



# Preserving medical correctness, readability and consistency in de-identified health records

**Kostas Pantazos\***, Soren Lauesen and  
**Soren Lippert**

IT-University of Copenhagen, Denmark

## Abstract

A health record database contains structured data fields that identify the patient, such as patient ID, patient name, e-mail and phone number. These data are fairly easy to de-identify, that is, replace with other identifiers. However, these data also occur in fields with doctors' free-text notes written in an abbreviated style that cannot be analyzed grammatically. If we replace a word that looks like a name, but isn't, we degrade readability and medical correctness. If we fail to replace it when we should, we degrade confidentiality. We de-identified an existing Danish electronic health record database, ending up with 323,122 patient health records. We had to invent many methods for de-identifying potential identifiers in the free-text notes. The de-identified health records should be used with caution for statistical purposes because we removed health records that were so special that they couldn't be de-identified. Furthermore, we distorted geography by replacing zip codes with random zip codes.

## Keywords

anonymity, consistency, correctness, de-identification, electronic health records, readability

## Introduction

Electronic health record (EHR) systems store large amounts of data and are essential for all clinical work. According to ANSI,<sup>1</sup> important qualities of an EHR are confidentiality and accessibility only by authorized persons. An EHR system must ensure confidentiality since exposing health records are against law and ethical principles. In order to create data for testing EHR systems, for presenting it to others and for teaching, access is needed to large amounts of EHR data, but it is hard to get the necessary permissions. Access to de-identified (anonymized) health records would in many cases be sufficient. However, the de-identified data should meet certain quality criteria: **IAQ: 1**

---

\*Died October 2015.

## Corresponding author:

Soren Lauesen, IT-University of Copenhagen, Rued Langgaards Vej 7, DK-2300 Copenhagen S, Denmark.

Email: slausen@itu.dk

- 1 1. *Medical correctness*: Each health record must show a true medical picture of a real patient.
- 2 2. *Anonymity*: It must not be possible to see who the real patient is.
- 3 3. *Readability*: The health record must look real. As an example, patient names and addresses
- 4 that have become *F274, XXXX* or *\*\*\*\** don't meet this criterion.
- 5 4. *Consistency*: The patient's identifiers must be consistent with the medical picture, for
- 6 instance, an age that is like the real person's age. If the real patient's name is Peter, but his
- 7 de-identified name is Jens, then Peter must be replaced by Jens also in the clinician's free-
- 8 text notes. Furthermore, his wife's health record may refer to him as Peter, and this Peter
- 9 must also be changed to Jens.

10  
11 A health record database contains fields with structured data that can identify the patient,  
12 such as patient ID, patient name and phone number. These data are fairly easy to replace with  
13 other identifiers, in that way ensuring anonymity. The database also contains fields with medi-  
14 cal data such as diagnosis codes, blood pressure and other measured values. They have to be  
15 preserved to ensure medical correctness. The problem is the unstructured data fields with doc-  
16 tors' free-text notes. They contain important medical information that has to be preserved, but  
17 may also contain phone numbers, patient names, name of the spouse and other identifying  
18 items. Furthermore, clinicians write notes in an abbreviated—often personal—style that cannot  
19 be analyzed grammatically.

20 Considerable work<sup>2–12</sup> has been done in developing de-identification algorithms using various  
21 techniques such as natural language processing (NLP), named entity recognition and machine  
22 learning. These approaches de-identify database records (e.g. pathology documents) that do not  
23 relate to other records. We will refer to this as a record-oriented de-identification approach. Recent  
24 work<sup>13–15</sup> has focused on utilizing a full database rather than records. The approach presented in  
25 this article was briefly presented by Pantazos et al.<sup>16</sup>

26 Previous research has not looked at quality attributes for database-oriented de-identifications. In  
27 this article, we focus on the four quality attributes above: medical correctness, anonymity, readabil-  
28 ity and consistency.

## 30 **Background**

31  
32 In 2010, we started work on an EHR system with a high degree of data visualization. We cooper-  
33 ated with a Danish software house that had delivered EHR systems to many clinics and small  
34 hospitals in Denmark. In order to get test data, we made a copy of the full database and de-identi-  
35 fied it. The database consisted of 437,164 patient health records. The work took place on their  
36 premises since no real health records could go outside the company.

37 The idea was to make a mapping table that translated all patient identifiers into patient identi-  
38 fiers for other patients. In this way, patient B got patient C's first name, patient D's last name,  
39 patient E's street name and so on. Patient B would get a randomized civil registration number  
40 (CPR) that preserved his year of birth and gender. In this way, we would ensure consistency across  
41 patients. Somebody looking at a full de-identified patient record would know that this was a real  
42 patient, but he or she was not called C, nor D and didn't have address E. We had outlined the con-  
43 version program and expected the whole thing to take a couple of days, but—alas—unexpected  
44 problems turned up. We spent 3 months.

