

PHI Scrubber: A Deep Learning Approach

Abhai Kollara Dilip* Kamal Raj K* Malaikannan Sankarasubbu
Saama Technologies AI Research Lab
Chennai, India

{a.kollara, kamal.raj, malaikannan.sankarasubbu}@saama.com

Abstract— Confidentiality of patient information is an essential part of Electronic Health Record System. Patient information, if exposed, can cause a serious damage to the privacy of individuals receiving healthcare. Hence it is important to remove such details from physician notes. A system is proposed which consists of a deep learning model where a de-convolutional neural network and bi-directional LSTM-CNN is used along with regular expressions to recognize and eliminate the individually identifiable information. This information is then removed from a medical practitioner’s data which further allows the fair usage of such information among researchers and in clinical trials.

I. INTRODUCTION

Patient identification for clinical trials is one of the major challenges in pharma industry since major percentage of clinical trials fail due to patient enrollment issues. The structured EHR data that are used to identify potential candidates for trials is designed to support billing processes with Health Insurance companies and may not include all the necessary information for identification. Doctor notes on the other hand contain a wealth of information but remain inaccessible due to the presence of protected health information that must remain confidential. The HIPAA act of 1996 sets forth national standards for the transaction of electronic healthcare information. As a direct result of this, any document containing information that can be used to trace the patient has to be treated as protected health information (PHI). Such documents must remain confidential unless there is a clear consent from the patient involved. HIPAA specifies a set of 18 identifiers whose presence make a document PHI (Refer methods section). De-identifying such documents involve removing the identifiers. The notes however are unstructured data and thus requires a system that can learn patterns in the language.

II. DATA

Authentic protected health information is not available for research (in large amounts) due to the very problem we are trying to solve. Since deep learning models require a large amount of data to perform well, we are forced to use a substitute dataset. The task at hand is an entity recognition problem and thus we use other datasets that are widely used for the problem.

Ontonotes corpus [1] is a large manually annotated corpus containing text from a variety of genres (news, talk shows,

newsgroups, conversational phone calls) in 3 languages (English, Chinese and Arabic). The dataset contains several layers of annotations for various natural language processing tasks. For the purpose of our experiments, we use the entity names layer of English language to train the models.

The dataset contains sentences and the corresponding entity label of each token in the sentence. The sentences are segmented into tokens using the Penn Treebank tokenizer. The tokenization splits the sentences based on whitespace as well as punctuation (eg: Robert’s friend cannot be there -> ”Robert”, ”s”, ”friend”, ”can”, ”not”, ”be”, ”there”).

The named entities in the dataset are classified into 18 possible categories which are labelled using BILOU labelling scheme. The BILOU scheme labels each word as either of Beginning (**B**), Inside (**I**) or Last (**L**) of an entity if the entity is multi-word. Unit length entities are marked as **U** while a non-entity is labelled as **O**.

PERSON	Person
NORP	Nationalities or religious or political groups
FACILITY	Buildings, airports, highways, bridges, etc
ORGANIZATION	Companies, agencies, institutions, etc.
GPE	Countries, cities, states
LOCATION	Non-GPE locations, mountain ranges, bodies of water
PRODUCT	Vehicles, weapons, foods, etc. (Not services)
EVENT	Named hurricanes, battles, wars, sports events, etc.
WORK OF ART	Titles of books, songs, etc.
LAW	Named documents made into laws
LANGUAGE	Any named language
DATE	Absolute or relative dates or periods
TIME	Times smaller than a day
PERCENT	Percentage
MONEY	Monetary values, including unit
QUANTITY	Measurements, as of weight or distance
ORDINAL	first, second etc
CARDINAL	Numerals that do not fall under another type

TABLE I
ENTITY TYPES IN ONTONOTES DATASET

III. METHODS

A PHI scrubber should remove the following identifiers from a given document.

