

Automatic recognition of handwritten medical forms for search engines

Robert Jay Milewski · Venu Govindaraju · Anurag Bhardwaj

Received: 26 June 2007 / Revised: 6 November 2008 / Accepted: 23 December 2008 / Published online: 12 February 2009
© Springer-Verlag 2009

Abstract A new paradigm, which models the relationships between handwriting and topic categories, in the context of medical forms, is presented. The ultimate goals are: (1) a robust method which categorizes medical forms into specified categories, and (2) the use of such information for practical applications such as an improved recognition of medical handwriting or retrieval of medical forms as in a search engine. Medical forms have diverse, complex and large lexicons consisting of English, Medical and Pharmacology corpus. Our technique shows that a few recognized characters, returned by handwriting recognition, can be used to construct a linguistic model capable of representing a medical topic category. This allows (1) a reduced lexicon to be constructed, thereby improving handwriting recognition performance, and (2) PCR (Pre-Hospital Care Report) forms to be tagged with a topic category and subsequently searched by information retrieval systems. We present an improvement of over 7% in raw recognition rate and a mean average precision of 0.28 over a set of 1,175 queries on a data set of unconstrained handwritten medical forms filled in emergency environments.

Keywords Handwriting analysis · Language models · Pattern matching · Retrieval models · Search process

This work was supported by the National Science Foundation.

R. J. Milewski (✉) · V. Govindaraju · A. Bhardwaj
Center of Excellence for Document Analysis and Recognition,
UB Commons, 520 Lee Entrance, Suite 202,
Amherst, NY 14228, USA
e-mail: jyoryoku@gmail.com; milewski@cedar.buffalo.edu

V. Govindaraju
e-mail: govind@cedar.buffalo.edu

A. Bhardwaj
e-mail: ab94@cedar.buffalo.edu

1 Introduction

This paper describes the first automatic recognition system for handwritten medical forms. In the United States, any pre-hospital emergency medical care provided must be documented. Departments of Health for each state provide a standard medical form to document all patient care from the beginning of the rescue effort until the patient is transported to the hospital. State laws require that emergency personnel fill out one form for each patient. Automatic recognition and retrieval of these forms is quite challenging for several reasons: (1) handwritten data in the form is unconstrained in terms of writing styles, variability in font type or size and choice of text due to emergency situations, (2) form images are noisy since they are obtained from carbon copies of the original forms, (3) dictionary of medical words is huge with over 40,000 words which leads to poor recognition results.

Figure 1 shows an example Pre-Hospital Care Report (PCR) [67] form which contains 16 information regions (see Table 1). Handwriting, from PCR regions 8, 9, 11, 13 and 14 are used for recognition and retrieval analysis. There are two phases to our research: (1) the recognition of handwriting on the medical form, and (2) a medical form query retrieval engine. Handwriting recognition is used to tag medical forms with a topic category to subsequently improve recognition performance. The medical forms reflect large lexicons containing Medical, Pharmacology and English corpus. While current state of the art recognizers report recognition performance between ~58–78%, on comparable lexicon sizes in the postal application [36,68,69], our experiments show ~25% raw match recognition performance on the medical forms. This underscores the extremely complicated nature of medical handwriting (Fig. 1). We have developed a method of automatically determining the topic category of a PCR form using machine learning and computational linguistics

Fig. 1 Pre-Hospital Care Report (PCR) labeled [67]

The image shows a Prehospital Care Report (PCR) form with handwritten entries and red numbered annotations (1-16) pointing to various sections. The form is titled "Prehospital Care Report" and includes the following sections:

- 1:** Call information: Date of call (01/02/00), Time (01:23:45), Agency (4-3881352), Mileage (12345), and Station (54321).
- 2:** Patient information: Name (JO SCHMOE), Address (123 EASY ST., JAMESTOWN NY 14701), Dispatch information (D2B), and Call location (123 EASY ST.).
- 3:** Care in progress on arrival: Medications (No PMD) and other details.
- 4:** Dispatch information: Dispatch information (D2B), Mileage (1720), and Station (1720).
- 5:** Call type as received: Emergency, Non-Emergency, Stand-by, and other options.
- 6:** Call record: CALL REC'D (0:34:7), ENROUTE (0:34:8), ARRIVED AT SCENE (0:35:1), FROM SCENE (0:40:8), AT DESTIN (0:41:2), IN SERVICE, and IN QUARTERS.
- 7:** Medications on arrival: MVA (seat belt used), Fall of feet, SSW, Machinery, Extractions, Seat belt used, and other options.
- 8:** Chief complaint: "I CAN'T BREATHE".
- 9:** Subjective assessment: "20 y/o ♀ PT FOUND SITTING ON OUTSIDE STEPS @ 020 E.H.A. ASTHMA. PT STATES SHE DOES NOT HAVE HER INHALERS. FRIEND STATES COUPLE".
- 10:** Presenting problem: Allergic reaction, Unconscious/Unresp, Shock, Major trauma, OB/GYN, Burns, Seizure, Head injury, Trauma-Blunt, Burns, Behavioral disorder, Spinal injury, Trauma-Penetrating, Environmental, General illness/Malaise, Substance abuse (Potential), Fracture/Dislocation, Soft tissue injury, Heat, Cold, Gastro-intestinal distress, Poisoning (Accidental), Amputation, Bleeding/Hemorrhage, Hazardous materials, Diabetic related (Potential), Cardiac arrest, Pain, Other, Obvious death.
- 11:** Past medical history: Allergy to, Hypertension, Stroke, Seizures, Diabetes, COPD, Cardiac, Asthma, Other (List), Current Medications (List) (ALBUTEROL).
- 12:** Vital signs table with columns for TIME, RESP, PULSE, B.P., LEVEL OF CONSCIOUSNESS, GCS, R, PUPILS, L, SKIN, and STATUS. Entries include 04:00 and 04:10.
- 13:** Objective physical assessment: "20 y/o ♀ PT. CACX3 @ HEENT @ JVD @ TRACHEAL SHIFT ↓ BREATH SOUNDS BILAT. @ C/P @ TRIMM. EXTREMITIES W/R @ D2B P70, NRB @ 12 LPM → NEURALGIC ALBUTEROL @ 6 LPM → IV NTS 20 CC".
- 14:** Comments: "MONITOR: NRB ST @ 102 BPM, IMPROVED RESPIRATIONS ON LOCATION DEF HOSPITAL ED. PT. TRANSFERRED TO CARE BY RN STAFF IN LS-1. 10/11/00 NRBNT-P".
- 15:** Treatment given: Medication administered (NS2), Cath. Gauge (20), and other options.
- 16:** Disposition: (See list) DEF HOSPITAL, DISP. CODE (1243), and other options.

techniques. We demonstrate the strategy for improving the raw word recognition rate by about 7% for a lexicon size of over 5,000 words.

2 Background

Though the task of efficient retrieval of text documents has been addressed by information retrieval community for several years [70], robust document search and retrieval has received some considerable attention lately [16]. The existing methods for document retrieval can be broadly classified into two categories: (1) OCR based methods [28,58,65], and (2) Word image matching based methods [2,54–56,64]. On one hand word image matching based methods rely heavily on the proper selection of image features [53] and similarity

methods [2,55], the OCR based methods depend on the word recognition accuracy. It has been shown that higher word recognition error rate adversely affects the document retrieval performance [14,40]. Therefore, an improved word recognition algorithm forms a basis for an efficient document retrieval system.

The basis for reducing the lexicon to improve recognition is a well researched strategy in handwriting recognition [26,68]. Although handwriting recognition and lexicon pruning/reduction [43] have been researched substantially over the years, many challenges still persist in the offline domain. Word recognition applications range from automated check recognition [35], postal recognition [20], historical documents recognition [18,21,25] and now emergency medical documents [45–47]. Strategic recognition techniques for

Table 1 PCR form description for Fig. 1

Tag	Form region
(1)	Form, agency, and ambulance vehicle identification
(2)	Patient and physician contact information
(3)	Care in progress on arrival and mechanism of injury
(4)	Dispatch information
(5)	Patient transfer information
(6)	Time duration between rescue and transport phases
(7)	Extrication and patient vehicle information
(8)	Chief complaint
(9)	Subjective assessment
(10)	Presenting problem
(11)	Past medical history
(12)	Vital/signs
(13)	Objective physical assessment
(14)	Physical assessment extension and/or comments
(15)	Treatment given
(16)	Ambulance crew identification

handwriting algorithms such as Hidden Markov Models (HMM) [11, 17, 18, 31, 37, 44, 48], Artificial Neural Networks (ANN) [6, 12, 13, 22, 50], Boosted Decision Trees [30] and Support Vector Machines (SVM) [1, 7] have been developed. Lexicon reduction has been shown to be critical to improvement of performance primarily because of the minimization of possible choices [26]. Even the systems dealing with a large vocabulary corpus have been successful [37, 38, 72].