45 We had to invent many methods for locating and anonymizing potential identifiers in the free-  
46 text notes. To our surprise, 3–4 percent of the words in free-text fields were potential patient identi-  
47 fiers. Consistency across patients turned up to be more important than we had expected: around  
90 percent of the patients had one or more relatives in the database.

1 We detected that many health records had been created for testing or were left uncompleted due  
2 to system errors (“corrupt” data). In 69,914 cases, we had to delete such patient records. In 43,119  
3 cases, data couldn’t be safely de-identified without manual intervention, so we made the program  
4 delete these patient records. As an example, we deleted patients that had a rare Danish name that  
5 was also the name of a disease, for instance, Aaron, which is also a medical term. The program  
6 couldn’t tell whether Aaron in a free text was a medical term that shouldn’t be changed or the name  
7 of a related patient that should. We ended up with 323,122 patient health records.

8 We manually compared 369 random, anonymized patient records against the original records,  
9 checking for medical correctness, readability and anonymity. These quality factors were preserved  
10 on an acceptable level for our purpose. Consistency was ensured by the algorithm, so we checked  
11 it in a few places only.

## 12 **Related work**

13 A good de-identification system must replace all data that are personal identifiers in structured  
14 data, as well as in free text.<sup>5</sup> One of the first de-identification systems for patient records was  
15 Scrub.<sup>2</sup> It was evaluated against 275 English patient records and 3198 letters to physicians from the  
16 pediatric department. External sources, predefined templates and rules (e.g. the format of a phone  
17 number and address) were included in the algorithm. This algorithm had a 99–100 percent success  
18 rate for de-identifying personal identifiers. Another system was developed by Ruch et al.<sup>3</sup> It resem-  
19 bles Scrub, but added NLP. NLP tools use a medical semantic dictionary with word-sense and  
20 morph-syntactic labeling. This system located 98–99 percent of all personal identifiers. To  
21 anonymize the data, the authors replaced all identifiers with XXX’s, which had a negative impact  
22 on readability.

23 Several systems<sup>4–6,17</sup> were developed in the last decade to de-identify pathology reports. Thomas  
24 et al.<sup>4</sup> developed an algorithm that scored 98.7 percent successful name replacements using English  
25 syntactic rules, prefixes, suffixes and names composed of first and last names. Gupta et al.<sup>6</sup> con-  
26 ducted an iterative evaluation of their system. At the end of the third iteration, the authors claimed  
27 that their method generated anonymized and readable reports. An algorithm designed by Berman<sup>5</sup>  
28 replaced words with codes from the Unified Medical Language System (UMLS) and asterisks. It  
29 produced hardly-readable documents. Beckwith et al.<sup>17</sup> evaluated an open source system which  
30 replaced identifiers with X’s in pathology reports.

31 Seven de-identification systems were evaluated in the “Challenge in NLP for clinical data”  
32 workshop, using medical discharge letters as input.<sup>18</sup> In this workshop, the systems were evaluated  
33 using three performance measures: precision, recall and f-measure. The highest f-measure was  
34 99.7534 achieved by a novel approach based on Named Entity Recognition combined with itera-  
35 tive machine learning.<sup>8</sup> This application finds personal identifiers in the structured data and uses  
36 them to locate identifiers in free-text data.

37 Hanauer et al.<sup>12</sup> introduced the iterative tag-a-little, learn-a-little approach for a particular docu-  
38 ment type. The authors used the MITRE Identification Scrubber Toolkit<sup>19</sup> to integrate their  
39 approach. They obtained an f-measure of 95.

40 Susilo and Win<sup>20</sup> present a new approach for patient confidentiality that utilizes searching through  
41 encrypted data. Huang et al.<sup>21</sup> focused on portable EHRs for privacy preservation. The authors stress  
42 the feasibility of the approach, which can meet patient confidentiality requirements.

43 Even though most of the research has been in an English context, there are some studies on de-  
44 identifying in other contexts. Tveit et al.<sup>22</sup> present their approach to de-identify Norwegian general  
45 practitioner medical records. Their approach consists of six steps: create dictionaries, find exact  
46 match and tag, identify approximate match and tag, replace tags, tackle untagged words and

1 generate the de-identified output. However, this approach was not evaluated empirically. A Swedish  
2 de-identification system was developed by Kokkinakis and Thurin,<sup>7</sup> using named entity recogni-  
3 tion. This approach de-identified 200 Swedish discharge letters with a precision of 96.97 percent,  
4 recall of 89.35 percent and f-measure of 93. Velupillai et al.<sup>10</sup> adjusted an English de-identification  
5 system for Swedish medical records. This transformation did not produce the expected results  
6 (f-measure in total=65, f-measure for names=80). Consequently, the authors reported that build-  
7 ing a Swedish system from scratch was more efficient. This phenomenon was also observed and  
8 confirmed by Grouin et al.<sup>9</sup> who adjusted a de-identification system from English to French and  
9 obtained poor results.

