



**HAL**  
open science

## De-Identification of French Unstructured Clinical Notes for Machine Learning Tasks

Yakini Tchouka, Jean-François Couchot, Maxime Coulmeau, David Laiymani,  
Azzedine Rahmani

► **To cite this version:**

Yakini Tchouka, Jean-François Couchot, Maxime Coulmeau, David Laiymani, Azzedine Rahmani.  
De-Identification of French Unstructured Clinical Notes for Machine Learning Tasks. 2022. hal-  
03720808

**HAL Id: hal-03720808**

**<https://hal.science/hal-03720808>**

Preprint submitted on 21 Jul 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# De-Identification of French Unstructured Clinical Notes for Machine Learning Tasks

Yakini Tchouka<sup>1</sup>, Jean-François Couchot<sup>1</sup>, Maxime Coulmeau<sup>1,3</sup>, David Laiymani<sup>1</sup>,  
Philippe Selles<sup>2</sup>, and Azzedine Rahmani<sup>2</sup>

<sup>1</sup>Femto-ST Institute, Univ. Bourg. Franche-Comté, CNRS, France

<sup>2</sup>Nord Franche-Comté Hospital, Trevenans, France

<sup>3</sup>University of Technology of Belfort-Montbéliard, France

July 13, 2022

## Abstract

Unstructured textual data are at the heart of health systems: liaison letters between doctors, operating reports, coding of procedures according to the ICD-10 standard, etc. The details included in these documents make it possible to get to know the patient better, to better manage him or her, to better study the pathologies, to accurately remunerate the associated medical acts. . . All this seems to be (at least partially) within reach of today by artificial intelligence techniques. However, for obvious reasons of privacy protection, the designers of these AIs do not have the legal right to access these documents as long as they contain identifying data. De-identifying these documents, i.e. detecting and deleting all identifying information present in them, is a legally necessary step for sharing this data between two complementary worlds. Over the last decade, several proposals have been made to de-identify documents, mainly in English. While the detection scores are often high, the substitution methods are often not very robust to attack. In French, very few methods are based on arbitrary detection and/or substitution rules. In this paper, we propose a new comprehensive de-identification method dedicated to French-language medical documents. Both the approach for the detection of identifying elements (based on deep learning) and their substitution (based on differential privacy) are based on the most proven existing approaches. The result is an approach that effectively protects the privacy of the patients at the heart of these medical documents. The whole approach has been evaluated on a French language medical dataset of a French public hospital and the results are very encouraging.

## 1 Introduction

Omnipresent in the fields of finance, transport, information - the list is necessarily incomplete - Artificial Intelligence (AI) governs our lives. The field of health is no exception, and even on unstructured data (e.g. textual), which is reputed to be the most difficult to manipulate. Problems that were inaccessible a short time ago are becoming soluble, such as the search for similar patient files, ICD-10 classification [AND<sup>+</sup>19, NRG<sup>+</sup>18], hospital readmission prediction [HMS<sup>+</sup>10], patient clustering [HSQ<sup>+</sup>19] . . . .

However, these processes can be carried out by computer specialists in deep learning from the data science and big data professions, but not yet by doctors. It therefore appears necessary to "share" the data between medical actors and data science specialists. Because of the critical nature of the data involved, this sharing implies

a de-identification process which can only take place within a legal framework that governs the actors in the medical world. Institutional official rules are for instance U.S. Health Insurance Portability and Accountability Act (HIPAA)[CM18] and the European General Data Protection Regulation (GDPR) [gdp16]. Does a technical implementation of this de-identification that respects these rules is possible? Does it allows documents to be sufficiently rich and not too degraded to be used by A.I. algorithms for automatic ICD-10 classification for example?

This work focuses on the ability to share textual medical documents, often written by doctors and can take the form of operating reports, clinical notes or biological examination results. To facilitate privacy protection, de-identification methods [Lev03, Laf, HHD<sup>+</sup>20, BAC<sup>+</sup>21, DLUS16] have been proposed as a process to remove or mask any type of Protected Health Information (PHI) from a patient, so that it becomes difficult to link an individual to its data. The type of information that constitutes PHI is defined in part by the privacy laws of the relevant jurisdiction. For example, HIPAA regulations define 18 categories of PHI including names, geographic locations and telephone numbers. In Europe, since the GDPR does not provide such PHI definitions, most research uses the HIPAA definitions.

A de-identification process can thus be summarized as an algorithm with two main phases. The former is detecting all compromising information (names, addresses, ages, dates, numbers) or equivalently as a Named Entity Recognition (NER) phase. The latter consists in replacing these elements by simple substitute data or more complex context-specific labels, classically denoted as Named Entity Substitution (NES) phase. The difficulty here is that the de-identification process (both NER and NES) must be balanced between too much removing (limiting the data usefulness for downstream tasks such as ICD10 classification or clustering) and not enough removing (allowing the public releasing of PHI information).

In this paper we propose a complete de-identification tool (NER and NES) able to label/delete/surrogate names, locations, organisations, dates, telephone numbers, emails and urls for French Unstructured health Record. For the NER phase, since deep learning approaches represent state-of-art, the availability of labeled datasets is mandatory. Unfortunately, the first challenging part of our work is the lack of french labeled datasets and more precisely the lack of medical french labeled datasets. We propose an hybrid process involving a Rule-based technique and a Bert-based deep learning approach to overcome this problem. For the NES phase, the foundation is Differential Privacy (DP) as introduced in [DMNS06] and particularized in Local Differential Privacy (LDP) [DJW13] when acquiring consecutive individual data.

The second challenging part of these work is the evaluation of our process. We decided to look at how the ICD-10 classification can be impacted by the de-identification. For this, we collaborate with the Nord Franche-Comté Hospital (HNFC) <sup>1</sup>, a mid-size hospital at the east of France. This hospital employs a team of 13 people dedicated to the ICD-10 classification which can be viewed as the assignation of some diagnosis codes (also known as ICD-10 codes, which stands for the 10th version of the International Classification of Diseases [Wik21]) to a clinical note. We then asked some of the ICD-10 coders to classify a dataset of non de-identified clinical notes and then to do the same thing with the de-identified version of the dataset.

The results we obtain are very encouraging and show an accuracy of ... which can be seen as a very good result since our models for the NER part does not train on medical data but mainly on general purpose data (i.e. Wikipedia).

The following of the paper is structured as follows. Section 2 details the context of our study and related work on medical document de-identification. Section 3 described the NER part of our process. We detail the datasets, the models and the overall architecture we used in our approach. Experimental results are also presented. Section 4 focuses on the NES phase and the substitution strategies of each HIPAA are formally described. The evaluation of our process through the ICD-10 classification task is presented in Section 5. We end in section 6 by some concluding remarks and future works.

---

<sup>1</sup><https://www.hnfc.fr/>

## 2 Context and Related work

This section first starts with the legal context of privacy for medical documents and then presents the state of the art of medical textual document de-identification.

### 2.1 Legal Context

Prior to any handling of medical records by a party external to the institution that disposes of them, it is necessary to ensure that the confidential health information is protected. Quoting GDPR [gdp16, (Recital 35)], "personal data concerning health should include all data pertaining to the health status of a data subject which reveal information relating to the past, current or future physical or mental health status of the data subject". This regulation however allows to outsource such kind of health data, but in the restricted context of public health as precised in Recital 54: "the processing of special categories of personal data may be necessary for reasons of public interest in the areas of public health without consent of the data subject".

But analyzing medical documents to extract codes (as in the specific context of this work) is not a public health issue. It is a global approach to optimizing hospital resources. Thus, this restrictive framework authorizing the use of raw health data cannot be applied here as in many other situations.

However GDPR "does not apply to anonymous information, namely information which does not relate to an identified or identifiable natural person or to personal data rendered anonymous in such a manner that the data subject is not or no longer identifiable"[gdp16, (Recital 26)]. The ability to work with health data is therefore based on the fact that the medical documents can be fully de-identified. State of the art of de-Identification methods for french clinical notes are recalled in the next section.

### 2.2 Related Work on De-Identification of (French) Clinical Notes

Automatic de-identification is a challenging problem whose first works date back to the late 1990s [Swe96, GSG04].

Nevertheless, in the medical field and to our knowledge, there is no commercial tool or large-scale deployment. One explanation of this is the difficulty of dealing with any unstructured texts and so to guarantee the complete de-identification of all PHI. As previously stated, a too strong de-identification may lead to an information loss that can be detrimental to the analysis tasks that may follow. For example: "Charcot disease" must not be de-identified while "M. Charcot suffers from vertigo..." must be de-identified. Another example is the term "PSA" which stands for a french car company and also for a blood medical exam. These few examples underline the importance of the context and the challenging nature of the de-identification problem.