- 1) Name
- 2) Address (all geographic subdivisions smaller than state, including street address, city county, and zip code)
- 3) All elements (except years) of dates related to an individual (including birthdate, admission date, discharge date, date of death, and exact age if over 89)

*Both authors had equal contribution

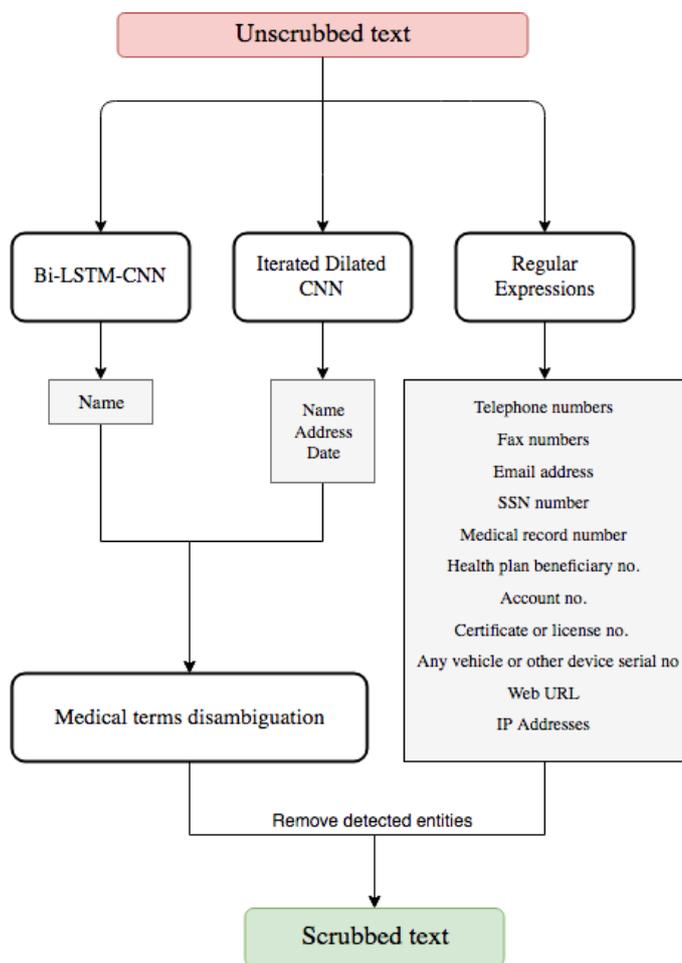


Fig. 1. Scrubber architecture

- 4) Telephone numbers
- 5) Fax number
- 6) Email address
- 7) Social Security Number
- 8) Medical record number
- 9) Health plan beneficiary number
- 10) Account number
- 11) Certificate or license number
- 12) Any vehicle or other device serial number
- 13) Device identifiers and serial numbers
- 14) Web URL
- 15) Internet Protocol (IP) Address
- 16) Finger or voice print
- 17) Photographic image - Photographic images are not limited to images of the face.
- 18) Any other characteristic that could uniquely identify the individual

For text data we ignore 16, 17 and 18. The dataset we chose contains 1, 2 and 3 thus allowing us to use deep learning models to identify them. The models can capture variations of the entities that cannot be captured using traditional methods. For the rest of the identifiers regular expressions are sufficient since all of them follow

a predictable and limited set of pattern.

Deep learning models

A combination of two deep learning models are applied here - A bidirectional variation of long short-term memory combined with character level convolutional neural networks [2] and an iterated dilated convolutions model (ID-CNN) [3]

Since recurrent neural networks are known for its performance in sequence labeling tasks, initially we applied the bi-LSTM CNN model for this task. But due to the extremely high number of false positives, it produced we used it solely for identifying names. An iterated dilated convolutions model was then added to the model to improve its performance. This model was chosen for its state of the art results in named entity recognition. Both models are trained separately to identify all the entity types available in the Ontonotes dataset. But for inference, Bi-LSTM detects "NAME" entity while ID-CNN model captures "NAME", "ADDRESS" and "DATE" entities.

