

# IBM Research Report

## Creating a Corpus of Clinical Notes Manually Tagged for Part-of-Speech Information

**Serguei V. Pakhomov\*, Anni Coden, Christopher G. Chute\***

IBM Research Division  
Thomas J. Watson Research Center  
P.O. Box 704  
Yorktown Heights, NY 10598

\*Division of Medical Informatics Research  
Mayo Clinic  
Rochester, MN



Research Division

Almaden - Austin - Beijing - Haifa - India - T. J. Watson - Tokyo - Zurich

# Creating a Corpus of Clinical Notes Manually Tagged for Part-of-Speech Information

**Serguei V. PAKHOMOV**  
Division of Medical Informatics  
Research, Mayo Clinic  
Rochester, MN  
[Pakhomov.Serguei@mayo.edu](mailto:Pakhomov.Serguei@mayo.edu)

**Anni CODEN**  
IBM, T.J. Watson Research  
Center,  
Hawthorne, NY 10532  
[anni@us.ibm.com](mailto:anni@us.ibm.com)

**Christopher G. CHUTE**  
Division of Medical  
Informatics Research, Mayo  
Clinic Rochester, MN  
[Chute.Christopher@mayo.edu](mailto:Chute.Christopher@mayo.edu)

## Abstract

This paper presents a pilot project aimed at generating a corpus of linguistically annotated clinical data to be used for training and testing of NLP and other text processing applications and techniques in the medical domain. We describe and discuss the process of training three medical index experts to perform linguistic annotation. We list some of the challenges as well as encouraging results pertaining to inter-rater agreement and consistency of annotation for part-of-speech information and lay out future directions for this initiative. We also present preliminary experimental results indicating the necessity for adapting state-of-the-art POS taggers to the sublanguage domain of clinical notes.

## 1 Introduction

Having reliable part-of-speech (POS) information is critical to successful implementation of NLP techniques for processing unrestricted text in the biomedical domain. State-of-the-art automated POS taggers achieve accuracy of 93% - 98% and the most successful implementations are based on statistical approaches to POS tagging. Taggers based on Hidden Markoff Model (HMM) technology currently appear to be in the lead. The prime examples of such implementations include the Trigrams'n'Tags tagger (Brandts 2000), Xerox tagger (Cutting et al. 1992) and LT POS tagger (Mikheev 1997). Maximum Entropy (MaxEnt) based taggers also seem to perform very well (Ratnaparkhi 1996, Jason Baldrige, Tom Morton, and Gann Bierner <http://maxent.sourceforge.net>). One of the issues with statistical POS taggers is that most of them need a representative amount of hand-labeled training data either in the form of a comprehensive lexicon and a corpus of untagged data or a large corpus of text annotated for POS or a combination of the two. Currently, most of the POS tagger accuracy reports are based on the experiments involving Penn Treebank data that consists of Department of Energy abstracts,

completely re-tagged with the Treebank tagset Brown corpus, Department of Agriculture bulletins, Library of America texts, MUC-3 texts, sentences from IBM computer manuals and the ATIS corpus of spontaneous speech transcriptions of Air Travel Information Systems project (Marcus, 1993). The purpose of the selection of sources in Treebank is to represent the general English sublanguage domain. It is not entirely clear how representative the general English language vocabulary and structure are of a specialized sub-domain such as clinical reports.

A well-recognized problem is that the accuracy of all current POS taggers drops dramatically on unknown words. For example, while the TnT tagger performs at 97% accuracy on known words in the Treebank, the accuracy drops to 89% on unknown words (Brandts, 2000). The LT POS tagger is reported to perform at 93.6-94.3% accuracy on known words and at 87.7-88.7% on unknown words using a cascading guesser (Mikheev, 1997). The overall results for both of these taggers are much closer to the high end of the spectrum because the rate of the unknown words in the tests performed on the Penn Treebank corpus is generally relatively low – 2.9% (Brandts, 2000). From these results, we can conclude that the higher the rate of unknown vocabulary, the lower the overall accuracy will be, necessitating the adaptation of the taggers trained on Penn Treebank to sublanguage domains with vocabulary that is substantially different from the one represented by the Penn Treebank corpus.

Based on the observable differences between the clinical and the general English discourse and POS tagging accuracy results on unknown vocabulary, it is reasonable to assume that a tagger trained on general English may not perform as well on clinical notes, where the percentage of unknown words will increase. However, in order to test this assumption, a “gold standard” corpus of clinical notes needs to be manually annotated for POS information. The issues with the annotation process constitute the primary focus of this paper.

