

An Automated System for Conversion of Clinical Notes into SNOMED Clinical Terminology

Jon Patrick, Yefeng Wang and Peter Budd

School of Information Technologies
University of Sydney
New South Wales 2006, Australia

{jonpat,ywang1,pbud3427}@it.usyd.edu.au

Abstract

The automatic conversion of free text into a medical ontology can allow computational access to important information currently locked within clinical notes and patient reports. This system introduces a new method for automatically identifying medical concepts from the SNOMED Clinical Terminology in free text in near real time. The system presented consists of 3 modules; an Augmented Lexicon, term compositor and negation detector. The Augmented Lexicon indexes the SNOMED-CT terms, the term compositor finds qualification relationships between concepts and the negation detector identifies negative concepts. The system delivers the services through a variety of interfaces including direct programming access and web-based access. It is currently in use in a hospital environment to capture patient data response with SNOMED-CT codes in real time at the point of care. No strict evaluation has been done on the system to date, however preliminary results indicate performance within acceptable time and accuracy limits.

Keywords: Concept Indexing, SNOMED-CT, Concept Mapping, Medical Terminology, Information Retrieval.

1 Introduction

A substantial amount of clinical data is locked away in a non-standardised form of clinical language which could be usefully mined to gain greater understanding of patient care and the progression of diseases if standardised. Unlike well written texts, such as scientific papers and formal medical reports, which generally conform to conventions of structure and readability, the clinical notes about patients written by a general practitioners, are in a less structured and often minimal grammatical form. However there is increasing interest in the automatic

extraction of the contents of clinical text and perform data analysis and retrieval on it.

In principle, clinical texts could be recorded in a coded form such as SNOMED CT (SNOMED International, 2006) or UMLS (Lindberg et al., 1993), in practice notes are written and stored in a free text representation. It is believed that the encoding of notes will provide better information for document retrieval and research into clinical practice (Brown and Sönksen, 2000). Many clinical information systems enforce standard semantics by mandating structured data entry. Transforming findings, diseases, medication procedures in clinical notes into structured, coded form is essential for clinician research and decision support system. Converting free text in clinical notes to terminology is a fundamental problem in many advanced medical information systems.

However, there is widely varying terminology for the same phenomena in the medical field. Different hospitals and clinics have their own information systems and each information system use their own ad hoc terminologies to code medical data. In order to manage the quality and cost of care cross these organizations data aggression is needed form different information systems. Lack of consistency in medical terminologies used in information system reduced the interoperability between information systems cross the entire health care organization. To ensure the interoperability, the use of standard terminologies for data representation is critical.

SNOMED CT is the most comprehensive reference medical terminology in the world. It consists of a set of concepts and relationships that provides a common reference point for comparison and aggregation of data about the entire health care process, recorded by multiple different individuals, systems, or institutions. Such reference terminology is useful for clinical data retrieval and analysis of data relating to the causes of disease, the treatment of patients, and the outcomes of the overall health care process. (Spackman et, al. 1997). SNOMED CT has currently been adopted by the Australia government to encode clinical disease and patient reports. We have been interested in a system to develop a standard terminology for reporting medical complaints so that their information is exchangeable and semantically consistent for other practitioners, and permit automatic extraction of the con-tents of clinical texts to compile statistics about diseases and their treatment. We are

developing a system that automatically transforms medical concepts in free text clinical notes into SNOMED CT terminology.

There are many researcher who have been working on mapping text to UMLS (The Unified Medical Language System), however, there is only a little work done on this topic for the SNOMED CT terminology. The present work proposes a natural language processing algorithm to automatically recognise medical terms in free text clinical notes and map them into SNOMED CT terminology. The algorithm is able to identify core medical terms in clinical notes in real-time as well as negation terms and qualifiers. In some circles SNOMED CT is termed an ontology, however this paper only covers its role as a terminology so we will use that descriptor only.