Lexicon reduction schemes in general, rely upon finding a specific topic of the document and then using a fixed smaller vocabulary of the chosen category as the reduced lexicon. This is usually achieved by performing categorization of the OCR'd document text which is noisy. Bayer et al. [3] in their work learn the noise model of the OCR using word substrings extracted with an iterative procedure. Taghva et al. [63] study the performance of a naive bayes classifier applied to 400 recognized documents with an OCR error rate of nearly 14%. In this experiment, 6 categories out of 52 are analyzed and the highest rate of correct classification achieved is 83.3%. However, both of these strategies have been applied to machine print OCR'd text where the noise level is not as high as the handwritten documents. In the context of medical forms, where the word recognition rate is very low (~25%) and only few characters are recognized with high confidence scores, such methods are not applicable. Vinciarelli et al. [66] study noisy text categorization over synthetic handwritten data. In this research, noisy data is obtained by changing a certain percentage of characters obtained from the OCR. However this method only handles the case when the character is changed to another list of known characters, whereas

in the text obtained from medical forms, there are slots for potentially unknown or human unreadable characters.

Additionally, some lexicon reduction strategies have used the extraction of character information for lexicon reduction, such as that by Guillevic et al. [27]. However, such strategies reduce the lexicon for a single homogeneous category, namely cities within the country of Finland. In addition, usage of word length estimates for a smaller lexicon are available [27]. Caesar et al. [8] also state that prior reduction techniques [51, 60, 61] are unsuitable since they can only operate on very small lexicons due to enormous computational burdens [8]. Caesar [8] further indicates that Suen's [62] approach of n-gram combinatorics is sensitive to segmentation issues, a common problem with medical form handwriting [8]. However, Caesar's method [8] and those which are dependent on using the character information, and/or the character information of only one word to directly reduce the lexicon, suffer if one of the characters is selected incorrectly [8]. This is observable in the cursive or mixed-cursive handwriting types.

Many existing schemes, such as that of Zimmermann [71], assume that some characters can be extracted. However, in the medical handwriting domain this task is error prone. Therefore, operating a reduction scheme which can be robust to incorrectly chosen characters is necessary. We use sequences of characters to determine the medical topic category which has a lexicon of its own, thereby reducing the issues of using the character information directly. Similar to the study by Zimmermann et al. [71], the length of words are used with phrases.

Kaufmann et al. [34] present another HMM strategy which is primarily a distance-based method and uses model assumptions which are not applicable in the medical environment. For example, Kaufmann [34] assumes that "...people generally write more cooperatively at the beginning of the word, while the variability increases in the middle of the word." In the medical environment, variability is apparent when multiple health care professionals enter data on the same form. The medical environment also has exaggerated and/or extremely compressed word lengths due to erratic movement in a vehicle and limited paper space. Kaufmann [34] only provides a reduction of 25% of the lexicon size with little to no improvement in error rate, and the experiments are run only on a small sample of words.

3 Lexicon reduction

This research proposes the following hypothesis which is verified experimentally: a sequence of confidently recognized characters, extracted from an image of handwritten medical text, can be used to represent a topic category. A reduced lexicon can then be constructed specifically for a medical

form based on its classified categories to improve recognition performance on a second feedback loop.

A medical form training and test set have been created from actual PCR data. A software data entry system has been developed which allows human truthers to segment all PCR form regions and words, and provide a human interpretation for the word, denoted as the truth. Truthing is done in two phases: (1) the digital transcription of medical form text and (2) the classification of forms into topic categories. The distribution of PCR forms under each category is approximately equal in both the training and test set. The task has been supervised and performed by a health care professional with several years of field emergency medical services (EMS) experience. This emergency medical data set is the first of its kind.

A PCR can be tagged with multiple categories from Table 2. In our data set, no form had more than five category tags. The subjectivity involved in determining the categories makes the construction of a hierarchical chart representing all patient scenarios with respective prioritized anatomical regions a difficult task and exceeds the scope of this research. The following are some examples for classifying medical form text into categories (see Table 2):

Example 1 A patient treated for an emergency related to her pregnancy would be included in the *Reproductive System* category (see Table 2).

Example 2 A conscious and breathing patient treated for gun shot wounds to the abdominal region would fall into the *Circulatory/Cardiovascular System* due to potential loss of blood, as well as being categorized for *Abdominal, Back, and Pelvic* categories (see Table 2).

We take characters with the highest recognition as an input and produce higher level topic categories. A knowledge base is constructed during the *training phase* from a set of PCR forms. The knowledge base contains the relationships between terms and categories and essentially describes the features for topic categorization. The *testing phase* takes an unknown form, and reduces the lexicon using the knowledge base. This phase is evaluated using a separate testing set. Finally, after all content of the PCR form has been recognized, a search can take place by entering in a query. This phase is tested by querying the system with a set of phrase inputs. The forms are then ranked accordingly and returned to the user. The complete architecture of the proposed algorithm is also shown in Fig. 2.

In the training phase, a mechanism for relating uni-grams and bi-grams (henceforth uni/bi-grams) as well as categories from a PCR training set are constructed. The testing phase then evaluates the algorithm's ability to determine the categories from a test form by using a lexicon driven word recognizer (LDWR) [36] to extract the top-choice uni/bi-gram

Table 2 Categories are denoted by these anatomical positions

Ten body systems	Circulatory/cardiovascular, digestive, endocrine, excretory, immune, integumentary, musculoskeletal, nervous, reproductive, respiratory
Six body range locations	Abdomen, back/thoracic/lumbar, chest, head, neck/cervical, pelvic/sacrum/coccyx
Four extremity locations	Arms/shoulders/elbows, feet/ankles/toes, hands/wrists/fingers, legs/knees
Four general	Fluid/chemical imbalance, full body, hospital transfer/transport, senses

characters from all words. Since the recognizer output is very noisy given the unconstrained handwriting, very few characters are correctly recognized per word image. In our setup, we consider a maximum of 2 characters per word since LDWR [36] successfully extracts a bi-gram with spatial encoding information 40% of the time. If ≥ 3 characters are selected, then LDWR [36] successfully extracts a character $\leq 1\%$ of the time owing to the noisy recognition output. Hence the maximum value of n in the n -grams is taken to be 2 (see Fig. 4). A list of about 400 stopwords provided by PubMed are omitted from text analysis [29,49].

3.1 Training

The training stage involves a series of steps to construct a matrix that represents relationships between terms and categories. The training phase is divided into the following steps: (1) cohesive phrase generation, (2) spatial encoding strategy, (3) normalization, (4) term discrimination ability, and (5) reduced singular value decomposition [9,10]. The process begins with the extraction of cohesive phrases from the medical forms. These phrases are then converted to ESI encoding terms (ESI denotes "Exact Spatial Information" used as the encoding procedure for the uni/bi-gram terms; see definitions later in this section). A matrix is then constructed utilizing the ESI terms for the rows and the categories in the columns. The matrix is then normalized, weighted, and prepared in Singular Value Decomposition format.

Step 1 (cohesive phrase generation). A cohesive phrase is defined as the frequency of two words co-occurring versus occurring independently (see Eq. 1). Figure 3 shows a high cohesive phrase extraction example from a PCR. A passage