In the remainder of the paper, we describe an effort to train three medical coding experts to mark

the text of clinical notes for part-of-speech information. The motivation for using medical coders rather than trained linguists is threefold. First of all, due to confidentiality restrictions, in order to develop a corpus of hand labeled data from clinical notes one can only use personnel authorized to access patient information. The only way to avoid it is to anonymize the notes prior to POS tagging which in itself is a difficult and expensive process (Ruch et al. 2000). Second, medical coding experts are well familiar with clinical discourse, which helps especially with annotating medical specific vocabulary. Third, the fact that POS tagging can be viewed as a classification task makes the medical coding experts highly suitable because their primary occupation and expertise is in classifying patient records for subsequent retrieval.

We discuss the training process, various issues that surfaced due to certain peculiarities of clinical notes discourse, and the results of a pilot study that evaluates the level of inter-rater agreement between the three annotators. We show that given a good set of guidelines, medical coding experts can be trained in a limited amount of time to perform a linguistic task such as POS annotation at a high level of agreement on both clinical notes and Penn Treebank data. Finally, we report on a set of training experiments performed with the TnT tagger (Brandts, 2000).

## 2 Annotator Training

Prior to this study, the three annotators who participated in it had a substantial experience in coding clinical diagnoses but virtually no experience in POS markup. The training process consisted of a general and rather superficial introduction to the issues in linguistics as well as some formal training using the POS tagging guidelines developed by Santorini (1991) for tagging Penn Treebank data. The formal training was followed by informal discussions of the data and difficult cases pertinent to the clinical notes domain which often resulted in slight modifications to the Penn Treebank guidelines. Below is a list of such cases and the modifications:

### 2.1 Drug Names

Drug names are tagged as personal pronouns (NNP) regardless of whether they are capitalized. Numeric and other types of attributes of drug names are considered to be part of the drug name and are tagged as NNP as well (Ex. Extra/NNP Strength/NNP Tylenol/NNP #4/NNP). At a later point #4 in this example and other examples of a symbol followed by a numeral, were converted to

two tokens “#” and the numeral and re-tagged as a SYM and CD respectively.

### 2.2 Dosages

Dosages actually have a recognizable structure which, in most cases, consists of the numeric magnitude of the dose (Ex. 500), the measurement unit (Ex. mg), method of delivery (Ex. orally) and the interval of delivery (Ex. b.i.d.). The amount is tagged as a cardinal number CD, the measurement unit is tagged as NNS or NN depending on whether the measurements is of a plural or singular entity, the method of delivery is tagged the same as normal text. A more problematic case is the interval of delivery which may include actual phrases “twice daily” or their Latin equivalents “b.i.d.” The former is tagged as if it were normal text, the latter is tagged as a foreign word - FW.

### 2.3 Foreign Words

The medical domain is permeated with words of Latin and Greek origin and making a clear distinction between a medical word of foreign origin and a medical word that is a foreign word is rather difficult even in cases where the word has retained its foreign pronunciation. – “polymyalgia rheumatica”. In the majority of such cases, the words in question are somehow related to a condition, procedure or some other clinically related phenomenon and normally not used on their own. Based on that, we tag potential candidates for foreign words as if they formed a unit. By this token, “polymyalgia rheumatica” is tagged as a noun followed by another noun and form a compound noun - polymyalgia/NN rheumatica/NN.

### 2.4 Special Symbols

Medical transcriptionists often use a one keystroke shorthand for words that occur relatively frequently. For example, “+” is often used for “positive” as in “positive throat culture”, “-“ – for “negative throat culture.” The pound sign “#” is often used to mean “pounds” or “fracture.” “x” is often used to mean “scar” as is “x 2” – two scars. It is not entirely clear at this point if it would be more beneficial to treat these symbols as actual symbols (SYM) or as special kinds of abbreviations. In the latter case, they would be tagged as if the actual word was used in place of the symbol. For the time being we rather arbitrarily decided to mark them as symbols SYM – “x/SYM 2/CD.”

## 3 Annotation structure and format

Each clinical note is represented as an XML (Extensible Markup Language) document with the underlying schema shown in Figure 1. The raw

clinical notes go through automatic POS tagging whose results are rendered in XML format and are then presented in a graphical XML editor for correction.