2 Related Work

2.1 Medical Term Mapping

There has been a large effort spent on automatic recognition of medical and biomedical concepts and mapping them to medical terminology. The Unified Medical Language System Meta-thesaurus (UMLS) is the world's largest medical knowledge source and it has been the focus of much research. Some prominent systems to map free text to UMLS include SAPHIRE (Hersh et al., 1995), MetaMap (Aronson, 2001), IndexFinder (Zou et al., 2003), and NIP (Huang et al., 2005). The SAPHIRE system automatically maps text to UMLS terms using a simple lexical approach. IndexFinder added syntactic and semantic filtering to improve performance on top of lexical mapping. These two systems are computationally fast and suitable for real-time processing. Most of the other researchers used advanced Natural Language Processing Techniques combined with lexical techniques. For example, NIP used sentence boundary detection, noun phrase identification and parsing. However, such sophisticated systems are computationally expensive and not suitable for mapping concepts in real time.

MetaMap has the capacity to code free text to a controlled terminology of UMLS. The MetaMap program uses a three step process started by parsing free-text into simple noun phrases using the Specialist minimal commitment parser. Then the phrase variants are generated and mapping candidates are generated by looking at the UMLS source vocabulary. Then a scoring mechanism is used to evaluate the fit of each term from the source vocabulary, to reduce the potential matches. The MetaMap program is used to detect UMLS concepts in e-mails to improve consumer health information retrieval (Brennan and Aronson, 2003).

The work done by (Hazelhurst et al., 2005) is on taking free text and mapping it into the classification system UMLS (Unified Medical Language System). The basic structure of the algorithm is to take each word in the input, generate all synonyms for those words and find the best combination of those words which matches a concept from the classification system. This research is not directly applicable to our work as it does not run in real

time, averaging 1 concept matched every 20 seconds or longer.

2.2 Negation Identification

In medical notes, the presence of a concept in a document does not mean that the finding or disease is related to that concept. The concept may refer to a finding that is absent or a suggested procedure that is not chosen. The term pneumonia and no pneumonia should be encoded differently and the system should not retrieve the latter in retrieving the patient records involving the condition of pneumonia.

Negation in medical domains is important, however, in most information retrieval systems negation terms are treated as stop words and are removed before any processing. UMLS is able to identify propositions or concepts but it does not incorporate explicit distinctions between positive and negative terms. Only a few works have reported negation identification (Mutalik et al., 2001; Chapman et al., 2001; Elkin et al., 2005).

Negation identification in natural languages is complex and has a long history. However, the language used in medical domains is more restricted and so negation is believed to be much more direct and straightforward. Mutalik et al (2001) demonstrated that negations in medical reports are simple in structure and syntactic methods are able to identify most occurrences. In their work, they used a lexical scanner with regular expressions and a parser that uses a restricted context-free grammar to identify pertinent negatives in discharge summaries. They identify the negation phrase first then identify the term being negated.

3 SNOMED CT Terminology

The Systematized Nomenclature of Medicine Clinical Terminology (SNOMED CT) is developed and maintained by College of American Pathologists. It is a comprehensive clinical reference terminology which contains more than 360,000 concepts and over 1 million relationships. The concepts in SNOMED CT are organised into a hierarchy and classified into 18 top categories, such as Clinical Finding, Procedure, Body Part, Qualifier etc. Each concept in SNOMED CT has at least three descriptions including 1 preferred term, 1 fully specified name and 1 or more synonyms. The synonyms provide rich information about the spelling variations of a term, and naming variants used in different countries. The concepts are connected by complex relationship networks that provide generalisation, specialisation and attribute relationships, for example, "focal pneumonia" is a specialisation of "pneumonia". It has been proposed for coding patient information in many countries.

4 Methods

In this section, we describe in detail the core algorithm TokenMatcher for mapping text to SNOMED CT terminology. We first present the token matcher system architecture. We then describe the token matching algorithm to map text to SNOMED CT concepts. Finally, we present the negation identification method.