10 Meystre et al.<sup>23</sup> reviewed recent de-identification algorithms and found that the majority of the  
11 algorithms focus on de-identifying structured data and not free text. However, in accordance with  
12 Dalianis and Velupillai<sup>24</sup> and Hanauer et al.,<sup>12</sup> there is immense valuable information in the free  
13 text. We found the same in our data.

## 14 **Quality factors**

15  
16 An EHR contains database fields with structured data that can identify the patient, for instance,  
17 CPR, patient name and phone number. Other structured data fields contain medical data such as  
18 diagnosis codes and blood pressure. The EHR also contains free-text fields, for instance, doctor's  
19 notes and discharge letters. It may also contain pictures of body parts, X-ray and so on, usually  
20 with a patient ID embedded in the picture. We have not dealt with pictures in this project.

21  
22 Some data are *quasi-identifiers* because they can narrow down the set of patients that might  
23 have this health record. Examples are street name, zip code, birth date, hospital or clinician who  
24 treated the patient. Two or more quasi-identifiers in combination may identify the patient.<sup>25</sup>

## 25 **Anonymity**

26  
27 In order to ensure anonymity, all patient identifiers and quasi-identifiers must be de-identified, that  
28 is, replaced with something else. It is fairly easy to do this for structured data, but very hard for  
29 free-text data. Often the computer has no way to tell whether a free-text word is an everyday word,  
30 a medical term or part of a patient name. As an example, Aaron's sign is a medical term, but it  
31 might also be the name of a person.

## 32 **Readability**

33  
34 In order to ensure readability, we have to replace the patient name with a new name that looks real.  
35 Inside this patient's record, we have to be consistent so that we replace with the same name for all  
36 occurrences.

## 37 **Consistency**

38  
39 In the database we worked with, 90 percent of the patients had one or more relatives in the data-  
40 base. Most likely, the patient's name and/or CPR will occur in one of these related health records  
41 in free-text fields. To ensure consistency, we have to replace also these identifiers with the same  
42 new identifier.

43  
44 There are other aspects of consistency, for instance, that the distribution of names should remain  
45 much the same. If rare names suddenly turn up for a large number of patients, the health record  
46 database will not look real.

## 1 *Medical correctness*

2  
3 If we replace a medical term that looks like a person name, with the new person name, the health  
4 record will look odd. We have lost medical correctness and readability. In many cases, a clinician  
5 can guess what the medical term was and in that way get to know the original name of all patients  
6 that have this new name.

7 Another aspect of medical correctness is age. If birth dates are transformed in a way that makes  
8 the patient have a very different age, it will not match the patient's diagnosis pattern.

## 10 **Solution**

11 We will first give an overview of the solution and then explain the details and where the data came  
12 from.

### 15 *Permutation tables*

16 For some identifiers, we made a *permutation table* that mapped existing identifiers to new ones.  
17 We picked the new identifier at random from the same table, avoiding reuse of identifiers. Any  
18 occurrence of an identifier from this table would be translated into the new one.

19 As an example, we created a permutation table of all last names. The last name Jensen would be  
20 translated into Petersen wherever it occurred. Petersen was another last name in the table, with a  
21 similar frequency. This ensured readability and consistency across all patient records.

22 We made permutation tables for these identifiers and quasi-identifiers: first male names, first  
23 female names, last names, street names, zip codes, hospital and clinic names.

### 26 *Distorted identifier table*

27 For the CPR, we made a mapping table from existing CPR to a distorted CPR in this way: The  
28 Danish CPR format is: DDMMYY-CSSG where DDMMYY is the birth date. The day (DD) and  
29 month (MM) were changed to a random valid day and month. The year (YY) was kept. C indicates  
30 birth century (1900 or 2000). This was not changed. SS (serial number) was changed, while G  
31 indicates the gender and wasn't changed.

32 This ensured readability (clinical users see lots of CPR numbers and can easily spot wrong  
33 ones) and medical correctness (because age and gender were kept).

### 36 *Randomized identifier*

37 For other identifiers, we randomized the identifier without caring about readability or consistency.  
38 This applied to phone numbers, e-mail addresses and URLs.

### 41 *Ambiguous words*

42 Ambiguous words could be part of a person's name or something else, for instance, a medical term  
43 or a common word. Through many sources, we created a list of ambiguous words. When the de-  
44 identification program meets an ambiguous word B in a free-text field, it has three choices:

- 46 1. Replace the word B with its corresponding new name, C. If the word B actually is part of a  
47 person's name, everything is fine. But if B actually is a medical term or a common word,

1 the clinician can see from the context that C probably means B. If he knows the replace-  
2 ment rules, he now knows that everybody in the database with the name C is actually B. We  
3 lose not only medical correctness and readability but also some anonymity. For this reason,  
4 we never replace ambiguous words.

