

OXFORD
UNIVERSITY PRESSJAMIA: Journal of the
American Medical Informatics Association**De-identification of Patient Notes with Recurrent Neural Networks**

| | |
|---------------|---|
| Journal: | <i>Journal of the American Medical Informatics Association</i> |
| Manuscript ID | amiajnl-2016-005114.R2 |
| Article Type: | Research and Applications |
| Keywords: | De-identification, Natural language processing, Protected Health Information (PHI), HIPAA, Artificial Neural Networks |
| | |

SCHOLARONE™
Manuscripts

De-identification of Patient Notes with Recurrent Neural Networks

Franck Deroncourt^{1*}, Ji Young Lee^{1*}, Ozlem Uzuner², Peter Szolovits¹

¹ Massachusetts Institute of Technology, Cambridge, MA.

² University at Albany, SUNY, Albany, NY.

* These authors contributed equally to this work.

| | | | |
|--------------------|----------------|------------------|-----------------|
| Franck Deroncourt* | Ji Young Lee* | Ozlem Uzuner | Peter Szolovits |
| 32 Vassar St., | 32 Vassar St., | 1400 Washington | 32 Vassar St., |
| 32-293, | 32-253, | Ave., | 32-254, |
| Cambridge, MA | Cambridge, MA | Draper 114A, | Cambridge, MA |
| 02139 | 02139 | Albany, NY 12222 | 02139 |

Corresponding Author: Franck Deroncourt
Email: francky@mit.edu
Tel: +1-443-637-2659

Keywords: Medical Language Processing, De-identification, Neural Networks

Word count: 3,946 words

De-identification of Patient Notes with Recurrent Neural Networks

Abstract

Objective Patient notes in electronic health records (EHRs) may contain critical information for medical investigations. However, the vast majority of medical investigators can only access de-identified notes, in order to protect the confidentiality of patients. In the United States, the Health Insurance Portability and Accountability Act (HIPAA) defines 18 types of protected health information (PHI) that needs to be removed to de-identify patient notes. Manual de-identification is impractical given the size of EHR databases, the limited number of researchers with access to the non-de-identified notes, and the frequent mistakes of human annotators. A reliable automated de-identification system would consequently be of high value.

Materials and Methods We introduce the first de-identification system based on artificial neural networks (ANNs), which requires no handcrafted features or rules, unlike existing systems. We compare the performance of the system with state-of-the-art systems on two datasets: the i2b2 2014 de-identification challenge dataset, which is the largest publicly available de-identification dataset, and the MIMIC de-identification dataset, which we assembled and is twice as large as the i2b2 2014 dataset.

Results Our ANN model outperforms the state-of-the-art systems. It yields an F1-score of 97.85 on the i2b2 2014 dataset, with a recall of 97.38 and a precision of 98.32, and an F1-score of 99.23 on the MIMIC de-identification dataset, with a recall of 99.25 and a precision of 99.21.

Conclusion Our findings support the use of ANNs for de-identification of patient notes, as they show better performance than previously published systems while requiring no manual feature engineering.

1 INTRODUCTION AND RELATED WORK

In many countries such as the United States, medical professionals are strongly encouraged to adopt electronic health records (EHRs) and may face financial penalties if they fail to do so [1,2]. The Centers for Medicare & Medicaid Services have paid out more than \$30 billion in EHR incentive payments to hospitals and providers who have attested to meaningful use as of March 2015. Medical investigations may greatly benefit from the resulting increasingly large EHR datasets. One of the key components of EHRs is patient notes: the information they contain can be critical for a medical investigation because much information present in texts cannot be found in the other elements of the EHR. However, before patient notes can be shared with medical investigators, some types of information, referred to as protected health information (PHI), must be removed in order to preserve patient confidentiality. In the United States, the Health Insurance Portability and Accountability Act (HIPAA) [3] defines 18 different types of PHI, ranging from patient names to phone numbers. Table 1 presents the exhaustive list of PHI types as defined by HIPAA.

The task of removing PHI from a patient note is referred to as de-identification, since the patient cannot be identified once PHI is removed. De-identification can be either manual or automated. Manual de-identification means that the PHI is labeled by human annotators. There are three main shortcomings of this approach. First, only a restricted set of individuals is allowed to access the identified patient notes, thus the task cannot be crowdsourced. Second, humans are prone to mistakes. Neamatullah et al. [4] asked 14 clinicians to detect PHI in approximately 130 patient notes: the results of the manual de-identification varied from clinician to clinician, with recall ranging from 0.63 to 0.94. Third, human annotation is costly. Douglass et al. [5,6] reported that annotators were paid US\$50 per hour and read 20,000 words per hour at best.

Table 1 PHI types as defined by HIPAA, i2b2, and MIMIC. Classification of PHI into categories and types are as defined in the i2b2 dataset. During training, the PHI types are used as the labels to predict. The mark “-” denotes that 2 or fewer instances of the corresponding PHI types are present in the whole dataset, and no instance is present in the test set. In the MIMIC dataset, some PHI types are mapped to a different PHI type due to data ambiguity or sparsity issues: these PHI types are marked with the specific PHI type that it is mapped to instead of the mark “x”.

| PHI categories | PHI types | Descriptions | HIPAA | i2b2 | MIMIC |
|----------------|----------------|--|-------|------|----------------|
| AGE | AGE | Ages \geq 90 | x | x | x |
| | | Ages < 90 | | x | |
| CONTACT | PHONE | Telephone numbers | x | x | x |
| | FAX | Fax numbers | x | x | PHONE |
| | EMAIL | Electronic mail addresses | x | x | |
| | URL | URLs | x | - | |
| | IPADDRESS | IP addresses | x | - | |
| DATE | DATE | Dates (month and day parts) | x | x | x |
| | | Year | | x | x |
| | | Holidays | | x | x |
| | | Day of the week | | x | |
| ID | IDNUM | Social security numbers | x | x | x |
| | | Account numbers | x | x | x |
| | | Certificate or license numbers | x | x | x |
| | MEDICALRECORD | Medical record numbers | x | x | IDNUM |
| | DEVICE | Vehicle or device identifiers | x | x | IDNUM |
| | HEALTHPLAN | Health plan numbers | x | - | IDNUM |
| | BIOID | Biometric identifiers or full face photographs | x | - | |
| LOCATION | STREET | Street address | x | x | x |
| | CITY | City | x | x | LOCATION-OTHER |
| | ZIP | Zip | x | x | x |
| | STATE | State | | x | x |
| | COUNTRY | Country | | x | x |
| | LOCATION-OTHER | Other identifiable locations such as landmarks | | x | x |
| | ORGANIZATION | Employers | x | x | |
| | HOSPITAL | Hospital name | | x | x |
| | | Ward name | | | x |
| NAME | PATIENT | Names of patients and family members | x | x | x |
| | DOCTOR | Provider name | | x | x |
| | USERNAME | User ID of providers | | x | |
| PROFESSION | PROFESSION | Profession | | x | |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

As a matter of comparison, the MIMIC dataset [7,8], which contains data from 50,000 intensive care unit (ICU) stays, consists of 100 million words. This would require 5,000 hours of annotation, which would cost US\$250,000 at the same pay rate. Given the annotators’ spotty performance, each patient note would have to be annotated by at least two different annotators: it would therefore cost at least US\$500,000 to de-identify the notes in the MIMIC dataset.

In order to reduce the cost of annotating, many studies investigate the use of machine pre-annotation, where human annotators are provided with machine-annotated data to reduce the annotation time. Lingret et al. [9] show that using pre-annotation resulted in 13.85-21.5% of time savings for developing a clinical named entity recognition corpus. However, another study by South et al. [10] show that using a machine pre-annotation along with an interactive annotation tool neither improved the quality nor decreased the time investment when annotating clinical text de-identification corpus.

Instead of annotating all documents at the same time either from raw or pre-annotated texts, Hanauer et al. [11] took a novel approach where annotations are performed alternately between humans and machine. More specifically, the clinical notes are divided into multiple batches of 10, 20, or 40 notes and each batch is annotated sequentially by human annotators after being pre-annotated by a de-identifier trained on previously annotated batches. They show that the annotation time for each instance decreased in later batches as the de-identifier’s performance improved, achieving an F1-score of 0.95 after just over 8 hours of annotation time (after 20 batches of 10 notes each). Similarly, Gobbel et al. [12] present a tool called RapTAT to assist human annotators by pre-annotating the documents interactively while the annotators are working on them, resulting in up to 50% of reduction in annotation time.

Automated de-identification systems can be classified into two categories: rule-based systems and machine-learning-based systems. Rule-based systems typically rely on patterns, expressed as regular

expressions and gazetteers, defined and tuned by humans. They do not require any labeled data (aside from labels required for evaluating the system), and are easy to implement, interpret, maintain, and improve, which explains their large presence in the industry [13]. However, they need to be fine-tuned for each new dataset, are not robust to language changes (e.g., variations in word forms, typographical errors, or infrequently used abbreviations), and cannot easily take into account the context (e.g., “Mr. Parkinson” is PHI, while “Parkinson’s disease” is not PHI). Rule-based systems are described in [4,14–22]. To alleviate some downsides of the rule-based systems, there have been many attempts to use supervised machine learning algorithms to de-identify patient notes. These algorithms are used to train a classifier to label each word as PHI or not PHI, sometimes distinguishing between different PHI types. Common statistical methods include decision trees [23], log-linear models, support vector machines [24–26], and conditional random fields [27]. The latter is employed in most of the state-of-the-art systems. For a thorough review of existing systems, see [28,29]. All these methods share two downsides: they require a decent sized labeled dataset and much feature engineering. As with rules, quality features are challenging and time-consuming to develop.

Recent approaches to natural language processing based on artificial neural networks (ANNs) do not require handcrafted rules or features. Instead, ANNs can automatically learn effective features by performing composition over tokens which are represented as vectors, often called token embeddings. The token embeddings are jointly learned with the other parameters of the ANN. They can be initialized randomly, but can be pre-trained using large unlabeled datasets typically based on token co-occurrences [30–32]. The latter often performs better, since the pre-trained token embeddings explicitly encode many linguistic regularities and patterns. As a result, methods based on ANNs have shown promising results for various tasks in natural language processing, such as language modeling [33], text classification [34–37], question answering [38,39], machine translation [40–42], as well as named entity

recognition [31,43,44]. A few methods also use vector representations of characters as inputs in order to either replace or augment token embeddings [43-45].

Inspired by the performance of ANNs for various other NLP tasks, this article introduces the first de-identification system based on ANNs. Unlike other machine learning based systems, ANNs do not require manually-curated features, such as those based on regular expressions and gazetteers. We show that ANNs achieve state-of-the-art results on de-identification of two different datasets for patient notes, the i2b2 2014 challenge dataset and the MIMIC dataset. To the best of our knowledge, this is the first paper to introduce ANN-based approaches using token and character embeddings to clinical de-identification task.

There have been a few related publications that apply ANNs and word embeddings for clinical NLP tasks. Wu et al. [46] investigate the use of deep neural networks to learn word embeddings and perform named entity recognition of four types of clinical entities – problems, lab tests, procedures, and medications – on Chinese clinical text. Two submissions [47,48] to a recent SemEval-2016 Task 12: Clinical TempEval challenge also reports ANN-based methods for information extraction from clinical notes and pathology reports. Li and Huang [47] use convolutional neural network and Fries [48] compares the performance of recurrent neural network and DeepDive [49] for the task.

2 METHODS AND MATERIALS

We first present a de-identifier we developed based on a conditional random field (CRF) model in Section 2.1. This de-identifier yields state-of-the-art results on the i2b2 2014 dataset, which is the reference dataset for comparing de-identification systems. This system will be used as a challenging baseline for the ANN model that we will present in Section 2.2. The ANN model outperforms the CRF model, as outlined in Section 3.4.