The main goal was to take free text and map it into SNOMED CT concepts. At first, the text tokens are matched to SNOMED CT concept tokens, then we use more sophisticated natural language processing to break the input up into “chunks” based on their meaning and then run the standard algorithm over the chunks to allow for identification of negations and qualification. The token matching system consists of Pre-processing, SNOMED CT concept matching, qualification identification and negation identification

4.1 Pre-processing

4.1.1 Pre-processing of Clinical Notes

The clinical notes were processed at sentence level, because it is believed that the medical terms and negations do not often cross sentence boundaries. A maximum entropy model based sentence boundary detection algorithm (Reynar, and Ratnaparkhi, 1996) was implemented and trained on medical case report sentences. The sentence boundary detector reports an accuracy of 99.1% on test data. Since there is a large variation in vocabulary written in clinical notes compared to the vocabulary in terminology, normalisation of each term is necessary. The term normalisation process includes stemming, converting the term to lower case, tokenising the text into tokens and spelling variation generation (haemocyte vs. hemocyte). After term normalisation, the sentence then is tagged with POS tag and chunked into chunks using the GENIA tagger (Tsuruoka et al., 2005). We did not remove stop words because some stop words are important for negation identification.

4.1.2 Administrative Entity Recognition

Entities such as Date, Dosage and Duration are useful in clinical notes, which are called administration entities. A regular expression based named entity recognizer was built to identify administration units in the text, as well as quantities such as 5 kilogram. SNOMED CT defined a set of standard units used in clinical terminology in the subcategory of unit (258666001). We extracted all such units and integrated them into the recognizer. The identified quantities are then assigned the SNOMED CT codes according to their units. Table 1 shows the administration entity classes and examples.

Entity Class	Examples
Dosage	40 mg/day
Blood Pressure	105mm of Hg
Demography	69 year-old man
Duration	3 weeks
Date	May 1993
Quantity	55 mm

Table 1: Administration Entities and Examples

4.2 SNOMED CT Concept Matcher

The slowness of previous approaches is mainly due to the large size of the terminology. UMLS has over 2 million

concepts and SNOMED CT has over 1 million descriptions. Using naïve approach to search concepts in such big terminology requires millions of comparison. For example, one naïve approach is to check each description in the terminology to see if the description is a substring of the sentence. The complexity of this approach is $O(mn)$ where m is average length of the sentence and n is the size of the terminology. Since n is very big, this approach is very high time complexity. Thus, naïve approach is not suitable for a real time application. In this section, we develop an indexing algorithm that can map free text into SNOMED CT terminology at real-time or pseudo real-time (response within few seconds).

4.2.1 Augmented Lexicon

The Augmented Lexicon is a data structure developed by the researchers to keep track of the words that appear and which concepts contain them in the SNOMED CT terminology. The Augmented Lexicon is built from the Description table of SNOMED CT. In SNOMED CT each concept has at least three descriptions, preferred term, synonym term and fully specified name. The fully specified name has the top level hierarchy element appended which is removed. The description is then broken up into its atomic terms, i.e. the words that make up the description. For example, *Myocardial Infarction* has the atomic word *Myocardial* and *Infarction*. The UMLS Specialist Lexicon was used to normalise the term. The normalisation process includes removal of stop words, stemming, and spelling variation generation. For each atomic word, a list of the Description IDs that contain that word is stored as a linked list in the Augmented Lexicon. An additional field is stored alongside the augmented lexicon, called the "Atomic term count" to record the number of atomic terms that comprise each description. The table is used in determining the accuracy of a match by informing the number of tokens needed for a match. Figure 1 contains a graphical representation of the Augmented SNOMED CT Lexicon.