Fig. 2 Proposed algorithm road map

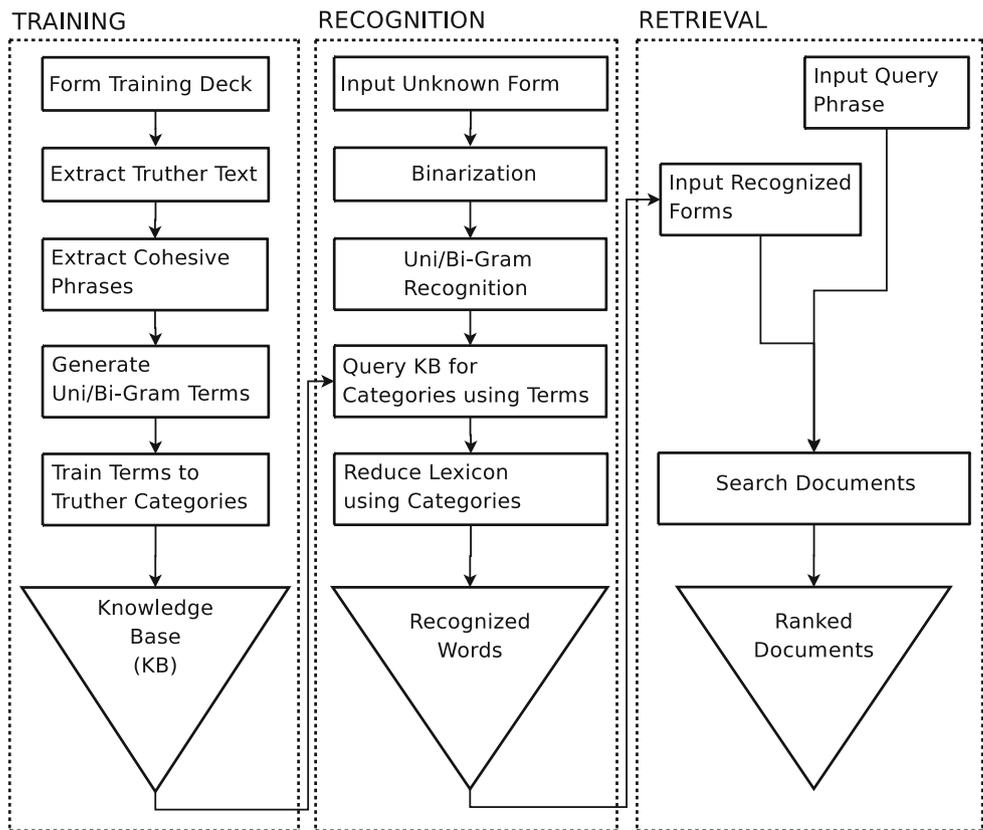
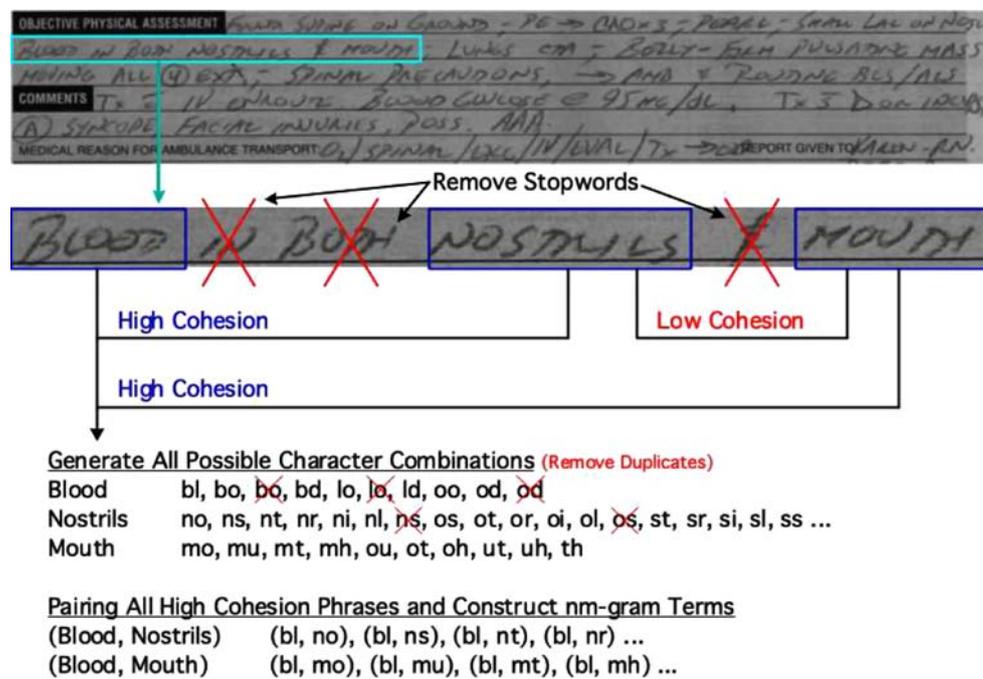


Fig. 3 Term extraction from high cohesive phrases



P is the set of all words w for a PCR form under a category C treated as a single string. For each C, every pair of passages, denoted P_1 and P_2 , is compared. A phrase is defined as a sequence of cohesive non-stopwords [19]. Here we denote w_x as a word located at position x within a passage P. Let a, b

and a', b' denote the index of words in an ordered passage P_1 and P_2 respectively ($w_a \in P_1, w_{a'} \in P_2, w_b \in P_1, w_{b'} \in P_2$ such that $b' > a'$ and $b > a$) then a potential phrase consisting of exactly two words is constructed. The cohesion of phrases under each C is then computed. If the cohesion is

Fig. 4 NSI encodings example (blue letters LDWR [36] successfully extracted)

(ID: 342-2)		*C*S* *A*N*
(ID: 407-1)		*C*M* *P*A*
(ID: 407-2)		*H*O* *A*W*
(ID: 473)		*C*F* *L*A*
(ID: 643)		*C*H* *A*W*
(ID: 695-1)		*L*S* *L*R*
(ID: 98)		*C*S* *S*S*
(ID: 606)		*C*W* *D*M*
(ID: 695-2)		*C*F* *W*L*

above a threshold, then that phrase represents that category C. Thus a category C is represented by a sequence of high cohesion phrases using only those PCR passages manually categorized under C. An additional list of about 50 words (e.g. male, female, etc.) found in most PCR's, which have little bearing on the category are omitted from the cohesion analysis but retained in the final lexicon.

$$\text{cohesion}(w_a, w_b) = z \cdot \frac{f(w_a, w_b)}{\sqrt{f(w_a)f(w_b)}} \quad (1)$$

The cohesion between any two words w_a and w_b is computed by the frequency that w_a and w_b occur together versus existing independently. The top 40 cohesive phrases are retained for each category (see Eq. 1). In the given equation, z is a constant weight which can be used if external information on the relationship of w_a and w_b is available ($z = 2$ in this research). The idea here is to analyze relationships between two words based on their correlations. If the two words are related to a category in some way, a higher correlation measure would reflect it accordingly.

Consider the following two unfiltered strings of words S_1 and S_2 under the category *legs*:

S_1 : “right femur fracture”

S_2 : “broken right tibia and femur”

The candidate phrases CP_1 and CP_2 after the filtering step are:

CP_1 : “right femur” ... “right fracture” ... “femur fracture”

CP_2 : “broken right” ... “right femur” ...

The phrase “right femur” is computed from CP_1 and CP_2 , given that w_a and $w'_a =$ “right”, w_b and $w'_b =$ “femur”, and the conditions $b > a$ and $b' > a'$ have been met. If the cohesion for “right femur” is above the threshold across all PCR forms under the *legs* category, then this phrase is retained as a representative of the category *legs*.

Tables 3 and 4 illustrate some top choice cohesive phrases generated. Digestive system and pelvic region are anatomically *close*. However, different information is reported in these two cases resulting in mostly different cohesive phrases. Those which are the same, such as *CHEST PAIN* have different cohesion values. This implies that it is likely that the term frequencies will also be different and therefore commonly occurring terms need to be weighted appropriately to their categories (this will be discussed in more detail later). Phrases sometimes may not make sense by themselves, however, this is the result of using a cohesive phrase formula in which words may not be adjacent.

Step 2 (spatial encoding strategy). Select one of three possible term representation strategies: NSI, ESI and ASI. These terms will later be modeled to an anatomical category and used as the essential criterion for lexicon reduction. The notation c denotes that a single character (uni-gram) is extracted from a word whereas c_1 and c_2 denote two ordered characters (bi-gram) are extracted from a word.

Table 3 Top cohesive phrases for the category: *pelvis*

Frequency	Cohesion	Phrase
2251	3.01	HIP PAIN
1860	2.49	PAIN HIP
390	0.83	PAIN CHEST
275	0.81	HIP FX
6	0.67	DCAP BTLS
288	0.59	HIP CHEST
213	0.55	PAIN LEG
202	0.50	PAIN JVD
205	0.42	CHEST HIP
163	0.40	JVD PAIN
106	0.40	CAOX3 PAIN
112	0.39	PAIN CHANGE
91	0.38	PAIN 0
82	0.38	PAIN 10
110	0.37	HIP CHANGE
118	0.36	PAIN FX
166	0.35	CHEST PAIN
144	0.34	HIP JVD
121	0.33	FELL HIP
3	0.33	PERPENDICULAR DECREASE

Table 4 Top cohesive phrases for the category: *digestive system*

Frequency	Cohesion	Phrase
52	0.81	STOMACH PAIN
30	0.72	PAIN INCIDENT
42	0.54	PAIN CHEST
25	0.54	PAIN SBM
39	0.51	CHEST PAIN
31	0.47	PAIN JVD
25	0.44	PAIN PMSX4
6	0.43	VOMITING ILLNESS
31	0.37	PAIN X4
11	0.34	PAIN EDEMA
5	0.31	PAIN TRANSPORTED
11	0.25	PAIN LEFT
9	0.25	HOME PAIN
4	0.24	CHEST SOFT
6	0.21	PAIN SOFT
3	0.21	SBM INCIDENT

No Spatial Information (NSI). An asterisk (*) indicates that zero or more characters are between or outside of c , c_1 and c_2 . NSI encodings are the most simple form of encoding (see Fig. 4 examples).

UNI-GRAM ENCODING: *c*

BI-GRAM ENCODING: *c₁*c₂*

BI-GRAM ENCODING EXAMPLE: BLOOD → *L*D*

Exact Spatial Information (ESI): The integers (x, y, z) represent the precise number of characters between or outside of c , c_1 and c_2 . ESI encodings are an extension of the NSI encodings with the inclusion of precise spatial information. In other words, the number of characters before, after and between the highest confidence c_1 and c_2 characters are part of the encoding. These encodings are the most successful in our experiments since there are fewer term collisions involved. Hence the ESI encodings are preferred.