2.1 CRF model

In the CRF model, each patient note is tokenized using the Stanford CoreNLP tokenizer [50], and features are extracted for each token. During the training phase, the CRF's parameters are optimized to maximize the likelihood of the gold standard labels. During the test phase, the CRF predicts the labels. The performance of a CRF model depends mostly on the quality of its features. We used a combination of lexical, morphological, temporal, semantic, gazetteer, and regular expression features. [Table 2](#) lists some of the features used in the CRF model. The regular expressions were written mostly based on the best-performing CRF-based competitors in the i2b2 challenge [51]. The gazetteers were compiled using common resources from the web, and most other features were from [52]. See [51,52] for details about the relevant features.

In order to effectively incorporate context when predicting a label, all the features for a given token are computed based on that token and on the four surrounding tokens.

Table 2 Examples of features used in the CRF model.

| Feature types | Features |
|---------------------|---|
| Lexical/Syntactic | Token, lemma, tense, part-of-speech |
| Morphological | Ends with s, contains a digit, is numeric, is alphabetic, is alphanumeric, is title case, is all lower case, prefix, suffix |
| Temporal | Season, month, weekday, time of the day |
| Semantic/Wordnet | Hypernyms, senses, lemma names |
| Gazetteers | First names, last names, medical titles, medical specialties, cities, states (including abbreviations), countries, organizations, professions, holidays |
| Regular expressions | Email, age, date, phone, zip code, id number, medical record number |

2.2 ANN model

The main components of the ANN model are recurrent neural networks (RNNs). In particular, we use a type of RNN called Long Short Term Memory (LSTM) [53], as discussed in Section 2.2.1. The system is composed of three layers:

- Character-enhanced token embedding layer (Section 2.2.2),
- Label prediction layer (Section 2.2.3),
- Label sequence optimization layer (Section 2.2.4).

As in the CRF model, the patient notes are first tokenized using the Stanford CoreNLP tokenizer. The character-enhanced token embedding layer maps each token into a vector representation. The sequence of vector representations corresponding to a sequence of tokens are input to the label prediction layer, which outputs the sequence of vectors containing the probability of each label for each corresponding token. Lastly, the sequence optimization layer outputs the most likely sequence of predicted labels based on the sequence of probability vectors from the previous layer. All layers are learned jointly. Figure 1 shows the ANN architecture.

In the following, we denote scalars in italic lowercase (e.g., k , b_f), vectors in bold lowercase (e.g., \mathbf{s} , \mathbf{x}_i), and matrices in italic uppercase (e.g., W_f) symbols. We use the colon notations $\mathbf{x}_{i:j}$ and $\mathbf{v}_{i:j}$ to denote the sequence of scalars x_i, \dots, x_j , and vectors $\mathbf{v}_i, \mathbf{v}_{i+1}, \dots, \mathbf{v}_j$, respectively.

2.2.1 Bidirectional LSTM

An RNN is a neural network architecture designed to handle input sequences of variable sizes, but it fails to model long term dependencies. An LSTM is a type of RNN that mitigates this issue by keeping a memory cell that serves as a summary of the preceding elements of an input sequence. More specifically, given a sequence of vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, at each step $t = 1, \dots, n$, an LSTM takes as input

$\mathbf{x}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}$ and produces the hidden state \mathbf{h}_t and the memory cell \mathbf{c}_t based on the following formulas:

$$\mathbf{i}_t = \sigma(W_i [\mathbf{x}_t; \mathbf{h}_{t-1}; \mathbf{c}_{t-1}] + \mathbf{b}_i)$$

$$\mathbf{c}_t = (1 - \mathbf{i}_t) \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tanh(W_c [\mathbf{x}_t; \mathbf{h}_{t-1}] + \mathbf{b}_c)$$

$$\mathbf{o}_t = \sigma(W_o [\mathbf{x}_t; \mathbf{h}_{t-1}; \mathbf{h}_{t-1}] + \mathbf{b}_o)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

where W_i, W_c, W_o are weight matrices and $\mathbf{b}_i, \mathbf{b}_c, \mathbf{b}_o$ are bias vectors used in the input gate, memory cell, and output gate calculations, respectively. The symbols $\sigma(\cdot)$ and $\tanh(\cdot)$ refer to the element-wise sigmoid and hyperbolic tangent functions, and \odot is the element-wise multiplication. $\mathbf{h}_0 = \mathbf{c}_0 = \mathbf{0}$.

A bidirectional LSTM consists of a forward LSTM and a backward LSTM, where the forward LSTM calculates the forward hidden states $(\vec{\mathbf{h}}_1, \vec{\mathbf{h}}_2, \dots, \vec{\mathbf{h}}_n)$, and the backward LSTM calculates the backward hidden states $(\overleftarrow{\mathbf{h}}_1, \overleftarrow{\mathbf{h}}_2, \dots, \overleftarrow{\mathbf{h}}_n)$ by feeding the input sequence in the backward order, from \mathbf{x}_n to \mathbf{x}_1 .

Depending on the application of the LSTM, one might need an output sequence corresponding to each element in the sequence, or a single output that summarizes the whole sequence. In the former case, the output sequence $\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n$ of the LSTM is obtained by concatenating the hidden states of the forward and the backward LSTMs for each element i.e., $\vec{\mathbf{h}}_t = (\vec{\mathbf{h}}_t; \overleftarrow{\mathbf{h}}_t)$ for $t = 1, \dots, n$. In the latter case, the output is obtained by concatenating the last hidden states of the forward and the backward LSTMs i.e., $\vec{\mathbf{h}} = (\vec{\mathbf{h}}_n; \overleftarrow{\mathbf{h}}_n)$.

2.2.2 Character-enhanced token embedding layer

The character-enhanced token embedding layer takes a token as input and outputs its vector representation. The latter results from the concatenation of two different types of embeddings: the first one directly maps a token to a vector, while the second one comes from the output of a character-level

token encoder. The direct mapping $\mathcal{V}_T(\cdot)$ from token to vector, often called a token (or word) embedding, can be pre-trained on large unlabeled datasets using programs such as word2vec [30,54,55] or GloVe [32], and can be learned jointly with the rest of the model. Token embeddings, often learned by sampling token co-occurrence distributions, have desirable properties such as locating semantically similar words closely in the vector space, hence leading to state-of-the-art performance for various tasks.

While the token embeddings capture the semantics of tokens to some degree, they may still suffer from data sparsity. For example, they cannot account for out-of-vocabulary tokens, misspellings, and different noun forms or verb endings. One solution to remediate some of these issues would be to lemmatize tokens before training, but this approach may fail to retain some useful information such as the distinction between some verb and noun forms.

We address this issue by using character-based token embeddings, which incorporate each individual character of a token to generate its vector representation. This approach enables the model to learn sub-token patterns such as morphemes (e.g., suffix or prefix) and roots, thereby capturing out-of-vocabulary tokens, different surface forms, and other information not contained in the token embeddings.

Let $x_{i,1}, \dots, x_{i,\ell(i)}$ be the sequence of characters that comprise the i^{th} token x_i , where $\ell(i)$ is the number of characters in x_i . The character-level token encoder generates the character-based token embedding of x_i by first mapping each character $x_{i,j}$ to a vector $\mathcal{V}_C(x_{i,j})$, called a character embedding, via the mapping $\mathcal{V}_C(\cdot)$. Then the sequence $\mathcal{V}_C(x_{i,1}), \dots, \mathcal{V}_C(x_{i,\ell(i)})$ is passed to a bidirectional LSTM, which outputs the character-based token embedding \vec{b}_i .

As a result, the final output e_i of the character-enhanced token embedding layer for i^{th} token x_i is the concatenation of the token embedding $\mathcal{V}_T(x_i)$ and the character-based token embedding \vec{b}_i . In summary, when the character-enhanced token embedding layer receives a sequence of tokens $x_{1:n}$ as input, it will output the sequence of token embeddings $e_{1:n}$.

2.2.3 Label prediction layer

The label prediction layer takes as input the sequence of vectors $e_{1:n}$, i.e., the outputs of the character-enhanced token embedding layer, and outputs $a_{1:n}$, where the t^{th} element of a_n is the probability that the n^{th} token has the label t . The labels are either one of the PHI types or non-PHI. For example, if one aims to predict all 18 HIPAA-defined PHI types, there would be 19 different labels.

The label prediction layer contains a bidirectional LSTM that takes the input sequence $e_{1:n}$ and generates the corresponding output sequence $\vec{d}_{1:n}$. Each output \vec{d}_i of the LSTM is given to a feed-forward neural network with one hidden layer, which outputs the corresponding probability vector a_i .

2.2.4 Label sequence optimization layer

The label sequence optimization layer takes the sequence of probability vectors $a_{1:n}$ from the label prediction layer as input, and outputs a sequence of labels $y_{1:n}$, where y_i is the label assigned to the token x_i .

The simplest strategy to select the label y_i would be to choose the label that has the highest probability in a_i , i.e. $y_i = \text{argmax}_k a_i[k]$. However, this greedy approach fails to take into account the dependencies between subsequent labels. For example, it may be more likely to have a token with the PHI type STATE followed by a token with the PHI type ZIP than any other PHI type. Even though the label prediction layer has the capacity to capture such dependencies to a certain degree, it may be preferable to allow the model to directly learn these dependencies in the last layer of the model.

One way to model such dependencies is to incorporate a matrix T that contains the transition probabilities between two subsequent labels. $T[i, j]$ is the probability that a token with label i is followed by a token with the label j . The score of a label sequence $y_{1:n}$ is defined as the sum of the probabilities of individual labels and the transition probabilities:

$$s(y_{1:n}) = \sum_{i=1}^n \mathbf{a}_i [y_i] + \sum_{i=2}^n T [y_{i-1}, y_i].$$

These scores can be turned into probabilities of the label sequences by taking a softmax function over all possible label sequences. During the training phase, the objective is to maximize the log probability of the gold label sequence. In the testing phase, given an input sequence of tokens, the corresponding sequence of predicted labels is chosen as the one that maximizes the score.

3 EXPERIMENTS AND RESULTS

3.1 Datasets

We evaluate our two models on two datasets: the i2b2 2014 and MIMIC de-identification datasets. The i2b2 2014 dataset was released as part of the 2014 i2b2/UTHealth shared task Track 1 [29]. It is the largest publicly available dataset for de-identification. Ten teams participated in this shared task, and 22 systems were submitted. As a result, we used the i2b2 2014 dataset to compare our models against state-of-the-art systems.

The MIMIC de-identification dataset was created for this work as follows. The MIMIC-III dataset [7,8,56] contains data for 61,532 ICU stays over 58,976 hospital admissions for 46,520 patients, including 2 million patient notes. In order to make the notes publicly available, a rule-based de-identification system [5,6,57] was written for the specific purpose of de-identifying patient notes in MIMIC, leveraging

dataset-specific information such as the list of patient names or addresses. The system favors recall over precision: there are virtually no false negatives, while there are numerous false positives. To create the gold standard MIMIC de-identification dataset, we selected 1,635 discharge summaries, each belonging to a different patient, containing a total of 60.7k PHI instances. We then annotated the PHI instances detected by the rule-based system as true positives or false positives. We found that 15% of the PHI instances detected by the rule-based system were false positives.

Table 1 introduces the PHI types used as labels for training and Table 3 presents the datasets' sizes. For the test set, we used the official test set for the i2b2 dataset, which is 40% of the dataset; we randomly selected 20% of the MIMIC dataset as the test set for this dataset.

Table 3 Overview of the i2b2 and MIMIC datasets.

| | i2b2 | MIMIC |
|-------------------------|---------|-----------|
| Vocabulary size | 46,803 | 69,525 |
| Number of notes | 1,304 | 1,635 |
| Number of tokens | 984,723 | 2,945,228 |
| Number of PHI instances | 28,867 | 60,725 |
| Number of PHI tokens | 41,355 | 78,633 |

3.2 Evaluation metrics

To assess the performance of the two models, we computed the precision, recall, and F1-score. Let TP be the number of true positives, FP the number of false positives, and FN the number of false negatives. Precision, recall, and F1-score are defined as follows:

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

$$\text{precision} = \frac{TP}{TP+FP}, \text{ recall} = \frac{TP}{TP+FN}, \text{ and F1-score} = \frac{2*\text{precision}*\text{recall}}{\text{precision}+\text{recall}}.$$

Intuitively, precision is the proportion of the predicted PHI labels that are gold labels, recall is the proportion of the gold PHI labels that are correctly predicted, and F1-score is the harmonic mean of precision and recall.