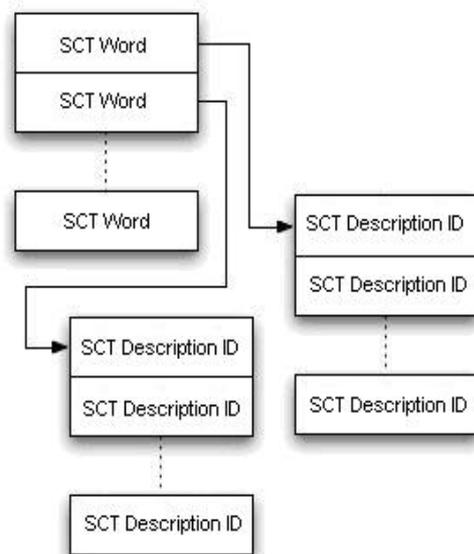


Figure 1: Augmented SNOMED CT Lexicon

4.2.2 Token Matching Algorithm

The token matching algorithm takes unstructured text and pre-processes it using the same techniques as are applied to the concepts when generating the augmented lexicon. It then attempts to find each SNOMED CT Description which is contained in the input sentence. For each word, the algorithm looks up the Augmented Lexicon, retrieving a list of the descriptions which contain the word. Figure 2 gives a graphical representation of the data structure used in the algorithm. The elements of the matrix are n-grams from the input sentence with the diagonal line sequence runs of two words. The remainder of the matrix is the cell to the left of it with the next word appended onto it. In this way the algorithm covers every possible sequence of sequential tokens.

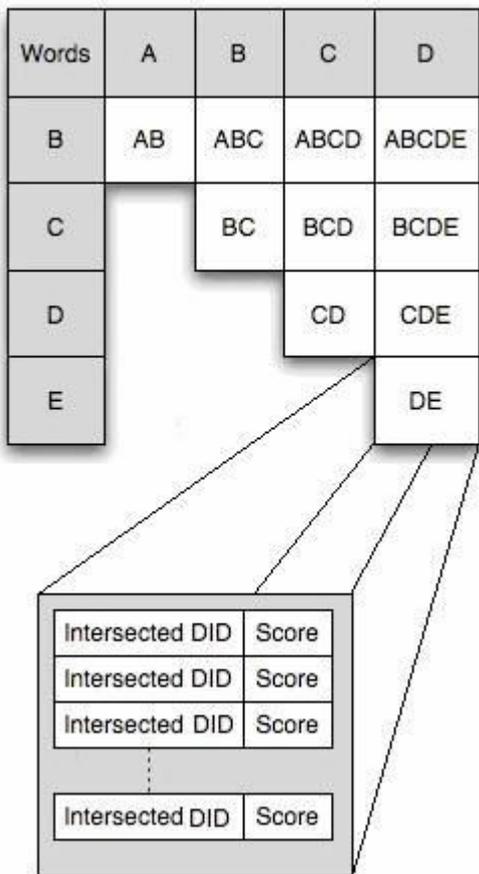


Figure 2: Matching Matrix

The data stored in each cell is a list of Description IDs (DID) that are in all the tokens that comprise the cell, i.e. the intersection of each set of DID of each word. The score is then calculated using the "atomic term count", which stores the number of tokens that make up that description. The score is the number of tokens in the current cell that have the DID in common divided by the number of tokens in the full description, i.e.:

$$\text{Score} = \frac{\text{\#of Tokens in Sequence}}{\text{\#of Tokens in Full Description}}$$

The algorithm itself is shown here in Figure 3 as pseudo-code. Step 1 is building the Matching Matrix. Step 2 is using the Matching Matrix to find the best combination of sequences that gives the highest score. This final score is dependant on the number of tokens used to make the match divided by the total number of tokens in the input stream, i.e.:

$$\text{Score} = \frac{\text{\#of Tokens used in all matches}}{\text{\#of Tokens in total input stream}}$$

```

STEP 1
for each word in list:
    add entry to the Matching Matrix
    for new column:
        Intersect new word with
        cell from matching table
Sort the matching array in descending order based off the
scores
for each row in the matrix:
    start at the right most cell

STEP 2
if the top score for the cell is 1.0
    add cell details to current best match list, update
    current match score.
    recursively call STEP 3 on cell (row=column+2,
    column=right)
else:
    move one column left to the next cell
    or
    the right-most cell of the next row if left cell
    empty
repeat STEP 3 until visited all cells

FINISH
return match with highest score.
    
```

Figure 3: Matching Algorithm Pseudo-code

An example of this process would be using the input sentence "Patient presents with bacterial pneumonia". The pre-processing step of the algorithm tokenises the sentence by the removal of stop words and stemming:

patient present bacterial pneumonia

The next step involves generating the matching matrix, and intersecting the description IDs for each concept in the matrix. "Patient" exists in the augmented lexicon as a single term concept so matches 100%. "bacterial" and

“pneumonia” intersect with the same concept “bacterial pneumonia” which is a 2 term concept, therefore also matching 100%. Pneumonia also exists in the augmented lexicon as a single term concept, therefore also matching 100%. No other intersections between the individual words are found in the augmented lexicon with a 100% matching.

