

# Automatic annotation of French medical narratives with SNOMED CT concepts

Gaudet-Blavignac Christophe<sup>a</sup>, Foufi Vasiliki<sup>a</sup>, Wehrli Eric<sup>b</sup>, Lovis Christian<sup>a</sup>

<sup>a</sup> Division of Medical Information Sciences, Geneva University Hospitals and University of Geneva, Switzerland

<sup>b</sup> Laboratoire d'Analyse et de Technologie du Langage, University of Geneva, Switzerland

## Summary

**RESEARCH QUESTION:** Medical data are multimodal. In particular, they are composed of both structured data and narrative data (free text). Narrative data are a type of unstructured data that, although containing valuable semantic and conceptual information, is rarely reused.

**METHODS:** In order to assure interoperability of medical data, automatic annotation of free text with SNOMED CT concepts via Natural Language Processing (NLP) tools is proposed. This task is performed using a hybrid multilingual syntactic parser.

**RESULTS:** A preliminary evaluation was completed on a small corpus of five discharge summaries chosen randomly. The corpus was at first manually annotated with SCT concepts by one expert. The same corpus was then processed by the parser and 421 medical terms were automatically annotated. A manual comparison of the two outputs was performed in order to evaluate the system. Precision of 0.7173 and recall of 0.517 were achieved.

**CONCLUSION:** The preliminary annotation results are encouraging and confirm that semantic enrichment of patient-related narratives can be accomplished by hybrid NLP systems, heavily based on syntax and lexicosemantic resources.

**Keywords:** interoperability, narrative data, SNOMED CT, NLP

## Introduction

Medical data are composed of both structured and unstructured data. Narratives are a type of unstructured data that contains crucial semantic information but is not readily usable by computers. Moreover, narratives constitute a challenge for interoperability between healthcare systems, hospitals and departments [1]. SNOMED CT (henceforth SCT), created in 2002, constitutes a terminology organised as a directed graph with concepts as nodes and relationships as edges. Currently, SCT contains more than 350,000 concepts. Property rights and developments are held by SNOMED International (London, UK). The goal of SNOMED International is to develop SCT and to ensure that it becomes the “most comprehensive and precise com-

mon global language for health terms in the world” [2]. Since SCT supports post-coordination, i.e., a formal grammar that can associate existing concepts, qualifiers and predicates, it has properties similar to those of a natural language. In this paper, a method for automatic annotation of French medical narratives with SCT codes is proposed.

Processing medical data with various terminologies, and recently SCT, has been a research focus of other studies as well [3, 4]. Those studies have pursued two kinds of goals. The first goal was the classification of documents, such as pathology or radiology reports [5, 6] in categories related to the disease mentioned in the text. The second one was a more general information retrieval task that aimed at extracting codes or annotating free text with concepts [7]. Commonly used terminologies are the International Classification of Diseases (ICD)-10 [8], the Unified Medical Language System (UMLS) [9] or SCT [10]. In some cases, preliminary work involves creating a subset of those terminologies in relation to a specific goal. It is the case with the UMLS because its metathesaurus contains more than 100 terminologies, classifications and thesauri.

The methods used to annotate or classify free-text documents vary. Rule-based methods need to be manually or semi-manually developed but require no training corpus and can produce very satisfying results when combined in a pipeline [11, 12]. On the other hand, machine learning and statistical methods, such as Naïve Bayes or Support Vector Machine, do not require the manual creation of rules. However, access to large gold standard corpora used as training sets is essential [13]. Hybrid NLP systems integrating both statistical and linguistic approaches have also been proven to be very efficient at NLP tasks targeting medical language [13]. The work presented in this article differs in several ways from the studies previously mentioned. First, the language of the free-text documents used in those references is mostly English. Working with another language requires translation of the terms and adaptation of the rules to the specificities of the target language. Second and most important, the absence of syntactic-semantic parsing of the text to detect terms in various morphological or syntactic structures makes the method presented in this paper innovative. Our system performs analysis of French medical texts on the morphological, syntactic and semantic level

## Correspondence:

Christophe Gaudet-Blavignac, Division of Medical Information Sciences, University Hospitals of Geneva, University of Geneva, Rue Gabrielle Perret-Gentil 4, CH-1211 Genève, Christophe.Gaudet-Blavignac[at]hcuge.ch

and annotates the recognised terms with SCT concepts simultaneously.

## Method

In this research, SCT was approached as a natural language. Automatic annotation of narratives with SCT concepts therefore required the processing of texts using NLP tools.

### Tool

The tool used for this goal was the hybrid multilingual syntactic parser Fips [14]. It relies on generative grammar concepts and is made of a generic parsing module that can be refined to suit the specific needs of a particular language or sublanguage. The lexicon is one of the key components of the parser. It contains detailed morphosyntactic and semantic information, selectional properties, valency information and syntactic-semantic features that influence the syntactic analysis. To achieve automatic annotation of medical narratives, modifications were needed to correctly process the specificities of the French medical language such as abbreviations or technical terms.

