

# Automated Reading and Mining of Pre-Hospital Care Reports

Venu Govindaraju and Robert Milewski

Center of Excellence for Document Analysis and Recognition (CEDAR)

! Department of Computer Science & Engineering

CEDAR, UB Commons, Suite 202

University at Buffalo, The State University of New York

Amherst, NY 14228

*email:* govind@cedar.buffalo.edu

## Abstract

*The lack of use of high technology in the health care delivery system is especially apparent in the Emergency Medical Information Systems (EMIS) area. For example, in New York State all patients who enter the EMS are tracked through their pre-hospital care to the emergency room using the PreHospital Care Report (PCR). Our goal is to automate the collection of data from the PCR and enable efficient maintenance, and dissemination of information. The task involves the automatic extraction and transliteration of the handwritten text in the response boxes provided on the form. The objective is to produce for each form image, an ASCII transcription of the handwritten contents of the response boxes. These responses could be then used to populate a database. The database itself will then emerge as a valuable resource for enabling data mining and knowledge discovery for the entire medical community.*

## 1 Introduction

The lack of use of high technology in the health care delivery system is especially apparent in the Emergency Medical Information Systems (EMIS) area. For example, in New York State all patients who enter the EMS are tracked through their pre-hospital care to the emergency room using the PreHospital Care Report (PCR). The PCR (Figure 1) 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.

Figure 2 illustrates the flow of PCR data, describing a patients condition, from EMS to the WREMS server. The PCR as completed by the treating EMS personnel is inserted in the patient's record in the emergency room. The PCR data is keystroked into a database format, often several months after the emergency run was completed. A copy of the database is forwarded to the New York State Health department for macro-analysis. Completed analysis reports are returned to the regional offices after a few years!

Our goal is to automate the collection of data from the PCR and enable efficient maintenance, and dissemination of information. The task involves the automatic extraction and transliteration of the handwritten text in the response boxes provided on the form. The objective is to produce for each form image, an ASCII transcription of the handwritten contents of the response boxes. These responses could be then used to populate a database. The database itself will then emerge

Figure 1: Pre-Hospital Care Report Form

as a valuable resource for enabling data mining and knowledge discovery for the entire medical community.

The inputs are to include a NYS PCR form image (Figure 1) and a knowledge base pertaining to the medical field in general. This knowledge base will be generated from the *Medline* online database accessible from the National Library of Medicine website <http://www.nlm.nih.gov/>. The database can be eventually augmented with the patient's medical history if available from other sources.

There are 10 different fields on the PCR form, with responses recorded as marked checkboxes, or writing a few words. Consider the field for past medical history. A paramedic can check a few boxes pertaining to medical history and write down the medications being taken by the patient. A database can be populated based on this field, where data is collected on which medications are patients using for different ailments. These entries can be then correlated with other data gathered from the other fields on the form. Such a database would serve as a rich source of research data.

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) writing often extends beyond the response box provided and sometimes overlaps the printed text of the questionnaire, (iv) in some cases, the response wraps around into more than one line, (v) when a respondent is nearing the right margin of the box and has more to write, the writing gets squeezed in the space available, (vi) some input fields contain words that have been subsequently struck out, marks for insertion of additional text etc. Generally speaking, these subtasks can be divided into two broad categories. First, the task of segmentation of the image of the handwritten word or phrase of interest from its surroundings. Sec-

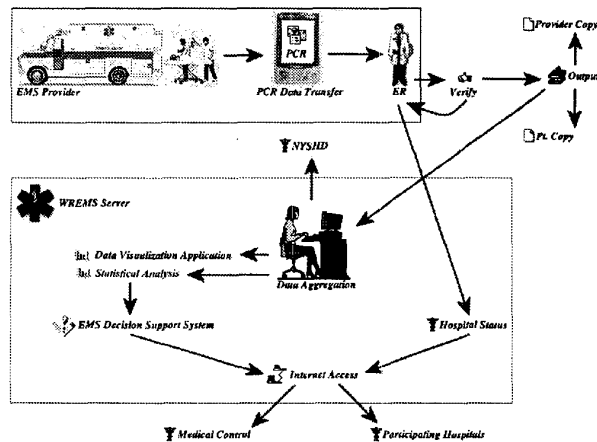


Figure 2: Project Schematic

ond, the task of recognition of the handwritten text. Each of these tasks is nontrivial and a research issue in its own right. The focus of this paper is on the latter task of recognition.

The success of Handwritten Word Recognition (HWR) technology primarily depends on how well the domain knowledge is captured in a lexicon of possible words associated with each word image. Needless to say, smaller the lexicon greater the chance of succeeding in the recognition task. Our aim is to constrain the size of the lexicon. To this purpose we are adopting several novel research techniques.