There are 3 concepts matched in this stage of the algorithm, 2 of which overlap. There are 2 different possibilities which do not overlap; “patient”, “pneumonia” and “patient”, “bacterial pneumonia”. The former uses 2 of the 4 available words, scoring 50%, while the latter uses 3 out of the 4 available words scoring 75%. Therefore the latter is chosen as having the best match, returning:

Patient presents with Bacterial pneumonia

4.2.3 Abbreviations

Different sub-domains have different definitions of abbreviations. In medical domain, the abbreviations are highly ambiguous, as (Liu et al., 2002) show that 33% of abbreviations in UMLS are ambiguous. In different hospitals, they have their own convention of abbreviations, and the abbreviations used are not the same cross the sections in the same sub-domain. This creates difficulties for resolving the abbreviation problem. As we are processing clinical data in the RPAH (Royal Prince Alfred Hospital) ICU (Intensive Care Unit), we believe that the abbreviations used in their reports are restricted to a sub-domain and not that ambiguous. We use a list of abbreviations provided by the ICU department, and integrated them into the Augmented Lexicon. The abbreviations are manually mapped to SNOMED CT concepts by two experts in RPAH. The list consists of 1,254 abbreviations, 57 of them are ambiguous (4.5%). We decided not to disambiguate the abbreviations in the token matching, but return a list of all possible candidates and leave it for later stage to resolve the ambiguity.

4.3 Negation Identification

In our system, we aim to identify the negation phrases and the scope of the negation. Two kinds of negations are identified, the pre-coordinated SNOMED CT concepts and concepts that are explicitly asserted as negative by negation phrases. A pre-coordinated phrase is a term that exists in SNOMED CT terminology that represents a negative term, for example no headache.

SNOMED CT contains a set of pre-coordinated medical terms that represent negative meaning under the Clinical Finding Absent (373572006) category that indicate the absence of findings and diseases. However, SNOMED CT is not able to capture all negated medical terms. Moreover clinical notes have many negation forms other than absence, such as "denial of procedures". For a negative term that has a pre-coordinated mapping in SNOMED CT, we mark up this term using the SNOMED CT concept id (CID), for other negations, we identify the negation phrases and the SNOMED CT concepts that the

negation applies on. The following examples show the two different negations:

no headache (negative concept in SNOMED CT)

"absent of"

CID: 162298006

no headache (context-dependent category)

no evidence of neoplasm malignant

(Explicitly asserted negation)

negation phrase: "no evidence of"

CID: 363346000

malignant neoplastic disease (disorder)

Figure 4: Example of Negation

To identify explicitly asserted negation, we implemented a simple-rule based negation identifier similar to (Chapman et al, 2001; Elkin et al, 2005). At first the SNOMED CT concept id is assigned to each medical term, the negation phrases then are identified using a list of negation phrases in (Chapman et al, 2001). Then a rule base is applied on the negation phrase to check at its left and right contexts to see if any surrounding concepts have been negated. The algorithm is able to identify the negation of the form:

negation phrase ... (SNOMED CT phrase)*

(SNOMED CT phrase)* ... negation phrase

The contexts can up to 5 non-stopwords long, which allow identification of negation of coordination structure, for example in the following sentence segment:

... and pelvis **did not** reveal **retroperitoneal lymphadenopathy** or **mediastinal lymphadenopathy** ...

The negation stops propagation when another negation phrase is encountered or some phrases such as “*other than*” that usually stop negation propagation are encountered. For example in the following sentence segment:

She had **no significant surgical history** other than having undergone cholecystectomy ...

Whenever there is a overlapping between pre-coordinated negation and explicitly asserted negation, we identify the term as pre-coordinated negation. For example, no headache (162298006) will not be identified as *no + headache* (25064002).