UNI-GRAM ENCODING: xcy

BI-GRAM ENCODING: xc_1yc_2z BI-GRAM ENCODING EXAMPLE: BLOOD → 1L2D0

Approximate Spatial Information (ASI): The integers (x_a , y_a , z_a), denoted as length codes, represent an estimated range of characters between or outside of c , c_1 and c_2 . A '0' indicates no characters, a '1' indicates between one and two characters, and a '2' represents greater than 2 characters. The ASI encodings are an approximation of ESI encodings designed to handle cases when the precise number of characters is not known with high confidence.

UNI-GRAM ENCODING: $x_a c y_a$

BI-GRAM ENCODING: $x_a c_1 y_a c_2 z_a$

BI-GRAM ENCODING EXAMPLE: BLOOD → 1L1D0

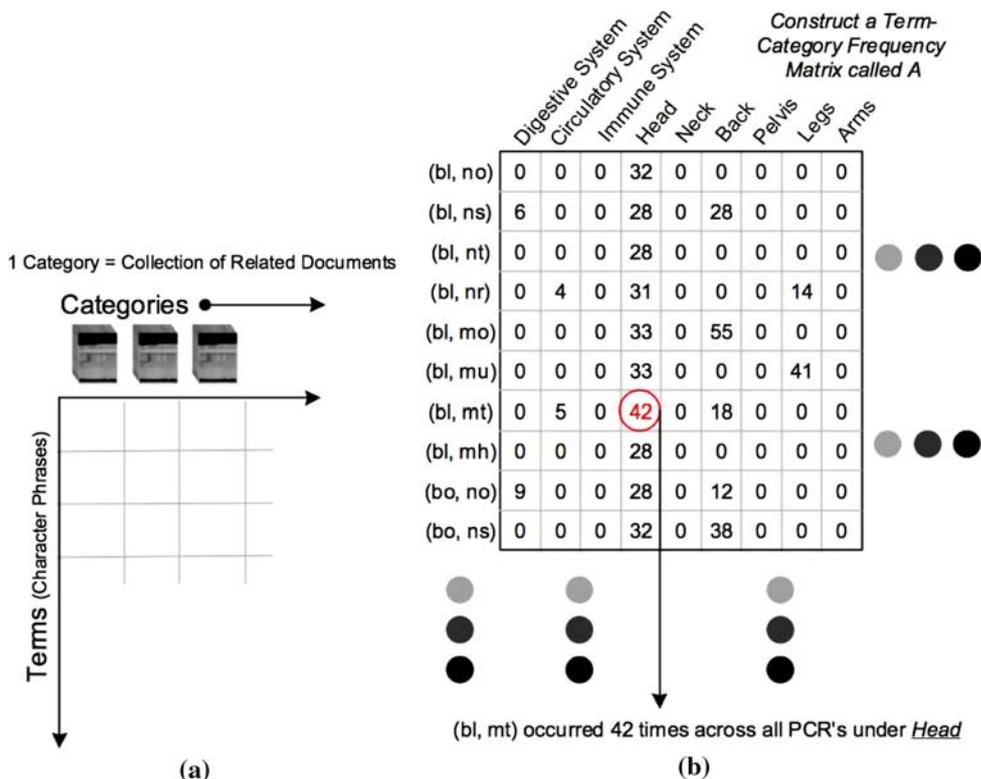
Combinatorial analysis

The quantity of all possible NSI, ESI and ASI uni-gram and bi-gram combinations, for a given word of character length n , such that $n \geq 1$, is represented by Eq. 2. Regardless of the encoding, the same quantity of combinations exist since there is only one encoding slot between or outside of the selected characters c , c_1 and c_2 . This is helpful in measuring the computational complexity of the encoding.

$$\mathcal{F}(n) = \left(\left(\sum_{i=1}^{n-1} (n-i) \right) + n \right) = \left(\left(\binom{n}{2} (n-1) \right) + n \right) \quad (2)$$

However, the function \mathcal{F} only considers the combinations of an individual word. The combination inflation of a uni/bi-gram phrase is shown by Eq. 3. The equation parameters a and b represent the string lengths of the words considered in a phrase. The total number of possible uni/bi-gram combinations resulting from a phrase P containing two words of length a and b is the product of the possible combinations of

Fig. 5 **a** Term Category Matrix (TCM) overview and **b** TCM frequency construction example



each word denoted as $\mathcal{F}(a)$ and $\mathcal{F}(b)$ respectively.

$$\mathcal{P}(a, b) = \mathcal{F}(a) \cdot \mathcal{F}(b) \tag{3}$$

For example

- Let the phrase for uni/bi-gram combinations be *PULMONARY DISEASE*
- Let $n = \text{length}(\text{"PULMONARY"}) = 9$
- Let $m = \text{length}(\text{"DISEASE"}) = 7$
- $\mathcal{F}(n) = 45$ uni-gram + bi-gram combinations for "PULMONARY"
- $\mathcal{F}(m) = 28$ uni-gram + bi-gram combinations for "DISEASE"
- $\mathcal{P}(n, m) = 1,260$ uni/bi-gram combinations for *PULMONARY DISEASE*

Each of these encodings has its advantages and disadvantages. The choice is ultimately based on the quality of the handwriting recognizer's (LDWR) ability to extract characters. If the handwriting recognizer cannot successfully extract positional information, then NSI is the best approach. If extraction of positional information is reliable, then the ESI is the best approach. However, NSI and ASI create more possibilities for recognizer confusion since distances are either approximated or omitted. ESI is more restrictive on the possibilities as the precise spacing is used leading to lesser confusion among terms.

Using the ESI protocol, all possible uni/bi-gram terms are synthetically extracted from each cohesive phrase under each category. For example, BLOOD can be encoded to the uni-gram 0B4 (zero characters before 'B' and four characters after 'B') and the bi-gram 0B3D0 (zero characters before 'B', three characters between 'B' and 'D' and zero characters following 'D'). All possible synthetic positional encodings are generated for each phrase and chained together (a '\$' is used to denote a chained phrase). For example, CHEST PAIN encodes to: 0C4\$0P0A2 ... 0C4\$1A2 ... 0C0H3\$0P1I1 ... 0C0H3\$0P2N0, etc. To improve readability, the notation (W_1, W_2) is used to represent an ESI encoding of a two-word phrase (e.g. Myocardial Infarction: (my, in), (my, if), (my, ia), etc ...). Therefore, each category now has a list of encoded phrases consisting of positional encoded uni/bi-grams. These terms are the most primitive representative links to the category used throughout the training process. In the training phase, the synthetic information can be extracted since the text is known. However, in the testing phase, a recognizer will be used to automatically produce an ESI encoding since the test text is not known.

A matrix *A*, of size $|T|$ by $|C|$, is constructed such that the rows of the matrix represent the set of terms *T*, and the columns of the matrix represent the set of category *C* as shown in Fig. 5a. The value at matrix coordinate (t, c) is the frequency that each term is associated with the category. The term frequency corresponds to the phrasal frequency from which it was derived. It is the same value as the numerator in the

cohesion formula (refer to Eq. 1): $f(w_a, w_b)$. For example, if the frequency of CHEST PAIN is 50, then all term encodings generated from CHEST PAIN, such as (ch, pa), will also receive a frequency of 50 in the matrix. An example of term frequency construction is shown in Fig. 5b.

Step 3 (normalization). Compute the normalized matrix B from A using Eq. 4 [9, 10], where normalization for a term is done over all possible categories.

$$B_{t,c} = \frac{A_{t,c}}{\sqrt{\sum_{e=1}^n A_{t,e}^2}} \quad (4)$$

Matrix A is the input matrix containing raw frequencies, Matrix B is the output matrix with normalized frequencies, and (t, c) is a (term, category) coordinate within a matrix. The normalization equation is used to normalize the frequency count of a term in a given category by the frequency of the same term in all possible categories, which reflects how representative the term is with respect to the given category.

Step 4 (term discrimination ability). The Term Frequency times Inverse Document Frequency (TF \times IDF) are used to favor those terms which occur frequently with a small number of categories as opposed to their existence in all categories [41, 59]. While Luhn [41] asserts that medium frequency terms would best resolve a document, it precludes classification of rare medical words. Salton's [59] theory, stating that terms with the most discriminatory power are associated with fewer documents, allows a rare word to resolve the document.

Step 4A. Compute the weighted matrix X from B using Eq. 5 [9, 10] [29]. IDF gives the inverse-document-frequency on term t , where $c(t)$ is the number of categories containing term t .

$$\text{IDF}(t) = \log_2 \frac{n}{c(t)} \quad (5)$$

Step 4B. Weight the normalized matrix by IDF values using Eq. 6 [9, 10, 29, 32]. Matrix B is the normalized matrix from Step 3, IDF is the computational step defined in Step 4A, and Matrix X is a normalized and weighted matrix.

$$X_{t,c} = \text{IDF}(t) \cdot B_{t,c} \quad (6)$$

Step 5 (reduced singular value decomposition). The normalized and weighted term-category matrix can now be used as the knowledge base for subsequent classification. A singular value decomposition variant, which incorporates a dimensionality reduction step allows a large term-category matrix to represent the PCR training set (see Eq. 7). This facilitates a category query from an unknown PCR using the LDWR [36] determined terms [9, 10, 15].

$$X = U \bullet S \bullet V^T \quad (7)$$

Matrix X is decomposed into three matrices: U is a $(T \times k)$ matrix representing term vectors, S is a $(k \times k)$ matrix, and V is a $(k \times C)$ matrix representing the category vectors. The value k represents the number of dimensions to be finally retained. If k equals the targeted number of categories to model, then SVD is performed without the reduction step. Therefore, in order to reduce the dimensionality, the condition $k < |C|$ is necessary to reduce noise [15].