For the NER phase (and for the English language), several works have focused on the use of machine learning models such as support vector machine or decision tress [GGR<sup>+</sup>06, USLS08, LCT<sup>+</sup>15]. Recent advances in the neural approach and deep learning have led to important advances. In 2016 Lample et al. in [LBS<sup>+</sup>16] and Deroncourt et al. in [DLUS16] proposed the first architectures based on Neural Networks for the de-identification of medical unstructured texts. Deroncourt used a Recurrent Neural Networks (RNN) trained on two medical datasets, namely i2b2 [DcrftDoBIDitBIaHMS] and MIMIC [JPM16]. The obtained results represent state-of-the-art with F1-scores reaching 97.85% and 99.23% depending on the dataset they used. In [HHD<sup>+</sup>20], the authors propose a comparison studies of deep learning systems ranging from off-the-shelf to fully customized. The authors use an hybrid system based on RNNs coupled with a CRF (similar to [DLUS16, LTWC17]). The customization levels depend on the dataset and the embedding layer they used. Unsurprisingly, their custom approaches are able to deliver the most accurate results with a F1-score ranging from 97 to 99% on par with Denoncourt et al's results.

The most consistent studies on de-identification of medical documents in French are mainly those carried out by C. Grouin [GN14, GGN15]. However, all of the implementations are based either on CRF or on regular expression rules. More recently, [BAC<sup>+</sup>21] focuses on the de-identification of french emergency medical records. The

approach consists in two steps: first the authors use Flaubert [LVF<sup>+</sup>20] (see later) to classify the notes containing data to de-identified. They then compare different approaches combining rules-based techniques and LSTM (via Flair [ABB<sup>+</sup>19]). Note that Flair was trained on WikiNER (see next section). For the evaluation phase, they have also manually annotated a relatively small number of notes (414) where only the persons names were detected.

### 3 French Named Entity Recognition

In this section, we describe our approach for the NER phase of the de-identification problem of French clinical notes.

#### 3.1 Named Entities

Named Entity Recognition is the task of identifying and categorizing key information (entities) in a text. An entity can be any word or series of words. So one must first identify which word classes may have content that could reveal personal information. Unfortunately, due to the richness of natural language and the uniqueness of many human behaviors, there is no definitive answer to this question. Some combinations of even innocuous keywords can uniquely identify a person and can thus be seen as quasi-identifiers. An acceptable answer may be to rely on what is accepted as identifiers for a specific search domain.

For example, all the decisions published in France on the Cour de Cassation (legal area) website have had the first and last names of individuals mentioned in the decisions removed and replaced by letters. Additional deletions of other elements that allow the identification of persons (address, telephone number, email address...) and whose disclosure would be likely to undermine their security (or that of their entourage or the respect of their private life or that of their entourage), were also carried out before the decisions were put online.

In the field of health, the Health Insurance Portability and Accountability Act (HIPAA) [CM18] provides safe harbor guidelines that define what information that can be considered as private: Private Health Information (PHI). The HIPAA categories form an acceptable consensus [PKSK17, FM08, BM10], even outside their field of application, which is the USA. For the sake of completeness, Table 1 recalls these categories.

1. Names;
2. All geographic subdivisions smaller than a state, including street address, city, county, precinct, zip code, and their equivalent geocodes;
3. All date elements (except year) for dates directly related to an individual including, birth date, admission date, discharge date, death date, etc. date of birth, date of admission, date of discharge, date of death; and all ages greater than 89 years and all date elements (including year) indicative of that age, except that such ages and elements may be aggregated into a single age category of 90 years or older ;
4. Telephone numbers;
5. Fax numbers;
6. E-mail addresses;
7. Social security numbers;
8. Medical record numbers;
9. Health plan beneficiary numbers;
10. Account numbers;
11. Certificate/license numbers;
12. Vehicle identifiers and serial numbers, including license plate numbers;
13. Device identifiers and serial numbers;
14. Web universal resource locators (URLs);
15. Internet Protocol (IP) address numbers;
16. Biometric identifiers, including fingerprints and voice prints;
17. Full face photographic images and any comparable images;
18. Any other unique identifying number, feature, or code.

Table 1: HIPAA categories [CM18]

## 3.2 NER Related Work

The first NER works of the 1990s were all based on rule-based techniques which typically use regular expressions manually defined. One of the most popular rule-based technique is the Conditional Random Field (CRF) technique [LMP01]. These approaches are not machine learning approaches and do not need any labeled data to be trained on. But they are unable to differentiate "M. Charcot" from "Charcot disease" i.e to take the context into account. The recent introduction of the Transformers [VSP<sup>+</sup>17] and Bert [DCLT18] models has allowed a new evolution of the field. The literature is abundant and we can cite here two comparative studies of different transformer models: [Han21] and [PdGS21].

All previous works focused on the English language where big labeled clinical datasets exists, i2b2 [DcrftDoBIDitBIaHMS] and MIMIC [JPM16] for example. This allows some very specific training of the different deep learning models. Unfortunately, no such datasets exist for the French language. This constrains us to train our models on a generalist dataset which will be less sensitive to the medical context of a clinical note. We can cite, the French WikiNER (see next section) dataset which is derived from Wikipedia. Nevertheless, several works exist in the general field of French NER. For example in [SDM<sup>+</sup>20], the authors present a new state-of-the-art for French Named Entity Recognition in a general purpose context (not medical). They use the recent Bert models in their French version i.e. Camembert [MMOS<sup>+</sup>20] and achieve an F1-score of 90.25% in the task of Named Entity Recognition. As specified by the authors, a major difficulty is to obtain French labelled datasets in order to train and evaluate their models (transfer learning). In this way, to carry out their work they had to manually label their own (relatively small) dataset.

## 3.3 Non Private Training Monolingual Datasets

As we will see in the remainder of this section, our approach relies on the use of transformers deep learning models such as FlauBERT. Note that, we choose to use monolingual models since several studies [PdGS21] show that they outperform multilingual approaches. In order to specialized FlauBERT on the NER task, we must trained it with a dedicated dataset. Unfortunately and as stated above, there are very few french NER datasets of sufficiently large size. Among them, the WikiNER dataset [NRR<sup>+</sup>13], created by Nothman et al. contains manually-labelled Wikipedia articles across 9 languages, namely English, German, French, Polish,... In French the size of the dataset consists in more than 61 000 pages and more than 3 000 000 words. The annotations are of 4 main types: LOCation, PERson, ORGanisation and MISCellaneous. Since this dataset is based on many Wikipedia pages and is unfortunately very general, it seems natural to augment this dataset with elements specific to the field studied (here the medical field) and the country (here France). Adding specificities of the domain allows indeed to hope to treat them automatically and more precisely (cf "Charcot Disease"). To this goal, the QUAERO [NGL<sup>+</sup>14] french medical corpus has been added to the former WikiNER dataset. It is not a clinical notes dataset but a selection of MEDLINE titles and EMEA <sup>2</sup> documents that were manually annotated. The annotations are of ten types: Anatomy, Chemical and Drugs, Devices, Disorders, Geographic Areas, Living Beings, Objects, Phenomena, Physiology and Procedures. It can be noticed that the QUAERO dataset does not contain any personal information. The main characteristics (content, quantitative attributes and number of labelled items) of these datasets are presented in Table 2. Notice that the HNFC dataset will be introduced in section 3.5.1.

## 3.4 System Architecture

The approach we proposed is an hybrid one combining state of the art approaches of Transformers and CRF. Indeed, an approach based only on Transformers does not seem to be adapted because of the lack of a consistent french

---

<sup>2</sup><http://opus.lingfil.uu.se/EMEA.php>

Dataset	Content	Nb of characters	Nb of words	Nb of sentences	Nb of labels
WikiNER	Generalist	27 926 161	3 054 130	135 276	415 088
Quaero	Medical	180 919	17 164	1839	2562
HNFC	Medical	775 035	156 423	9993	23829

Table 2: Main characteristics of the different french datasets

dataset in which one can find almost all PHI. Since some entities (email addresses, phone numbers, for example) are built from regular expressions, it seems that an approach based on rules would obtain higher prediction scores on these than another tool based on learning. Combining rule-based approaches with supervised or unsupervised learning approaches seems to be a necessity when the objective is to increase the accuracy of the overall approach. Figure 2 presents the general architecture of our proposal and is detailed in the following sections.

### 3.4.1 Deep Learning based approach

As stated previously, we use some models based on BERT (Bidirectional Encoder Representations from Transformers) [DCLT18]. It is a Transformer network composed of a suite of encoders only ( $N = 12$  or  $24$  depending on the version: base with 110 millions parameters or large with 340 millions parameters). BERT was pre-trained on a large corpus of unlabelled text including the entire Wikipedia and Book Corpus. BERT is a bidirectional model meaning that BERT learns information from both the left and the right side of a word’s context. Since its introduction, BERT and all its "descendants" have been widely used and proved to be very efficient. In our case the choice of an architecture based on BERT was guided by the availability of pre-trained French models i.e CamemBERT and FlauBERT described hereafter.

#### CamemBERT and FlauBERT

CamemBERT [MMOS<sup>+</sup>20] is a BERT type model developed by Facebook and the INRIA in France. It has been pre-trained on a 138Gb French corpus. More precisely it is based on the RoBERTa architecture [LOG<sup>+</sup>19].

FlauBERT [LVF<sup>+</sup>20] is another French BERT developed by the CNRS in France. It has been pre-trained on a large heterogeneous French corpus and its performances compared to CamemBERT are very close. More generally, the results obtained with both models show that a specific French language model improves the results compared to similar multilingual BERT models [MMOS<sup>+</sup>20].