4.4 Qualification

In medical terminology a term may contain an atomic concept or composition of multiple concepts, for example the term pain is an atomic concept and “back pain” represents composition two atomic concepts back and pain. Some composite concepts appear as single concepts in medical terminology, for example “back pain” is a single concept in SNOMED CT. Such concept is called

pre-coordinated concept. However, the medical terms can be composed by adding adjective modifiers to form new terms, for example, the add qualifiers to the concept “pain” can have “back pain”, “chronic back pain”, “chronic low back pain” etc. It is impossible to pre-coordinate combinations of all qualifiers into a terminology, because it will lead to term explosion. Term composition allows user to create new composite concepts using two or more single or composite concept. It is a solution to so called content completeness problem.

The SNOMED CT terminology has a subclass of terms called qualifier values. The qualifier values are used to qualify core concepts. The SNOMED CT defined qualifying relationship adds additional information about a concept without changing its meaning. In most cases, the qualifier is an adjective. There are also some nouns classified as qualifiers, such as fractions (278277004).

The purpose of the qualifier matching is to perform term composition. We separate the qualifiers apart from the Augmented Lexicon when performing concept matching, and build another lexicon that contains only qualifiers. Another reason for treating the qualifier differently is that the qualifier values always conflict with commonly used English words, for example, the unit qualifier day (258703001), side qualifier left (7771000), technique qualifier test (272394005). Such qualifiers cause noise when mapping text to concepts, and they should be refined by looking at their context.

The Concept Matchers runs at first to identify any SNOMED CT concepts and qualifiers. A search then is run to look at the qualifiers’ surroundings using the following rules to identify the scope of qualification. A concept can have multiple qualifiers to modify it.

(Qualifier / JJ|NN)* ... (Concept / NN)*

(Concept / NN)* ... (Qualifier / JJ|NN)*

The first rule aims identify left hand side qualifications, for example in the following sentence segment:

... She had **severe lethargy** and **intermittent right upper abdominal discomfort** ...

The second rule aims to identify right hand side qualification, for example:

... **autoimmune screening** were **normal** ...

If no concepts are found with in a context window, the qualifier then is not considered as a modifier to any medical concepts, thus removed to reduce noise.

5 Results and Discussion

The token matching algorithm has been implemented as a web-based service named TTSCCT (Text to SNOMED CT) that provides web interfaces for users to submit clinical notes and respond with SNOMED CT codes in real-time. The system is able to encode SNOMED CT concepts, qualifiers, negations, abbreviations as well as administration entities. It has been developed as the first step to the analysis and deep understanding of clinical notes and patient data. The system has been installed in RPAH (Royal Prince Alfred Hospital) ICU (Intensive Care Unit) aiming to collect bedside patient data. The web interface has been implemented in several clinical form templates the RPAH, allowing data to be captured as the doctors fill in these forms. A feedback form has been implemented allowing clinicians to submit comments, identify terms that are missed by the system and submit corrections to incorrectly labelled terms. Figure 5 shows the concepts that have been identified by the TTSCCT system and Figure 6 shows the responding SNOMED CT codes.

We are currently collecting test data and evaluating the accuracy of our method. We plan to collect patient reports and cooperate with the clinicians in the RPAH to identify correct mappings, missing mappings and incorrect mappings. Although the algorithm hasn’t been comprehensively evaluated on real data, we have collected some sample patient reports and a few feedback from some clinicians. Preliminary results demonstrate that the algorithm is able to capture most of the terms within acceptable accuracy and response time.

By observation, missing terms and partially identified terms are mainly due to the incompleteness in SNOMED CT. In the above example, the *atypical urothelial cells* is only partially matched, because neither *atypical urothelial cell* is present in SNOMED CT as a single term nor *urothelial* can be found as a qualifier in SNOMED CT. However the qualified term *moderate urothelial cell atypia* can be found in SNOMED CT. This raises the question of term composition and decomposition because the terms in the terminology have different levels of composition and the qualification can be written in a different order with morphological transformation (urothelia cell atypia vs. atypical urothelial cell). The qualifier ontology and term relationships must be addressed to make sure term composition is done in a reliable manner.