3.2 Testing

Given an unknown PCR form, the task is to determine the form categories, and construct a reduced lexicon from those classified categories to drive the word recognizer, LDWR [36]. In addition, the categories determined can be used to tag the form which can be subsequently used for information retrieval. The testing phase is divided into the following steps: (1) term extraction, (2) pseudo-category generation, (3) candidate category selection, and (4) reduced lexicon recognition [9, 10].

Step 1 (term extraction). Given a new PCR image, all image words are extracted from the form, and LDWR [36] is used to produce a list of confidently recognized characters for each word. These are used to encode the positional uni/bi-grams consistent with the format during training. All combinations of uni/bi-phrases in the PCR form are constructed. Each word has exactly one uni-gram and exactly one bi-gram. A phrase consists of exactly two unknown words. Therefore it is represented by precisely four uni/bi-phrases (BI-BI, BI-UNI, UNI-BI and UNI-UNI).

Step 2 (pseudo-category generation). A $(m \times 1)$ query vector Q is derived, which is then populated with the term frequencies for the generated sequences from the term extraction step. If a term is not encountered in the training set, then it is not considered. Positional bi-grams are generated to yield the trained terms 37% of the time, and similarly positional uni-grams 57% of the time. The experiments here illustrate this to be a sufficient number of terms. A scaled vector representation of Q is then produced by multiplying Q^T and U .

Once the pseudo-category is derived, R-SVD is applied for the following reasons: (1) it converts the query into a vector space compatible input and (2) the dimensional reduction can help reduce noise [15]. Since the relationship between terms and categories is scaled by variance, the reduction allows parametric removal of less significant term-category relationships.

Step 3 (candidate category selection). The task is now to compare the pseudo-category vector Q with each vector in $V_r \bullet S_r$ (from the training phase) using a scoring mechanism. The cosine rule is used for matching the query [9, 10]. Both x and y are dimensional vectors used to compute the cosine in

Eq. 8. Vectors x and y in the equations represent the comparison of the vectors: pseudo-category Q with every column vector in $V_r \bullet S_r$.

$$z = \cos(x, y) = \frac{x \cdot y^T}{\sqrt{\sum_{i=1}^n x_i^2 \cdot \sum_{i=1}^n y_i^2}} \quad (8)$$

Each cosine score is mapped onto a sigmoid function using the least square fitting method, thereby producing a more accurate confidence score [9, 10]. The least squares regression line used to satisfy the equation $f(x) = ax + b$ are shown in Eqs. 9 and 10 [39]:

$$a = \frac{n \sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \quad (9)$$

$$b = \frac{1}{n} \left(\sum_{i=1}^n y_i - a \sum_{i=1}^n x_i \right) \quad (10)$$

The fitted sigmoid confidence is produced using the cosine score and the regression line, using Eq. 9:

$$\text{confidence}(a, b, z) = \frac{1}{1 + e^{-(az+b)}} \quad (11)$$

The confidence scores are then used to rank the categories. If a category is above an empirically chosen threshold, then that category is retained for the PCR. Multiple categories may be thus retained.

Step 4 (reduced lexicon recognition). All words corresponding to the selected categories are then used to construct a new reduced lexicon which is submitted to the LDWR recognizer [36] on its second round (i.e. bootstrapping/feedback loop). Note that the first LDWR execution round occurred during *Step 1 Term Extraction*. Given a test PCR form, and the reduced lexicon, the LDWR [36] converts the handwritten medical words to ASCII. Each word which is recognized is compared with the truth. However, a simple string comparison is insufficient due to spelling mistakes and root variations of word forms which are semantically identical. This occurs 20% of the time within the test set words. Therefore, Porter stemming [33, 52, 57] and Levenshtein String Edit Distance [4] of 1 allowable penalty are performed on both the truth and the recognizer result before they are compared. Levenshtein is only applied to a word that is believed to be ≥ 4 characters in length. For example, PAIN and PAINS are identical. However, this also results in an improper comparison in about 11% of the corrections. These corrections would only affect the performance measurements of match rate and do not affect the systems recognition ability.

Table 5 Handwriting recognition system environment

Environment item	Value
Training set PCR size	750
Testing set PCR size	62
Training set lexicon size	5,628
Testing set lexicon size	2,528
Training + testing set lexicon size	8,156
Training set words for modeling	42,226
Testing set words to recognize	3,089
Modeled categories/RSVD dimensions	24
Category selection threshold	0.55
Maximum categories per form	5
Average categories per form	1.40
Max phrases per category	50
Apple OS X memory usage	520 MB
Apple OS X G4 1 GHZ train time	15–20 min/exp
Apple OS X G4 1 GHZ test time	3 h/exp

4 Recognition experiments

Our training data consists of 750 PCR forms and the test data consists of a separate blind set of 62 PCR forms. In all experiments it is assumed that the word segmentation and extraction has been performed by a person. Also, forms in which 50% of the content is indecipherable by a human being are omitted. This occurs 15% of the time. A description of the training and test sets can be found in Table 5.

4.1 Performance measures

Table 6 contains seven columns corresponding to performance measure in recognition performance. These fields are EXP, ACCEPT, ERROR, RAW, LEX, ABSENT and ILLEG which are explained as follows:

EXP: the experiment that the performance values refer to.

ACCEPT (accept recognition rate): number of words the word recognizer accepts above an empirically decided threshold.

ERROR (error recognition rate): number of words incorrectly recognized among the accepted words.

RAW (raw recognition rate): top choice word recognition rate without use of thresholds.

LEX (lexicon size): the lexicon size for the experiment after any reductions.

ABSENT (truth word not present in the lexicon): percentage of words (for a specific experiment) not in the lexicon as a result of incorrectly chosen categories or due to the absence of that word in the training set.

ILLEG (word is illegible to a human being): percentage of the ABSENT set in which even human beings could not

Table 6 Handwriting recognition performance

EXP	ACCEPT (%)	ERROR (%)	RAW (%)	LEX	ABSENT (%)	ILLEG (%)
CL	76.34	71.93	23.31	5,628	–	–
CLT	76.92	69.65	25.32	8,156	–	–
AL	63.52	57.24	32.31	1,193	23.89	33.33
ALT	66.59	47.12	41.73	1,246	8.02	97.98
SL	70.51	62.26	30.30	2,514	16.06	48.19
SLT	71.51	59.44	32.73	2,620	10.46	73.99
RL	70.70	62.04	30.62	2,401	16.61	46.59
RLT	71.06	59.45	32.63	2,463	12.23	62.96

Table 7 Comparison and improvements illustrated by experiments

	CLT versus RLT (%)	CL versus RL (%)	CLT versus ALT (%)	CLT versus SLT (%)
RAW	↑ 7.48	↑ 7.42	↑ 17.58	↑ 7.42
ERROR	↓ 10.78	↓ 10.88	↓ 24.53	↓ 10.21

reliably decipher all or some of the characters in the word (given the context).

Table 7 contains conclusions in raw recognition and error rate based on the experiments in Table 6. These fields are RAW and ERROR which are explained as follows:

RAW: shows the improvement (denoted by an upward arrow in Table 7) in raw recognition performance between experiments.

ERROR: shows the reduction (denoted by a downward arrow in Table 7) in the incorrect accept rate between experiments.

4.2 Experiments

This section describes several kinds of experiments which correspond to Table 6. The purpose of these experiments is to compare and contrast the theoretical maximum recognition performance with the actual recognition performance. There are 4 major types of experiments: (C)omplete, (A)ssumed, (R)educed, and (S)ynthetic. The complete experiment means the recognizer was executed with the full lexicon. The assumed experiment means that a theoretically reduced lexicon is constructed under the assumption that the medical form categories are supplied by an oracle. The reduced experiment means that the actual latent semantic analysis in this paper is used to extract a reduced lexicon from recognized medical form categories. The synthetic experiment means that the uni/bi-grams were theoretically known (i.e. the handwriting recognizer always extracted 2 characters with 100% accuracy). However, since all words in a test set may not have been seen in a training set, the 4 experiments are executed in two modes: (1) with just words from the training set lexicon (L) and (2) words merged from both the training lexicon and

testing sets (LT). These two modes allow us to compare the performance in situations of known versus unseen words in a form. To indicate in the charts the different of each of 4 experiments in 2 modes, we use acronyms: CL and CLT for complete lexicon analysis in mode 1 and 2 respectively, and similarly AL versus ALT, SL versus SLT, and finally RL versus RLT. The experimental results can be found in Tables 6 and 7 with discussion that follows.

