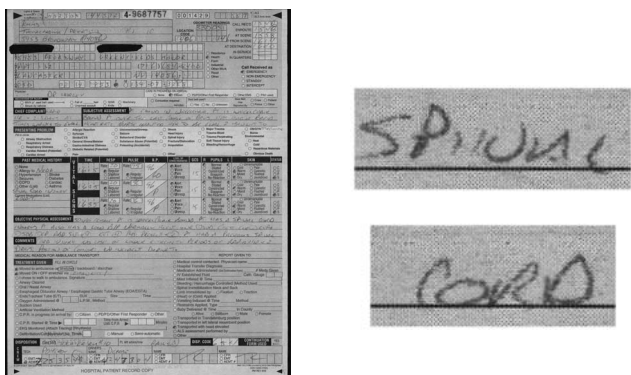
If we only consider text retrieval, given a set of documents  $\{d_j\}$  and a set of terms  $\{t_i\}$ , the term frequency of term  $t_i$  in document image  $i$  is defined by formula

$$tf_{i,j} = \frac{freq_{i,j}}{L} \quad (1)$$

where  $freq_{i,j}$  is the number of occurrences of term  $t_i$  in document  $d_j$  and  $L$  is the total number of occurrences of all terms in document  $d_j$ , i.e., the length of the  $j$ -th document. The inverted document frequency is defined by formula

$$idf_i = \log \frac{\|\{d_j\}\|}{\|\{d_j | freq_{i,j} > 0\}\|} \quad (2)$$

where  $\|\{\cdot\}\|$  denotes the number of elements in set  $\{\cdot\}$  [1].



(a)

(b)

(c)

**Figure 2. An example of PCR forms. (a) A entire PCR form. (b) A small local region showing obscure text and background noise array. (c) Fields of interest in the PCR form.**

In the modified vector model, the term frequency of term  $t_i$  in document image  $d_j$ , denoted by  $tf_{i,j}^{\text{ocr}}$ , and the inverted document frequency of term  $t_i$  in the collection of document images, denoted by  $idf_i^{\text{ocr}}$ , are defined in the following equations:

$$tf_{i,j}^{\text{ocr}} = \frac{E\{freq_{i,j}\}}{L}, \quad (3)$$

and

$$idf_i^{\text{ocr}} = \log \frac{\|\{d_j\}\|}{\|\{d_j | E\{freq_{i,j}\} > 0.5\}\|} \quad (4)$$

where  $E\{freq_{i,j}\}$  denotes the expectation, or mean value, of  $freq_{i,j}$ . The query term frequency is defined by

$$tf_{i,q} = \begin{cases} 1 & \text{if term } t_i \text{ is in query phrase } q, \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

The similarity between any document  $d_j$  and the query  $q$  is defined by

$$\text{sim}(d_j, q) = \sum_i tf_{i,j}^{\text{ocr}} \cdot idf_i^{\text{ocr}} \cdot tf_{i,q}. \quad (6)$$

$E\{freq_{i,j}\}$  in equations (3) and (4) can be calculated using formula

$$E\{freq_{i,j}\} = \sum_{I_w \in d_j} \Pr(V_{(I_w)} = \text{true}) \cdot \Pr(t_i | I_w) + \epsilon, \quad (7)$$

where  $I_w$  runs over all segmentations of word images in document  $d_j$ , boolean function  $V_{(I_w)} \in \{\text{true}, \text{false}\}$  indicates if  $I_w$  is a valid segmentation of word image,  $\Pr(t_i | I_w)$  is the posterior probability of word recognition ignoring stop words filtering and trimming, and  $\epsilon$  is the average number of lost terms  $t_i$  due to incorrect segmentation.

The value of  $\epsilon$  is very small for the short content of each PCR and is omitted.  $\Pr(V_{(I_w)} = \text{true})$  is estimated in the following way: suppose a hypothesis  $I_w$  overlaps  $m$  other hypotheses, then

$$\hat{\Pr}(V_{(I_w)} = \text{true}) = \frac{1}{1 + m}. \quad (8)$$

Using a Distance Modeling Method, the posterior probability  $\Pr(t_i | I_w)$  in equation (7) is estimated from distances produced by the word recognition algorithm.

### Distance-Modeling Method

The observation density  $\Pr(I_w | t_i)$  can be modeled as a gaussian function of the distance  $d_i$  returned by the recognizer, i.e.[15],

$$\Pr(I_w | t_i) \sim e^{-\frac{d_i^2}{2\sigma^2}} \quad (9)$$

In our test, the parameter  $\sigma$  is empirically selected as  $\sigma = 0.707$ . Assuming that all terms have identical prior probability, then the posterior probability is estimated by

$$\hat{\Pr}(t_i | I_w) = \frac{e^{-\frac{d_i^2}{2\sigma^2}}}{\sum_i e^{-\frac{d_i^2}{2\sigma^2}}}. \quad (10)$$

There are two alternative methods for estimating the probability: Top-Rank Method and Weighted Top-Rank Method[4].