Bidirectional LSTM-CNN

The Bi-LSTM-CNN model has the ability to capture context from the entire sequence as well as incorporate character level features. For the model, we create 3 different vector representations for each word in a sentence

- Word embedding
- Character level embedding
- Additional word features

Word embeddings are obtained using a lookup table. The pre-trained Glove embeddings [4] are used for this purpose.

For character level embeddings, we represent each character in a word by a vector, which are then concatenated. A 1D convolution is applied over each word and is followed by max pooling to obtain character level representations.(Fig.3)

Finally, we use additional word features by categorizing each word into one of the following classes, which are then mapped to a random vector.

- Numeric
- All lower case
- All upper case
- Initial upper case
- Mostly numeric
- Contains digit
- Other

The three different representations are then concatenated and fed into a stacked bidirectional LSTM[5][6] layer. The outputs at each timestep are passed through a softmax layer to get the score for each word. Each word is classified into one of 5 varieties (we use BILOU labeling scheme) of the 18 different categories available in the Ontonotes dataset. Even though we require only 3 of the 18 entity types, we train the model using all entity types.

We used bidirectional LSTM-CNN model since the CNNs are able to capture character level features. The model combines character level features with word-level features thus making it tolerant to minute spelling variations.

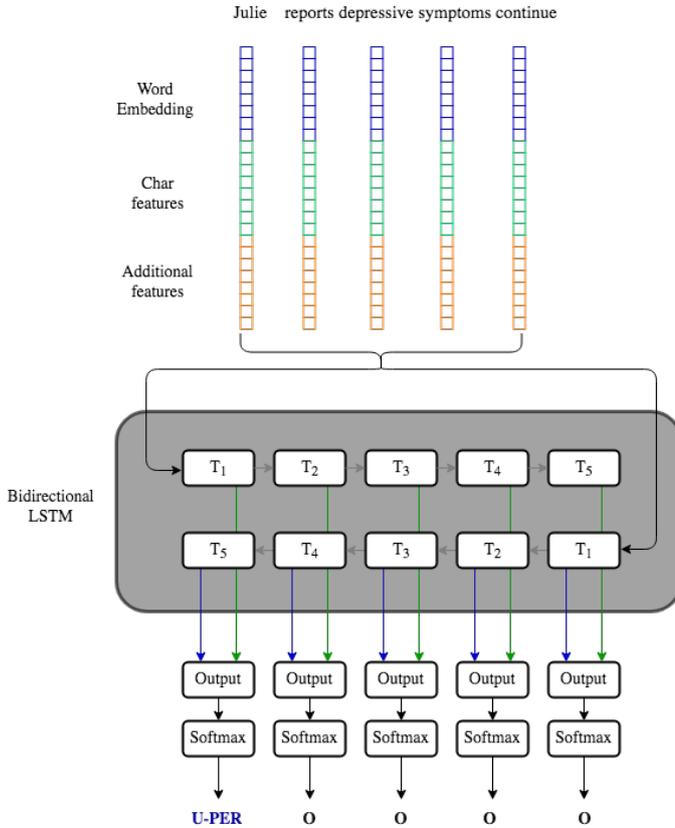


Fig. 2. Bidirectional LSTM

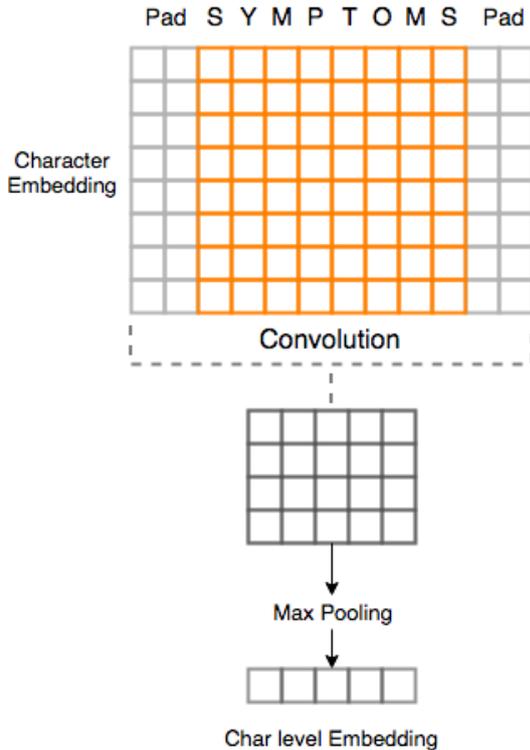


Fig. 3. Character CNN

Word embedding size	300
Character embedding size	95
LSTM cell size	200
LSTM Layers	2
Learning rate	0.001
Optimizer	Nadam
Epochs	65

TABLE II
TRAINING PARAMETERS FOR BI-LSTM-CNN NETWORK

Iterated Dilated Convolution

Here, we replicate the model created by Strubell et al.[3]. Similar to the previous model, this network takes in a sequence of words and outputs for each word a probability distribution over all possible labels. Besides the superior results, the use of convolutions allows us to exploit parallel computation to the maximum.