### Creation of electronic dictionaries

Specific lexicons have been developed and incorporated in the parser:

A French medical language dictionary was created by extracting simple words and collocations from a corpus of discussions of 11,000 discharge summaries from the internal medicine division of the University Hospitals of Geneva during 2012 to 2014. In its current version, the lexicon comprises 4454 simple words and 5640 collocations (groups of words) manually processed.

A SCT dictionary. To perform automatic annotation of French narratives with SCT codes, the SCT terminology was added as a new language in the parser. 173,067 SCT concepts and their equivalent code were entered in this dictionary.

A bilingual French-SCT dictionary. In the aim of automatic annotation, the target language (SCT) must be linked to the source language (French medical language) in a bilingual dictionary. In the current version of the system, 5842 medical terms have been mapped to SCT concepts.

### Automatic annotation

In this research, the automatic annotation procedure consisted of parsing the initial text and recognising medical terms. Then, the system looked up the dictionaries (both monolingual and bilingual) and proceeded to the SCT code attribution. Terms in medical terminologies can be affected by syntagmatic and paradigmatic variation to different degrees, or may be too precise or complex to actually be used in electronic health records [15]. By providing syntactic analysis and a proper recognition of collocations, the parser can detect concepts regardless of the specific morphological or syntactic form under which they appear in the text. Table 1 shows an example of a sentence annotated with SCT concepts:

We can observe that the system is capable of recognising structures in various forms, e.g., *iv*, the abbreviated form of *intraveineux* “intravenous”. It can also identify complex

structures even if their constituents do not follow the canonical order and are found in different positions, such as the verbal collocation *poursuivre un traitement* “continue a treatment”, *les traitements ... sont poursuivis* “the treatments ... are being continued”.

## Results

### Automatic annotation

Automatic annotation using the syntactic parser was performed on a corpus of 11,000 discharge summaries. Table 2 displays the results of the automatic annotation procedure.

### Preliminary evaluation

A preliminary evaluation was completed on a small corpus of five randomly chosen discharge summaries (1820 words) written by four different clinicians. The corpus was first de-identified (protected health information [PHI] was removed) and then manually annotated with SCT concepts by one expert. The concepts used for the annotation were selected from the set of codes that are incorporated in the parser’s SCT dictionary. The same corpus was processed by the parser and 421 medical terms were automatically annotated. Then, a manual comparison of the two outputs was performed in order to evaluate the system. The performance of the system is very encouraging since precision of 0.7173 and recall of 0.517 were achieved. However, an evaluation on a bigger corpus would allow a more precise measurement of the efficiency of the method.

## Discussion

### Annotation procedure

The rules used to annotate a narrative with SCT concepts are subject to debate. Since the terminology is structured as a graph with a treelike disposition, there are various levels of granularity for each concept. For instance, *douleur abdominale* “abdominal pain” could be annotated with a unique SCT code (21522001, cf. table 1) or could be annotated with several more specific concepts (22253000 | *douleur* “pain” |, 277112006 | *abdominal* “abdominal” |). At the current stage of the research, the annotation was the concept that corresponded to the largest text structure.

Table 1: Example of SCT annotation.

Initial phrase	SCT Annotation
En raison des douleurs abdominales, un traitement de morphine iv est débuté et les traitements habituels du patient sont poursuivis.	{21522001   douleur abdominale  }, {373529000   morphine  }, {255560000   intraveineux  }, {40451002   habituel  }, {116154003   patient  }, {266714009   poursuivre le traitement  }

Table 2: Automatic annotation of a corpus of 11,000 discharge summaries.

Words	4481,191
Annotated terms	892,787
Unique SCT concepts	7569
Annotated terms per sentence	4.17

### Limitations and future work

Medical documents contain sensitive information and as a consequence access to corpora, and in particular annotated corpora, is a well-known challenge in this field. This is especially true for languages other than English. The size of the evaluation corpus is one of the major limitations of this paper. In addition, evaluation of medical free-text annotation must be performed in a specific setting to ensure that the results are reliable. The manual annotation task, in particular, should be performed by at least two annotators not directly involved in the development of the automatic annotation tool in order to avoid bias. Having more than one annotator is important to compute the inter-annotator agreement and set an upper bound on the annotation task. The annotation of French narratives with SCT concepts is a first step toward the ultimate goal, which is the complete representation of patient-related narratives in a formal language. The next step in this research will be the processing of post-coordinated concepts according to the SCT compositional grammar. Post-coordination will enable the storage of the full information contained in the text into SCT post-coordinated sentences.