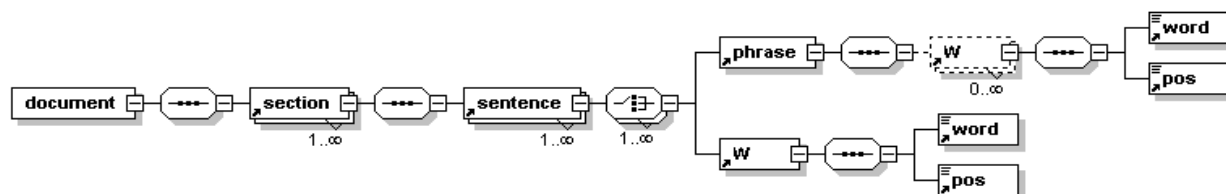


Figure 1. XML Schema for a Clinical NLP document

The top node in the XML schema represents the whole document and can have a number of “section” nodes under it that represent various subsections of a clinical note such as History of Present Illness (HPI), Chief Complaint (CC) and other standard sections compliant with HL-7 Clinical Document Architecture specification<sup>1</sup>. The section nodes branch out into sentence nodes and each sentence node can branch out directly into W nodes, which are combinations of word orthographic representation and its part-of-speech information. The W nodes can be arranged into phrases with various types of attributes. These can be either traditional linguistic phrases such as NP, VP, etc. or specialized medical phrases such as “drug mention.” The “phrase” nodes group W nodes under it into meaningful chunks. The “meaningfulness” of the chunks as well as their boundaries is left up to the trained medical indexers to determine. The “phrase” node also allows simultaneous anonymization of data by introducing “phrase” nodes with special attributes to group patient identifying information.

The “dual residence” of the W node under both the “phrase” node and the “sentence” node in the schema reflects the fact that we are aiming at working with incomplete shallow parses rather than full parses. In the former case, not all phrases that comprise a sentence have to be identified and, therefore, not every word has to be part of a phrase.

All data presented to the annotators is preprocessed before annotation. The pre-processing includes sentence boundary detection, tokenization and priming with part-of-speech tags generated by a MaxEnt tagger (Maxent 1.2.4

<sup>1</sup> Health Level 7 is a medical standards organization part of whose purpose is to establish and maintain various standards applicable to medical information management.

package (Baldrige et al.)) trained on Penn Treebank data. To expedite the annotation process, an in-house Java based editor was developed.

## 4 Annotator agreement

In order for any large scale annotation project based on the efforts of the three annotators in this study to be successful, we need to ensure internal as well as external consistency of the annotation. First of all, we need to make sure that the annotators agree amongst themselves (internal consistency) on how they mark up text for part-of-speech information. Second, we need to find out how closely the annotators generating data for this study agree with the annotators of an established project such as Penn Treebank (external consistency). If both tests show high levels of agreement, then we can safely assume that the annotators in this study are able to generate part-of-speech tags for biomedical data that will be consistent with a widely recognized standard and can work independently of each other.

### 4.1 Methods

Two types of measures of consistency were collected – absolute agreement and Kappa coefficient. The absolute agreement was calculated according to the following formula in (1).

$$(1) \quad Abs\ agr = 100 * \left( \frac{M}{T} \right)$$

where M is the total number of time all annotators agreed on a tag and T is the total number of tags.

Kappa coefficient is given in (2) (Carletta 1996)

$$(2) \quad Kappa = \frac{P(A) - P(E)}{1 - P(E)}$$

P(A) is the proportion of times the annotators actually agree and P(E) is the proportion of times the annotators are expected to agree due to chance.

The Absolute Agreement is most informative when computed over several sets of labels and where one of the sets represents the “authoritative” set. In this case, the ratio of matches among all the sets including the “authoritative” set to the total

number of labels shows how close the other sets are to the “authoritative” one. The Kappa statistic is useful in measuring how consistent the annotators are compared to each other as opposed to an authority standard.

#### 4.2 Internal consistency

In order to test for internal consistency, we analyzed inter-annotator agreement where the three annotators tagged the same small corpus of clinical dictations.

File ID	Abs agr.	Kappa	N Samples
1137689	93.24%	0.9527	755
1165875	94.59%	0.9622	795
1283904	89.79%	0.9302	392
1284881	90.42%	0.9328	397
1307526	84.43%	0.8943	347
<b>Total</b>			<b>2686</b>
<b>Average</b>	<b>90.49%</b>	<b>0.9344</b>	

Table 1. Annotator agreement results based on 5 clinical notes

The results were compared and the Kappa-statistic was used to calculate the inter-annotator agreement. The results of this experiment are summarized in Table 1. For the absolute agreement, we computed the ratio of how many times all three annotators agreed on a tag for a given token to the total number of tags.

Based on the small pilot sample of 5 clinical notes (2686 words), the Kappa test showed a very high agreement coefficient – 0.93. An acceptable agreement for most NLP classification tasks lies between 0.7 and 0.8 (Carletta 1996, Poessio and Vieira 1988). Absolute agreement numbers are consistent with high Kappa as they show an average of 90% of all tags in the test documents assigned exactly the same way by all three annotators.