3.3 Training and hyperparameters

The model is trained using stochastic gradient descent, updating all parameters, i.e., token embeddings, character embeddings, parameters of bidirectional LSTMs, and transition probabilities, at each gradient step. For regularization, dropout is applied to the character-enhanced token embeddings before the label prediction layer. Training the model takes approximately 2 days on an Nvidia Titan X GPU. The actual running time depends on the choice of hyperparameters, the weight initialization, and the size of the dataset.

Below are the choices of hyperparameters and token embeddings, optimized using a subset of the training set:

- character embedding dimension: 25
- character-based token embedding LSTM dimension: 25
- token embedding dimension: 100
- label prediction LSTM dimension: 100
- dropout probability: 0.5

As mentioned in Section 2.2.2, token embeddings can be pre-trained, and during training the token mapping $\mathcal{V}_T(\cdot)$ is initialized with the pre-trained token embeddings. We tried pre-training token

embeddings on the i2b2 2014 dataset and the MIMIC dataset¹ using word2vec and GloVe. Both word2vec and GloVe were trained using a window size of 10, a minimum vocabulary count of 5, and 15 iterations. Additional parameters of word2vec were the negative sampling and the model type, which were set to 10 and skip-gram, respectively. We also experimented with the publicly available² token embeddings such as GloVe trained on Wikipedia and Gigaword 5 [58]. The results were quite robust to the choice of the pre-trained token embeddings. The GloVe embeddings trained on Wikipedia articles yielded slightly better results, and we chose them for the rest of this work.

3.4 Results

All results were computed using the official evaluation script from the i2b2 2014 de-identification challenge. Table 4 presents the main results, based on binary token-based precision, recall, and F1-score for HIPAA-defined PHI only. These PHI types are the most important since only those are required to be removed by law. The results for each PHI type, dataset, and system are presented in Appendix 1, Tables A1 and A2.

On the i2b2 dataset, our ANN model has a higher F1-score and recall than our CRF model as well as the best system from the i2b2 2014 de-identification challenge, which was the Nottingham system [51]. The only freely available, off-the-shelf program for de-identification, called the MITRE Identification Scrubber Toolkit (MIST) [27], performed the worst. The outputs of our ANN and CRF models can be combined by considering a token to be PHI if it is identified as such by either model. This further increases the performance in terms of F1-score and recall. It should be noted that the Nottingham system was specifically fine-tuned for the i2b2 dataset as well as the i2b2 evaluation script. For example,

¹ For MIMIC, we used the entire dataset containing 2 million notes and 800 million tokens.

² <http://nlp.stanford.edu/projects/glove/>

the Nottingham system post-processes the detected PHI terms in order to match the offset of the gold PHI tokens, such as modifying “MR:6746781” to “6746782” and “MWFS” to “M”, “W”, “F”, “S”.

Table 4 Performance (%) on the PHI as defined in the HIPAA. We evaluated the systems based on the detection of PHI tokens versus non-PHI tokens (i.e., binary HIPAA token-based evaluation). The best performance for each metric on each dataset is highlighted in bold. Nottingham is the best performing system from the 2014 i2b2/UTHealth shared task Track 1. MIST, the MITRE Identification Scrubber Toolkit, is a freely available de-identification program. CRF is the model based on Conditional Random Field, ANN is the model based on Artificial Neural Network, and CRF+ANN is the result obtained by combining the outputs of the CRF model and the ANN model. The tagsets used for training the CRF and ANN models are the same as in Table 1, and the configuration of MIST is presented in Appendix 2. The Nottingham system could not be run on the MIMIC dataset, as it is not publicly available.

| Model | i2b2 | | | MIMIC | | |
|------------|---------------|---------------|---------------|---------------|---------------|---------------|
| | Precision | Recall | F1 | Precision | Recall | F1 |
| Nottingham | 99.000 | 96.400 | 97.680 | - | - | - |
| MIST | 91.445 | 92.745 | 92.090 | 95.867 | 98.346 | 97.091 |
| CRF | 98.560 | 96.528 | 97.533 | 99.060 | 98.987 | 99.023 |
| ANN | 98.320 | 97.380 | 97.848 | 99.208 | 99.251 | 99.229 |
| CRF + ANN | 97.920 | 97.835 | 97.877 | 98.820 | 99.398 | 99.108 |

On the MIMIC dataset, our ANN model also has a higher F1-score and recall than our CRF model. Interestingly, combining the outputs of our ANN and CRF models did not increase the F1-score, because precision was negatively impacted. However, the recall did benefit from combining the two models. MIST was much more competitive on this dataset.

We calculated the statistical significance of the differences in precision, recall, and F1-score between the CRF and ANN models using approximate randomization with 9999 shuffles. The significance levels of the differences in precision, recall, and F1-score are 0.37, 0.02, 0.22 for the i2b2 dataset, and 0.08, 0.00, 0.00 for the MIMIC dataset, respectively.

3.5 Error analysis

Figure 2 shows the binary token-based F1-scores for each PHI category. The ANN model outperforms the CRF model on all categories for both datasets, with the exception of the ID category (which mostly contains medical record numbers) in the i2b2 dataset. This is due to the fact that the CRF model uses sophisticated regular expression features that are tailored to detect ID patterns such as “38:Z8912708G”.

Another interesting difference between the ANN and the CRF results is the PROFESSION category: the ANN significantly outperforms the CRF. The reason behind this result is that the embeddings of the tokens that represent a profession tend to be close in the token embedding space, which allows the ANN model to generalize well. We tried assembling various gazetteers for the PROFESSION category, but all of them were performing significantly worse than the ANN model.

Table 5 presents some examples of gold PHI instances correctly predicted by the ANN model that the CRF model failed to predict, and conversely. This illustrates that the ANN model efficiently copes with the diversity of the contexts in which tokens appear, whereas the CRF model can only address the contexts that are manually encoded as features. In other words, the ANN model’s intrinsic flexibility allows it to better capture the variances in human languages than the CRF model. For example, it would be challenging and time-consuming to engineer features for all possible contexts such as “had a stroke at 80”, “quit smoking in 08”, “on the 29th of this month”, and “his friend Epstein”. The ANN model is also very robust to variations in surface forms, such as misspellings (e.g., “in teh late 60s”, “Khazakhstan”, “01/19/:0”), tokenizations (e.g., “Results02/20/2087”, “MC # 0937884Date”), and different phrases referring to the same semantic meaning (e.g., “San Rafael Mount Hospital”, “Rafael Mount”, “Rafael Hospital”). Furthermore, the ANN model is able to detect many PHI instances despite not having explicit gazetteers, as examples in the LOCATION and PROFESSION categories illustrate. We

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Table 5 Examples of correctly detected PHI instances (in bold) by the ANN and CRF models for the i2b2 dataset. The examples in the ANN column are only predicted by the ANN model and not predicted by the CRF model, and conversely. Typographical errors are from the original text.

| PHI category | ANN | CRF |
|--------------|--|--|
| AGE | Father had a stroke at 80 and died of ?another stroke at age PERSONAL DATA AND OVERALL HEALTH: Now 63 , despite his FH: Father: Died @ 52 from EtOH abuse (unclear exact etiology) Tobacco: smoked from age 7 to 15 , has not smoked since 15. History of Present Illness 86F reports worsening b/l leg pain. | HPI: 53RHM who going to bed Wednesday was in usoh, but Tobacco: Quit at 38 y/o; ETOH: 1-2 beers/week; Caffeine: |
| CONTACT | by phone, Dr. Ivan Guy. Call w/ questions 86383 . Keith Gilbert, H/O paroxysmal afib VNA 171-311-7974 ===== Medications | |
| DATE | During his May hospitalization he had dysphagia Social history: divorced, quit smoking in 08 , sober x 10 yrs, She is to see him on the 29th of this month at 1:00 p.m. He did have a renal biopsy in teh late 60s adn thus will look for results, Results 02/20/2087 NA 135, K 3.2 (L), CL 96 (L), CO2 30.6, BUN 1 Jose Church, M.D. /ray DD: 01/18/20 DT: 01/19/:0 DV: 01/18/20 | She is looking forward to a good Christmas . She is here today |
| ID | placed 3/23 for bradycardia. P/G model # 5435 , serial # 4712198, Consult NotePt: Ulysses Ogrady MC # 0937884 Date: 10/07/69 | DD:05/05/2095 DT:05/05/2095 WK:65255 :4653 NO GROWTH TO DATE Specimen: 38:Z8912708G Collected |
| LOCATION | Works in programming at Audiovox . Formerly at BrightPoint. He has remote travel hx to the Rockefeller Centre , more recent global History of Present Illness: Pt is a 59 yo Khazakhstani male, with who was admitted to San Rafael Mount Hospital following a syncopal nauseas and was brought to Rafael Mount ED. Five weeks ago prior Anemia: On admission to Rafael Hospital , Hb/Hct: 11.6/35.5. | 2nd set biomarkers (WPH): Creatine Kinase Isoenzymes Hospitalized 2115 TCH for ROMI 2120 TCH new onset |
| NAME | ATCH: 655-75-45 Dear Harry and Yair : My thanks for your kind Patient lives in Flint with his friend Epstein . He has 3 children. Health care proxy-Yes, son (West) Allergies DUTASTERIDE - cough, | Lab Tests Amador : the lab results show good levels of 10MG PO qd : 05/10/2066 - 04/15/2068 ACT : rosenberg 128 Williams Ct M OSCAR, JOHNNY Hyderabad, WI 62297 |
| PROFESSION | Social history: Married, glazier , 3 grown adult children Has VNA. Former civil engineer, supervisor , consultant. He was formerly self-employed as a CPA and would often travel Communications senior manager, marketing , worked for Brinker and Concrete Finisher (25yrs). He is a veteran . Former tobacco user, works part time in securities . | Social history: He is retried Motor Vehicle Body Repairer . |

Table 6 Examples of PHI instances undetected by CRF+ANN, i.e., undetected by both the CRF and the ANN, for the i2b2 dataset. Each row present one or two false negatives (marked in bold letters, and underlined). The “Reason” column specifies what we believe to be the main factor that caused the CRF+ANN to fail to detect the token(s) as a PHI instance. Ab: abbreviation; Am: ambiguity; D: debatable annotation; S: data sparsity. The “FN” column indicates how many tokens of a given PHI type are false negatives. The “Support” column indicates the number of tokens of a given PHI type in the test set.