4.3 Discussion

In reference to Table 7 which is computed from the most relevant changes of Table 6: The theoretical RLT (i.e. comparing RLT to CLT) improves the RAW match rate by 7.48% and drops the error rate 10.78% with a *degree of reduction* $\rho = 61.59\%$. The practical RL (i.e. comparing RL to CL) improves the RAW match rate by 7.42% and drops the error rate by 10.88%. The RLT and RL numbers are close due to the difference in the initial lexicon sizes: CLT/RLT starts with 6,561 words (i.e. training set and testing set lexicons) whereas the CL/RL starts with 5,029 words (i.e. training set lexicon only). The RLT lexicon is more complete, but the lexicon is larger. The RL lexicon is less complete, but the lexicon is smaller. Thus, RLT gives the advantage that the recognizer has a greater chance of the word being a possible selection and RL gives the advantage of the lexicon being smaller. The ALT shows the theoretical upper bound for the paradigm: (1) the categories are correctly determined 100% and (2) the lexicon is complete. The ALT (i.e. comparing ALT to CLT) improves the RAW match rate by 17.58% and drops the error rate 24.53% with a *degree of reduction* $\rho = 83.01\%$. The synthetic experiments (SL and SLT) also do not offer much improvement which shows perfect

Table 8 Lexicon reduction performance between the complete lexicon (CL) and the reduced lexicon (RL)

Lexicon analysis metric	Value
Accuracy of reduction (α)	0.33
Degree of reduction (ρ)	0.83
Reduction efficacy (η)	0.06
Lexicon density (ϱ')	1.07 \rightarrow 0.87
Lexicon density (ϱ'')	0.50 \rightarrow 0.78

character extraction does not guarantee recognition improvement. This is due to two reasons: (1) a form is a representation of many characters and so some incorrectly recognized characters are tolerated and (2) the remaining words on the form to be recognized are difficult to determine even when the lexicon is constructed with only the words of known uni/bi-gram terms.

Table 8 provides insight into the effectiveness of the lexicon reduction from the complete lexicon (CL) to the reduced lexicon (RL) experiments. The performance measures for lexicon reduction as described by Madhvanath [42] and Govindaraju et al. [26] are used with alteration to the definition of reduction efficacy. The *Accuracy of Reduction* $\alpha = E(\mathcal{A})$, such that $\alpha \in [0, 1]$ [42], and \mathcal{A} is a random variable [5], indicates the existence of the truth in the lexicon. The function E computes the expectation [5]. The *Degree of Reduction* $\rho = E(\mathcal{R})$, such that $\rho \in [0, 1]$ [42], represents the mean size of the reduced lexicon. The *Reduction Efficacy* $\eta = \Delta_{\text{LDWR}} \times \alpha^{1-\rho}$, such that $\Delta_{\text{LDWR}}, \eta, \alpha, \rho \in [0, 1]$, is a measure of the effectiveness of a lexicon with respect to a lexicon driven recognizer. This formula is defined differently in this research to weigh the importance of accuracy over the reduction and include the reductions effect on the recognizer. The larger the efficacy value is, the better is the effectiveness of the reduction for one recognizer versus another. The larger the *Lexicon Density* $\varrho_{\text{LDWR}}(\mathcal{L}) = (\nu_{\text{LDWR}}(\mathcal{L}))(f_{\text{LDWR}}(n) + \delta_{\text{LDWR}})$ value (such that $\nu_{\text{LDWR}}(\mathcal{L}) = \frac{n(n-1)}{\sum_{i \neq j} d_{\text{LDWR}}(\omega_i, \omega_j)}$ and $d_{\text{LDWR}}(\omega_i, \omega_j)$ is a recognizer dependent computation used to denote a distance metric between two supplied words) the more *similar* or *close* the lexicon words are [26]. A supplemental distance measure denoted by the *N-Gram Lexicon Distance Metric* $d_{\text{LDWR}}(\omega_i, \omega_j) = \gamma(\omega_i, \omega_j) / \Gamma(\omega_i, \omega_j)$, introduced in this research and substituted into the lexicon density equation ϱ , provides a measure of uni/bi-grams existing within the lexicon. The value γ represents the number of uni/bi-gram terms that are *not* common between ω_i and ω_j . Γ denotes the total number of uni/bi-gram term combinations between ω_i and ω_j . In order to distinguish between the *lexicon density distance metric* and the *n-gram lexicon distance metric* equations, the values ϱ' and ϱ'' will be respectively used. The *lexicon density distance metric* ϱ' shows less

confusion among lexicon words considering all the characters are equally important. This implies that the reduced lexicon will be less confusing to the recognizer. The *n-gram lexicon distance metric* shows an increase in the quantity of words with common NSI encodings. This implies the recognizer has a greater chance of selecting a word using the confidently selected characters.

5 Search experiments

The ability to query a set of PCR medical forms which match a user supplied input phrase is important for Health Surveillance applications. Searching text in digital format is easily accomplished but this is much harder to do for scanned handwritten documents. While searching handwriting has only been demonstrated in certain areas [56]. The experiments in this section illustrate search effectiveness even when words are incorrectly recognized. Both the original LDWR (CL) and the reduced lexicon LDWR (RL) PCR medical form data sets are compared.

In order to have a query set of sufficient size, the test set is constructed using a leave-1-out strategy. There are eight rounds of recognition such that each round of the 800 PCR's are divided into two different groups of 100 and 700. During each of the rounds, the content of the 100 PCR's is recognized using the 700 PCR's as the training data. This allows the full set to be evaluated with no bias. Finally, a set of 1175 phrases, constructed from adjacent non-stopwords, are extracted from a blind set of 200 PCR forms (i.e. these 200 forms are not a subset of the 800 set) such that each phrase is found in at least one form in the 800 set. Each of the query phrase in the query set consists of exactly two words. Different experiments are conducted which search the PCR forms for at least one of the words or both the words from the input query phrase.

$$d(a_i, b_j) = w_{ij} * \frac{1}{|a_i - b_j|} \quad (12)$$

A query is performed by scanning the forms in the 800 test set for recognized words that match a two-word input query phrase. Any LDWR recognized form which contains the occurrence of both query words independently in the document are considered matched results. Relevancy is determined if the input query words, for example *CHEST* and *PAIN*, are actually found on that form according to the human truth. A two-step ranking algorithm is then performed on all matching documents. First documents are ranked according to the frequencies of the occurring words. Second, those documents with the same word frequency are ranked using the distance measurement in Eq. 12. Let $d(a_i, b_j)$ be a function which computes the distance between two matched words, a_i and b_j such that i and j respectively represent the word

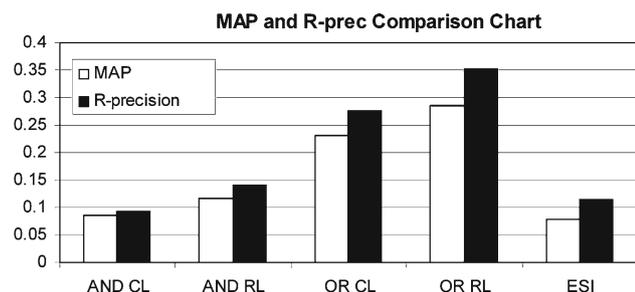


Fig. 6 Mean average precision and R-precision comparison for experiments

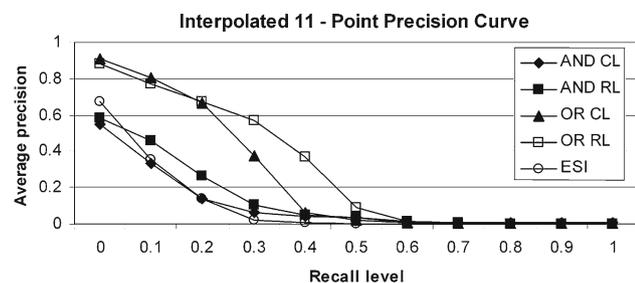


Fig. 7 Interpolated 11-point precision curve

position in the document. w_{ij} here is a weight based on the frequency of occurrences of words a and b in the document. This is especially necessary in situations where word a exists and b does not, and vice versa. Documents with closer proximity words are given a higher rank. Discussion on proximity based metrics can be found here [23]. Finally, the search methods are evaluated using the standard trec_eval system. To account for cases, where the system improperly returns no documents for a given query, -c option of trec_eval is used to include the relevance count of these queries in the final calculation.

5.1 Performance measures

MAP (mean average precision) is the mean of the average precision of all individual queries in the set. Average precision of a single query is defined as the mean of the precision after every relevant document retrieved. This performance measure emphasizes on retrieving relevant documents earlier.

R-prec (R-precision) is the precision at R , where R denotes the total number of relevant documents for the given query. This measure emphasizes on retrieving more relevant documents.

5.2 Experiments

AND CL. Given a query phrase of two words, both words are found in a PCR form during the search process using a complete training lexicon.

AND RL. Given a query phrase of two words, both words are found in a PCR form during the search process using a reduced training lexicon.

OR CL. Given a query phrase of two words, at least one of the words is found in a PCR form during the search process using a complete training lexicon.