#### 4.3 External consistency

The external consistency with the Penn Treebank annotation was computed using a small sample of 939 words from the Penn Treebank-2 WSJ Corpus annotated for POS information.

Annotator	Abs agr
A1	88.17%
A2	87.85%
A3	87.85%
<b>Average</b>	<b>87.95%</b>

Table 2. Absolute agreement results based on 5 clinical notes with an “authority” label set.

The sample had not been seen by the annotators prior to the test. The annotators were asked to make POS judgments on the WSJ corpus sample just as they would on the clinical notes.

The labels were compared to the Penn Treebank annotation individually by annotator and the results of the comparison are presented in Table 2. The results indicate that the three annotators who participated in this project are on average 88% consistent with the annotators of the Penn Treebank corpus.

#### 4.4 Descriptive statistics for the corpus of clinical notes

The annotation process resulted in a corpus of 273 clinical notes annotated with POS tags. The corpus contains 100650 tokens from 8702 types distributed across 7299 sentences. Table 3 displays frequency counts for the top most frequent syntactic categories.

Add table number

Category	Count	% total
NN	18372	18%
IN	8963	9%
JJ	8851	9%
DT	6796	7%
NNP	4794	5%

Table 3. Syntactic category distribution in the corpus of clinical notes.

The distribution of syntactic categories suggests the predominance of nominal categories, which is consistent with the nature of clinical notes reporting on various patient characteristics such as disorders, signs and symptoms.

Another important descriptive characteristic of this corpus is that the average sentence length is 13.79 tokens per sentence, which is relatively short as compared to the Treebank corpus where the sentence length is 24.16 tokens per sentence. This supports our informal observation of the clinical notes data to contain multiple sentence fragments and short diagnostic statements. Shorter sentence length implies greater number of inter-sentential transitions and therefore is likely to present a challenge for a stochastic process.

#### 5 Training a POS tagger on medical data

In order to test some of our assumptions regarding how the differences between general English language and the language of clinical notes may affect POS tagging, we have trained the HMM-based TnT tagger with default parameters at the tri-gram level both on Penn Treebank and the clinical notes data. The clinical notes data was split

at random 10 times in 80/20 fashion where 80% of the sentences were used for training and 20% were used for testing. This technique is a variation on the classic 10-fold validation and appears to be more suitable for smaller amounts of data.

We conducted two experiments. First, we computed the correctness of the Treebank model on each fold of the clinical notes data. We tested the Treebank model on the 10 folds rather than the whole corpus of clinical notes in order to produce correctness results on exactly the same test data as would be used for validation tests of models build from the clinical notes data. Then, we computed the correctness of each of the 10 models trained on each training fold of the clinical notes data using the corresponding testing fold of the same data for testing.

Correctness was computed simply as the percentage of correct tag assignments of the POS tagger (hits) to the total number of tokens in the test set.

$$(3) \quad \text{Correctness} = 100 \cdot \frac{\text{Hits}}{\text{Total}}$$

Table 4 summarizes the results of testing the Treebank model.

Split	Hits	Total	Correctness
1	21560	23872	90.32%
2	22122	24665	89.69%
3	21417	23923	89.52%
4	21970	24461	89.82%
5	22079	24665	89.52%
6	21649	24049	90.02%
7	21598	24040	89.84%
8	21379	23882	89.52%
9	22131	24610	89.93%
10	22358	24923	89.71%
<b>Average</b>	<b>21826.3</b>	<b>24309</b>	<b>89.79%</b>

Table 4 Correctness results for the Treebank model.

Table 5 summarizes the testing results for the models trained on the clinical notes.

The average correctness of the Treebank model tested on clinical notes is ~88%, which is considerably lower than the state-of-the-art performance of the TnT tagger - ~96%. Training the tagger on a relatively small amount of clinical notes data brings the performance much closer to the state-of-the-art - ~95%.

Split	Hits	Total	Correctness
1	22654	23872	94.90%
2	23332	24665	94.60%
3	22645	23923	94.66%
4	23206	24461	94.87%
5	23326	24665	94.57%
6	22732	24049	94.52%
7	22807	24040	94.87%
8	22603	23882	94.64%
9	23316	24610	94.74%
10	23563	24923	94.54%
<b>Average</b>	<b>23018.4</b>	<b>24309</b>	<b>94.69%</b>

Table 5 Correctness results for the clinical notes model.