| PHI categorie | PHI type | Examples | Reason | FN | Support |
|---------------|---------------|---|--------------------------------|----|---------|
| AGE | AGE | ASSESSMENT AND PLAN: A <u>seventy-one</u> -year-old woman with multiple medical Both parents died of sudden death in their <u>82nd</u> year. Brother had SCD at <u>66</u> . smoked from age 7 to 15, has not smoked since <u>15</u> . (+)(-)prostate/colonCa/CADPGP d 80s ?cause, MGF d <u>90</u> age, MGM d <u>73</u> CVAM d 73 stomach Ca, OA, obeseF d <u>84</u> multi-infarct dementiaS b66 DMS b74 | S S S S S | 19 | 790 |
| CONTACT | PHONE | Wheatland Manor: 154-734-1487, x <u>557</u> (4th floor) | S | 1 | 410 |
| | FAX | Phone: (091)920-5569 Fax: <u>(251)628-xxxx</u> | S | 3 | 6 |
| | EMAIL | E-Mail: iparedes@oachosp.org | S | 3 | 3 |
| DATE | DATE | PARONYCHIAL INFECTION : LEFT HAND <u>78</u> Ectopic pregnancy : <u>74</u> alb 4.2 fe 50, tbc 204, ferritin 878 <u>8/27</u> inr 1.1 pth 115 8/27 lipitor20 mg and lopid 600 bid. Prior HDL 19. <u>8/67</u> TC 170, TG 162, H40, L98, ratio 4.3diabetes Referral submitted to GI <u>6/65</u> : saw GI - going for scope to eval pancreas multi-infarct dementiaS b66 DMS b74 DMSon b93D b94 GC due <u>22</u> D Fran b03 Abn pap24 Nephropathy Patient was last seen in clinic in <u>11-70</u> after which time she left for Ghana for the past | Am Am Am Am S S | 60 | 12534 |
| ID | IDNUM | Influenza vaccine Received 11/95 <u>MLL</u> disp #100 order number <u>38/48</u> ALLERGY NKDA | Am S | 9 | 382 |
| | MEDICALRECORD | Patient: Vincent Ware (71417347 <u>2Y</u>) | S | 1 | 732 |
| | DEVICE | Interrogation today of his Medtronic Kappa <u>QQ 626</u> pacemaker reveals that his underlying rhythm | S | 4 | 12 |
| LOCATION | STREET | - | | 0 | 416 |
| | CITY | Oriented to "LCC" in " <u>Galena</u> ," "March 2095." Speech fluent in Dutch. | S | 8 | 344 |
| | ZIP | - | | 0 | 144 |
| | STATE | BP has been well-controlled in <u>VA</u> , usually in the 128 systolic range. | Ab/Am | 9 | 205 |
| | COUNTRY | is here with her husband who is translating from <u>columbian</u> . | S | 13 | 130 |

| | | | | | |
|------------|----------------|---|-----------------------------|----|------|
| | LOCATION-OTHER | travel hx to the Rockefeller Centre, more recent <u>global</u> travel (Fernley, Cartersville, Iceland) and has infrequently visited <u>Storting</u> and <u>Acropolis</u> . | D S | 12 | 20 |
| | ORGANIZATION | diabetes diet - he enjoys a blueberry muffin from <u>RR Donnelley</u> daily. his level of fatigue. He continues to go to <u>the library</u> daily. He continues | S D | 42 | 147 |
| | HOSPITAL | were placed at Pomeroy Care Center (Big Rapids, <u>AC</u>) and also he had evaluation at the Corcoran Medication List for QUICK,ISABELLE Y 6557545 (<u>ATCH</u>) 52 F 2. DM, stable, Glyburide increased at <u>MS</u> . Dietary rec's reviewed. | Ab/Am Ab Ab/Am | 44 | 1595 |
| NAME | PATIENT | ct dementiaS b66 DMS b74 DMSon b93D b94 GC due22D <u>Fran</u> b03 Abn pap24 Nephropathy 3/25 TACT: Gracen Logan (HCP, daughter) 625-248-3647; <u>Flowers</u> (son) 705-690-8475 Patient Name: JIMENEZ,YOUSSEF I [0554733(LCH)] | Am Am Ab/Am | 6 | 1450 |
| | DOCTOR | Snyder/Ophthalmology - Insley/Endocrinology - End 6 <u>Lane</u> /Neurology - NEU 265 Script: Amt: 30 Refill: 3 Date: 03/11/2074 : <u>um</u> If the latter, will change it. <u>Q</u> Plasma Sodium 138 | Am Am Ab/Am | 35 | 3297 |
| | USERNAME | - | | 0 | 92 |
| PROFESSION | PROFESSION | however he would like to try to <u>intern</u> , when he feels up to it. Patient lives in Lake Pocotopaug with wife. <u>Justice of the peace</u> . 2ppd x 40y, denies etoh. On disability. Volunteers - <u>animal rescue</u> . No current or previous tobacco Social History <u>NP</u> in Laplace - waiting for researcher job. He has continued actively <u>managing production</u> and is planning a trip to Italy next | D S S Ab/Am S/D | 69 | 340 |

conjecture that the character-enhanced token embeddings contain rich enough information to effectively function as gazetteers, as tokens with similar semantics are closely located in the vector representation [26,27,41].

On the other hand, CRF is good at rarely occurring patterns that are written in highly specialized regular expression patterns (e.g., “38:Z8912708G”, “53RHM”) or tokens that are included in the gazetteers (e.g., “Christmas”, “WPH”, “rosenberg”, “Motor Vehicle Body Repairer”). For example, the PHI token “Christmas” only occurs in the test set, and unless the context gives a strong indication, the ANN model cannot detect it, whereas the CRF model could, as long as it is included in the gazetteers.

Table 6 presents examples of PHI instances that are false negatives in the system that combines the CRF and ANN outputs. In other words, these PHI instances are detected by neither the CRF nor the ANN. The sources of errors may be classified into four main categories:

- Abbreviations: some PHI instances are abbreviations, which are sometimes challenging to detect, especially when they are short and ambiguous.
- Ambiguities: a human reader may not be able to tell whether a token is PHI. Examples include names involving common words, or numbers that could be date or test result. Ambiguities may stem from the token itself as well as its context.
- Data sparsity: the training samples do not contain enough PHI instances similar to the ones that are missed in the test set. Also, some PHI instances are more difficult to detect than others and subsequently require more training samples.
- Debatable annotations: some tokens are questionably marked as PHI instances.

Abbreviations and ambiguities are among the most challenging sources of errors to address in order to further improve the performance. We anticipate that the data sparsity issues may partly be resolved by increasing the size of the training set to contain more instances of difficult PHI types.

3.6 Effect of training set size

Figure 3 shows the impact of the training set size on the performance of the models on the MIMIC dataset. When the training set size is very limited, the CRF performs slightly better than the ANN model, since the CRF model can leverage handcrafted features without much training data. As the training set size increases, the ANN model starts to significantly outperform the CRF model, since the parameters including the embeddings are automatically fine-tuned with more data, and therefore the features learned by the ANN model become increasingly more refined than the manually handcrafted features. As a result, combining the outputs of the CRF and ANN models increases the F1-score over the ANN model only for small training set size and yields a less competitive F1-score than the ANN model for bigger training set size.

Figure 4 details the impact of the number of labeled PHI instances in the training set on the model’s performance for a given PHI type, in the i2b2 dataset. As expected, PHI types with a large number of labeled PHI instances tend to be detected more accurately than rarer PHI types. However, the correlation is far from perfect: some PHI types with a lower number of labeled PHI instances are detected more accurately than some PHI types with a higher number of labeled PHI instances. This indicates that some PHI types are harder to detect than others. For example, although the PHI type “PHONE” has fewer labeled PHI instances than the PHI type “PROFESSION” (310 vs. 425 instances), the former is much more accurately detected than the latter (F1-score of 99.272 vs. F1-score of 86.642): this

result is expected since tokens containing a phone number are typically very similar, whereas professions can appear in many different forms.

3.7 Ablation analysis

In order to quantify the importance of various elements of the ANN model, we tried 4 variations of the model, eliminating different elements one at a time. Figure 5 presents the results of the ablation tests. Removing either the label sequence optimization layer, pre-trained token embeddings, or token embeddings slightly decreased the performance. Surprisingly, the ANN performed pretty well with only character embeddings and without the token embeddings, and eliminating the character embeddings was more detrimental than eliminating the token embeddings. This suggests that the character-based token embeddings may be capturing not only the sub-token level features, but also the semantics of the tokens themselves.

4 CONCLUSION

We proposed the first system based on ANN for patient note de-identification. It outperforms state-of-the-art systems based on CRF on two datasets, while requiring no handcrafted features. Utilizing both the token and character embeddings, the system can automatically learn effective features from data by fine-tuning the parameters. It jointly learns the parameters for the embeddings, the bidirectional LSTMs as well as the label sequence optimization, and can make use of token embeddings pre-trained on large unlabeled datasets. Quantitative and qualitative analysis of the ANN and CRF models indicates that the ANN model better incorporates context and is more flexible to variations inherent in human languages than the CRF model.

From the viewpoint of deploying an off-the-shelf de-identification system, our results in Table 4 demonstrate recall on the MIMIC discharge summaries over 99%, which is quite encouraging. Figure 2,

however, shows that the F1-score on the NAME category, probably the most sensitive PHI type, falls just below 98% for the ANN model. We anticipate that adding gazetteer features based on the local institution's patient and staff census should improve this result, which will be explored in future work.

5 ACKNOWLEDGMENTS

We warmly thank Michele Filannino, Alistair Johnson, and Tom Pollard for their helpful suggestions and technical assistance.

6 FUNDING STATEMENT

The project was supported by Philips Research. The content is solely the responsibility of the authors and does not necessarily represent the official views of Philips Research.

7 COMPETING INTERESTS STATEMENT

The authors have no competing interests to declare.

8 CONTRIBUTORSHIP STATEMENT

Franck Dernoncourt and Ji Young Lee contributed equally to this work. They designed and implemented the CRF and ANN models, annotated the MIMIC de-identification dataset, evaluated the systems' performances, created the figures and wrote the paper. Ozlem Uzuner and Peter Szolovits formulated the original problem, provided direction and guidance, and gave helpful feedback on the paper.

9 REFERENCES

1. DesRoches CM, Worzala C, Bates S. Some hospitals are falling behind in meeting “meaningful use” criteria and could be vulnerable to penalties in 2015. Health Affairs. Health Affairs; 2013;32:1355–60.
2. Wright A, Henkin S, Feblowitz , et al. Early results of the meaningful use program for electronic health records. New England Journal of Medicine. Mass Medical Soc; 2013;368:779–80.
3. Office for Civil Rights H. Standards for privacy of individually identifiable health information. final rule. Federal Register. 2002;67:53181.
4. Neamatullah I, Douglass MM, Li-wei HL, et al. Automated de-identification of free-text medical records. BMC Med Inform Decis Mak. BioMed Central; 2008;8:1.
5. Douglass M, Clifford G, Reisner A, et al. De-identification algorithm for free-text nursing notes. Computers in cardiology, 2005. IEEE; 2005:331–4.
6. Douglas M, Clifford G, Reisner A, et al. Computer-assisted de-identification of free text in the mIMIC ii database. Computers in cardiology, 2004. IEEE; 2004:341–4.
7. Goldberger AL, Amaral LA, Glass L, et al. Physiobank, physiotoolkit, and physionet components of a new research resource for complex physiologic signals. Circulation. Am Heart Assoc; 2000;101:e215–20.

8. Saeed M, Villarroel M, Reisner AT, et al. Multiparameter intelligent monitoring in intensive care II (MIMIC-II): A public-access intensive care unit database. *Crit Care Med*. NIH Public Access; 2011;39:952.

9. Lingren T, Deleger L, Molnar K, et al. Pre-annotating clinical notes and clinical trial announcements for gold standard corpus development: Evaluating the impact on annotation speed and potential bias. *Proc - 2012 IEEE 2nd Conf Healthc Informatics, Imaging Syst Biol HISB 2012*. 2012:108. doi:10.1109/HISB.2012.33.

10. South BR, Mowery D, Suo Y, et al. Evaluating the effects of machine pre-annotation and an interactive annotation interface on manual de-identification of clinical text. *J Biomed Inform*. 2014;50:162-172. doi:10.1016/j.jbi.2014.05.002.

11. Hanauer D, Aberdeen J, Bayer S, et al. Bootstrapping a de-identification system for narrative patient records: Cost-performance tradeoffs. *Int J Med Inform*. 2013;82(9):821-831. doi:10.1016/j.ijmedinf.2013.03.005.

12. Gobbel GT, Garvin J, Reeves R, et al. Assisted annotation of medical free text using RapTAT. *J Am Med Inform Assoc*. 2014;21(5):833-841. doi:10.1136/amiainl-2013-002255.

13. Chiticariu L, Li Y, Reiss FR. Rule-based information extraction is dead! Long live rule-based information extraction systems! *EMNLP*. 2013:827–32.