OR RL. Given a query phrase of two words, at least one of the words is found in a PCR form during the search process using a reduced training lexicon.

ESI. An additional query expansion experiment was also performed in which a document was matched if at least one ESI encoding sequence was found in the document (i.e. the requirement for matching words was removed). For example, consider input query phrase *CHEST PAIN* where *CHEST* is decomposed into CH, CE, CS, CT, HE, HS, HT, ES, ET, C, H, E, S, and T., and *PAIN* is decomposed into PA, PI, PN, AI, AN, IN, P, A, I, and N. Since the input phrase is known, and hence the spatial encodings between characters are also known, the ESI encodings for the terms are known. The ESI encodings for *CHEST* are decomposed into: 0C0H3, 0C1E2, 0C2S1, 0C3T0, 1HE2, 1H1S1, 1H2T0, 2E0S1, 2E1T0, 0C4, 1H3, 2E2, 3S1, and 4T0. The ESI encodings for *PAIN* are decomposed into: 0P0A2, 0P1I1, 0P2N0, 1A0I1, 1A1N0, 2I0N0, 0P3, 1A2, 2I1, and 3N0. Finally, all possible ESI sequences from the input words are generated: 0C0H3 \$0P0A2, 0C0H3\$0P1I1, 0C0H3\$0P2N0, 0C0H3\$1A0I1, etc.

5.3 Discussion

The experimental results for each algorithm in terms of MAP and R-precision are shown in Fig. 6. As shown, retrieval based on reduced lexicon (RL) outperform retrieval based on complete lexicon (CL). This behavior is observed irrespective if the search is performed using both words from the query phrase (AND) or at least one of the words from the query phrase (OR). An interpolated 11 - point precision curve shown in Fig. 7 also supports this observation. As shown in the figure, after a recall level of 0.2, OR-RL method retrieves relevant documents earlier in the order as compared to OR-CL method. In the case of AND logic, RL based method performs better than CL based methods at all recall levels. The improvement in the search performance due to lexicon reduction algorithm used highlights the effectiveness of the proposed method. For the query expansion experiment (ESI) as intuitively expected, the uni/bi-grams match more terms in the test set due to the loss in word information. The precision chart in Fig. 7 illustrates this drop in retrieval effectiveness and shows that searches are more effective at the word level rather than raw encoding level. Similar drop in performance is observed for the query expansion technique in Figs. 6 and 8.

To study the effect of different methods on the total number of relevant documents retrieved, we also compute the

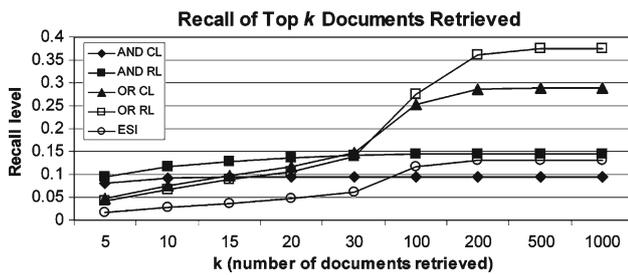


Fig. 8 Recall level of top k documents retrieved

value of recall and precision levels independently for the top k documents retrieved as shown in Fig. 8. The results from Fig. 8 suggest that reduced lexicon (RL)-based methods not only retrieve relevant documents earlier, but also retrieve more relevant documents overall as compared to their counterpart complete lexicon (CL)-based methods. The contribution of this research is that the lexicon reduction strategy (i.e. the RL experiment) improves both handwriting recognition and search effectiveness.

6 Conclusions

This paper defines a new paradigm for lexicon reduction and information retrieval in the complex situation of handwriting recognition of medical forms. An improvement in raw recognition rate from about 25% of the words on a PCR form to approximately about 33% has been shown with a reduction in false accepts by about 7%, a reduction in error rate by about 10–25%, and a lexicon reduction from 32–85%. The addition of a category driven query facilitates a mean average precision of 0.28 and R-prec of 0.35 for 1175 queries in a search engine experiment with medical forms. Additionally, using a reduced lexicon for searching medical form also enables retrieving more relevant number of documents overall, as compared to complete lexicon search.

Interestingly, certain computational elements of bootstrapping, described in our work, are consistent with the human interpretation of ambiguous handwriting using contextual cues. Our methodology accomplishes this by modeling character terms as a higher level semantic concept which mimics the human ability to recognize a word within context, when some characters are unknown.

Acknowledgments (1) This material is based upon work supported by the National Science Foundation under Grant IIS 0429358. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation. (2) Casey Czamara, of the Western Regional Emergency Medical Services (WREMS) program operating under the New York State Department of Health for providing necessary resources.

References

- Bahlmann, C., Haasdonk, B., Burkhardt, H.: On-line handwriting recognition with support vector machines—a kernel approach. *International Workshop On Frontiers in Handwriting Recognition* (2002)
- Balasubramanian, A., Meshesha, M., Jawahar, C.V.: Retrieval from document image collections. In: *Proceedings of Seventh IAPR Workshop on Document Analysis Systems*, pp. 1–12 (2006)
- Bayer, T., Kressel, U., Mogg-Schneider, H., Renz, I.: Categorizing paper documents. *Comput. Vis. Image Understand.* **70**(3), 299–306 (1998)
- Black, P.E. (ed.): *Levenshtein distance. Algorithms and Theory of Computation Handbook*; CRC Press LLC, dictionary of Algorithms and Data Structures, NIST (1999)
- Blum, J.R., Rosenblatt, J.I.: *Probability and statistics. Random Variables and Their Distributions*, chap. 4. Expectations, Moment Generating Functions, and Quantiles, chap. 6. W.B. Saunders Company, USA (1972)
- Blumenstein, M., Verma, S.: A neural based segmentation and recognition technique for handwritten words. *IEEE Int. Conf. Neural Netw.* (1998)
- Byun, H., Lee, S.W.: *Applications of support vector machines for pattern recognition: a survey. Lecture Notes in Computer Science*. Springer, Berlin (2002)
- Caesar, T., Gloger, J.M., Mandler, E.: Using lexical knowledge for the recognition of poorly written words. In: *Third International Conference on Document Analysis and Recognition*, vol. 2, pp. 915–918 (1995)
- Chu-Carroll, J., Carpenter, B.: Dialogue management in vector-based call routing. In: *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics*, pp. 256–262 (1999)
- Chu-Carroll, J., Carpenter, B.: Vector-based natural language call routing. *Comput. Linguist.* **25**(3), 361–388 (1999)
- Chen, M.Y., Jundu, A., Zhou, J.: Off-line handwritten word recognition using a hidden markov model type stochastic network. *IEEE Trans. Pattern Anal. Mach. Intell.* **16**(5), 481–496 (1994)
- Cho, S.B., Kim, J.H.: Applications of neural networks to character recognition. *Pattern Recognit.* (1991)
- Cho, S.B.: Neural-network classifiers for recognizing totally unconstrained handwritten numerals. *IEEE Trans. Neural Netw.* **8**(1), 43–53 (1997)
- Croft, B., Harding, S.M., Taghva, K., Borsack, J.: An evaluation of information retrieval accuracy with simulated OCR output. In: *Proceedings of Symposium on Document Analysis and Information Retrieval*, pp. 115–126 (1994)
- Deerwester, S., Dumais, S.T., Furnas, G.Q., Landauer, T.K., Harshman, R.: Indexing by latent semantic analysis. *J. Am. Soc. Inf. Sci.* **41**(6), 391–407 (1990)
- Doermann, D.: The indexing and retrieval of document images: a survey. *Comput. Vis. Image Understand.* **70**(3), 287–298 (1998)
- Edwards, J., Forsyth, D.: Searching for character models. In: *Proceedings of the 19th Annual Conference on Neural Information Processing Systems*, Vancouver, Canada, pp. 331–338 (2005)
- Edwards, J., Teh, Y.W., Forsyth, D., Bock, R., Maire, M., Vesom, G.: Making Latin manuscripts searchable using (gHMM)'s. In: *Proceedings of the 18th Annual Conference on Neural Information Processing Systems*, pp. 385–392 (2004)
- Fagan, J.: The effectiveness of a non-syntactic approach to automatic phrase indexing for document retrieval. *J. Am. Soc. Inf. Sci.* **40**, 115–132 (1989)
- Favata, J.T.: Offline general handwritten word recognition using an approximate BEAM matching algorithm. *IEEE Trans. Pattern Anal. Mach. Intell. (PAMI)* **23**(9), 1009–1021 (2001)