## 6 Discussion

This paper was intended to share the goals and challenges experienced during annotation of a small sample of clinical data for part-of-speech information. The annotation was performed by experts in the domain of indexing medical content who are minimally trained to label medical texts for part of speech. We have outlined some of the challenging issues in clinical note annotation as well as some of the more outstanding differences between clinical notes and other types of written and spoken discourse widely used in training NLP applications.

The results of this pilot project are very encouraging. It is clear that with appropriate supervision, people who are well familiar with medical content can be reliably trained to carry out some of the tasks traditionally done by trained linguists. We have shown that the three annotators in this study have been able to achieve relatively high levels of inter-rater agreement (Kappa ~ .93) as well as compliance with an authoritative Penn Treebank annotation (Abs agr. ~ 89%).

This study also indicates that an automatic POS tagger trained on data that does not include clinical documents may not perform as well as a tagger trained on data from the same domain. A comparison between the Treebank and the clinical notes data shows that the clinical notes corpus contains 3,239 lexical items that are not found in Treebank. The Treebank corpus contains over 40,000 lexical items that are not found in the corpus of clinical notes. 5,463 lexical items are found in both corpora. In addition to this 37% out-of-vocabulary rate (words in clinical notes but not the Treebank corpus), the picture is further complicated by the differences between the n-gram tag transitions within the two corpora. For example, the likelihood of a DT → NN bigram is 1 in Treebank and 0.75 in the clinical notes corpus. On the other hand, JJ → NN transition in the clinical notes is 1 but in the Treebank corpus it has

a likelihood of 0.73. This is just to illustrate the fact that not only the “unknown” out-of-vocabulary items may be responsible for the decreased accuracy of POS taggers trained on general English domain and tested on the clinical notes domain, but the actual n-gram statistics may be a major contributing factor.

Due to a relatively small size of the available hand-labeled data, it is hard to draw final conclusions. However, it is clear at this point that, at least for the domain of clinical notes, it is necessary to obtain domain specific data in order to train state-of-the-art POS taggers.

## 7 Conclusion

Several questions remain unresolved. First of all, it is unclear how much domain specific data is enough to achieve state-of-the-art performance on POS tagging. Second, given that it is somewhat easier to develop lexicons for POS tagging than to annotate corpora, we need to find out how important the corpus statistics are as opposed to a domain specific lexicon. In other words, can we achieve state-of-the-art performance in a specialized domain by simply adding the vocabulary from the domain to the POS tagger’s lexicon? We intend to address both of these questions with further experimentation.

## 8 Acknowledgements

Our thanks go to Barbara Abbot, Pauline Funk and Debora Albrecht for their persistent efforts in the difficult task of corpus annotation.

## References

- Baldrige, J., Morton, T., and Bierner, G URL: <http://maxent.sourceforge.net>
- Brandts, T (2000) “TnT – A Statistical Part-of-Speech Tagger.” In Proc. NAACL/ANLP-2000.
- Carletta, J. (1996). Assessing agreement on classification tasks: The Kappa statistic. *Computational Linguistics*, 22(2) pp. 249-254.
- Cutting, D., Kupiec, J., Pedersen, J, and Sibun, P. A (1992). Practical POS Tagger. In Proc. ANLP’92.
- Jurafski D. and Martin J. (2000). *Speech and Language Processing*. Prentice Hall, NJ.
- Manning, C. and Shute H. (1999). *Foundations of Statistical Natural Language Processing*. MIT Press, Cambridge, MA.
- Marcus, M., B. Santorini, and M. A. Marcinkiewicz (1993). Building a large

annotated corpus of English: the Penn Treebank. *Computational Linguistics* 19, 297-352.

Mikheev, A. (1997). Automatic Rule Induction for Unknown-Word Guessing. [Computational Linguistics](#) 23(3): 405-423

Poessio, M. and Vieira, R. (1988). “A corpus based investigation of definite description use” *Computational Linguistics*, pp 186-215.

Ratnaparkhi A. (1996). A maximum entropy part of speech tagger. In *Proceedings of the conference on empirical methods in natural language processing*, May 1996, University of Pennsylvania

Ruch P, Baud RH, Rassinoux AM, Bouillon P, Robert G. Medical document anonymization with a semantic lexicon. *Proc AMIA Symp.* 2000; 729-33.

Santorini B. (1991). *Part-of-Speech Tagging Guidelines for the Penn Treebank Project*. Technical Report. Department of Computer and Information Science, University of Pennsylvania.

UMLS. (2001). *UMLS Knowledge Sources (12th ed.)*. Bethesda (MD): National Library of Medicine.