The ability of recurrent neural networks to capture temporal information from the entire sequence of its inputs has made it the workhorse of NLP tasks in deep learning. However, due to the nature of LSTM computations, they are difficult to run parallelly. Convolutional networks on the other hand can easily be run in parallel but have fixed contexts ie they fail to incorporate the entire sequence as its context. Dilated convolutions provide a workaround for this obstacle[7]. Stacked dilated CNNs can easily incorporate global information from a whole sentence or document. They allow the effective input width to grow exponentially with the depth of the network. Dilated CNNs operate on similar principles of convolutional networks except that the dilated window skips over every dilation width 'd' inputs (See Fig 3.).

Similar to the Bi-LSTM-CNN model each word has multiple representations. A pretrained word embedding and additional word feature embedding. Word embeddings are obtained from pre-trained Lample embeddings[8]. The additional word features are similar to those used in Bi-LSTM CNN model. However, the following 4 categories are used instead

- All upper case
- Initial upper case
- Camel case
- Other

The two different representation are concatenated and fed into iterated dilated CNN layers. The ID-CNN architecture repeatedly applies the same block of dilated convolution to the input representations. The transition parameter and score of each token is passed through a Viterbi decoding to get the score for each word. Similar to the Bi-LSTM-CNN model, we use a BILOU labeling scheme here.

Regular expressions

Regular expressions are used to detect all the below-given identifiers. We observed regular expression were sufficient to detect the following identifiers.

- Telephone numbers

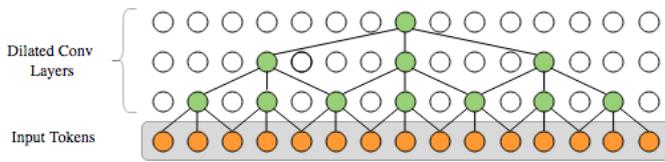


Fig. 4. Iterated Dilated Convolution

Word embedding size	100
ID-CNN layers	3
ID-CNN dilation	[1,2,1]
ID-CNN width	3
ID-CNN filters	400
ID-CNN blocks	1
Learning rate	0.0001
Optimizer	adam
Epochs	100

TABLE III

TRAINING PARAMETERS FOR ID-CNN NETWORK

- Fax numbers
- Email address
- SSN number
- Medical record number
- Health plan beneficiary number
- Account number
- Certificate or license number
- Any vehicle or other device serial number
- Web URL
- Internet Protocol (IP) Address