21. Feng, S.L., Manmatha, R.: Classification models for historic manuscript recognition. In: Proceedings of the Eighth International Conference on Document Analysis and Recognition (ICDAR) (2005)
22. Gader, P.D., Keller, J.M., Krishnapuram, R., Chiang, J.H., Mohamed, M.A.: Neural and fuzzy methods in handwriting recognition. *Computer* **30**(2), 79–86 (1997)
23. Goldman, R., Shivakumar, N., Venkatasubramanian, S., Garcia-Molina, H.: Proximity search in databases. *IEEE Proc. Int. Conf. Very Large Databases*, pp. 26–37 (1998)
24. Golub, G.B., Van Loan, C.E.: *Matrix Computations*, 2nd edn. John Hopkins University Press, Baltimore (1989)
25. Govindaraju, V., Xue, H.: Fast handwriting recognition for indexing historical documents. In: First International Workshop on Document Image Analysis for Libraries (DIAL) (2004)
26. Govindaraju, V., Slavik, P., Xue, H.: Use of lexicon density in evaluating word recognizers. *IEEE Trans. PAMI* **24**(6), 789–800 (2002)
27. Guillevic, D., Nishiwaki, D., Yamada, K.: Word lexicon reduction by character spotting. In: Seventh International Workshop on Frontiers in Handwriting Recognition, Amsterdam (2000)
28. Harding, S.M., Croft, W.B., Weir, C.: Probabilistic retrieval of OCR degraded text using n-grams. In: *Research and Advanced Technology for Digital Libraries*, pp. 345–359 (1997)
29. Hersh, W.R.: *Information Retrieval: A Health and Biomedical Perspective*, 2nd edn. Springer-Verlag, New York, Inc. USA (2003)
30. Howe, N.R., Rath, T.M., Manmatha, R.: Boosted decision trees for word recognition in handwritten document retrieval. In: Proceedings of the 28th Annual Int'l ACM SIGIR Conference, pp. 377–383 (2006)
31. Hu, J., Brown, M.K., Turin, W.: HMM based online handwriting recognition. *IEEE Trans. Pattern Anal. Mach. Intell. (PAMI)* **18**(10), 1039–1045 (1996)
32. Jones, K.S.: A statistical interpretation of term specificity and its application in retrieval. *J. Document.* **28**(1), 11–20 (1972)
33. Jones, K.S., Willet, P.: *Readings in Information Retrieval*. Morgan Kaufmann, San Francisco (1997)
34. Kaufmann, G., Bunke, H., Madom, M.: Lexicon reduction in an HMM-Framework based on quantized feature vectors. In: Proceedings of the 4th International Conference on Document Analysis and Recognition (ICDAR), vol. 2, pp. 1097–1101 (1997)
35. Kim, G., Govindaraju, V.: Bank check recognition using cross validation between legal and courtesy amounts. *IJPRAI* **11**(4), 657–674 (1997)
36. Kim, G., Govindaraju, V.: A lexicon driven approach to handwritten word recognition for real-time applications. *IEEE Trans. PAMI* **19**(4), 366–379 (1997)
37. Koerich, A.L., Sabourin, R., Suen, C.Y.: Fast two-level HMM decoding algorithm for large vocabulary handwriting recognition. Ninth International Workshop on Frontiers in Handwriting Recognition (IWFHR-9), pp. 232–237 (2004)
38. Koerich, A.L., Sabourin, R., Suen, C.Y.: Large vocabulary off-line handwriting recognition: a survey. *Pattern Anal. Appl.* **6**, 97–121 (2003)
39. Larson, R.E., Hostetler, R.P., Edwards, B.H.: *Calculus with Analytic Geometry*, chap. 13, sect. 13.9, 5th edn. D.C. Heath and Company, USA (1994)
40. Lopresti, D., Zhou, J.: Retrieval strategies for noisy text. In: Proceedings of Symposium on Document Analysis and Information Retrieval, pp. 255–270 (1996)
41. Luhn, H.: A statistical approach to mechanized encoding and searching of literary information. *IBM J. Res. Dev.* **1**, 309–317 (1957)
42. Madhvanath, S.: The holistic paradigm in handwritten word recognition and its application to large and dynamic lexicon scenarios. Ph.D. Dissertation, University at Buffalo Computer Science and Engineering (1997)
43. Madhvanath, S., Krpasundar, V., Govindaraju, V.: Syntactic methodology of pruning large lexicons in cursive script recognition. *J. Pattern Recognit. Soc. Pattern Recognition*, vol. 34. Elsevier Science, Amsterdam (2001)
44. Marti, U.V., Bunke, H.: Using a Statistical Language Model to Improve the Performance of an HMM-based Cursive Handwriting Recognition Systems. *World Scientific Series in Machine Perception and Artificial Intelligence Series* (2001)
45. Milewski, R., Govindaraju, V.: Medical word recognition using a computational semantic lexicon. In: Eighth International Workshop on Frontiers in Handwriting Recognition, Canada (2002)
46. Milewski, R., Govindaraju, V.: Handwriting analysis of pre-hospital care reports. In: *IEEE Proceedings of the Seventeenth IEEE Symposium on Computer-Based Medical Systems (CBMS)* (2004)
47. Milewski, R., Govindaraju, V.: Extraction of handwritten text from carbon copy medical forms. *Document Analysis Systems (DAS)*. Springer, Berlin (2006)
48. Nakai, M., Akira, N., Shimodaira, H., Sagayama, S.: Substroke approach to HMM-based on-line kanji handwriting recognition. In: Sixth International Conference on Document Analysis and Recognition (2001)
49. National Library of Medicine. PubMed Stop List
50. Oh, I.-S., Suen, C.Y.: Distance features for neural network-based recognition of handwritten characters. *Int. J. Doc. Anal. Recognit. (IJ DAR)* **1**(2), 73–88 (2004)
51. Okuda, T., Tanaka, E., Kasai, T.: A method for the correction of garbled words based on the levenshtein distance. *IEEE Trans. Comput.*, Col. C-25, No. 2 (1976)
52. Porter, M.F.: An algorithm for suffix stripping. *Program* **14**, 130–137 (1980)
53. Rath, T.M., Manmatha, R.: Features for word spotting in historical manuscripts. In: Proceedings of IEEE International Conference on Document Analysis and Recognition, pp. 218–222 (2003)
54. Rath, T.M., Manmatha, R.: Word spotting for historical documents. *IJDAR* **9**(2), 139–152 (2007)
55. Rath, T.M., Manmatha, R.: Word image matching using dynamic time warping. In: Proceedings of the Conference on Computer Vision and Pattern Recognition, vol.2, pp. 521–527, Madison, WI (2003)
56. Rath, T.M., Manmatha, R., Lavrenko, V.: A search engine for historical manuscript images. In: *ACM SIGR*, pp. 369–376 (2004)
57. Rijsbergen, C.J. van, Robertson, S.E., Porter, M.F.: *New models in probabilistic information retrieval*. British Library, London (1980)
58. Russell, G., Perrone, M.P., Chee, Y.M.: Handwritten document retrieval. In: Proceedings of International Workshop on Frontiers in Handwriting Recognition, pp. 233–238 (2002)
59. Salton, G.: *Introduction to Modern Information Retrieval*. McGraw-Hill, New York (1983)
60. Sinha, R.M.K., Prasada, B.: Visual text recognition through contextual processing. *Pattern Recognit.* **21**(5), 463–479 (1988)
61. Srihari, S.N., Hull, J.J., Choudhari, R.: Integrating diverse knowledge sources in text recognition. *ACM Trans. Office Inf. Syst.* **1**(1), 68–87 (1983)
62. Suen, C.Y.: N-gram statistics for natural language understanding and processing. *IEEE Trans. Pattern Anal. Mach. Intell.* **1**(2), 164–172 (1979)
63. Taghva, K., Narkter, T., Borsack, J., Lumos, S., Condit, A., Young, R.: Evaluating text categorization in the presence of OCR errors. In: Proceedings of IS&T SPIE 2001 International Symposium on Electronic Imaging Science and Technology, pp. 68–74 (2001)
64. Tan, C.L., Huang, W., Yu, Z., Xu, Y.: Imaged document text retrieval without OCR. *IEEE Trans. Pattern Anal. Mach. Intell.* **24**(6), 838–844 (2002)

65. Vinciarelli, A.: Application of information retrieval techniques to single writer documents. *Pattern Recognit. Lett.* **26**(14–15), 2262–2271 (2005)
66. Vinciarelli, A.: Noisy text categorization. *IEEE Trans. Pattern Anal. Mach. Intell.* **27**(12), 1882–1295 (2005)
67. Western Regional Emergency Medical Services. Bureau of Emergency Medical Services. New York State (NYS) Department of Health (DoH). Prehospital Care Report v4
68. Xue, H., Govindaraju, V.: Stochastic models combining discrete symbols and continuous attributes—application in handwriting recognition. In: *Proceedings of 5th IAPR International Workshop on Document Analysis Systems*, pp. 70–81 (2002)
69. Xue, H., Govindaraju, V.: On the dependence of handwritten word recognizers on lexicons. *IEEE Trans. PAMI* **24**(12), 1553–1564 (2002)
70. Yates, B.R., Ribeiro-Neto, B.: *Modern Information Retrieval*. Addison-Wesley, Reading (1999)
71. Zimmermann, M., Mao, J.: Lexicon reduction using key characters in cursive handwritten words. *Pattern Recognit. Lett.* **20**, 1297–1304 (1999)
72. Zobel, J., Dart, P.: FInding approximate matches in large lexicons. *Softw. Pract. Experience* **25**(3), 331–345 (1995)