14. Berman JJ. Concept-match medical data scrubbing: How pathology text can be used in research. *Arch Pathol Lab Med*. 2003;127:680–6.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
15. Beckwith BA, Mahaadevan R, Balis UJ, et al. Development and evaluation of an open source software tool for deidentification of pathology reports. *BMC Med Inform Decis Mak. BioMed Central*; 2006;6:1.
16. Fielstein E, Brown S, Speroff T. Algorithmic de-identification of VA medical exam text for HIPAA privacy compliance: Preliminary findings. *Medinfo*. 2004;1590.
17. Friedlin FJ, McDonald CJ. A software tool for removing patient identifying information from clinical documents. *Journal of the American Medical Informatics Association. The Oxford University Press*; 2008;15:601–10.
18. Gupta D, Saul M, Gilbertson J. Evaluation of a deidentification (De-Id) software engine to share pathology reports and clinical documents for research. *Am J Clin Pathol. The Oxford University Press*; 2004;121:176–86.
19. Morrison FP, Li L, Lai AM, et al. Repurposing the clinical record: Can an existing natural language processing system de-identify clinical notes? *Journal of the American Medical Informatics Association. The Oxford University Press*; 2009;16:37–9.
20. Ruch P, Baud RH, Rassinoux A-M, et al. Medical document anonymization with a semantic lexicon. *Proceedings of the AMIA symposium. American Medical Informatics Association*; 2000:729.
21. Sweeney L. Replacing personally-identifying information in medical records, the Scrub system. *Proceedings of the AMIA annual fall symposium. American Medical Informatics Association*; 1996:333.

22. Thomas SM, Mamlin B, Schadow G, et al. A successful technique for removing names in pathology reports using an augmented search and replace method. Proceedings of the AMIA symposium. American Medical Informatics Association; 2002:777.

23. Szarvas G, Farkas R, Kocsor A. A multilingual named entity recognition system using boosting and c4.5 decision tree learning algorithms. Discovery science. Springer; 2006:267–78.

24. Guo Y, Gaizauskas R, Roberts I, et al. Identifying personal health information using support vector machines. I2b2 workshop on challenges in natural language processing for clinical data. 2006:10–1.

25. Uzuner Ö, Sibanda TC, Luo Y, et al. P. A de-identifier for medical discharge summaries. Artif Intell Med. Elsevier; 2008;42:13–35.

26. Hara K. Applying a SVM based chunker and a text classifier to the deid challenge. I2b2 workshop on challenges in natural language processing for clinical data. Am Med Inform Assoc; 2006:10–1.

27. Aberdeen J, Bayer S, Yeniterzi R, et al. The MITRE Identification Scrubber Toolkit: Design, training, and assessment. Int J Med Inform. Elsevier; 2010;79:849–59.

28. Meystre SM, Friedlin FJ, South BR, et al. Automatic de-identification of textual documents in the electronic health record: A review of recent research. BMC Med Res Methodol. BioMed Central; 2010;10:1.

29. Stubbs A, Kotfila C, Uzuner Ö. Automated systems for the de-identification of longitudinal clinical narratives: Overview of 2014 i2b2/UTHealth shared task track 1. J Biomed Inform. Elsevier; 2015;58:S11–9.
30. Mikolov T, Sutskever I, Chen K, et al. Distributed representations of words and phrases and their compositionality. Adv Neural Inf Process Syst. 2013:3111–9.
31. Collobert R, Weston J, Bottou L, et al. Natural language processing (almost) from scratch. The Journal of Machine Learning Research. JMLR. org; 2011;12:2493–537.
32. Pennington J, Socher R, Manning CD. GloVe: Global vectors for word representation. Proceedings of the Empirical Methods in Natural Language Processing (EMNLP 2014). 2014;12:1532–43.
33. Mikolov T, Karafiát M, Burget L, et al. Recurrent neural network based language model. INTERSPEECH. 2010:3.
34. Socher R, Perelygin A, Wu JY, et al. Recursive deep models for semantic compositionality over a sentiment treebank. Proceedings of the conference on empirical methods in natural language processing (EMNLP). Citeseer; 2013:1642.
35. Kim Y. Convolutional neural networks for sentence classification. Proceedings of the 2014 conference on empirical methods in natural language processing. Association for Computational Linguistics; 2014:1746–51.
36. Blunsom P, Grefenstette E, Kalchbrenner N, et al. A convolutional neural network for modelling sentences. Proceedings of the 52nd annual meeting of the association for

computational linguistics. Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics; 2014.

37. Lee JY, Dernoncourt F. Sequential short-text classification with recurrent and convolutional neural networks. Human language technologies 2016: The conference of the North American Chapter of the Association for Computational Linguistics, NAACL HLT 2016.

38. Weston J, Bordes A, Chopra S, et al. Towards AI-complete question answering: A set of prerequisite toy tasks. arXiv preprint arXiv:1502.05698. 2015.

39. Wang D, Nyberg E. A long short-term memory model for answer sentence selection in question answering. Proceedings of the 53rd annual meeting of the association for computational linguistics and the 7th international joint conference on natural language processing (volume 2: Short papers) [Internet]. Beijing, China: Association for Computational Linguistics; 2015:707–12.

40. Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate. arXiv preprint arXiv:1409.0473. 2014.

41. Tamura A, Watanabe T, Sumita E. Recurrent neural networks for word alignment model. ACL (1). 2014:1470–80.

42. Sundermeyer M, Alkhoul T, Wuebker J, et al. Translation modeling with bidirectional recurrent neural networks. EMNLP. 2014:14–25.

43. Lample G, Ballesteros M, Subramanian S, et al. Neural architectures for named entity recognition. arXiv preprint arXiv:1603.01360. 2016;

44. Labeau M, Löser K, Allauzen A. Non-lexical neural architecture for fine-grained POS tagging. Proceedings of the 2015 conference on empirical methods in natural language processing [Internet]. Lisbon, Portugal: Association for Computational Linguistics; 2015:232–7.
45. Kim Y, Jernite Y, Sontag D, et al. Character-aware neural language models. arXiv preprint arXiv:1508.06615. 2015.
9. Lingren T, Deleger L, Molnar K, et al. Pre-annotating clinical notes and clinical trial announcements for gold standard corpus development: Evaluating the impact on annotation speed and potential bias. *Proc - 2012 IEEE 2nd Conf Healthc Informatics, Imaging Syst Biol HISB 2012*. 2012:108. doi:10.1109/HISB.2012.33.
10. South BR, Mowery D, Suo Y, et al. Evaluating the effects of machine pre-annotation and an interactive annotation interface on manual de-identification of clinical text. *J Biomed Inform*. 2014;50:162-172. doi:10.1016/j.jbi.2014.05.002.
11. Hanauer D, Aberdeen J, Bayer S, et al. Bootstrapping a de-identification system for narrative patient records: Cost-performance tradeoffs. *Int J Med Inform*. 2013;82(9):821-831. doi:10.1016/j.ijmedinf.2013.03.005.
12. Gobbel GT, Garvin J, Reeves R, et al. Assisted annotation of medical free text using RapTAT. *J Am Med Inform Assoc*. 2014;21(5):833-841. doi:10.1136/amiainl-2013-002255.
46. Wu Y, Jiang M, Lei J, et al. Named Entity Recognition in Chinese Clinical Text Using Deep Neural Network. *Stud Health Technol Inform*. 2015;216:624-628. doi:10.3233/978-1-61499-564-7-624.

47. Li P, Huang H. UTA DLNLP at SemEval-2016 Task 12: Deep Learning Based Natural Language Processing System for Clinical Information Identification from Clinical Notes and Pathology Reports. In: *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*. San Diego, California: Association for Computational Linguistics; 2016:1268-1273. <http://www.aclweb.org/anthology/S16-1197>.

48. Fries JA. Brundlefly at SemEval-2016 Task 12: Recurrent Neural Networks vs. Joint Inference for Clinical Temporal Information Extraction. In: *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*. San Diego, California: Association for Computational Linguistics; 2016:1274-1279. <http://www.aclweb.org/anthology/S16-1198>.

49. Zhang C. DeepDive: A Data Management System for Automatic Knowledge Base Construction. *Thesis*. 2015;53:1689-1699. doi:10.1017/CBO9781107415324.004.

50. Manning CD, Bauer J, Finkel J, et al. The Stanford CoreNLP Natural Language Processing Toolkit. *Proc 52nd Annu Meet Assoc Comput Linguist Syst Demonstr*. 2014:55-60. <http://aclweb.org/anthology/P14-5010>.

51. Yang H, Garibaldi JM. Automatic detection of protected health information from clinic narratives. *J Biomed Inform*. Elsevier; 2015;58:S30–8.

52. Filannino M, Brown G, Nenadic G. ManTIME: Temporal expression identification and normalization in the TempEval-3 challenge. *CoRR*. 2013;abs/1304.7(2005):5. <http://arxiv.org/abs/1304.7942>.

- 1
2
3 53. Hochreiter S, Schmidhuber J. Long short-term memory. Neural Comput. MIT Press;
4
5
6 1997;9:1735–80.
7
8
9 54. Mikolov T, Chen K, Corrado G, Dean J. Efficient estimation of word representations in vector
10
11 space. arXiv preprint arXiv:1301.3781. 2013.
12
13
14
15 55. Mikolov T, Yih W-t, Zweig G. Linguistic regularities in continuous space word
16
17 representations. HLT-nAACL. 2013:746–51.
18
19
20
21 56. Johnson AEW, Pollard TJ, Shen L, et al. MIMIC-III, a freely accessible critical care database.
22
23 Scientific Data. 2016.
24
25
26
27 57. Douglass M. Computer-assisted de-identification of free-text nursing notes [Master's
28
29 thesis]. Massachusetts Institute of Technology; 2005.
30
31
32
33 58. Parker R, Graff D, Kong J, et al. English Gigaword fifth edition, linguistic data consortium.
34
35 Technical Report. Linguistic Data Consortium, Philadelphia; 2011.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

10 LIST OF FIGURES

Figure 1: Architecture of the artificial neural network (ANN) model. RNN stands for recurrent neural network. The type of RNN used in this model is Long Short Term Memory (LSTM). n is the number of tokens, and x_i is the i^{th} token. \mathcal{V}_T is the mapping from tokens to token embeddings. $\ell(i)$ is the number of characters and $x_{i,j}$ is the j^{th} character in the i^{th} token. \mathcal{V}_C is the mapping from characters to character embeddings. e_i is the character-enhanced token embeddings of the i^{th} token. \vec{d}_i is the output of the LSTM of label prediction layer, a_i is the probability vector over labels, y_i is the predicted label of the i^{th} token.

Figure 2: Binary token-based F1-scores for each PHI category. The evaluation is based on PHI types that are defined by HIPAA as well as additional PHI types specific to each dataset. Each PHI category and the corresponding PHI types are defined in Table 1. The “All” category refers to the F1-score micro-averaged over all PHI categories. The PROFESSION category exists only in the i2b2 dataset, and was plotted separately to avoid distorting the y-axis. For the same reason, the AGE category in MIMIC was drawn separately.