- 5 2. Keep the old word B. This ensures medical correctness, readability and consistency. If the  
6 word actually is a person's name, the clinician can see it from the context. As an example,  
7 assume that the program finds Aaron in a free-text note. Since it is a medical term, it keeps  
8 it. However, a clinician can see that Aaron in this context is the name of a person. If he  
9 knows the rule of replacement, he now knows that the person referred to is really called  
10 Aaron, although this is not his name in the de-identified database. The clinician gets no clue  
11 to where Aaron's health record is. If there are only a few Aarons in real life, he might guess  
12 whom it is. If there are many Aarons, he cannot know. We decided that 200 occurrences  
13 was a safe limit. If the ambiguous name occurs more than 200 times, we keep it in the  
14 database.
- 15 3. Delete all patient records with name B. We do this when the ambiguous name occurs less  
16 than 200 times. This ensures all four quality factors, but we lose data. If a free text for  
17 another patient refers to patient B, the reference will now be to a deleted patient. We have  
18 lost a bit of consistency, but such data could exist anyway in the database.

## 20 **The database and the mapping tables**

21  
22 The EHR we de-identified is built on Microsoft Axapta, which is an ERP system that can be  
23 extended in many ways. It contained data from 79 clinics and hospitals (including a few in  
24 Greenland and the Faroe Islands) and contained 437,164 patient records in total. The entire data-  
25 base was 12GB. There were 65 health-related tables:

- 26  
27 1. 43 tables had no fields that could expose the patient identity. They included reference tables  
28 of drug codes, treatment codes and diagnosis codes.
- 29 2. 9 tables had fields that only contained personal identifiers in structured form, for instance,  
30 the patient table that contained patient ID, first name, last name, address, zip code, five  
31 phone numbers, birth date and date of death. Another example is a table of family relations,  
32 that is, relations between two patients. Clinicians, hospitals and clinics had their own tables  
33 with name, address and so on.
- 34 3. 13 tables had fields with free text. The largest one was Medical Record Lines, which occu-  
35 pied 7 of the 12GB in the database.

## 37 **Mapping tables**

38  
39 To be able to replace existing identifiers with new identifiers, we created the following mapping  
40 tables.

41  
42 *CPR.* We collected the CPR numbers from the patient table and gave each number a partially ran-  
43 dom new number according to the rules above. If the new number was already used as a new  
44 number, we randomized it once more.

45  
46 *Last names.* We used three sources to collect last names: the database's patient table, Danmarks  
47 Statistik's website<sup>26</sup> and a study of Danish names at University of Copenhagen, 2005.<sup>27</sup> We merged

1 these sources and obtained 56,339 last names. We counted how often each last name occurred in  
2 the patient table. Many names didn't occur at all in the patient table, but might occur in free-text  
3 notes. It was important to catch them too and de-identify them.

4 For each name, we assigned a new name from the table with a frequency similar to the old  
5 name. We used this approach: we divided the names into groups according to frequency. The group  
6 of most frequent last names consisted of 20 names with frequencies from 14,712 to 5319. We  
7 rotated these names a random number of steps to obtain the new names (a cyclical permutation).  
8 We used the same approach for groups of 30 names with decreasing frequencies. Rare names (fre-  
9 quency < 200) were randomly replaced with another name in the frequent part of the list. This also  
10 took care of the names that didn't occur at all in the patient table.

11  
12 *Male first names.* We used the same approach to collect and de-identify male first names. For  
13 names occurring in the patient table, we got the gender from the CPR number. Our external sources  
14 had separate lists for male and female first names. In total, we got 11,415 male first names.

15  
16 *Female first names.* We treated them in the same way as male names. In total, we got 13,044 female  
17 first names.

18  
19 *Street names.* We collected street names from the patient table's address field. The address field  
20 included also floor numbers and entrance letters. In Denmark, the street name is first, so we simply  
21 extracted the first real name from the address field. We also included street names from the CPR  
22 website. In total, we got 25,429 street names. We assigned a random street name as the new name  
23 without caring about frequencies or consistency with zip codes.

24  
25 *Zip codes.* We collected zip codes and related city names from Post Danmark<sup>28</sup> and assigned a  
26 random zip code and city name as the new name. In total, we got 1396 zip codes.

27  
28 *Hospital names and clinic names.* We collected hospital names from Region Hovedstaden, Region  
29 Sjælland, Region Syddanmark, Region Midtjylland, Region Nordjylland, Queen Ingrid's Hospital  
30 in Greenland, Faroe Islands website and our own EHR Database. We used Sygehusvalg,<sup>29</sup> Branche-  
31 foreningen for Privathospitaler og Klinikker (the trade association) and our own EHR database to  
32 extract names of clinics. In total, we got a list of 219 clinic names and 93 hospital names. We did  
33 not randomly assign new names to the clinics and hospitals. This would reduce medical correctness  
34 because clinicians know which clinics do what. On the other hand, being treated in a specific clinic  
35 is a quasi-identifier. We manually selected 41 hospital names and 92 clinic names and used them  
36 as new names. In many cases, the new name was simply "Hospital" or "Clinic." This was a reduc-  
37 tion in readability and to some extent in medical correctness.

