

# How to de-identify a large clinical corpus in 10 days

C. Grouin, PhD<sup>1</sup>, L. Deléger, PhD<sup>1</sup>, J.B. Escudié, MD<sup>2</sup>, G. Groisy, MD<sup>2</sup>,  
A.S. Jannot, MD, PhD<sup>2,3</sup>, B. Rance, PhD<sup>2,3</sup>, X. Tannier, PhD<sup>1,4</sup>, A. Névéol, PhD<sup>1</sup>

<sup>1</sup>CNRS UPR 3251 LIMSI, Orsay, France; <sup>2</sup>AP-HP, University Hospital HEGP, Biomedical Informatics and Public Health Department, Paris, France; <sup>3</sup>INSERM U1138, Université Paris Descartes, Sorbonne Paris Cité, Faculté de médecine, Paris; <sup>4</sup>Université Paris Sud, Orsay, France

## Abstract

*To de-identify a large corpus of clinical documents in French supplied by two different health care institutions, we apply a protocol built from previous work. We show that the protocol can be installed and executed by outside collaborators with little guidance from the authors. The automatic de-identification method used reaches 0.94 F-measure and human validation requires about 1 minute per document.*

## Introduction

Clinical corpora are useful for scientists to develop natural language processing methods for the clinical narrative. In order to ensure the robustness of those methods, access to real data is a critical point. As clinical records contain personal health information, de-identification tools have been designed to help protecting privacy.<sup>1</sup>

## Methods

Herein, we apply the protocol we designed to rapidly de-identify a new set of clinical records:<sup>2</sup> new documents to be de-identified are automatically pre-annotated with a statistical model built on similar resources. For new documents from the same hospital as our previous study<sup>2</sup> (corpus 1), we directly apply a model built with 100 training documents. For documents coming from a distinct hospital (corpus 2) we iteratively build a model as follows:

1. 20 documents were de-identified automatically, using a rule-based method<sup>3</sup> for the first iteration and a CRF (Conditional Random Fields) model for subsequent iterations, and manually validated by scientists and physicians
2. a statistical (CRF) model was built using all the manually validated de-identified documents available from all iterations. The latest model is then used for the next iteration of step 1. The features used to build the CRF model include surface features (*case, punctuation, digit, token length, etc.*), deep features (*part-of-speech, lexical look-up, etc.*), and external features (*position of the token in the record, cluster id*).

For the initial seed corpus of 20 documents, six different types of clinical notes were chosen to ensure robustness: 3 medical certificates, 3 consultation reports, 3 exam certificates, 3 hospitalization reports, 5 follow-up care letters, and 3 staff consultation reports.

## Results and Conclusion

Corpus 1: The automatic de-identification method reached 0.94 F-measure and human validation required less than 1 minute per document for 800 documents. Corpus 2: About 1 day was needed to install the tools to use in the de-identification protocol. The automatic de-identification method reached 0.92 F-measure and human validation required as little as 1 minute per document for 100 documents. Validation time decreased with annotator experience and improvement in the performance of the de-identification method (i.e. after more training documents become available as we progress through more iterations). The protocol can be executed by outside collaborators with no prior experience of de-identification. It provides adequate support for the de-identification of a large corpus requiring little time and guidance.

## References

- <sup>1</sup>Meystre SM, Friedlin FJ, South BR, Shen S, Samore MH. Automatic de-identification of textual documents in the electronic health record: a review of recent research. *BMC Med Res Methodol.* 2010;10.
- <sup>2</sup>Grouin C, Névéol A. De-identification of clinical notes in French: towards a protocol for reference corpus development. *J Biomed Inform.* 2014. In press.
- <sup>3</sup>Grouin C, Zweigenbaum P. Automatic de-identification of French clinical records: comparison of rule-based and machine-learning approaches. *Stud Health Technol Inform.* 2013;192(Part 1):476–80.