One of the main advantages of these models, and of BERT models in general, is their efficiency in case of transfer learning. The idea is to use a pre-trained model such as CamemBERT or FlauBERT and fine-tune it on a smaller and more specialized dataset. We then obtained and new specialized model. In the medical research area ClinicalBERT [HAR20] and BioBERT [LYK<sup>+</sup>19] are such models, fine-tuned on a medical corpus. Unfortunately all these models are English language models. To our knowledge, there is no French medical fine-tuned BERT model.

#### NERDA

Given a fine-tuned BERT model, it is now possible to add it some layers (typically a dense layer and classification layer) to perform some NLP tasks such as text classification or Name Entity Recognition as shown in the Figure 1. In this study we used the NERDA API [KN21]. 'NERDA' originally stands for 'Named Entity Recognition for DANish'. However, this is somewhat misleading, since the functionality is no longer limited to Danish. This Python

package offers an interface for fine-tuning pre-trained transformers for NER tasks. This architecture is presented in Figure 1. In the remainder we will use NERDA and FlauBERT indistinctly.

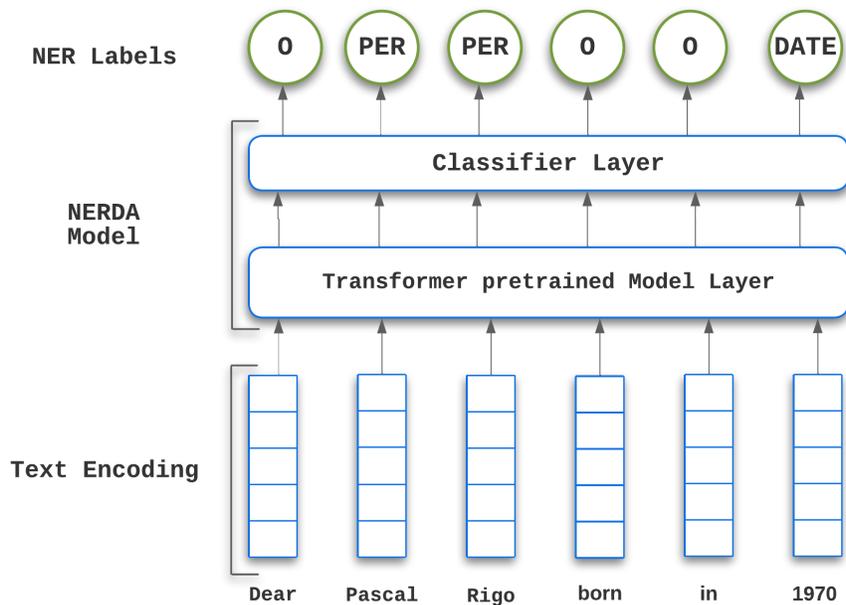


Figure 1: Deep Learning Model Architecture for NER

### 3.4.2 CRF based approach: MEDINA

MEDINA is machine learning based approach using Linear chain Condition Random Field (CRF) system [GZ13] as implemented in wapiti [LCY10] for French document de-identification.

In the context of medical notes, medical reports (written by doctors) often do not respect any typographic rule. For instance, in the sentence "**jean** habite **bermont 3 rue pierre dole**", (jean lives in bermont 3, pierre dole Street)many information can be missed by a CRF tool: "jean" is the name of a person normally in lower case, not preceded by "Dr" or "M./Mr" and may be classified without context as a classic pant.

### 3.4.3 Hybrid System

The hybrid approach we proposed is the following:

- A neural approach based on BERT and more precisely one of its french version, FlauBERT for very contextual entities such as: persons, locations, organisations and medical terms.
- A french CRF-based (MEDINA) to label entities such as: dates, telephone numbers, email and url.

Figure 2 illustrates this approach. The input of both models is a clinical note with sensitive information. The MEDINA model will tries to detect all attributes (person, location, age, date, phone number ...) present in the file. Since the deep learning model (FlauBERT) is trained only on the attributes PERson, ORGanization and LOCation (due to the WikiNER training dataset), it tries to detect only these 3attributes. The output of each of the two models is forwarded to a decision procedure whose objective is to make the most appropriate choice according to the detected entities. For a given word or group of words, let  $T_M$  and  $T_N$  be the tags proposed by MEDINA and by NERDA respectively. The decision procedure is based on the following rules.

- **Case 1** when  $T_M = T_N$ . When both tools associate the same tag, this one is simply considered as final and is returned.
- **Case 2** when  $T_M \neq T_N$  and  $T_M$  or  $T_N$  is equal to "O", the Outside value, i.e., no tag has been associated. In this case, the final returned tag is the tag that is not "O". This case illustrates the fact that an approach suggests a tag, which possibly identifies the patient (the ZIP code for instance), whereas the other one does not detect anything. If the associated final tag was "O", no substitution would be performed later on the corresponding word. A consequence would be that we would have forgotten to clean this word and the de-identification would not be robust. Conversely, if the word is common, if a wording is wrongly associated with it and then replaced because of its wording, the usefulness of the de-identification is reduced, but this does not affect the patient privacy. We have favored this second scenario in the present case.
- **Case 3** when " $O$ "  $\neq T_M \neq T_N \neq$  " $O$ ". The situation when the two approaches have associated two different tags both distinct of "O" to the same sequence occurs particularly when NERDA (FlauBERT) associates one tag in {PER, LOC, ORG}. The training dataset of this approach is indeed only based on this set of entities. It has been shown [SDM<sup>+</sup>20] that deep learning approaches are more accurate than CRF ones for very contextual entities such as persons, organizations and locations. So in this case, the final returned tag is  $T_N$ .

For example, in the sentence "M. **Jean** habite à **Bermont 90400**" (Mr. Jean lives in Bermont, 90400), MEDINA associates "PER" to Jean, "PER" to Bermont and "LOC" to 90400 whereas NERDA associates "PER" to Jean, "LOC" to Bermont and "O" to 90400. The final association returned by the decision procedure will be Jean is PERson (case 1), Bermont is LOCation (case 3) and 90400 to LOCation (case 2).

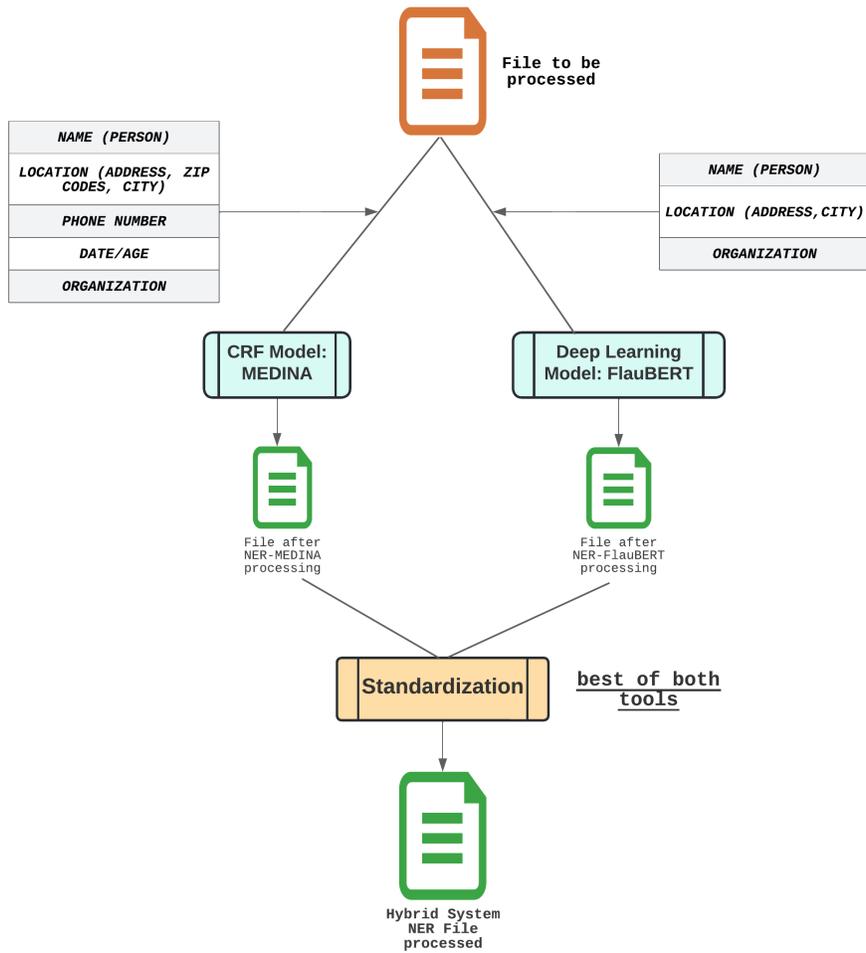


Figure 2: Our proposed hybrid approach for the NER part

## 3.5 Evaluating the NER approach

This section describes the evaluation methodology we carried out for the NER phase. figure 3 summarizes the whole evaluation approach.