<p>No neoplasm malignant <small>negation</small> seen.</p> <p>Sections confirm CRANIOPHARYNGIOMA <small>concept</small> with small <small>qualifier</small> fragments <small>qualifier</small> of adjacent <small>qualifier</small> brain tissue <small>concept</small>.</p> <p>The slides <small>concept</small> show degenerate atypical <small>qualifier</small> urothelial cells <small>concept</small> occurring in sheets <small>qualifier</small> and singly with hyperchromatic <small>qualifier</small> enlarged <small>qualifier</small> irregular <small>qualifier</small> nuclei <small>concept</small>.</p>
--

Figure 5: A Sample Clinical Note

SNOMED CT Concept	SCT Concept ID	SCT Fully Specified Name	
CRANIOPHARYNGIOMA	40009002	Craniopharyngioma (morphologic abnormality)	
	189179009	Craniopharyngioma (disorder)	
Brain tissue	256865009	Brain tissue (substance)	
Cells	4421005	Cell structure (cell structure)	
	362837007	Entire cell (cell)	
hyperchromatic	9767008	Hyperchromatism (morphologic abnormality)	

Qualifiers	SCT Concept ID	SCT Fully Specified Name	Scope of Qualification
Small	255507004	Small (qualifier value)	
	263796003	Lesser (qualifier value)	
Fragments	29140007	Fragment of (qualifier value)	
Adjacent	18769003	Juxta-posed (qualifier value)	brain tissue
Atypical	112231000	Atypical (qualifier value)	Cells
Sheets	255292000	Sheets (qualifier value)	
Enlarged	260376009	Enlarged (qualifier value)	
Irregular	49608001	Irregular (qualifier value)	
Fragments	29140007	Fragment of (qualifier value)	Tissue

Negation	Negation Phrase	Negative Term
No neoplasm malignant	No	neoplasm malignant (86049000)

Figure 6: Concepts, Qualifiers and Negations Identified From the Sample Note

Any terminology that can perform term composition or term decomposition will have to encounter multiple ways to express the same concepts (Spackman et al., 1997). In the above example, atypical urothelial cell can be composed of the qualifier *atypical*, concept *urothelium* and concept *cell*, or it can be composed of the qualifier *atypical* and the composite term *urothelial cell*. Since the redundancy is inevitable, there is a need to use description logic to automatically determine semantic equivalent concepts, to avoid same concept is represented by multiple concepts, for example *third degree burn, elbow* and *third degree burn of elbow*.

Restricting the concept mapping to noun phrase chunks can rule out many false positives and also increase the speed of processing, however many pre-coordinated terms and qualifications cross noun phrase boundaries, for example the term “*Third degree burn of elbow (87559001)*” will be broken into two terms “*Third degree burn (403192003)*” and “*elbow (76248009)*” and their relationship not preserved.

A limitation of this system is that the identification of SNOMED CT concepts is based on exact lexical matches in SNOMED CT descriptions. The system is likely fail to match long terms that have lexical variations or synonyms at token level. Such strict exact matching will definitely increase precision, but it also will decrease the coverage. A modification to the token matching algorithm that using a string similarity functions to match the concepts is a possible solution to this problem, but it may also introduce more incorrectly identified concepts.

6 Conclusions

In conclusion, we propose a system to translate free text clinical notes into medical terminology and perform simple term composition. We have implemented the system as a web-service system. The algorithm uses an augmented lexicon to index concept descriptors in SNOMED CT, which allow a much faster mapping of longest concepts in system than naïve searching approach. A qualification identifier and negation identifier have been implemented for recognising composite terms and negative concepts, which can then create more effective information retrieval and information extraction. The system is yet to be fully evaluated, nevertheless the test on sample data shows it is already meeting expectations. In the future, we will perform comprehensive evaluation for the algorithm on real clinical data, and compare the system with some well known term indexing algorithms.