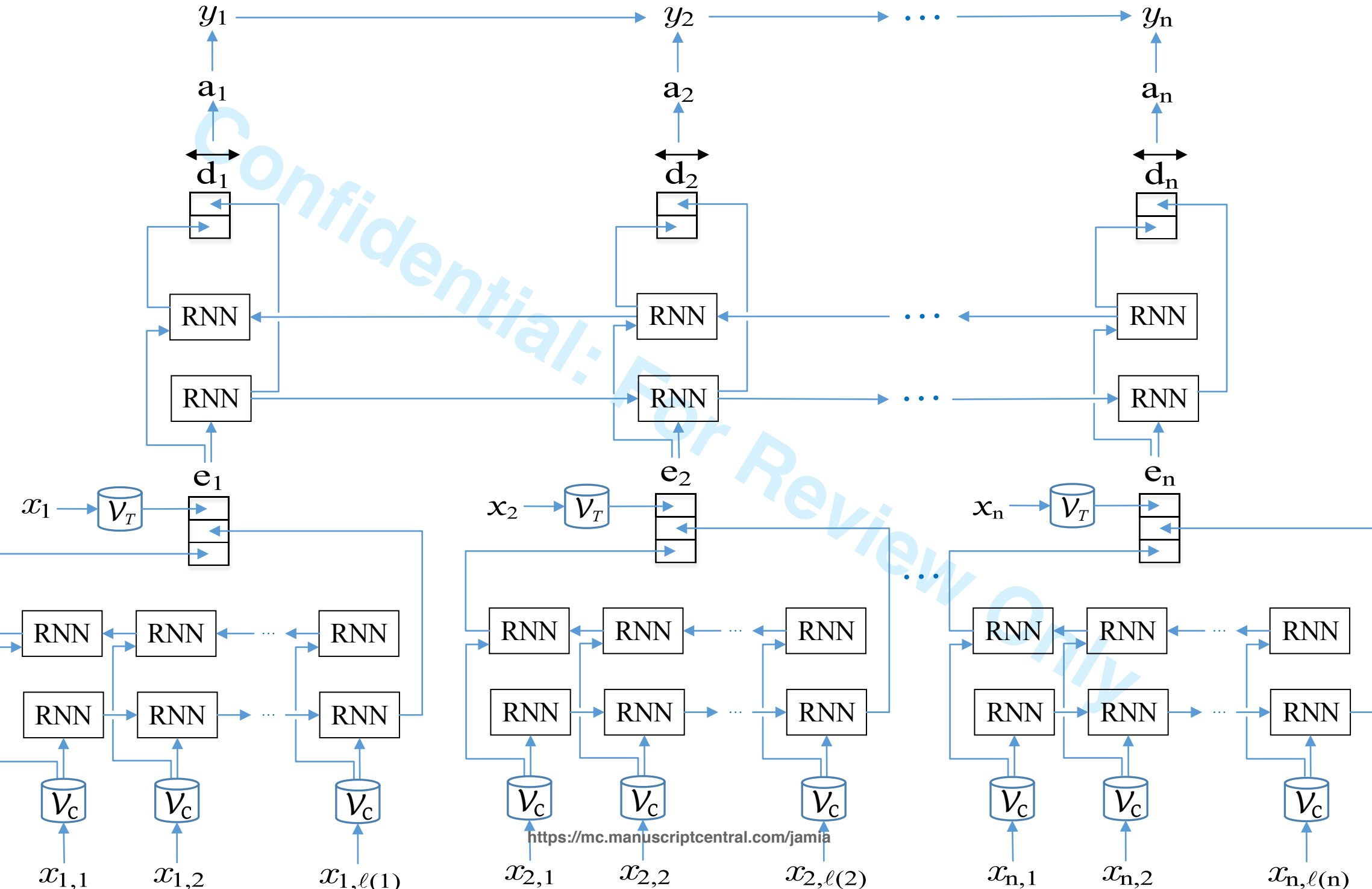
Figure 3: Impact of the training set size on the binary HIPAA token-based F1-scores on the MIMIC dataset. 100% training set size refers to using all of the dataset minus the test set, which amounts to 2,046,488 tokens and 42,531 PHI instances.

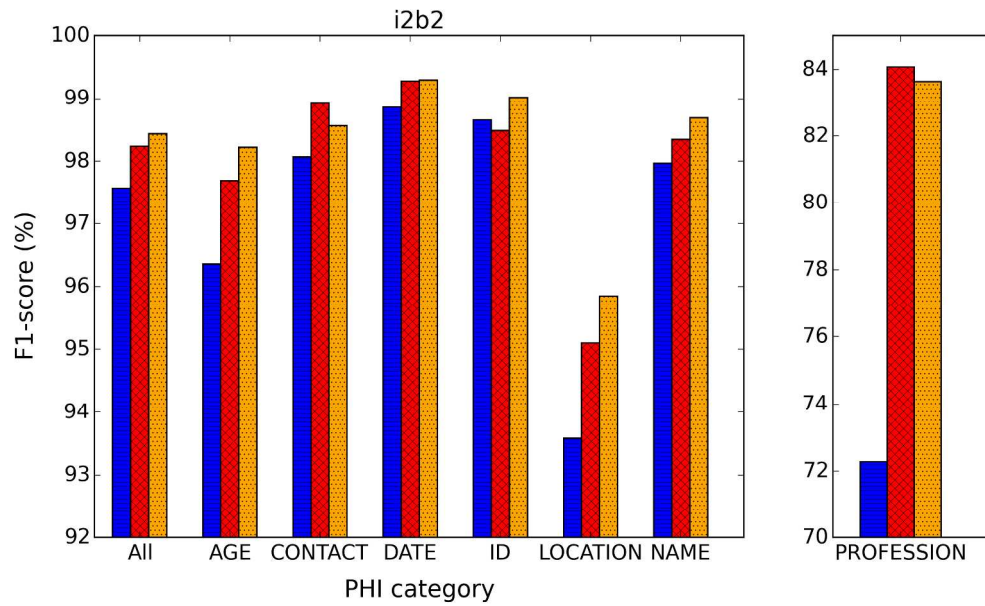
Figure 4: Impact of the number of labeled PHI instances in the training set on the model’s performance for each PHI type, in the i2b2 dataset. Figure (a) presents all PHI types, and Figure

(b) focuses on the most commonly occurring PHI types. More PHI instances in the training set helps increase F1-score, but some PHI types are harder to detect than others.

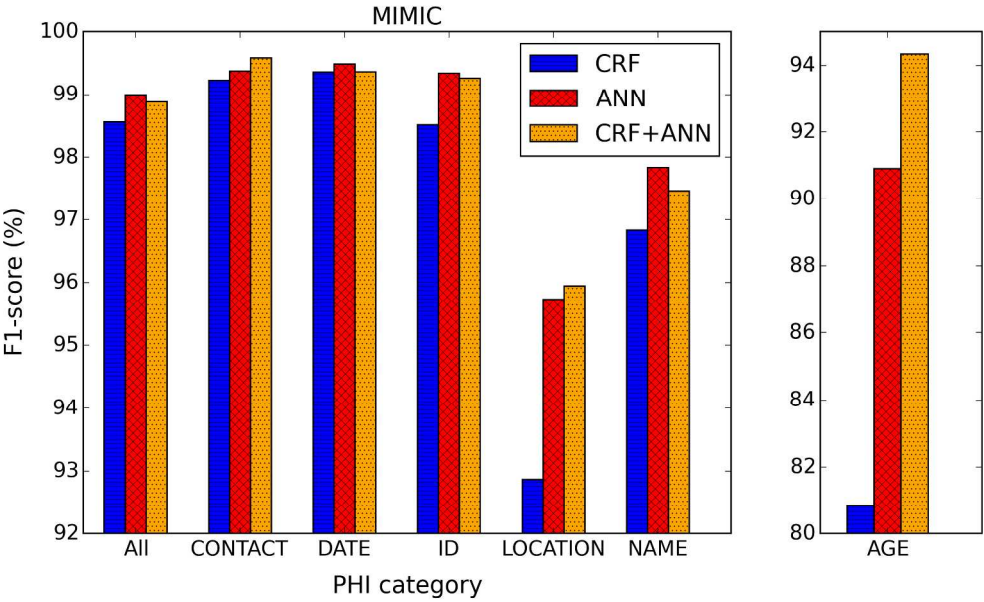
Figure 5: Ablation test performance based on binary HIPAA token-based evaluation. ANN is the model based on Artificial Neural Network. “— seq opt” is the ANN model without the label sequence optimization layer. “— pre-train” is the ANN model where token embeddings are initialized with random values instead of pre-trained embeddings. “— token emb” is the ANN model using only character-based token embeddings, without token embeddings. “— character emb” is the ANN model using only token embeddings, without character-based token embeddings.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46