### 38 39 *Ambiguous names*

40  
41 Ambiguous names in our context are first or last names of persons that happen to mean something  
42 else too. We need a table of them to decide how to treat such a name when it occurs in free text. As  
43 explained above, we have to delete patients with rare names if they appear in free text. If they are  
44 frequent names, we leave them as they are.

45 In healthcare, it is common that diseases, signs, symptoms and so forth are named after a person,  
46 most likely the one who discovered it. These names are called medical eponyms and may cause  
47 ambiguity. For example, according to Statistics Denmark in 2010,<sup>26</sup> there were 88 males using the

1 name Aaron. At the same time, Aaron is part of a medical eponym (Aaron sign). The algorithm  
2 knows too little about the context to decide whether to de-identify this name or not.

3 A similar ambiguity exists also with common words in a language. Each language contains  
4 several words whose meaning depends on the context. For example, in Danish, the word “hans”  
5 can be a pronoun or a male name. Another ambiguous case is abbreviations used by clinicians. For  
6 instance, instead of writing “kirurgisk” (in English: “surgical”) they use the abbreviation “kir,”  
7 which can be a last name as well. As another example, it is common that a city, hospital, clinic or  
8 street name is used as a first or last name. For instance, Aalborg is a city in Denmark, but it is a last  
9 name as well.

10 We derived the table of ambiguous names from several sources. We checked our lists of first and  
11 last names against the Danish Dictionary from Microsoft Office Word 2010. This created a list of  
12 3557 potentially ambiguous names. That a name exists in the dictionary doesn’t mean that it also  
13 has another meaning that can occur in health records. So, the medical specialist in our team  
14 (Lippert) scrutinized the list and came up with 1952 ambiguous names.

15 To the best of our knowledge, there is no official source that contains medical eponyms. So, we  
16 used the website “Who named it”<sup>30</sup> and extracted 3246 medical eponymous names from it. They  
17 too entered the list of ambiguous names.

## 18 19 **Applying the mappings**

20  
21 The mappings must be applied to the structured fields as well as to the free-text fields. We applied  
22 the mappings to the structured fields according to Table 1. Notice the last rule: remove all patients  
23 above 90 years. It came from the US Health Insurance Portability and Accountability Act guide-  
24 lines, HIPAA.<sup>31</sup> There are so few patients above 90 that their age exposes them.

25 It was harder to apply the mappings in free-text fields because we don’t know whether an identi-  
26 fier is a first name, a street name and so on. The program analyzed the free-text token by token and  
27 applied these rules.

### 28 29 *Name tokens*

- 30  
31 1. If the name is in one of the person name tables and also in the table of ambiguous words,  
32 do nothing or delete the related patient records depending on the name frequency.
- 33 2. If the name is in the last name table, replace it with the new name in the table.
- 34 3. If the name is in the first male name table, replace it with the new name.
- 35 4. If the name is in the first female name table, replace it with the new name.
- 36 5. If the name is in the table of ambiguous words, leave it as it is.
- 37 6. If the name is in the table of street names, replace it with the new name.
- 38 7. If the name is in the table of zip codes and city names, replace it with the new city name.
- 39 8. If the name is in the table of hospitals and clinics, replace it with the new name.
- 40 9. Otherwise, leave it as it is.

### 41 42 *Number tokens*

- 43  
44 10. If the number is in the table of CPR numbers, replace it with the new CPR number.
- 45 11. If the number looks like a CPR number (10 digits starting with a date), randomize it as other  
46 CPR numbers.
- 47 12. If the number has eight digits and is next to a word like tlf, tel and fax, randomize the  
number.



**Table 1.** Replacement rules for structured data fields.**Identifying fields**

Civil registration number (CPR)	Replace it with the new CPR in the CPR mapping table
First name	Select the first male or first female mapping table according to the gender code in CPR. Replace first name with the new name in the mapping table
Last name	Replace it with a new name according to the mapping table
Address	An address contains a street name, a house number and sometimes a floor number and entrance position (e.g. Byevej 21, 2tv). Replace the street name according to the street mapping table. Replace numbers randomly with a number that has the same number of digits
Phone numbers (up to five per patient)	Alter each phone number to a random number with the same number of digits
E-mail	Alter the address with random characters before the letter @ and change the domain name to <i>email.dk</i>

**Quasi-identifiers**

Zip code	Replace it according to the zip mapping table
City	Replace it with the city name in the zip mapping table
Country	Change it to <i>Denmark</i>
Date of birth	Set it from the new CPR
Date of death	Randomly change the day and month
Hospital name	Replace it with a new name according to the mapping table
Clinic name	Replace it with a new name according to the mapping table
Clinician first name	Replace it with a new name according to the mapping table for first names
Clinician last name	Replace it with a new name according to the mapping table for last names
Clinician alias	Replace it with the new first name of the clinician
Age	Remove all patients older than 90 years due to high anonymity risks

13. If the number is in the table of zip codes and next to a city name, replace it with the new zip code.

14. Otherwise, leave it as it is. (It may be a measured value, a lab-test number (eight digits), a house number and so on.)