### 3.5.1 HNFC dataset

Evaluating our hybrid approach requires an annotated set of french clinical notes. The only way to get a such dataset (especially in French) is to manually annotate some existing notes. We then use 375 files of deceased patients from the Nord Franche-Comté Hospital. These files were pre-annotated by the hybrid system previously presented and then were manually annotated by the hospital staff with Doccano [NKK<sup>+</sup>18] manual annotation tool. Given the automatically pre-annotated files, the annotator is responsible for checking, correcting and completing the possible errors produced by the model. HIPAA attributes were formalized to avoid ambiguity and multiple annotation criteria. For example, "Ehpad de Bermont" (Retirement Home in Bermont) can be annotated in several ways: Ehpad as O, Bermont as LOcation or Ehpad de Bermont as ORGanisation. Similarly, "dans 3 jours" (in the 3 days) is a date while "3 x par jour" (3 times per day) is not a date. To minimize the risks, two annotators worked in parallel on the same files and each annotation pair is then manually analyzed and merged into an unique one. This work was performed by 6 people during 6 hours. The approximations encountered around the attributes, like the examples given above, are the result of this experimentation. At the end of this process, we have obtained a reference dataset with which we can use to evaluate our model. In the following, this dataset is referred as HNFC dataset. Note that, for all our tests (annotation and evaluation) we worked on site and no note came out of the hospital.

### 3.5.2 Comparison of NER approaches on HNFC dataset

For evaluation, we use the classical metrics: precision, recall and F1-score. Let  $TP$  be the number of true positive annotations,  $FP$  the number of false positive annotations, and  $FN$  the number of false negative annotations. Then, the recall  $R$  is given by  $R = TP/(TP + FN)$ , and the precision  $p$  is given by  $P = TP/(TP + FP)$ . Recall and precision answer two questions about a named entity recognition tool, respectively "did we find everything we were looking for?" and "did we only label what we were looking for?". The F1-score metric combines precision and recall, usually by taking the harmonic mean of the two. To get a sense of the overall performance of the system, we use the micro-average of precision, recall, and F1-score. To compute the micro-average, a confusion matrix is created for all categories, and then precision and recall are computed from this table, giving the same weight to each PHI instance regardless of its category.

To establish some baselines for our evaluations, we recall that our model has been trained on the French WikiNER dataset. The french NER pipeline from the Python Spacy library [HM17] has also been trained on WikiNER (to which has been added the small dataset Sequoia). The model used is based on a Convolutional Neural Network (CNN) which provides a F1-score of 0.84 for predicting 4 labels: Person, Location, Organisation and Miscellaneous.

J-B. Polle<sup>1</sup> proposes CamemBERT-ner, a pre-trained model for french NER based on CamemBERT and trained with WikiNER. This model exhibits a F1-score of 0.89 but the evaluation was not performed by the author on a sub-set of WikiNER but on a personal crafted dataset composed of emails and chats.

Table 3 summarizes results of the different NER approaches applied to the HNFC dataset.

For all approaches PERson, ORGanisation and LOcation are searched in documents, since the training step of each of them is based on this set of tags. MEDINA, which is a CRF model, is able to perform research for Date, Age and Phone Number in addition to. For each tag, the higher results for Precision, Recall and  $F_1$ -score are emphasized in bold.

---

<sup>1</sup>Jean-Baptiste Polle (2020). CamemBERT model for Named Entity Recognition task

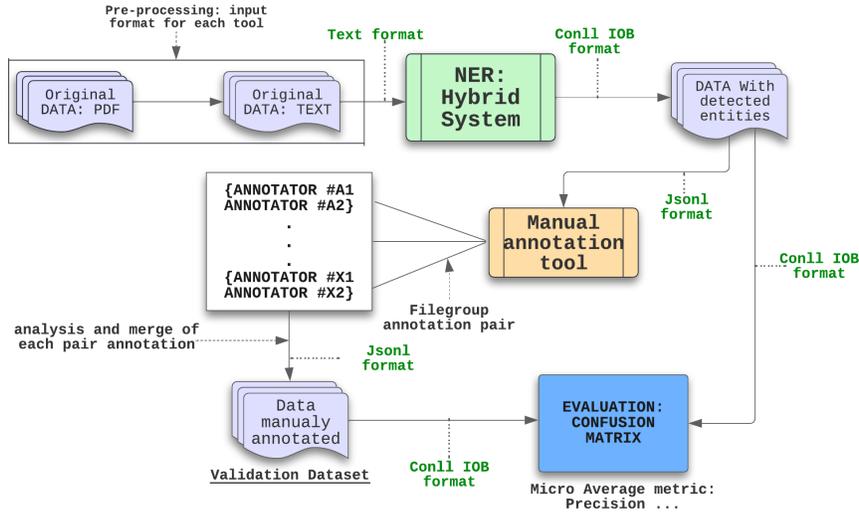


Figure 3: HNFC data set creation and evaluation process

Labels	Spacy			Camembert NER			MEDINA			NERDA			PROPOSAL		
	P	R	F <sub>1</sub>	P	R	F <sub>1</sub>	P	R	F <sub>1</sub>	P	R	F <sub>1</sub>	P	R	F <sub>1</sub>
PERson	59	76.8	67	89	99	93.8	<b>98.2</b>	97.7	<b>98.2</b>	91.8	97.6	94.6	96.3	<b>99.8</b>	98
ORGanizat.	2.2	10.9	3.6	7	21.8	11.1	32.6	24.8	28.1	16.9	34.1	22.6	<b>41.1</b>	<b>57.3</b>	<b>47.8</b>
LOCation	40.3	11.9	18.4	46	67.2	54.6	<b>98.8</b>	81.1	89.1	75.7	66.3	70.7	88.4	<b>95.8</b>	<b>92</b>
Date							97.7	86.6	91.9				<b>97.7</b>	<b>86.7</b>	<b>91.9</b>
Age							91.5	66.9	77.3				<b>91.5</b>	<b>66.9</b>	<b>77.3</b>
Phone N.							99.5	97.9	98.7				<b>99.5</b>	<b>97.9</b>	<b>98.7</b>
Micro Av.	54.9	54.8	54.9	70.8	51.5	59.6	<b>98.2</b>	91.2	94.5	85.8	86.7	86.3	94.6	<b>94.9</b>	<b>94.7</b>

Table 3: NER results on the HNFC dataset

Let us first focus on the common tags. For all recalls, the highest recall score is the one provided by our hybrid proposal. Notice that MEDINA provides more precise results in PERson and LOCation categories

As far as the detection of dates, ages and telephone numbers is concerned, it is stated that these tags are absent from all learning-based approaches. Only MEDINA detects them (partially). Due to rule 2 of the decision procedure of our approach, the entities detected in this set by MEDINA will systematically be returned as is. The results are therefore those of MEDINA.

The last line Micro Average summarises the contribution of this hybrid entity detection approach: our proposal is globally the one that detects entities most often (highest recall), without neglecting precision (with the highest F1-score). Notice that the aggregated micro average recall score is 94.9%. The drawbacks to this picture are clearly the detection of temporal aspects (dates, ages) and of organisations. From a privacy point of view, the incomplete identification of dates is problematic as these temporal elements are quasi-identifying (see section 4.6). As far as organisations are concerned, the hospital partner thinks that this data is less sensitive.

As detailed in the architecture section our model is a hybrid system that combines a machine learning method

(MEDINA) and a deep learning model (NERDA). MEDINA as illustrated in the Figure is in general less accurate in recall than in precision i.e. the model misses the attributes that need to be detected more than it detects the words that are not attributes (Location  $\Rightarrow$  0.811, Organization  $\Rightarrow$  0.248 ). It is very important to have a better recall in a privacy context. it is better to detect all the confidential information that needs to be detected than to detect the ones that are not. Our deep learning based system (NERDA) allows us to improve this. It outperforms MEDINA in recall, which improves the results in these hybrid model cases. The improvement is due to our hybrid system (decision procedure) which allows us to correct the errors of one of the methods by the other. The aggregation of the associations increases our accuracy of detection on the one hand and of association on the other hand.

Let's take for instance the sentence "Mr. **Jean** lives in **Bermont, 90400**". MEDINA associates Jean to PERson, Bermont to PERson and 90400 to LOCation whereas NERDA associates Jean to PERson, Bermont to LOCation and 90400 to Outside. Due to the decision procedure, our hybrid system will finally associate Jean to "PER", Bermont to "LOC" and 90400 to "LOC". It will go from 66% of precision for each tool to 100%. This hybrid system not only improves the results but also allows to take into account all the main attributes (HIPAA) present in the medical documents.

## 4 Surrogate Generation Strategies

This section describes the Named Entity Substitution step of our approach for the de-identification problem of French clinical notes. It starts with related works of the field and continue by showing that PHI can be divided into two types of categories (Section 4.2). Random based substitution are detailed in Section 4.3 whereas Local Differential Privacy based approaches are defined in Section 4.4.

### 4.1 Related Work

The most straightforward way to substitute entities is to replace each entity value with the name of the entity itself (LOC for example). It has been shown [Laf] that such replacement of all the 18 HIPAA entities efficiently protects the privacy of the patients/doctors: in this article, it is shown that only 2 over 32500 health notes have been indeed re-identified. However, this kind of coarse substitution drastically reduces the usefulness of the data. For example, it may become difficult to distinguish whether the subject of certain sentences is the patient or the doctor. More problematically, the chronology of dates is in this case very difficult to re-establish since they are completely hidden. Other more accurate substitution methods have been designed for this purpose.