1. Use medical information from the *Medline* website. For example, if the form has the check box for *Diabetes* marked and the medication is handwritten. The possible medications form the following lexicon of 5 words selected from a list of about 9,000 prescription and over-the-counter drugs ([http://www.nlm.nih.gov/medline/drug information.html](http://www.nlm.nih.gov/medline/drug%20information.html)): Sulfonyleureas Metformin Thiazolidinediones Alpha-glucosidase Repaglinide
2. Use multimodal input. For example, the recognizer when reading the location of injury could look at the corresponding human body diagram present in the full PCR form to see which area has been marked. This could help in constraining the lexicon. If the head is indicated in the diagram, the lexicon for the corresponding handwritten words can be obtained from the *Visual Human Project* website.
3. Use commonly used terms in forms such as *chronic*, *distal*, etc. from the health information dictionaries website ([health-dictionaries.html](http://www.health-dictionaries.com)). Also use common abbreviations list appropriately, such as CHF - Congestive Heart Failure.
4. Medical words are often put together, cobbled from two or more building blocks. Among these building blocks are the prefixes. The lexicon should make use of these. For example, *dextro-*: from the Latin "dexter" meaning on the right side. The opposition of *dextro-* is *levo-* from the Latin "laevus" meaning on the left side. This could be again combined with the information obtained from the markings on the human body diagram.

Many of the handwriting recognition tools, such as underline removal, line separation, word separation, and lexicon generation are available and have been successfully used in other applications. We are currently in the process of applying the technology to the task of reading and understanding PCR forms. We are working on data sets provided by WREMs. We expect to have preliminary results in a few months.

## 2 Example

Our word recognizer, WMR [4] accepts two inputs: a HIPS image and an ASCII flatfile lexicon. This lexicon sample contains 295 sample words. All images shown are binarized 300 DPI.

### 2.1 Sample lexicon

abdomen	abdominal	accident	admitted	advanced	after
afternoon	alcohol	allergies	amount	ampule	anterior
arrival	arterial	as	at	attach	axillary
basic	be	before	bilateral	birth	black
bladder	blood	body	both	bowel	breath
by	cancer	capsule	carbon	carcinoma	care
centigrade	centimeter	central	cerebral	cerebrovascular	cesarean
change	check	chest	chief	chloride	chronic
complaining	complaint	compound	confinement	congestive	consciousness
continue	coronary	cubic	daily	date	day
dead	decrease	deep	degree	delirium	department
desired	dextrose	diagnosis	dilute	dioxide	discharge
discontinue	disease	dram	dressing	drop	drug
dry	due	each	ear	ears	electocardiogram
emergency	equal	estimated	every	everyday	exam
eye	eyes	fahrenheit	failure	female	fetal
fever	fire	first	fluid	foreign	four
fracture	full	gall	gallon	gases	gastric
gastro	general	genito	glucose	gonococcal	gonorrhea
grain	gram	greater	gunshot	heart	heartbeat
height	history	hour	hours	hypodermic	if
immediately	in	include	incomplete	incontinent	increase
infarct	infection	inflammatory	inguinal	inspiration	intake
intensive	internal	intestinal	intramuscular	intravenous	involuntary
keep	kilogram	kilometer	known	last	lateral
left	less	licensed	life	light	limits
liquids	liter	loving	lower	male	maximum
medical	medicines	menstrual	microgram	milliequivalent	milligram
minute	more	morning	motion	mouth	movement
myocardial	nasal	nauseau	necessary	negative	nervous
no	noon	normal	nothing	number	nurse
obstetrics	obstructive	of	ointment	old	on
once	open	operative	or	orally	order
origin	ounce	out	output	partum	patient
pediatric	pelvic	percentage	percussion	period	physical
positive	possible	post	posterior	pound	practical
preparation	prescription	pulmonary	pulse	pupils	quadrant
questionable	range	reactive	rectal	reflex	regular
rehabilitation	respiration	respiratory	responder	right	rule
saline	saturated	second	secondary	section	separated
service	short	signs	sleep	sodium	solution
splinal	staphylococcus	sterile	streptococcus	stroke	subcutaneous
sublingual	sugar	support	surgical	system	tablespoon
tablet	teaspoon	technician	telephone	temperature	tender
tendon	than	than loss	three	times	to
tremens	twice	undetermined	unit	upper	urinary
vaginal	vascular	vein	verbal	vital	volume
vomiting	wall	water	way	week	weight
whenever	white	with	within	without	wound
years					

<i>CHEST</i>	1	chest	0.608320	<i>MVED</i>	1	limits	0.533817
	2	cancer	0.434726		2	infarct	0.518858
	3	carbon	0.419975		3	motion	0.515454
	4	last	0.414112		4	lower	0.507230
	5	post	0.410222		5	medical	0.483264
	6	chief	0.406720		6	more	0.462608
	7	after	0.402420		7	times	0.459020
	8	cesarean	0.396076		8	increase	0.457098
	9	care	0.393596		9	twice	0.450752
	10	first	0.393272		10	male	0.448123