These rules give priority to anonymity rather than medical correctness. As an example, a lab-test number or a date-time that looks like a CPR number will be de-identified and thus reduce the medical correctness.

## Evaluation of the quality factors

### *Anonymity, readability and medical correctness*

In order to evaluate the actual anonymity, readability and medical correctness, we need to know how many words were replaced incorrectly.

We selected a random sample of 369 full patient records. A clinician manually compared all the free-text fields in the old and the new version, in total 73,150 words. The result is shown in Table 2.

**Table 2.** Correct and incorrect replacements.

Number of words	Should be de-identified	Should not	Total
Was de-identified	1313	109	1422
Was not	7	71,721	71,728
Total	1320	71,830	73,150

Seven words should have been de-identified but wasn't. Only one of them was a person name. It was ambiguous and frequent (frequency >200) and consequently preserved according to our rules. Since it was frequent, we consider it a quasi-identifier. The other words were quasi-identifiers such as department names and misspelled street names that were not in our translation tables. In total, out of 73,150 words, we had seven anonymity leaks on quasi-identifiers and none on full identifiers.

A total of 109 words were replaced, but shouldn't. They were ambiguous, but not in our table of ambiguous words. One example was the word "Uno," which was the name of a drug, but also a male first name. These cases decreased medical correctness and readability. It also revealed a general weakness: also drug names should be considered a source of ambiguity.

Measured in the traditional way with *recall* and *precision*, the algorithm scored 99.5 percent for recall (the seven anonymity leaks) and 92.3 percent for precision (the 109 leaks in medical correctness). The f-measure was 95.7 percent. Our database-oriented approach compares favorably with previous work on record-based de-identification approaches.<sup>4-6,17</sup>

It would be interesting to compare how other de-identification approaches would handle our data. However, this is impossible because the approaches are very dependent on the language. Furthermore, we are not allowed to move our original data out of the company where it is hosted. We have not found publications about de-identification that discuss ambiguity. Most likely, they don't pay attention to it. It will probably cause some leaks of confidentiality that isn't detected.

### Consistency

The database can record family relations and other relations between patients. Around 90 percent of the patients have one or more recorded relatives. When a person name is de-identified in the structured patient table, it is important that the same name is de-identified in the same way in the rest of the patient's records and in records of related patients, also for free-text fields. This is solely a matter of correct programming. We checked it for a couple of patients in Table 3. Since the translation tables are used for all patients, consistency is also preserved for relatives who are not recorded as relatives.

### Results

Table 3 shows a (non-random) sample of patients with two or more relatives. It gives an impression of the variety and complexity of patient records. Several patients have eight relatives in the database, many have more than 100 measurements with notes (Clinical Data), many have more than 10 diagnoses and several hundred prescriptions.

Table 4 shows the results for the Medical Record Line and Clinical Data tables, which contain most of the free text in the database. In total, 3-4 percent of the words are personal identifiers.

This study is the first de-identification algorithm that focuses on anonymity, medical correctness, readability and consistency. Other approaches are limited to a few types of documents, while

**Table 3.** Sample of 20 patients showing the variety of patient records.

CPR	Relatives	Clinical data	Medical records	Diagnoses	Prescriptions	Total
2905931069	6	54	4	4	6	68
2904220702	2	335	9	3	678	1025
2812620120	4	37	2	42	177	258
2812351528	2	36	1	54	517	608
2811831753	2	30	1	1	18	50
2810291211	2	22	2	13	68	105
2809711115	4	15	4	9	15	43
2809550048	6	151	6	2	32	191
2808972414	8	50	6	9	7	72
2806492477	4	603	10	22	412	1047
2805832168	4	42	1	11	62	116
2805620030	4	176	5	1	29	211
2803961559	2	6	3	2	5	16
2801981465	8	76	4	1	1	82
2801460257	6	29	3	1	22	55
2712742278	6	186	8	7	64	265
2711743812	4	77	9	14	51	151
2711440133	2	98	4	11	100	213
2710592476	8	238	3	10	38	289
2709530059	4	22	1	1	9	33

CPR: civil registration number.

**Table 4.** Number of identifiers in free text.

	Medical record line	Clinical data
E-mails	18,858	727
Phone numbers	43,051	62,461
Clinics	114,318	17,213
CPRs	455,946	121,036
Zip codes	599,566	668
Hospitals	787,055	117,369
Cities	994,125	7557
Last names	2,675,386	254,915
Street names	3,156,356	125,470
First names	4,331,593	330,679
Total identifiers	(4%) 13,176,254	(3%) 1,038,095
Non-identifiers	322,734,954	32,052,044

CPR: civil registration number.