Prior to this work, the Scrub [Swe96] systems generated surrogate texts that match the format of the original ones. More precisely, for dates, Scrub associates to each detected date an approximate timestamp (the nearest month, for example). For names, a lookup table is used to ensure that the same name in the original file is always replaced by the same substitute.

In [DCR<sup>+</sup>04, Lev03] annotated PHIs from a corpus by hand and replacement carried out semi-automatically. More precisely, dates are shifted by a random number of weeks and years, while preserving the day of the week. As in Scrub, each person name is replaced by its associated one, here from a list (Boston residents). Locations are replaced by randomly selected small towns. Uzuner et al.[ULS07, USLS08] extend the works of Douglass et al[DCR<sup>+</sup>04] by replacing strings such as identification numbers and phone numbers with randomly generated digits and letters. For dates they maintained internal temporal relationships by shifting all dates in a document by the same number of days and ensuring that the surrogate dates were properly formatted.

In the de-identification approach of Deleger et al [DLN<sup>+</sup>14], names are replaced with surrogates by randomly selecting male, female, unisex and family names from pre-compiled lists. Then, all documents are parsed to store all 455 detected places (street names, city names, state names...), resulting in a corpus of places. Dates are replaced

by random dates, while respecting the format, and places by places from the above corpus, while respecting the format.

One of the systems used in recent research is the system developed by Stubbs et al.[SKU15]. The authors use, for names, numbers and letters the system described by Deleger et al.[DLN<sup>+</sup>14] For geographic locations, they use a pre-compiled list of different types of geographic locations, and a random choice to generate the surrogates. For the dates a uniform date shift with a random number of days was applied. This is the system used in medical corpora available today for research such as the 2014 i2b2/UTHealth[KSSU15] corpus. The French document de-identification tool (MEDINA[GZ13]) also uses the recommendations of Stubbs system. This paper uses a method that is based for some categories on the Stubbs system and makes contributions in the generation of dates, ages and geographical locations.

## 4.2 Splitting Categories

For some categories, generating consistent surrogates is straightforward. For example let us first consider phone numbers, URLs, email addresses, are alphanumeric strings of numbers, letters and special characters. All these elements are of course strongly identifying, but are not clearly linked to health data. Random substitutions can be applied on all these entities to ensure privacy without any consequence on utility provided the text format is respected.

Other categories are more problematic. Dates and ages, which are temporal data, clearly explain the chronology of medical developments. Locations can affect pathologies: for instance, some cities have high radon levels which can significantly increase the risk of cancer. Randomly substituting these data does not make sense as they directly affect health.

The approaches developed in the literature concerning dates can be summarised as generalisations (as in an approach based on  $k$ -anonymity) of a shift or addition of a bounded noise. In the first case, it is known that these generalisations are not robust to attacks by additional knowledge [MKGv07]. In the second case, it has been shown that this kind of noise is not robust [DN03]. Introduced by Dwork [DMNS06], differential confidentiality is a mathematical context that allows the publication of an individual’s information while respecting the latter’s privacy. It therefore guarantees the confidentiality of the individual during the process of disseminating information by means of queries on a database containing or not containing the data of the latter. In our case, we want to clean up all the dates of a document in order to be able to use it as many times as we want. This amounts to considering that the patient has a mechanism that he applies locally to his documents. This is known as local differential confidentiality (LDP) [DJW13], the definition of which is given below

The next two sections detail these two faces of sanitization and are summarized in Figure 4.

## 4.3 Random Substitutions

Among categories detected by the NER step, most of them contains values which can be replaced by random values without affecting further processing. For example, generating substitutes for categories such as emails, URLs, phone numbers ... corresponds to random substitution: a phone number can be replaced by any random phone number.

The names are a little different because it is interesting to preserve the affiliation within the documents. The algorithm is illustrated in the Figure 5. First, a lookup table is created for each file. We start by checking if the current full name is in the dictionary. If it is not, we check if the corresponding surname is in the dictionary, if it is not, it means that we have not processed it yet. Its processing consists in generating its substitute (last name & first name) and in registering it in the dictionary. Moreover, we only register the surname that corresponds with its substitute. If on the other hand its surname is in the dictionary, we recover its substitute and we generate only

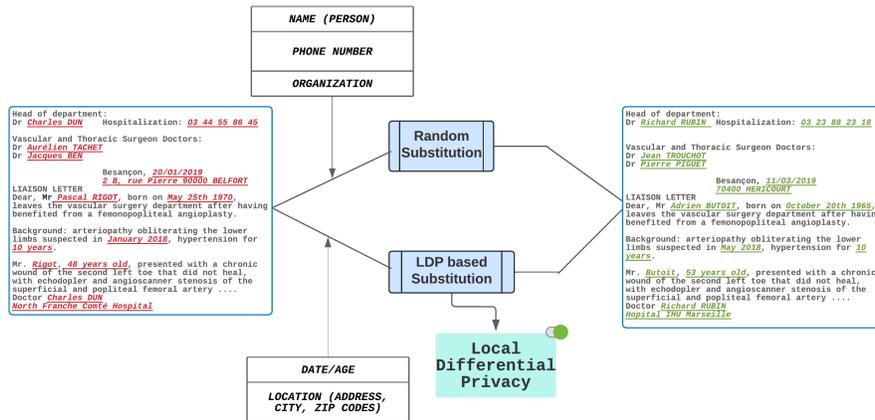


Figure 4: Substitution process

the first name and we record the couple (last name & first name) in the dictionary. Finally, if the full name is in the dictionary, we simply retrieve its substitute.

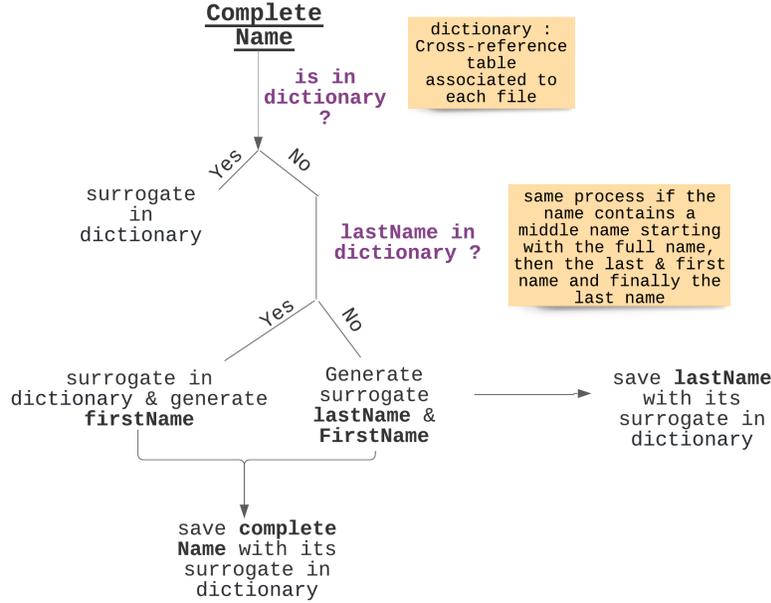


Figure 5: Name Surrogate Strategy

#### 4.4 Local Differential Privacy Context

Initially formalized in [DJW13], local differential privacy (LDP) ensures individual’s privacy during the data collection process. A formal definition of LDP is given in the following:

**Definition 1 ( $\epsilon$ -Local Differential Privacy)** *A randomized algorithm  $\mathcal{A}$  satisfies  $\epsilon$ -LDP if, for any pair of input values  $v_1$  and  $v_2 \in \text{Domain}(\mathcal{A})$  and any possible output  $y$  of  $\mathcal{A}$ :*

$$\Pr[\mathcal{A}(v_1) = y] \leq e^\epsilon \cdot \Pr[\mathcal{A}(v_2) = y].$$

This property is a strengthening of the original centralized model of differential privacy (DP) [DMNS06] since it applies on each records whereas the original one applies on each query on the whole dataset. Roughly speaking, a small-scale noise should suffice for a weak privacy constraint (corresponding to a large value of  $\epsilon$ ), while a greater level of noise would provide a greater degree of uncertainty in what was the original input (corresponding to a small value of  $\epsilon$ ).

Many mechanisms verify this local differential confidentiality. They can be classified according to the type of considered data (categorical, real, integer), and according to their usefulness with respect to a question. We recall below the Laplacian mechanism (introduced by Dwork [DMNS06]) because of the simplicity of its implementation on real data.

**Definition 2 (Laplacian mechanism in an interval of amplitude  $\Delta$ )** *In the Laplacian mechanism, a numerical value  $v$  is transformed into a numerical value  $\mathcal{M}_{\text{Lap}}(v, \Delta, \epsilon)$  such that*

$$\mathcal{M}_{\text{Lap}}(v, \Delta, \epsilon) = v + \text{Lap}\left(\frac{\Delta}{\epsilon}\right) \quad (1)$$

where  $\text{Lap}\left(\frac{\Delta}{\epsilon}\right)$  is the Laplace distribution centred in 0 and whose scale parameter is  $\frac{\Delta}{\epsilon}$ .