7 References

- Aronson, A. R. (2001). *Effective mapping of biomedical text to the UMLS Metathesaurus: the MetaMap program*. Proc AMIA Symp 17: 21.
- Brennan, P. F. and A. R. Aronson (2003). *Towards linking patients and clinical information: detecting UMLS concepts in e-mail*. Journal of Biomedical Informatics 36(4/5): 334-341.
- Brown, P. J. B. and P. Sönksen (2000). *Evaluation of the Quality of Information Retrieval of Clinical Findings*

- from a Computerized Patient Database Using a Semantic Terminological Model. *Journal of the American Medical Informatics Association* 7: 392-403.
- Chapman, W. W., W. Bridewell, et al. (2001). *A Simple Algorithm for Identifying Negated Findings and Diseases in Discharge Summaries*. *Journal of Biomedical Informatics* 34(5): 301-310.
- Chute, C. G. and P. L. Elkin (1997). *A clinically derived terminology: qualification to reduction*. *Proc AMIA Annu Fall Symp* 570: 4.
- Elkin, P. L., S. H. Brown, et al. (2005). *A controlled trial of automated classification of negation from clinical notes*. *BMC Medical Informatics and Decision Making* 2005(5): 13.
- Friedman, C., P. O. Alderson, et al. (1994). *A general natural-language text processor for clinical radiology*. *Journal of the American Medical Informatics Association* 1(2): 161-174.
- Friedman, C., L. Shagina, et al. (2004). *Automated Encoding of Clinical Documents Based on Natural Language Processing*. *J Am Med Inform Assoc* 11: 392-402.
- Hazlehurst, B., H. R. Frost, et al. (2005). *MediClass: A System for Detecting and Classifying Encounter-based Clinical Events in Any Electronic Medical Record*, American Medical Informatics Association.
- Hersh, W. R. and D. Hickam (1995). *Information retrieval in medicine: The SAPHIRE experience*. *Journal of the American Society for Information Science* 46(10): 743-747.
- Huang, Y., H. J. Lowe, et al. (2005). *Improved Identification of Noun Phrases in Clinical Radiology Reports Using a High-Performance Statistical Natural Language Parser Augmented with the UMLS Specialist Lexicon*, American Medical Informatics Association.
- Lindberg, D. A., B. L. Humphreys, et al. (1993). *The Unified Medical Language System*. *Methods Inf Med* 32(4): 281-91.
- Liu H., Johnson S. B. and Friedman C. (2002), *Automatic Resolution of Ambiguous Terms Based on Machine Learning and Conceptual Relations in the UMLS*. *Journal of the American Medical Informatics Association*, Vol. 9, No. 6, Pages 621-636.
- Mutalik, P. G., A. Deshpande, et al. (2001). *Use of General-purpose Negation Detection to Augment Concept Indexing of Medical Documents A Quantitative Study Using the UMLS*, *Journal of the American Medical Informatics Association* 2001.
- Nadkarni, P., R. Chen, et al. (2001). *UMLS Concept Indexing for Production Databases*. *Journal of the American Medical Informatics Association* 8: 80-91.
- Reynar, J. C. and A. Ratnaparkhi (1997). *A maximum entropy approach to identifying sentence boundaries*. *Proceedings of the Fifth Conference on Applied Natural Language Processing*: 16-19.
- SNOMED International. (2006). *SNOMED International: The Systematized Nomenclature of Medicine*. <http://www.snomed.org>. Last accessed 08-09-2006.
- Spackman, K. A., K. E. Campbell, et al. (1997). "SNOMED RT: a reference terminology for health care." *Proc AMIA Annu Fall Symp* 640(4). Stearns, M. Q., C. Price, et al. (2001). *SNOMED clinical terms: overview of the development process and project status*. *Proc AMIA Symp* 662(6). Association for Computing Machinery. 1983. *Computing Reviews*, 24(11):503-512.
- Tsuruoka, Y., Y. Tateishi, et al. *Developing a robust part-of-speech tagger for biomedical text*. *Lecture notes in computer science*: 382-392.
- Zou, Q., W. W. Chu, et al. (2003). *IndexFinder: A Method of Extracting Key Concepts from Clinical Texts for Indexing*. *Proc AMIA Symp* 763: 7.