#### Top-Rank Method

$$\hat{\Pr}(t_i|I_w) = \begin{cases} \frac{1}{S}, & 1 \leq \text{rank}(t_i) \leq S \\ 0, & \text{otherwise} \end{cases} \quad (11)$$

where  $S$  is a truncation threshold.

#### Weighted Top-Rank Method

$$\hat{\Pr}(t_i|I_w) = \text{top-R OCR rate} - \text{top-(R-1) OCR rate} \quad (12)$$

where  $R = \text{rank}(t_i)$ .

### 3 Experimental Results and Discussions

Our corpus contains 342 PCR forms with manually transcribed ground truth and coordinates of each word. The test bed involves 22 queries, and manual annotation of relevance of the 342 forms to these queries. An example of an entire PCR form and handwritten regions of interest in the PCR form are shown in figure 2. Of all the 22 queries, 14 contain one term each, 6 contain two terms each, and 2 contain three terms. The total number of relevant documents (with duplicates) is 386. Different from the tests on PCR Forms [4], only those words in certain fields (Chief Complaints, Subjective/Objective Physical Assessments, and Comments) of the form are involved in our IR test.

16303 word images were obtained after form layout analysis and word segmentation. Among all the word images, roughly 63% are correctly segmented word images, 32% are under or over segmented word images, and 5% are combinations of symbols or digits that are not in the lexicon. Word images were recognized and probabilities were estimated as discussed in section 2.2. The top-1, 2, 5, and 10 word recognition accuracies are 20%, 27%, 35%, and 42%.

We compared the performances of IR tests using different posterior probability estimation methods in equations (10)-(12) as well as IR tests on OCR'ed and ground-truth text using standard vector model. All above tests are run on human segmented word images and automatically segmented word images. For comparison, in addition to using equation (8), we also tested the effects of ignoring segmentation errors, that is, letting  $\hat{\Pr}(V_{I_w} = \text{true}) \equiv 1$  for any  $I_w$ .

The IR tests are evaluated in terms of Mean Average Precision (MAP) and R-Precision. The Mean Average Precision is the mean value of the average precisions computed

for each query. R-Precision is the mean value of precisions computed for each query when  $R$  documents are retrieved, where  $R$  is the number of relevant documents. Table 1 shows the MAP and R-Precision of every IR test. We can see that the use of the proposed OCR vector model produces remarkable improvement of MAP (from 6.7% to 18.7%) and R-Precision (from 7.4% to 22.6%) as compared to the standard method on ordinary OCR text. We can also see that the distance-modeling method in equation (10) outperforms two other probability estimation methods in equations (11) and (12). Finally, most tests, except the R-precision values of the group of tests using Weighted Top-Rank method, show that the use of word segmentation probability of equation (8) improves IR performance.

Figure 3 visually shows the improvement of the proposed method and the gap between it and the result on ground-truth text in terms of interpolated 11-point average precisions. In the standard method, the precision almost drops to zero when the recall is 40%, whereas in the proposed method, the precision remains above 10% when the recall is 60%. The proposed method improves IR performance in the sense that more relevant documents are returned while precision rate is maintained when using multiple OCR candidates under a probabilistic model.

### 4 Conclusion

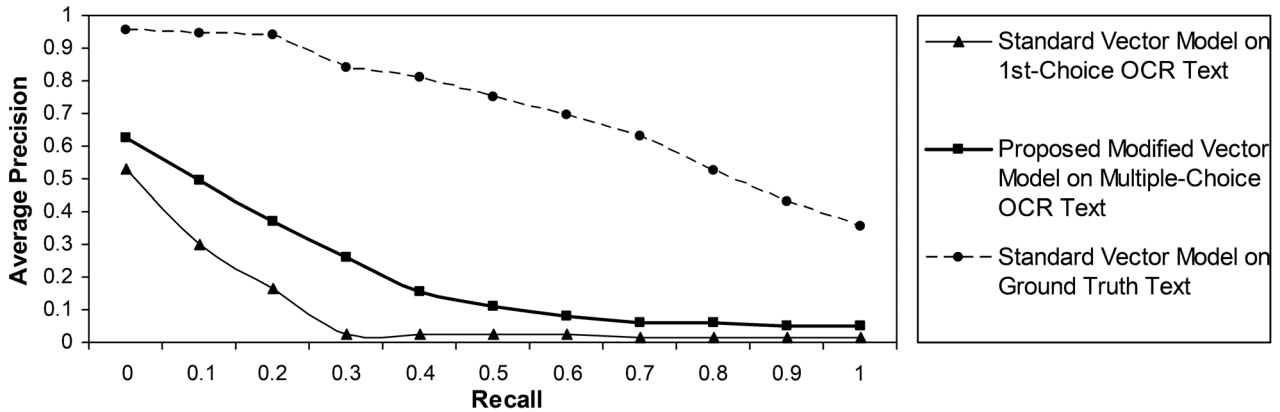
This paper proposed a vector model based method for retrieval of degraded medical form images with very low recognition rate. Our algorithm uses the posterior probabilities of word recognition results to estimate the term frequency and inverted document frequency. The proposed approach improved the performance of IR when compared to the standard Vector Model.