The addition of noise in the Laplace mechanism ensures local differential  $\epsilon$ -confidentiality. As defined, this noise depends on the amplitude of the input data interval and of  $\epsilon$ .

$\epsilon$ -LDP provides several important properties, e.g., immunity to post-processing (and thus robust against any supplementary knowledge based attack) and composition [DR<sup>+</sup>14]. That is, if mechanisms are sequentially applied to many elements inside a document of a person the whole budget  $\epsilon$  is equal to the sum of all the budgets of all mechanisms.

Generally, data sanitization mechanisms verifying  $\epsilon$ -LDP are optimised according to the nature of the data taken as input. For example, a mechanism adding noise according to a Laplace distribution to numerical data can be applied to integers, but its utility will be reduced compared to the exponential mechanism dedicated to discrete data. Thus, for each of the categories to be sanitized, a careful selection of the mechanism to be used must be made. In the context of medical notes de-identification, there remains 2 categories namely location and date elements, on which a  $\epsilon$ -LDP sanitization process will be applied on. Location is a spatial data whereas age/date may be seen as a number. Since the domains of both are not equal, a mechanism dedicated to each of those is thus developed. For any of them, the privacy budget  $\epsilon$  should be processed with the maximum care.

Defining the  $\epsilon$  leakage budget for each algorithm can be achieved by sharing the  $\epsilon$  budget among all of them. For instance, in our medical context, starting from the global  $\epsilon$  budget that is allocated to whole sanitization process,  $\frac{\epsilon}{4}$  can be associated to location, and  $\frac{3\epsilon}{4}$  to date&age, but any other partition of the  $\epsilon$  original budget would provide the same level of protection. Note, however, that some attributes do not have the same degree of criticality or the same discriminating ability. For example, in a set of textual documents where the whole cohort concerns retired people, age, which is in a rather low amplitude class, is not as critical nor as discriminating as in a set of adult people.

When the attributes have a similar level of criticality and therefore the choice of the partition is more open, it should be made with respect to the accuracy that can be obtained through predictions once the dataset is sanitized. However, knowing that these predictions can only be made once the dataset has been sanitized, the question arises of finding the best parameters, knowing that any re-reading of the personal data will itself reduce the leakage budget.

Notice finally that other combination are possible for instance share the  $\epsilon$  budget proportionally to the number of detected occurrences of each category, proportionally of the number of entities concerned (Date, Age, Location ...).

## 4.5 Sanitizing locations

To solve the problem of privacy location, we used the concept of Geo-Indistinguishability [ABCP13, CABP13] which is based on  $\epsilon$ -LDP. It consists briefly in retrieving, for a given location, a location at a certain distance from it by differential privacy. This process is detailed in the algorithm 1. The association  $(Z, Y)$  between the original location  $Z$  and the sanitized one  $Y$  is stored and used in the whole document if  $Z$  appears at least twice. This step, often denoted as memoization, allows to resist to a possible averaging attack.

---

### Algorithm 1 Sanitizing locations

---

1.  $F$  is a local list of cities  $City(long, lat)$  in the local area with there longitude and latitude data.
  2. Given a location  $Z$  to sanitize. Extract its  $(long, lat)$  data from locations  $F$
  3. From  $Z$ , generate of  $Z'(long, lat)$  by applying the **Geo-Indistinguishability**[ABCP13, CABP13] algorithm
  4. Let  $Y$  be the location in  $F$  closest to  $Z'$ .
  5. Save the mapping  $(Z, Y)$  in a correspondence table for this document.
-

## 4.6 Sanitizing dates and ages

The objective with dates is to preserve the temporality of events in the medical document to gain information during a second analysis while respecting the privacy of patients. In the public medical data sets available for research (i2b2 [DcrftDoBIDitBIaHMS] and MIMIC [JPM16]) a uniform shift of a randomly drawn number of days is performed on the dates. This process poses confidentiality problems. To evaluate this approach on the HNFC dataset (see Section 3.5.1) and for each document, all temporal elements (dates, ages converted to dates) are stored as an ordered sequence  $S_e = [e_0, e_1, e_2, \dots, e_n]$ , from to the current date  $e_0$ , the most recent in the document  $e_1$  and the oldest one  $e_n$ . A second sequence  $S_i$  of intervals is generated  $S_i = [e_0 - e_1, e_1 - e_2, \dots, e_{n-1} - e_n]$ , by computing the differences (expressed in days) between two consecutive temporal elements of  $S_e$ . Finally let  $S'_i$  be a copy of  $S_i$ , but the first element which is removed. Notice that applying a uniform shift (as done in [ULS07, USLS08]) on all the temporal dates will not modify the sequence  $S'_i$ . Note that, in the HNFC dataset we used for our evaluation, all the 375 documents provide distinct  $S_e$  and only 8 out of 375 documents did not contain unique sequences of intervals  $S'_i$ . We conclude that approximately 98% of the chronological intervals present in documents are unique and therefore very strongly identifiable.

As detailed above, for each document, sequences  $S_e$  and  $S_i$  are computed. Each interval in  $S_i$  is a number of days, *i.e.* a numerical value. Local differential privacy mechanism, like Laplacian one, can thus been applied on it. More precisely, the bounded Laplace mechanism [HABMA18] is applied here to avoid negative noise while preserving privacy. Concerning the budget for each interval, several questions arise: should the global budget allocated to the Date category be distributed uniformly? What are the most compromising dates? We see that the older the date, the more sensitive it is (Date of birth for example). How to deal with a huge interval amplitude  $\Delta$  (100 years) in this context?

To solve this problem, dates have been classified into categories (Short, Medium, Long term) and for each category, an interval amplitude  $\Delta$  is computed. More precisely,  $\Delta_S = 61$ ,  $\Delta_M = 660$  and  $\Delta_L = 36,000$ . Each temporal interval is associated to one of these 3 categories. In case of ambiguity (a date in the short term and a date in the medium one, for example), the smallest category is associated to each interval. The global budget  $\epsilon$  is split uniformly between all the dates. Then, a Bounded Laplacian Mechanism [HABMA18] with parameter  $\epsilon_i$  and  $\Delta_i$  is applied to each interval  $i$ , where  $\Delta_i$  an interval amplitude associated to  $i$ . This approach is detailed in the algorithm 2 and illustrated in Figure 6.

---

**Algorithm 2** Sanitizing dates and ages

---

1. Identification : Identify all the temporal elements  $e$  (dates, ages) of the document
  2. Normalization : Normalize each  $e$  in a standard format (ex. dd/mm/yyyy)
  3. Establish the chronology of the latter (classify from the first date to the last including the current date), *i.e.*, compute  $S_e$
  4. Define Date category (short, medium, long term)
  5. Compute the interval sequence  $S_i$  between consecutive dates in  $S_e$
  6. Apply to the intervals the local differential privacy with a Bounded Laplacian noise where  $\Delta$  is the category amplitude
  7. Reconstitute dates from the current date
-

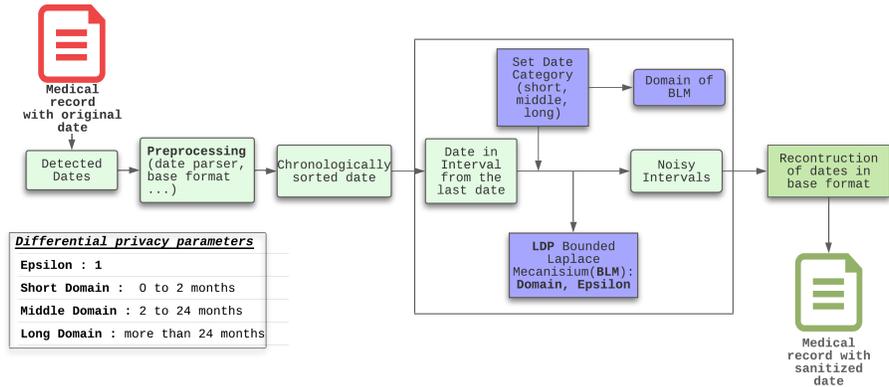


Figure 6: Date Surrogate Strategy

## 5 NES Evaluation

To evaluate the quality of our approach, we de-identified 64 documents from the HNFC and submitted them to an ICD-10 code association analysis (along with the original non de-identified documents). The objective is to code the two groups of documents : before and after the complete system (named entity recognition followed by entity substitution) and compare the results. In the absence of a system capable of automatically associating ICD-10 codes, this work was carried out by two hospital coding specialists. In this way, they manually went through the documents to associate the corresponding codes (a group of documents was assigned by coder). Note that, from a machine learning point of view, the ICD-10 classification is a multi-label text classification task as illustrated in Figure7.

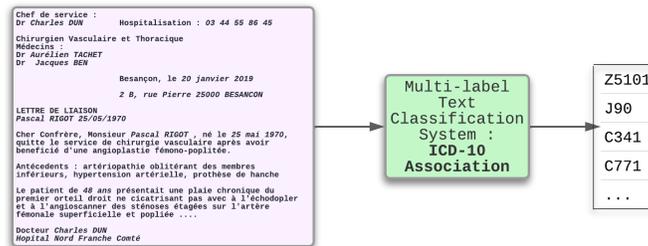


Figure 7: ICD-10 Code association

The performance measures used in this case were computed by considering the coding related to the original document as the truth and the coding related to the de-identified document as the prediction. As performance measures we use metrics: Precision, Recall, F1-score recalled hereafter.