Right align both columns

our approach deals with full EHR records from 79 hospitals and clinics. An important part of our approach was to collect ambiguous names from many sources.

We started out with 437,164 patient health records. We deleted 69,914 patient records because data were corrupted (old test data and records left after system failures). We deleted 43,119 patient

1 records because of rare ambiguous names or because the patient was older than 90. We ended up  
2 with 323,122 patient health records.

3 The distinction between frequent and rare names (fewer than 200 occurrences) is somewhat arbi-  
4 trary. The limit of 200 caused us to delete “only” 43,119 patient records because they had rare ambi-  
5 guous names. If all names were considered rare, we would have lost another 55,000 patient records.

6 We made a manual review of 369 patient records with 71,721 free-text words. It revealed seven  
7 words where a quasi-identifier hadn’t been de-identified. It revealed 109 words where it was de-  
8 identified, but shouldn’t because the word wasn’t in our list of ambiguous names. This reduced  
9 medical correctness and readability.

## 11 *Limitations and errors*

13 An EHR database contains also binary files (e.g. X-rays), scanned documents and Word docu-  
14 ments. They are not part of the database, but fields in the database contain the file names. Our  
15 approach is limited to structured data and free-text fields, and it doesn’t try to de-identify pictures  
16 and other files. The picture will usually contain patient identifiers such as CPR and name.  
17 De-identifying these would be a project of its own.

18 We have not tried to deal with spelling errors. It might have reduced the seven un-identified  
19 words above to around three. We could deal with spelling errors by looking at close matches of  
20 words instead of precise matches, but we don’t know how much it would have increased the num-  
21 ber of false de-identifications (the 109 words above).

22 We forgot to put also pharmaceutical names in the list of ambiguous words. This could have  
23 removed some of the 109 false de-identifications above.

24 We missed several clinical abbreviations as potential ambiguous names. A language analysis of  
25 the free-text notes might have revealed them.

26 The de-identified data should be used with caution for statistical purposes because of the way  
27 we had to remove health records that couldn’t be de-identified and also because we deleted patients  
28 older than 90 and distorted geography by replacing zip codes with random zip codes.

29 For statistical purposes, the de-identification should have been different. We shouldn’t care  
30 about readability or consistency, but simply replace all potential identifiers in free text with aster-  
31 isks or the like. We should only delete corrupted patient records. The mapping tables would still be  
32 needed, but only to detect what might be an identifier. We wouldn’t need to care about ambiguous  
33 words. The result would probably be similar to many other de-identification approaches.

## 35 **Funding**

36 The author(s) received no financial support for the research, authorship and/or publication of this article.

## 38 **References**

- 39 1. ISO/TR 20514:2005. Health informatics—electronic health record definition, scope and context.
- 40 2. Sweeney L. Replacing personally-identifying information in medical records, the scrub system. In:  
41 *Proceedings of the AMIA annual fall symposium*, Washington, DC, 1996, p. 333. Bethesda, MD:  
42 American Medical Informatics Association. **[AQ: 2]**
- 43 3. Ruch P, Baud RH, Rassinoux AM, et al. Medical document anonymization with a semantic lexicon. In:  
44 *Proceedings of the AMIA symposium*, Los Angeles, CA, 4–8 November 2000, p. 729. Bethesda, MD:  
45 American Medical Informatics Association. **[AQ: 3]**
- 46 4. Thomas SM, Mamlin B, Schadow G, et al. A successful technique for removing names in pathology  
47 reports using an augmented search and replace method. In: *Proceedings of the AMIA symposium*,  
San Antonio, TX, 9–13 November 2002, p. 777. Bethesda, MD: American Medical Informatics  
Association. **[AQ: 4]**