### References

- [1] R. A. Baeza-Yates and B. A. Ribeiro-Neto. *Modern Information Retrieval*. ACM Press / Addison-Wesley, 1999.
- [2] S. M. Beitzel, E. C. Jensen, and D. A. Grossman. A survey of retrieval strategies for ocr text collections. In *Proceedings of the Symposium on Document Image Understanding Technologies*, Greenbelt, Maryland, April 2003.
- [3] W. B. Croft, S. M. Harding, K. Taghva, and J. Borsack. An evaluation of information retrieval accuracy with simulated ocr output. In *Proceedings of the Symposium on Document Analysis and Information Retrieval*, 1994.
- [4] H.Cao, F. Farooq, and V. Govindaraju. Indexing and retrieval of degraded handwritten medical forms. In *Proceedings of the Workshop on Multimodal Information Retrieval at IJCAI-2007*, 2007.
- [5] N. R. Howe, T. M. Rath, and R. Manmatha. Boosted decision trees for word recognition in handwritten document retrievals. In *Proceedings of the SIGIR*, pages 377–383, 2005.

**Table 1. Performance evaluation of IR tests.**

Estimate of $\hat{P}_r(t_i I_w)$	OCR'ed text (no Estimate)	Top-Rank (S=5)	Top-Rank (S=7)	Weighted Top-Rank	Distance Modeling	Ground Truth (no Estimate)
MAP/R-Precision(%) (human word segmentation)	15.9/17.4	16.8/22.2	19.0/23.4	19.4/25.3	22.1/25.6	70.6/70.7
MAP/R-Precision(%) (not handling seg error)	<b>6.7/7.4</b>	12.8/16.9	14.0/18.8	15.9/20.3	17.6/20.8	N/A
MAP(%) (handling seg error)	7.8/7.7	13.9/18.6	15.5/21.2	16.3/20.1	<b>18.7/22.6</b>	N/A



**Figure 3. Precision-Recall curves of IR tests.**

- [6] H. Jing. Using hidden markov modeling to decompose human-written summaries. *Computational Linguistics*, 28(4):527–543, 2002.
- [7] G. Kim and V. Govindaraju. 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.
- [8] D.-R. Lee, W.-Y. Kim, and I.-S. Oh. Hangul document image retrieval system using rank-based recognition. In *Proceedings of the International Conference on Document Analysis and Recognition*, volume 2, pages 615–619, 2005.
- [9] R. Manmatha and W. Croft. Word spotting: Indexing handwritten manuscripts. In *Intelligent Multi-media Information Retrieval Collection*. AAAI press, May 1997.
- [10] R. Milewski and V. Govindaraju. Extraction of handwritten text from carbon copy medical form images. In *Document Analysis Systems*, pages 106–116, 2006.
- [11] E. Mittendorf, P. Schauble, and P. Sheridan. Applying probabilistic term weighting to OCR text in the case of a large alphabetic library catalogue. In *Research and Development in Information Retrieval*, pages 328–335, 1995.
- [12] M. Ohta, A. Takasu, and J. Adachi. Retrieval methods for english text with misrecognized ocr characters. In *Proceedings of the International Conference on Document Analysis and Recognition*, 1997.
- [13] T. M. Rath and R. Manmatha. Word spotting for historical documents. *International Journal on Document Analysis and Recognition*, 2006.
- [14] T. M. Rath, R. Manmatha, and V. Lavrenko. A search engine for historical manuscript images. In *Proceedings of the 27th annual international ACM SIGIR conference on research and development in information retrieval*, 2004.
- [15] S. Tulyakov and V. Govindaraju. Probabilistic model for segmentation based word recognition with lexicon. In *Proceedings. Sixth International Conference on Document Analysis and Recognition*, volume 2, pages 164–167, 2001.
- [16] B. Zhang, S. N. Srihari, and C. Huang. Word image retrieval using binary features. In *Document Recognition and Retrieval XI, SPIE*, pages 45–53. SPIE, 2004.