<i>VERBALLY</i>	1	infarct	0.485254	<i>PHYSICALLY</i>	1	vascular	0.456741
	2	regular	0.461089		2	kilometer	0.447316
	3	vascular	0.441192		3	milligram	0.446956
	4	allergies	0.411874		4	kilogram	0.424423
	5	reflex	0.410939		5	posterior	0.407979
	6	verbal	0.409466		6	axillary	0.406425
	7	operative	0.409145		7	allergies	0.400291
	8	greater	0.408897		8	emergency	0.396712
	9	intensive	0.405514		9	bilateral	0.391382
	10	necessary	0.401208		10	physical	0.384652

Figure 3: Top 10 results for some sample word images

### 3 Future work

Reducing the length of a lexicon greatly increases the performance of the WMR (Word Model Recognizer) algorithm, developed at CEDAR. However, the WMR algorithm, although proven effective in USPS postal recognition, does not have the luxury of such small lexicons in the medical world. There are at least two phases to modify WMR for successful lexicon reduction: (i) constructing the representation between words and morphemes (e.g. prefixes, suffixes and affixes) in the lexicon (ii) and pruning the lexicon for the more relevant data.

The derivation tree of words and morphemes can represent lexicon element relationships. The idea is to construct a tree for each element in the lexicon where the stem of a word could be the root node, and the child nodes could be a series of one or more morphemes [5]. For example, the words *moved*, *verbally*, and *physically*, in Figure 3, would have received higher rankings if a derivation tree contained the possible suffix endings *-ed* and *-ly*. The WMR algorithm could be applied twice: the first step, which has already been done, is to determine the scores for each lexicon entry with an image input. It is necessary to determine the most optimal top lexicon rankings for the second step (i.e. a pruning algorithm), based on PCR training data. The input, in the second step, would contain the top rankings with the addition and/or subtraction of morphemes. This new morphological lexicon may have great practical use in recognizing medical terms, which contain many repetitious morphemes [2].

Another data pruning method is to use the meanings of words and morphemes from foreign languages (e.g. Latin and Greek), to infer partial meanings of words [2]. Rather than using the partial meanings as a contextual vocabulary acquisition step, we could place words and/or morphemes into classes based on the probabilities that they are found together; common algorithms

are demonstrated in today's most popular search engines. Given an ordered set of words in a sentence, we could then minimize the lexicon for the current word, based on the frequency of the current and previous-words being adjacent [3]. Intuitively, more accurate results could be obtained by a weighted lexicon. For example, the word *fluid* in the phrase *amniotic fluid embolism* may commonly follow the word *amniotic* in medical documentation [1]. We could therefore have a derivational tree with *amniotic* as the root node, and the words *fluid* and *embolism* as one of several possible paths connected to the root node. If *fluid* is found in the lexicon, it should be given a higher weight. Furthermore, other words not in the same class as *fluid* could be either eliminated or be assigned a reduced weight. For example, the word *from*, may not appear frequently after *amniotic*, therefore, it is not in the same class as *fluid*, and its lexicon weight should be reduced.

For each element in the lexicon above, there also exists an abbreviation (e.g. *anterior* is abbreviated as *ant.*). In addition, various words in a derivational path have an abbreviation (e.g. *Chronic Obstructive Pulmonary Disease* is abbreviated as *COPD*). Some abbreviations may not necessarily be the prefix of a word (e.g. *check* is abbreviated *ck*). The abbreviations can also be used to assist WMR by minimizing the lexicon to the current words beginning and/or ending with combinatorial abbreviation substrings. Although the initials of a word are classified as affixes, the combination of morphemes (e.g. abbreviations may be the concatenation of a prefix and suffix) should be classified as an *n*-affix (take *n* to be a quantity of affixes). The design of an *n*-affix algorithm may also assist lexicon pruning. For example, if the word *abdominal* and its abbreviation *abdl.* are stored and related in a lexicon, we may want to run WMR on an input word. If the first three letters of the word is *abd* and the last letter is *l*, then we can increase the confidence that the input word is *abdominal* over *abdomen*.

## References

- [1] Beers, Mark H. et al. **The Merck Manual: Seventeenth Edition, Centennial Edition.** Merck Research Laboratories. USA. C1999. pp. 2063+.
- [2] Haubrich, William S. **Medical Meanings: A Glossary of Word Origins.** American College of Physicians. USA. C1997.
- [3] Jones, Daniel et al. **New Methods in Language Processing.** Centre for Computational Linguistics Press. London. C1997. pp. 18-27.
- [4] Kim, Gyengwhon and Govindaraju, Venu. **Lexicon Driven Handwritten Word Recognition for Real Time Applications.** *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(4), 1997, pp. 366-379.
- [5] Roche, Emmanuel and Clemenceau, David et al. **Finite State Language Processing.** *Finite-State Morphology: Inflections and Derivations in a Single Framework using Dictionaries and Rules.* MIT Press. USA. C1997. pp.68-96.