### Conclusion

In this paper, a method to annotate French medical texts with SCT concepts is proposed. This method relies on a syntactic-semantic parser specifically modified to meet the needs of this task. Lexicosemantic resources (monolingual and bilingual dictionaries as well as grammar rules) were constructed taking into consideration the specificities of the French medical language. A preliminary evaluation has shown encouraging results with a precision of 0.7173, a recall of 0.5171 and an F-score of 0.6009. Further research is needed to produce post-coordinated structures and full representation of medical narratives into SCT.

### Financial disclosure

This research has been financed by the “Réseau Thématique Langage & Communication” from the University of Geneva.

### Potential competing interests

No potential conflict of interest relevant to this article was reported.

### References

- 1 Meystre SM, Lovis C, Bürkle T, Tognola G, Budrionis A, Lehmann CU. Clinical Data Reuse or Secondary Use: Current Status and Potential Future Progress. *Yearb Med Inform.* 2017;26(1):38–52. doi: <http://dx.doi.org/10.15265/Y-2017-007>. PubMed.
- 2 Support: SNOMED International. [Online]. Available: <https://ihtsdo.freshdesk.com/support/home>. [Accessed: 20-Mar-2017].
- 3 Patrick J, Wang Y, Budd P. An Automated System for Conversion of Clinical Notes into SNOMED Clinical Terminology. Proceedings of the fifth Australasian symposium on ACSW frontiers; 2007.
- 4 Ruch P, Gobeill J, Lovis C, Geissbühler A. Automatic medical encoding with SNOMED categories. *BMC Med Inform Decis Mak.* 2008;8(Suppl 1):S6. doi: <http://dx.doi.org/10.1186/1472-6947-8-S1-S6>. PubMed.
- 5 Zuccon G, Waghlikar AS, Nguyen AN, Butt L, Chu K, Martin S, et al. Automatic Classification of Free-Text Radiology Reports to Identify Limb Fractures using Machine Learning and the SNOMED CT Ontology. *AMIA Jt Summits Transl Sci Proc.* 2013;2013:300–4. PubMed.
- 6 Nguyen A, Moore J, Zuccon G, Lawley M, Colquist S. Classification of pathology reports for cancer registry notifications. *Stud Health Technol Inform.* 2012;178:150–6. PubMed.
- 7 Torii M, Waghlikar K, Liu H. Using machine learning for concept extraction on clinical documents from multiple data sources. *J Am Med Inform Assoc.* 2011;18(5):580–7. doi: <http://dx.doi.org/10.1136/amia-jnl-2011-000155>. PubMed.
- 8 Koopman B, Zuccon G, Nguyen A, Bergheim A, Grayson N. Automatic ICD-10 classification of cancers from free-text death certificates. *Int J Med Inform.* 2015;84(11):956–65. doi: <http://dx.doi.org/10.1016/j.ijmedinf.2015.08.004>. PubMed.
- 9 Riedl B, Than N, Hogarth M. Using the UMLS and Simple Statistical Methods to Semantically Categorize Causes of Death on Death Certificates. *AMIA Annu Symp Proc.* 2010;2010:677–81. PubMed.
- 10 Ruch P, Gobeill J, Lovis C, Geissbühler A. Automatic medical encoding with SNOMED categories. *BMC Med Inform Decis Mak.* 2008;8(1, Suppl 1):S6. doi: <http://dx.doi.org/10.1186/1472-6947-8-S1-S6>. PubMed.
- 11 De Meyere D, Klein T, François T, Debongnie JC, Radulescu C, Fairon C, et al. Automatic annotation of medical reports using SNOMED-CT: a flexible approach based on medical knowledge databases. Proceedings of the 7th Language & Technology Conference; 2015 Nov 27–29; Poznań, Poland. Poznań, Poland: Fundacja Uniwersytetu im. A. Mickiewicza, 2015.
- 12 Aronson AR. Effective mapping of biomedical text to the UMLS Metathesaurus: the MetaMap program. *Proc AMIA Symp.* 2001:17–21. PubMed.
- 13 Kate RJ. Towards Converting Clinical Phrases into SNOMED CT Expressions. *Biomed Inform Insights.* 2013;6(Suppl 1):29–37. PubMed.
- 14 Wehrli E, Nerima L. The fips multilingual parser. In: Gala N, Rapp R, Bel-Enguix G, editors. *Language Production, Cognition, and the Lexicon*. Springer; 2015. p. 473–90.
- 15 Hansart C, De Meyere D, Watrin P, Bittar A, Fairon C. CENTAL at SemEval-2016 Task 12: a linguistically fed CRF model for medical and temporal information extraction. Proceedings of the 10th Int. Workshop Semantic Eval. *SemEval*; 2016 Jun 16–17; San Diego, California. Association for Computational Linguistics; 2016.