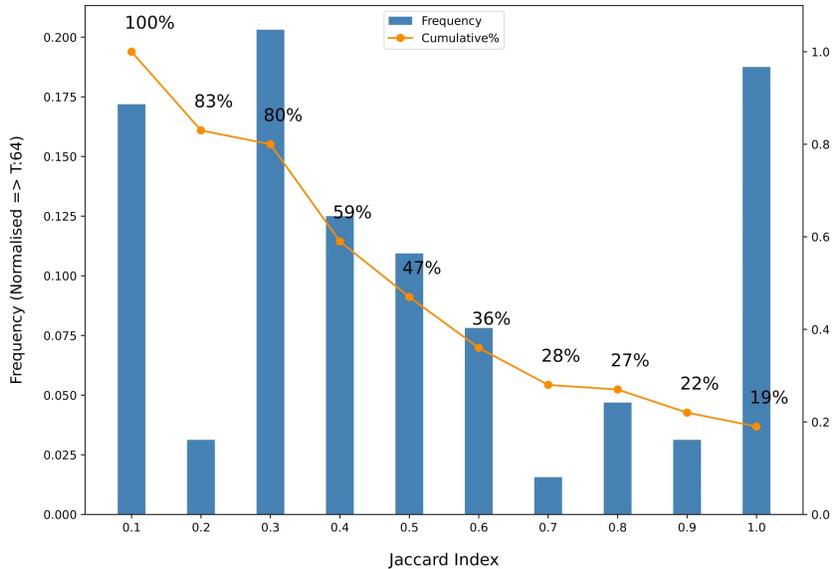
The performances are summarized in the table 4.

- **Precision:** the proportion of correct associations among all associations of a certain code. In other words, it is the proportion of true positives among all positive associations.

Labels	Precision	Recall	F1-score
246	0.396	0.588	0.473

Table 4: Performances after Coding

- **Recall:** the proportion of examples of a certain code that have been associated by the model as belonging to that code. In other words, it is the proportion of true positives among all true examples.
- **F1-score:** the harmonic mean of precision and recall, so it's an overall measure of the quality of a classifier's associations



Jaccard Index frequencies and there cumulative version computed with set of ICD-10 codes from original documents and the corresponding de-identified ones.

Figure 8: Jaccard Index frequency distribution of ICD-10 associations

For each original health document, let  $A$  be the set of original associated ICD-10 codes. Let then  $B$  be the set of ICD-10 codes associated to the sanitized version of this document. The Jaccard index  $\frac{|A \cap B|}{|A \cup B|}$  is a similarity measure between these sets of codes. On these 64 documents, 12 documents have a rate equal to 1, which means that sets  $A$  and  $B$  are equals. Indices are group by bids of size 0.1 and aggregated in a cumulative way. Figure 8 displays these Jaccard indices and there cumulative version. For instance, it can be noticed that for more than 27% (resp. 50%) of the documents more than 70% (resp. 40%) of classification codes can be deduced from the sanitized documents

This analysis (association of ICD-10 codes) on de-identified documents allowed us to see the limits of our method. The coders feedback shows us that the difficulties come mainly from the de-identification of dates where the chronology, the relationship between dates... are crucial for a medical analysis. On the other hand, it confirms the necessity to keep a maximum of information (date, age, location...) for a better analysis.

## 6 Conclusion and Future Work

This work presented a global method for de-identifying textual medical documents adapted to the French language. The HIPAA-defined entity recognition (NER) of medical documents is based on a combination of rule-based approaches implemented in MEDINA and deep learning-based approaches using FlauBERT transformer and NERDA. The substitution of these entities with plausible privacy features is either random when it has no impact on the document, or based on local differential privacy (LDP), a mathematical property accepted as a de facto standard in privacy.

The robustness of the approach was evaluated on an original medical dataset within a French public hospital. In terms of entity detection, the selected combination is the most efficient to date. In terms of privacy protection, the selected  $\epsilon$ -LDP mathematically guarantees that the addition of noise is largely sufficient to make re-identification impossible, if not very difficult. This is the first time that such a property has been established on medical documents.

The utility of the approach was globally observed on a dataset by exhaustively comparing the ICD-10 codes associated with and without this de-identification step. The association of these codes is of high quality, more accurate, but not systematically identical with and without de-identification, the treatment of dates being perhaps very/too protective for the patient.

Let us continue with future work. On manually annotated documents in French to be built, we first think of implementing machine learning to distinguish dates ("15 years ago") from ages ("she was 15 years old"), to allow a more accurate de-identification than the one currently implemented.

Regarding NES, we are thinking of implementing a metric-based LDP algorithm to substitute one location for another with similar epidemiological properties. Similarly, this context could also be applied to dates and avoid the medically sensible but non-linear division of short term, medium term, etc.

## References

- [ABB<sup>+</sup>19] A. Akbik, Tanja Bergmann, Duncan A. J. Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. Flair: An easy-to-use framework for state-of-the-art nlp. In *NAACL*, 2019.
- [ABCP13] Miguel E Andrés, Nicolás E Bordenabe, Konstantinos Chatzikoakolis, and Catuscia Palamidessi. Geo-indistinguishability: Differential privacy for location-based systems. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 901–914, 2013.
- [AND<sup>+</sup>19] Saadullah Amin, Günter Neumann, Katherine Dunfield, Anna Vechkaeva, Kathryn Annette Chapman, and Morgan Kelly Wixted. Mlt-dfki at clef ehealth 2019: Multi-label classification of icd-10 codes with bert. In *CLEF (Working Notes)*, pages 1–15, 2019.
- [BAC<sup>+</sup>21] Loick Bourdois, Marta Avalos, Gabrielle Chenais, Frantz Thiessard, Philippe Revel, Cédric Gil-Jardiné, and Emmanuel Lagarde. *De-identification of Emergency Medical Records in French: Survey and Comparison of State-of-the-Art Automated Systems*, volume 34. LibraryPress@UF, May 2021.
- [BM10] Kathleen Benitez and Bradley Malin. Evaluating re-identification risks with respect to the hipaa privacy rule. *Journal of the American Medical Informatics Association*, 17(2):169–177, 2010.
- [CABP13] Konstantinos Chatzikoakolis, Miguel E Andrés, Nicolás Emilio Bordenabe, and Catuscia Palamidessi. Broadening the scope of differential privacy using metrics. In *International Symposium on Privacy Enhancing Technologies Symposium*, pages 82–102. Springer, 2013.
- [CM18] I Glenn Cohen and Michelle M Mello. Hipaa and protecting health information in the 21st century. *Jama*, 320(3):231–232, 2018.

- [DCLT18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805, 2018.
- [DCR<sup>+</sup>04] Margaret Douglass, Gari D Clifford, Andrew Reisner, George B Moody, and Roger G Mark. Computer-assisted de-identification of free text in the mimic ii database. In *Computers in Cardiology, 2004*, pages 341–344. IEEE, 2004.
- [DcrftDoBIDitBIaHMS] Data and computing resources from the Department of Biomedical Informatics (DBMI) in the Blavatnik Institute at Harvard Medical School. Unstructured notes from the research patient data registry at partners healthcare (originally developed during the i2b2 project).
- [DJW13] John C Duchi, Michael I Jordan, and Martin J Wainwright. Local privacy and statistical minimax rates. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pages 429–438. IEEE, 2013.
- [DLN<sup>+</sup>14] Louise Deleger, Todd Lingren, Yizhao Ni, Megan Kaiser, Laura Stoutenborough, Keith Marsolo, Michal Kouril, Katalin Molnar, and Imre Solti. Preparing an annotated gold standard corpus to share with extramural investigators for de-identification research. *Journal of biomedical informatics*, 50:173–183, 2014.
- [DLUS16] Franck Dernoncourt, Ji Lee, Ozlem Uzuner, and Peter Szolovits. De-identification of patient notes with recurrent neural networks. *Journal of the American Medical Informatics Association : JAMIA*, 24, 06 2016.
- [DMNS06] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [DN03] Irit Dinur and Kobbi Nissim. Revealing information while preserving privacy. In *Proceedings of the twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 202–210, 2003.
- [DR<sup>+</sup>14] Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4):211–407, 2014.
- [FM08] F Jeff Friedlin and Clement J McDonald. A software tool for removing patient identifying information from clinical documents. *Journal of the American Medical Informatics Association*, 15(5):601–610, 2008.
- [gdp16] Regulation (eu) 2016/679 of the european parliament and of the council of 27 april 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing directive 95/46/ec (general data protection regulation) (text with eea relevance), May 2016.
- [GGN15] Cyril Grouin, Nicolas Griffon, and Aurélie Névéol. Is it possible to recover personal health information from an automatically de-identified corpus of french ehers? In *Proceedings of the Sixth International Workshop on Health Text Mining and Information Analysis*, pages 31–39, 2015.