- 1 5. Berman JJ. Concept-match medical data scrubbing: how pathology text can be used in research. *Arch*
- 2 *Pathol Lab Med* 2003; 127(6): 680–686.
- 3 6. Gupta D, Saul M and Gilbertson J. Evaluation of a de-identification (de-id) software engine to share
- 4 pathology reports and clinical documents for research. *Am J Clin Pathol* 2004; 121(2): 176–186.
- 5 7. Kokkinakis D and Thurin A. Identification of entity references in hospital discharge letters. In:
- 6 Proceedings of the 16th Nordic conference of computational linguistics, Tartu, 2007. **[AQ: 5]**
- 7 8. Szarvas G, Farkas R and Busa-Fekete R. State-of-the-art anonymization of medical records using an
- 8 iterative machine learning framework. *J Am Med Inform Assoc* 2007; 14(5): 574–580.
- 9 9. Grouin C, Rosier A, Dameron O, et al. Testing tactics to localize de-identification. In: **XXII European**
- 10 **congress of medical informatics, Sarajevo, 2009, pp. 735–739. [AQ: 6]**
- 11 10. Velupillai S, Dalianis H, Hassel M, et al. Developing a standard for de-identifying electronic patient
- 12 records written in Swedish: precision, recall and f-measure in a manual and computerized annotation
- 13 trial. *Int J Med Inform* 2009; 78(12): e19–e26.
- 14 11. Dalianis H, Nilsson G and Velupillai S. Is de-identification of electronic health records possible? Or
- 15 can we use health record corpora for research? In: AAAI fall symposium series, Arlington, VA, 5–7
- 16 November 2009.
- 17 12. Hanauer D, Aberdeen J, Bayer S, et al. Bootstrapping a de-identification system for narrative patient
- 18 records: cost-performance trade-offs. *Int J Med Inform* 2013; 82(9): 821–831.
- 19 13. Emam KE, Paton D, Dankar FK, et al. De-identifying a public use microdata file from the Canadian
- 20 national discharge abstract database. *BMC Med Inform Decis Mak* 2011; 11: 53.
- 21 14. Viangteeravat T, Huang EY and Wade G. Giving raw data a chance to talk: a demonstration of de-
- 22 identified pediatric research database (PRD) and exploratory analysis techniques for possible research
- 23 cohort discovery and identifiable high risk factors for re-admission. *BMC Bioinformatics* 2013; 14: A5.
- 24 15. Gordon JS. Altering the function of the electronic medical record: creating a de-identified database for
- 25 clinical researchers and educators. *Nurs Inform* 2012; 2012: 132.
- 26 16. Pantazos K, Lauesen S and Lippert S. De-identifying an EHR database—anonymity, correctness and
- 27 readability of the medical record. *Stud Health Technol Inform* 2011; 169: 862–866.
- 28 17. Beckwith BA, Mahaadevan R, Balis UJ, et al. Development and evaluation of an open source software
- 29 tool for deidentification of pathology reports. *BMC Med Inform Decis Mak* 2006; 6(1): 12.
- 30 18. Uzuner O, Luo Y and Szolovits P. Evaluating the state-of-the-art in automatic de-identification. *J Am*
- 31 *Med Inform Assoc* 2007; 14(5): 550–563.
- 32 19. Aberdeen J, Bayer S, Yeniterzi R, et al. The MITRE identification scrubber toolkit: design, training, and
- 33 assessment. *International J Med Inform* 2010; 79(12): 849–859.
- 34 20. Susilo W and Win K. Security and access of health research data. *J Med Syst* 2007; 31(2): 103–107,
- 35 <http://dx.doi.org/10.1007/s10916-006-9035-y> (accessed August 2010).
- 36 21. Huang LC, Chu HC, Lien CY, et al. Embedding a hiding function in a portable electronic health record
- 37 for privacy preservation. *J Med Syst* 2010; 34(3): 313–320, [http://dx.doi.org/10.1007/s10916-008-9243-](http://dx.doi.org/10.1007/s10916-008-9243-8)
- 38 [8](http://dx.doi.org/10.1007/s10916-008-9243-8) (accessed August 2010).
- 39 22. Tveit A, Edsberg O, Raast TB, et al. Anonymisation of general practitioner’s patient records. In:
- 40 **Proceedings of the HelsIT’04 conference, Trondheim, September 2004. [AQ: 7]**
- 41 23. Meystre SM, Friedlin FJ, South BR, et al. Automatic de-identification of textual documents in the elec-
- 42 tronic health record: a review of recent research. *BMC Medical Res Methodol* 2010; 10(1): 70.
- 43 24. Dalianis H and Velupillai S. De-identifying Swedish clinical text-refinement of a gold standard and
- 44 experiments with conditional random fields. *J Biomedical Semantics* 2010; 1: 6.
- 45 25. El Emam K, Jabbouri S, Sams S, et al. Evaluating common de-identification heuristics for personal
- 46 health information. *J Med Internet Res* 2006; 8(4): e28
- 47 26. Danmarks Statistik, <http://www.dst.dk/HomeUK/Statistics/Names.aspx> (accessed August 2010).
27. Copenhagen University, <http://danskernesnavne.navneforskning.ku.dk/TopNavne.asp> (accessed August
- 2010).
28. Post Denmark, <http://www.postdanmark.dk/> (accessed August 2010). **[AQ: 8]**
29. Sygehusvalg, <http://www.sygehusvalg.dk/> (accessed August 2010).
30. Who Named It, <http://www.whonamedit.com/azeponyms.cfm/> (accessed August 2010).
31. HIPAA. HIPAA privacy rules and public health, [http://www.cdc.gov/mmwr/preview/mmwrhtml/](http://www.cdc.gov/mmwr/preview/mmwrhtml/m2e411a1.htm)
- m2e411a1.htm (accessed August 2010).