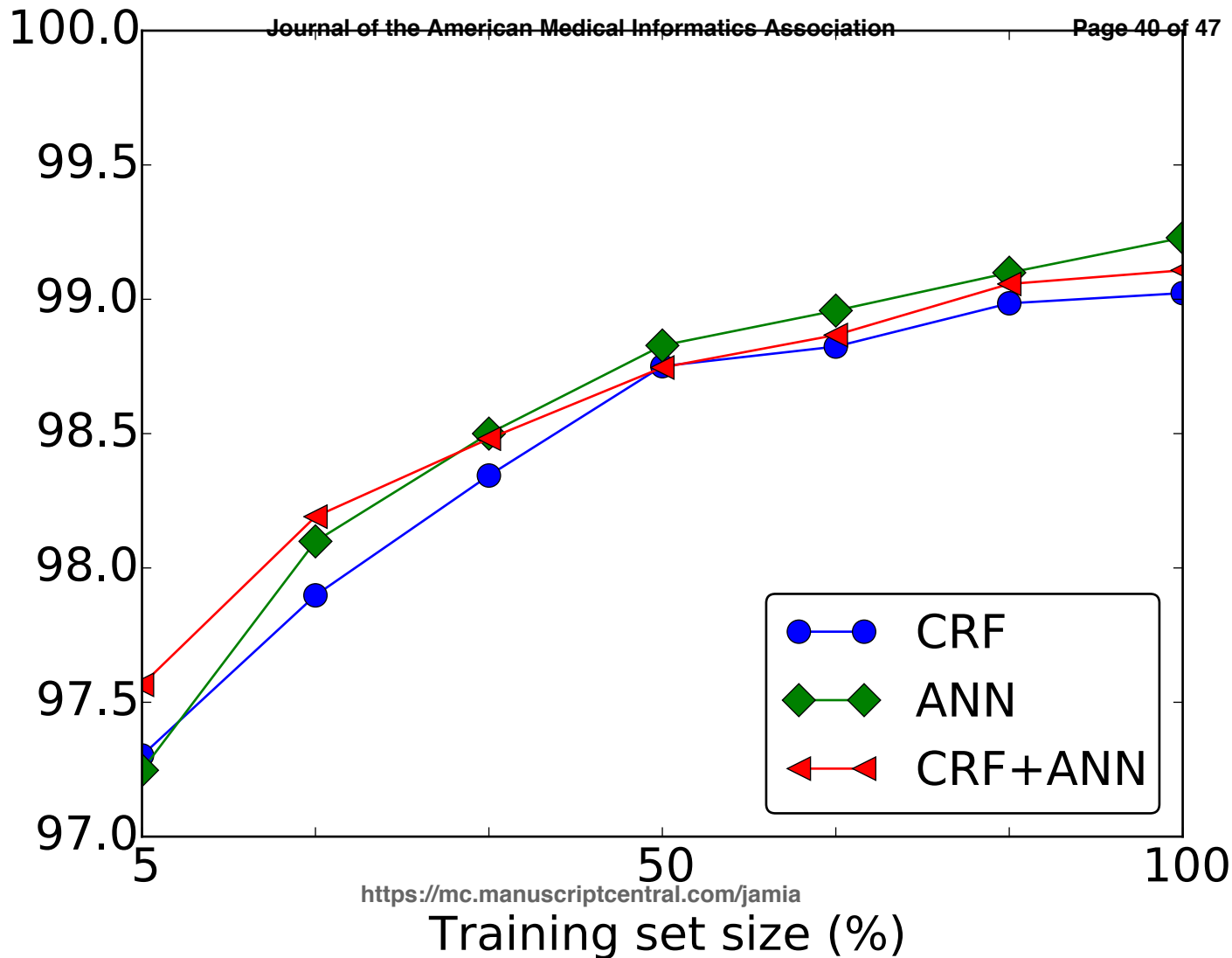


See submission.docx last page

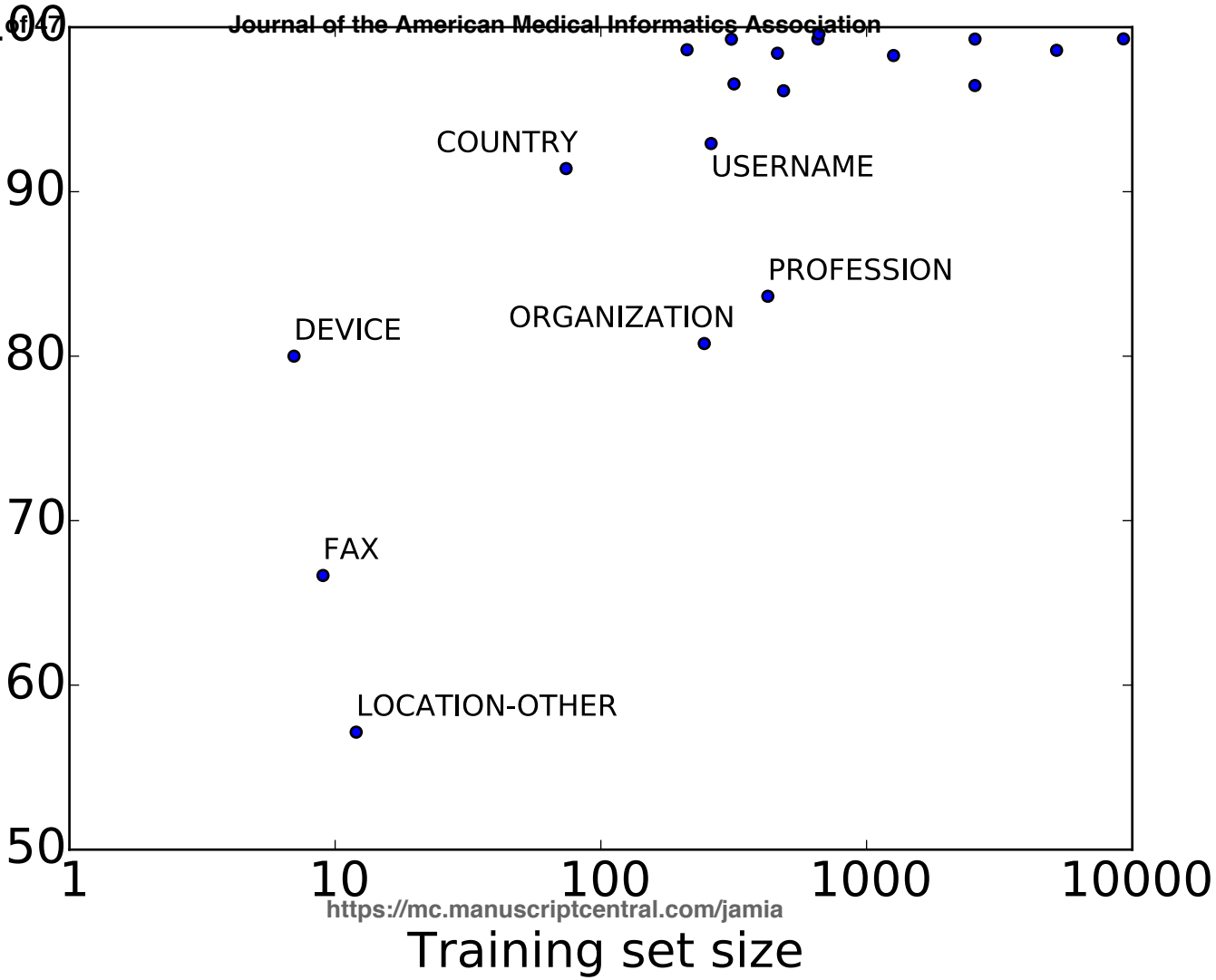


See submission.docx last page

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

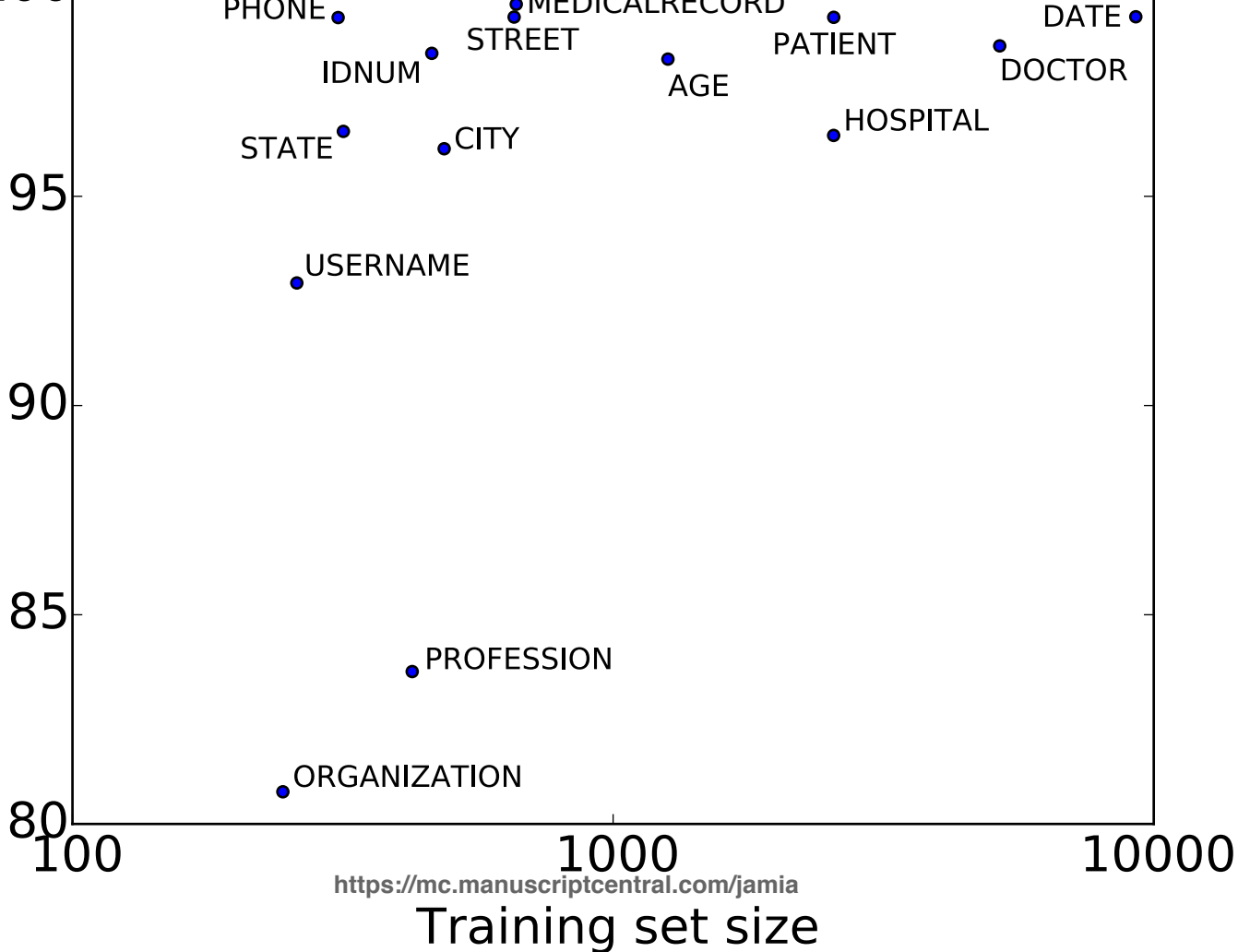


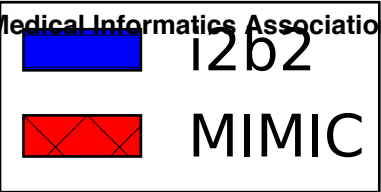
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32



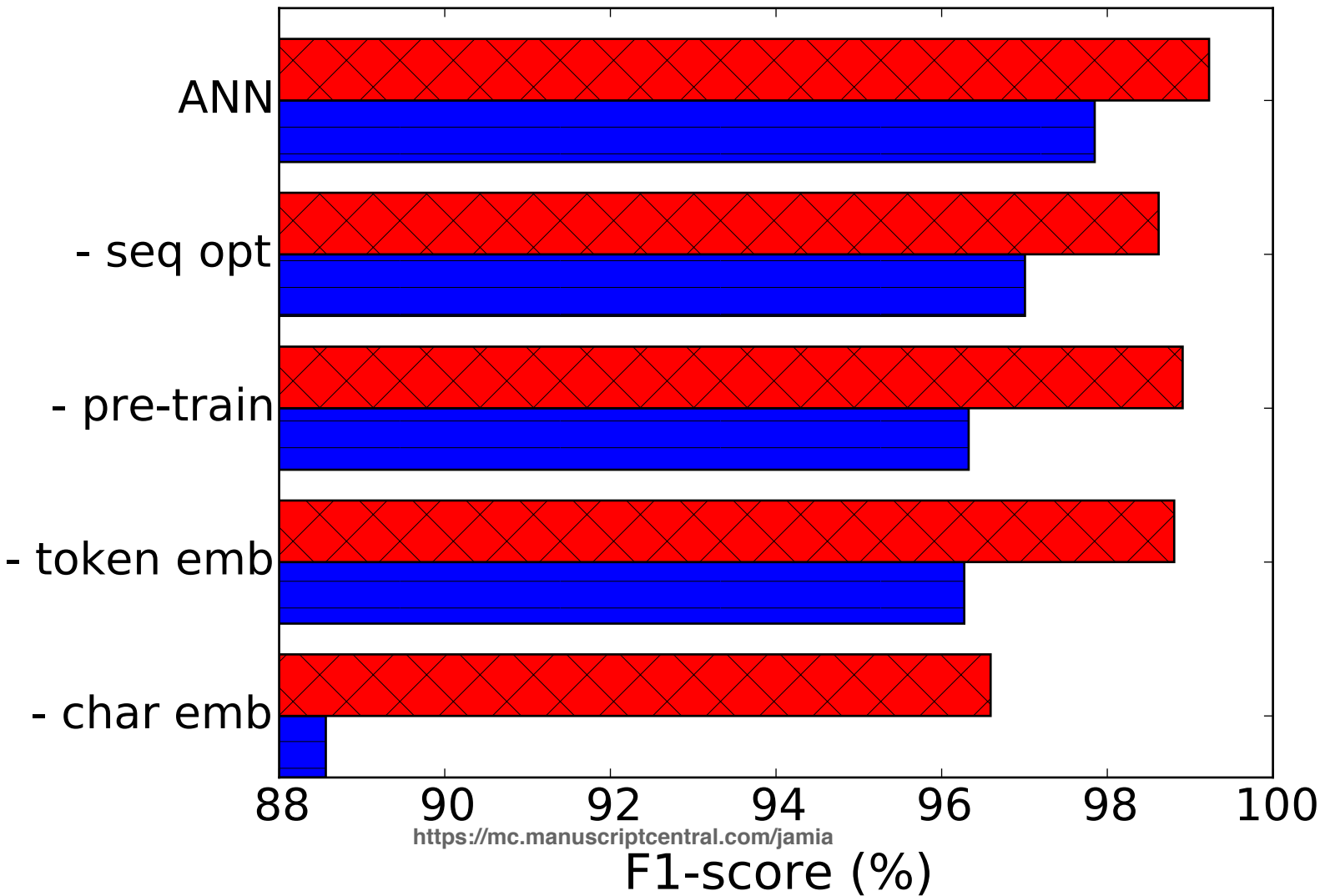
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

F1 score (%)





1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40



Appendix 1. Detailed result of all systems on i2b2 and MIMIC

Table A1 Performance (%) on i2b2 using binary, token-based evaluation for all PHI types. PHI types are grouped by categories and sorted in descending order of support. Since the official i2b2 evaluation script does not support evaluation at the PHI type label, we used our own code to compute the results.