- [GGR<sup>+</sup>06] Yikun Guo, Robert Gaizauskas, Ian Roberts, George Demetriou, Mark Hepple, et al. Identifying personal health information using support vector machines. In *i2b2 workshop on challenges in natural language processing for clinical data*, pages 10–11. Citeseer, 2006.
- [GN14] Cyril Grouin and Aurélie Névéol. De-identification of clinical notes in french: towards a protocol for reference corpus development. *Journal of Biomedical Informatics*, 50:151–161, 2014. Special Issue on Informatics Methods in Medical Privacy.
- [GSG04] Dilip Gupta, Melissa Saul, and John Gilbertson. Evaluation of a deidentification (de-id) software engine to share pathology reports and clinical documents for research. *American journal of clinical pathology*, 121:176–86, 03 2004.
- [GZ13] Cyril Grouin and Pierre Zweigenbaum. Automatic de-identification of french clinical records: comparison of rule-based and machine-learning approaches. In *MEDINFO 2013*, pages 476–480. IOS Press, 2013.
- [HABMA18] Naoise Holohan, Spiros Antonatos, Stefano Braghin, and Pól Mac Aonghusa. The bounded laplace mechanism in differential privacy. *arXiv preprint arXiv:1808.10410*, 2018.
- [Han21] Ridewaan Hanslo. Deep learning transformer architecture for named entity recognition on low resourced languages: State of the art results. *CoRR*, abs/2111.00830, 2021.
- [HAR20] Kexin Huang, Jaan Altonaar, and Rajesh Ranganath. Clinicalbert: Modeling clinical notes and predicting hospital readmission, 2020.
- [HHD<sup>+</sup>20] Tzvika Hartman, Michael Howell, Jeff Dean, Shlomo Hoory, Ronit Slyper, Itay Laish, Oren Gilon, Danny Vainstein, Greg Corrado, Katherine Chou, Ming Po, Jutta Williams, Scott Ellis, Gavin Bee, Avinatan Hassidim, Rony Amira, Genady Beryozkin, Idan Szpektor, and Yossi Matias. Customization scenarios for de-identification of clinical notes. *BMC Medical Informatics and Decision Making*, 20, 01 2020.
- [HM17] Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017.
- [HMS<sup>+</sup>10] Omar Hasan, David O Meltzer, Shimon A Shaykevich, Chaim M Bell, Peter J Kaboli, Andrew D Auerbach, Tosha B Wetterneck, Vineet M Arora, James Zhang, and Jeffrey L Schnipper. Hospital readmission in general medicine patients: a prediction model. *Journal of general internal medicine*, 25(3):211–219, 2010.
- [HSQ<sup>+</sup>19] Li Huang, Andrew L Shea, Huining Qian, Aditya Masurkar, Hao Deng, and Dianbo Liu. Patient clustering improves efficiency of federated machine learning to predict mortality and hospital stay time using distributed electronic medical records. *Journal of biomedical informatics*, 99:103291, 2019.
- [JPM16] A Johnson, T Pollard, and R. Mark. Mimic-iii clinical database (version 1.4). physionet, 2016. Available from: <https://doi.org/10.13026/C2XW26>.
- [KN21] Lars Kjeldgaard and Lukas Nielsen. Nerda, 2021.

- [KSSU15] Vishesh Kumar, Amber Stubbs, Stanley Shaw, and Özlem Uzuner. Creation of a new longitudinal corpus of clinical narratives. *Journal of biomedical informatics*, 58:S6–S10, 2015.
- [Laf] D Lafky. The safe harbor method of de-identification: an empirical test. fourth national hipaa summit west; 2010.
- [LBS<sup>+</sup>16] Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. Neural architectures for named entity recognition, 2016.
- [LCT<sup>+</sup>15] Zengjian Liu, Yangxin Chen, Buzhou Tang, Xiaolong Wang, Qingcai Chen, Haodi Li, Jingfeng Wang, Qiwen Deng, and Suisong Zhu. Automatic de-identification of electronic medical records using token-level and character-level conditional random fields. *Journal of Biomedical Informatics*, 58:S47–S52, 2015. Supplement: Proceedings of the 2014 i2b2/UTHealth Shared-Tasks and Workshop on Challenges in Natural Language Processing for Clinical Data.
- [LCY10] Thomas Lavergne, Olivier Cappé, and François Yvon. Practical very large scale crfs. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 504–513, 2010.
- [Lev03] Jason Michael Levine. *De-identification of ICU patient records*. PhD thesis, Massachusetts Institute of Technology, 2003.
- [LMP01] John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning, ICML '01*, pages 282–289, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc.
- [LOG<sup>+</sup>19] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.
- [LTWC17] Zengjian Liu, Buzhou Tang, Xiaolong Wang, and Qingcai Chen. De-identification of clinical notes via recurrent neural network and conditional random field. *Journal of Biomedical Informatics*, 75, 06 2017.
- [LVF<sup>+</sup>20] Hang Le, Loïc Vial, Jibril Frej, Vincent Segonne, Maximin Coavoux, Benjamin Lecouteux, Alexandre Allauzen, Benoît Crabbé, Laurent Besacier, and Didier Schwab. Flaubert: Unsupervised language model pre-training for french, 2020.
- [LYK<sup>+</sup>19] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. Biobert: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, Sep 2019.
- [MKGV07] Ashwin Machanavajjhala, Daniel Kifer, Johannes Gehrke, and Muthuramakrishnan Venkatasubramanian. l-diversity: Privacy beyond k-anonymity. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1):3–es, 2007.

- [MMOS<sup>+</sup>20] Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, and Benoît Sagot. Camembert: a tasty french language model. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.
- [NGL<sup>+</sup>14] Aurélie Névéol, Cyril Grouin, Jeremy Leixa, Sophie Rosset, and Pierre Zweigenbaum. The QUAERO French medical corpus: A resource for medical entity recognition and normalization. In *Proc of BioTextMining Work*, pages 24–30, 2014.
- [NKK<sup>+</sup>18] Hiroki Nakayama, Takahiro Kubo, Junya Kamura, Yasufumi Taniguchi, and Xu Liang. doccano: Text annotation tool for human, 2018. Software available from <https://github.com/doccano/doccano>.
- [NRG<sup>+</sup>18] Aurélie Névéol, Aude Robert, Francesco Grippio, Claire Morgand, Chiara Orsi, Laszlo Pelikan, Lionel Ramadier, Grégoire Rey, and Pierre Zweigenbaum. Clef ehealth 2018 multilingual information extraction task overview: Icd10 coding of death certificates in french, hungarian and italian. In *CLEF (Working Notes)*, pages 1–18, 2018.
- [NRR<sup>+</sup>13] Joel Nothman, Nicky Ringland, Will Radford, Tara Murphy, and James R. Curran. Learning multilingual named entity recognition from wikipedia. *Artificial Intelligence*, 194:151–175, 2013. Artificial Intelligence, Wikipedia and Semi-Structured Resources.
- [PdGS21] Marco Polignano, Marco de Gemmis, and Giovanni Semeraro. Comparing transformer-based NER approaches for analysing textual medical diagnoses. In Guglielmo Faggioli, Nicola Ferro, Alexis Joly, Maria Maistro, and Florina Piroi, editors, *Proceedings of the Working Notes of CLEF 2021 - Conference and Labs of the Evaluation Forum, Bucharest, Romania, September 21st - to - 24th, 2021*, volume 2936 of *CEUR Workshop Proceedings*, pages 818–833. CEUR-WS.org, 2021.
- [PKSK17] Fabian Prasser, Florian Kohlmayer, Helmut Spengler, and Klaus A Kuhn. A scalable and pragmatic method for the safe sharing of high-quality health data. *IEEE journal of biomedical and health informatics*, 22(2):611–622, 2017.
- [SDM<sup>+</sup>20] Pedro Javier Ortiz Suárez, Yoann Dupont, Benjamin Müller, Laurent Romary, and Benoît Sagot. Establishing a new state-of-the-art for french named entity recognition. *CoRR*, abs/2005.13236, 2020.
- [SKU15] Amber Stubbs, Christopher Kotfila, and Özlem Uzuner. Automated systems for the de-identification of longitudinal clinical narratives: Overview of 2014 i2b2/uthealth shared task track 1. *Journal of biomedical informatics*, 58:S11–S19, 2015.
- [Swe96] Latanya Sweeney. Replacing personally-identifying information in medical records, the scrub system. In *Proceedings of the AMIA annual fall symposium*, page 333. American Medical Informatics Association, 1996.
- [ULS07] Özlem Uzuner, Yuan Luo, and Peter Szolovits. Evaluating the state-of-the-art in automatic de-identification. *Journal of the American Medical Informatics Association*, 14(5):550–563, 2007.
- [USLS08] Özlem Uzuner, Tawanda C. Sibanda, Yuan Luo, and Peter Szolovits. A de-identifier for medical discharge summaries. *Artificial intelligence in medicine*, 42 1:13–35, 2008.

[VSP<sup>+</sup>17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017.

[Wik21] Wikipedia contributors. Icd-10 — Wikipedia, the free encyclopedia, 2021. [Online; accessed 2-January-2022].