Medical terms disambiguation

Due to the absence of medical terms in Ontonotes dataset several such terms are misclassified by the deep learning models. The presence of names in disease and drug names (eg: Parkinson’s disease) are also a cause for misclassification. The sentence structure at times causes body part names to be classified as a location. To circumvent this problem, we use a separate disambiguation module for the outputs of the named entity recognition module. The "PERSON", "LOCATION" and "DATE" entities that are detected by the deep learning models are passed into this disambiguation module. Here, a fuzzy search is done over a dictionary of medical terms. For non-abbreviated words, we consider the Levenshtein distance for matching. If the word has a Levenshtein distance less than 2 for any word in the dictionary we consider this as a match. For abbreviated words only exact matches are considered. Besides this, we also check if a word ends with drug name stem (eg.-dralazine, -pristone etc) to further check for drug names. Once all the matches have been found, their labels are removed thus preventing the entities from being removed.

IV. TRAINING

We trained the BiLSTM-CNN and ID-CNN model for 100 epochs using Adam optimizer. Both models are trained independently on the Ontonotes dataset to identify all types of entities in the dataset. While training only the best weights

Input

Mr. Stanford was previously seen on February 25, 2010 for evaluation of chronic GE reflux disease. Subsequently, an EGD was performed on March 1, 2010 concluding a small hiatal hernia and gastritis from which biopsies were obtained. There is no evidence of peptic ulcer disease or neoplasm. Pathology findings document no evidence of Barrett’s esophagus or H. pylori infection.

Fetch result

Output

Mr. PERSON was previously seen on DATE for evaluation of chronic GE reflux disease. Subsequently, an EGD was performed on DATE concluding a small hiatal hernia and gastritis from which biopsies were obtained. There is no evidence of peptic ulcer disease or neoplasm. Pathology findings document no evidence of Barrett’s esophagus or H. pylori infection.

Fig. 5. PHI De-Identifier

are saved. The training samples are batched into sentences of equal length. Parameters for training the models are given in tables II and III. A single Nvidia Geforce GTX 1080 took about 6 hours to complete the training for BiLSTM-CNN while the ID-CNN model took 1.5 hours.

V. RESULTS

Evaluation of the models was done on the Ontonotes dataset. The ID-CNN model achieves a segmented micro F1 score of 86.84 while Bi-LSTM model achieves 86.5. The models are evaluated on the performance on Ontonotes dataset. Even though both models produce near similar results, the ID-CNN model is significantly faster during training.

VI. CONCLUSION

We present a deep learning based approach for the removal of PHI from text documents. Currently, we use a regular expression to assist the removal of identifiers. In future, we hope a single unified deep learning model can identify all identifiers with higher accuracy.

ACKNOWLEDGMENT

The authors would like to thank Dr. Anand Dubey and V Archana for reviewing the manuscript and for providing their technical inputs.

REFERENCES

- [1] Eduard Hovy, Mitchell Marcus, Martha Palmer, Lance Ramshaw, and Ralph Weischedel. Ontonotes: The 90 In *Proceedings of the Human Language Technology Conference of the NAACL, Companion Volume: Short Papers*, NAACL-Short '06, pages 57–60, Stroudsburg, PA, USA, 2006. Association for Computational Linguistics.

- [2] Jason P. C. Chiu and Eric Nichols. Named entity recognition with bidirectional lstm-cnns. *CoRR*, abs/1511.08308, 2015.
- [3] Emma Strubell, Patrick Verga, David Belanger, and Andrew McCallum. Fast and accurate sequence labeling with iterated dilated convolutions. *CoRR*, abs/1702.02098, 2017.
- [4] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
- [5] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997.
- [6] Alex GRAVES and Jürgen SCHMIDHUBER. Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural networks*, 18(5-6):602–610, 2005.
- [7] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *CoRR*, abs/1511.07122, 2015.
- [8] Guillaume Lample, Miguel Ballesteros, Sandeep Subramanian, Kazuya Kawakami, and Chris Dyer. Neural architectures for named entity recognition. *CoRR*, abs/1603.01360, 2016.