| PHI type | MIST | | | CRF | | | ANN | | | CRF + ANN | | | Support |
|---------------------|-----------|--------|---------------|-----------|--------|--------------|-----------|--------|--------------|-----------|--------|--------------|---------|
| | Precision | Recall | F1-score | Precision | Recall | F1-score | Precision | Recall | F1-score | Precision | Recall | F1-score | |
| DATE category | 96.04 | 96.84 | 96.44 | 99.33 | 98.39 | 98.86 | 99.42 | 99.14 | 99.28 | 99.06 | 99.52 | 99.29 | 12532 |
| DATE | 96.04 | 96.84 | 96.44 | 99.33 | 98.39 | 98.86 | 99.42 | 99.14 | 99.28 | 99.06 | 99.52 | 99.29 | 12532 |
| NAME category | 97.60 | 93.18 | 95.34 | 99.11 | 96.84 | 97.96 | 98.67 | 98.02 | 98.34 | 98.22 | 99.15 | 98.68 | 4839 |
| DOCTOR | 97.01 | 93.60 | 95.28 | 98.98 | 96.60 | 97.77 | 98.65 | 97.54 | 98.09 | 98.25 | 98.94 | 98.60 | 3297 |
| PATIENT | 98.81 | 91.86 | 95.21 | 99.37 | 97.31 | 98.33 | 99.58 | 98.97 | 99.27 | 98.97 | 99.59 | 99.28 | 1450 |
| USERNAME | 100.00 | 98.91 | 99.45 | 100.00 | 97.83 | 98.90 | 86.79 | 100.00 | 92.93 | 86.79 | 100.00 | 92.93 | 92 |
| LOCATION category | 94.06 | 85.97 | 89.83 | 98.74 | 88.94 | 93.58 | 97.22 | 93.07 | 95.10 | 95.96 | 95.74 | 95.85 | 3001 |
| HOSPITAL | 94.10 | 89.97 | 91.99 | 99.32 | 91.35 | 95.17 | 96.17 | 94.48 | 95.32 | 95.68 | 97.24 | 96.46 | 1595 |
| STREET | 84.79 | 97.84 | 90.85 | 99.75 | 96.15 | 97.92 | 99.52 | 100.00 | 99.76 | 98.58 | 100.00 | 99.28 | 416 |
| CITY | 100.00 | 88.95 | 94.15 | 97.47 | 89.54 | 93.33 | 98.80 | 95.93 | 97.35 | 94.65 | 97.67 | 96.14 | 344 |
| STATE | 100.00 | 86.34 | 92.67 | 97.37 | 90.24 | 93.67 | 98.94 | 91.22 | 94.92 | 97.51 | 95.61 | 96.55 | 205 |
| ORGANIZATION | 100.00 | 34.01 | 50.76 | 92.21 | 48.30 | 63.39 | 97.17 | 70.07 | 81.42 | 92.92 | 71.43 | 80.77 | 147 |
| ZIP | 100.00 | 100.00 | 100.00 | 100.00 | 99.31 | 99.65 | 97.30 | 100.00 | 98.63 | 97.30 | 100.00 | 98.63 | 144 |
| COUNTRY | 100.00 | 43.85 | 60.96 | 96.26 | 79.23 | 86.92 | 95.15 | 75.39 | 84.12 | 92.86 | 90.00 | 91.41 | 130 |
| LOCATION-OTHER | 100.00 | 20.00 | 33.33 | 100.00 | 10.00 | 18.18 | 100.00 | 40.00 | 57.14 | 100.00 | 40.00 | 57.14 | 20 |
| ID category | 98.91 | 88.99 | 93.69 | 99.73 | 97.60 | 98.65 | 99.10 | 97.87 | 98.48 | 99.29 | 98.76 | 99.02 | 1126 |
| MEDICALRECORD | 100.00 | 89.89 | 94.68 | 100.00 | 99.73 | 99.86 | 99.32 | 99.45 | 99.39 | 99.32 | 99.86 | 99.59 | 732 |
| IDNUM | 96.87 | 89.01 | 92.77 | 99.19 | 95.81 | 97.47 | 98.65 | 95.81 | 97.21 | 99.20 | 97.64 | 98.42 | 382 |
| DEVICE | 100.00 | 33.33 | 50.00 | 100.00 | 25.00 | 40.00 | 100.00 | 66.67 | 80.00 | 100.00 | 66.67 | 80.00 | 12 |
| AGE category | 64.33 | 88.35 | 74.45 | 99.20 | 93.67 | 96.35 | 99.22 | 96.20 | 97.69 | 98.97 | 97.60 | 98.28 | 790 |
| AGE | 64.33 | 88.35 | 74.45 | 99.20 | 93.67 | 96.35 | 99.22 | 96.20 | 97.69 | 98.97 | 97.60 | 98.28 | 790 |
| CONTACT category | 99.27 | 97.61 | 98.44 | 99.27 | 96.90 | 98.07 | 99.52 | 98.33 | 98.92 | 98.80 | 98.33 | 98.57 | 419 |
| PHONE | 99.27 | 99.02 | 99.15 | 99.27 | 99.02 | 99.15 | 99.51 | 99.76 | 99.64 | 98.79 | 99.76 | 99.27 | 410 |
| FAX | 100.00 | 50.00 | 66.67 | 0.00 | 0.00 | 0.00 | 100.00 | 50.00 | 66.67 | 100.00 | 50.00 | 66.67 | 6 |
| EMAIL | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 3 |
| PROFESSION category | 100.00 | 0.88 | 1.75 | 92.24 | 59.41 | 72.27 | 91.67 | 77.65 | 84.08 | 87.99 | 79.71 | 83.64 | 340 |
| PROFESSION | 100.00 | 0.88 | 1.75 | 92.24 | 59.41 | 72.27 | 91.67 | 77.65 | 84.08 | 87.99 | 79.71 | 83.64 | 340 |
| All PHI types | 94.78 | 92.58 | 93.67 | 99.16 | 96.03 | 97.57 | 98.87 | 97.62 | 98.24 | 98.34 | 98.53 | 98.44 | 23047 |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Table A2 Performance (%) on MIMIC using binary, token-based evaluation for all PHI types. PHI types are grouped by categories and sorted in descending order of support. Since the official i2b2 evaluation script does not support evaluation at the PHI type label, we used our own code to compute the results.

| PHI type | MIST | | | CRF | | | ANN | | | CRF + ANN | | | Support |
|-------------------|-----------|--------|----------|-----------|--------|----------|-----------|--------|--------------|-----------|--------|--------------|---------|
| | Precision | Recall | F1-score | Precision | Recall | F1-score | Precision | Recall | F1-score | Precision | Recall | F1-score | |
| DATE category | 97.84 | 99.02 | 98.42 | 99.20 | 99.52 | 99.36 | 99.37 | 99.61 | 99.49 | 98.99 | 99.75 | 99.37 | 20627 |
| DATE | 97.84 | 99.02 | 98.42 | 99.20 | 99.52 | 99.36 | 99.37 | 99.61 | 99.49 | 98.99 | 99.75 | 99.37 | 20627 |
| NAME category | 97.17 | 94.82 | 95.98 | 97.72 | 95.95 | 96.83 | 98.01 | 97.66 | 97.83 | 96.87 | 98.06 | 97.46 | 3978 |
| DOCTOR | 96.98 | 95.32 | 96.15 | 97.63 | 96.16 | 96.89 | 98.03 | 97.63 | 97.83 | 96.78 | 98.07 | 97.42 | 3676 |
| PATIENT | 99.63 | 88.74 | 93.87 | 98.95 | 93.38 | 96.08 | 97.69 | 98.01 | 97.85 | 98.01 | 98.01 | 98.01 | 302 |
| LOCATION category | 92.80 | 88.67 | 90.69 | 97.17 | 88.94 | 92.87 | 97.22 | 94.28 | 95.73 | 96.07 | 95.82 | 95.95 | 1889 |
| HOSPITAL | 91.39 | 91.10 | 91.25 | 96.60 | 90.31 | 93.35 | 96.99 | 94.68 | 95.82 | 95.57 | 95.87 | 95.72 | 1259 |
| LOCATION-OTHER | 100.00 | 85.50 | 92.18 | 98.54 | 87.45 | 92.66 | 97.10 | 94.16 | 95.60 | 96.75 | 96.75 | 96.75 | 462 |
| STATE | 100.00 | 77.61 | 87.40 | 100.00 | 86.57 | 92.80 | 100.00 | 95.52 | 97.71 | 100.00 | 97.02 | 98.49 | 67 |
| STREET | 70.67 | 86.89 | 77.94 | 100.00 | 78.69 | 88.07 | 100.00 | 85.25 | 92.04 | 100.00 | 86.89 | 92.98 | 61 |
| ZIP | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 24 |
| COUNTRY | 100.00 | 25.00 | 40.00 | 75.00 | 56.25 | 64.29 | 93.33 | 87.50 | 90.32 | 82.35 | 87.50 | 84.85 | 16 |
| CONTACT category | 99.36 | 97.36 | 98.35 | 99.86 | 98.61 | 99.23 | 99.51 | 99.24 | 99.37 | 99.58 | 99.58 | 99.58 | 1438 |
| PHONE | 99.36 | 97.36 | 98.35 | 99.86 | 98.61 | 99.23 | 99.51 | 99.24 | 99.37 | 99.58 | 99.58 | 99.58 | 1438 |
| ID category | 98.33 | 95.92 | 97.11 | 99.67 | 97.39 | 98.51 | 99.67 | 99.02 | 99.34 | 99.35 | 99.18 | 99.26 | 612 |
| IDNUM | 98.33 | 95.92 | 97.11 | 99.67 | 97.39 | 98.51 | 99.67 | 99.02 | 99.34 | 99.35 | 99.18 | 99.26 | 612 |
| AGE category | 100.00 | 32.14 | 48.65 | 100.00 | 67.86 | 80.85 | 92.59 | 89.29 | 90.91 | 100.00 | 89.29 | 94.34 | 28 |
| AGE | 100.00 | 32.14 | 48.65 | 100.00 | 67.86 | 80.85 | 92.59 | 89.29 | 90.91 | 100.00 | 89.29 | 94.34 | 28 |
| All PHI types | 97.51 | 97.54 | 97.52 | 98.91 | 98.20 | 98.56 | 99.05 | 98.94 | 99.00 | 98.54 | 99.22 | 98.88 | 28572 |

Appendix 2. MIST configuration

For the MIST system, only 8 coarse-grained labels are supported when using the AMIA de-identification mode: HOSPITAL, PATIENT, DOCTOR, DATE, LOCATION, ID, PHONE, and AGE. Therefore, all PHI types were mapped to the most appropriate label in order to match the labels supported by the system. Table A3 presents the mapping. All PHI types supported by the MIST labels (DATE, DOCTOR, PATIENT, HOSPITAL, PHONE and AGE) were mapped to themselves. Each unsupported PHI type was mapped to its category, if the category is among the MIST labels. The PHI types USERNAME and FAX were mapped to the semantically closest labels, ID and PHONE, respectively. The two remaining PHI types EMAIL and PROFESSION were mapped to non-PHI, as there was no appropriate label. The PHI type EMAIL is negligible since it has only one instance (comprising 3 tokens) in the test set of only the i2b2 dataset, and the PHI type PROFESSION is non-HIPPA.

Table A3 Mapping from PHI types to MIST labels.

| PHI category | PHI type | MIST label |
|--------------|----------------|------------|
| DATE | DATE | DATE |
| NAME | DOCTOR | DOCTOR |
| | PATIENT | PATIENT |
| | USERNAME | ID |
| LOCATION | HOSPITAL | HOSPITAL |
| | STREET | LOCATION |
| | CITY | |
| | STATE | |
| | ORGANIZATION | |
| | ZIP | |
| | COUNTRY | |
| ID | LOCATION-OTHER | |
| | MEDICALRECORD | ID |
| | IDNUM | |
| AGE | DEVICE | |
| | AGE | AGE |
| CONTACT | PHONE | PHONE |
| | FAX | non-PHI |
| | EMAIL | |
| PROFESSION | PROFESSION | non-PHI |

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

The following code was used to train and tag with MIST. Note that MIST was configured to utilize the same gazetteers as the other systems.

- Train: bin/MATModelBuilder --task "AMIA Deidentification" --save_as_default_model --nthreads=20 --max_iterations=15 --lexicon_dir="\$PWD/sample/mist/gazetteers" --input_files "\$PWD/sample/mist/dataset/train/*.json"
- Tag: bin/MATEngine --task "AMIA Deidentification" --workflow Demo --input_dir "\$PWD/sample/mist/dataset/test" --input_file_type mat-json --output_dir "\$PWD/sample/mist/dataset/test/test_out" --output_file_type mat-json --tagger_local --steps "tag"