

RESEARCH ARTICLE

Open Access



# Recognizing chemicals in patents: a comparative analysis

Maryam Habibi<sup>1\*</sup> , David Luis Wiegandt<sup>1</sup>, Florian Schmedding<sup>2</sup> and Ulf Leser<sup>1</sup>

## Abstract

Recently, methods for Chemical Named Entity Recognition (NER) have gained substantial interest, driven by the need for automatically analyzing today's ever growing collections of biomedical text. Chemical NER for patents is particularly essential due to the high economic importance of pharmaceutical findings. However, NER on patents has essentially been neglected by the research community for long, mostly because of the lack of enough annotated corpora. A recent international competition specifically targeted this task, but evaluated tools only on gold standard patent abstracts instead of full patents; furthermore, results from such competitions are often difficult to extrapolate to real-life settings due to the relatively high homogeneity of training and test data. Here, we evaluate the two state-of-the-art chemical NER tools, tmChem and ChemSpot, on four different annotated patent corpora, two of which consist of full texts. We study the overall performance of the tools, compare their results at the instance level, report on high-recall and high-precision ensembles, and perform cross-corpus and intra-corpus evaluations. Our findings indicate that full patents are considerably harder to analyze than patent abstracts and clearly confirm the common wisdom that using the same text genre (patent vs. scientific) and text type (abstract vs. full text) for training and testing is a pre-requisite for achieving high quality text mining results.

**Keywords:** Chemical named entity recognition, Patent mining, Ensemble approach, Simple chemical elements, Performance measurements

## Background

Patents are an economically important type of text directly related to the commercial exploitation of research results. They are particularly essential for the pharmaceutical industry, where novel findings, such as new therapeutics or medicinal procedures, result from extremely cost-intensive, long-running research projects, but often are relatively easy to copy or reproduce [1]. Accordingly, a large number of commercial services exists regarding the formulation and retrieval of patents [2], and large companies devote entire departments to the creation, the licensing, and the defense of their patent portfolio. Such services must be supported by proper computational tools, as the number of patents is increasing rapidly. For instance, the European Patent Office granted 614,850

patents since 2006 [3]; the size of the United States Patent and Trademark Office corpus currently is 6,718,054 patents with a yearly increase of roughly 300,000 over the last 5 years [4]. However, current tools for patent management mostly support keyword search [5–9], whereas only few projects exist that target the extraction of specific facts from patents [10, 11].

In this work, we study the identification and extraction of chemical names<sup>1</sup> from patents. By extraction, we mean the identification of left and right borders of mentions in patents, a task usually referred to as Named Entity Recognition (NER). Extracting chemicals from scientific articles has been a topic of ample research over the last 15 years, leading to the creation of high quality tools like OSCAR [12] or ChemSpot [13] which focus on the particularities of chemical names when compared to other entities, such as genes or species [13–16]. However, the extraction of chemicals from patents has been neglected

\*Correspondence: [habibima@informatik.hu-berlin.de](mailto:habibima@informatik.hu-berlin.de)

<sup>1</sup> Knowledge Management in Bioinformatics, Humboldt-Universität zu Berlin, 12489 Berlin, Germany

Full list of author information is available at the end of the article

<sup>1</sup> Abbreviated as “chemical” from now on.

by the research community for long, mostly due to the difficulties in obtaining computer-readable patents at large—compared to the simple procedures necessary to download scientific articles from sources like PubMed-Central<sup>2</sup>—and the lack of properly annotated patents, i.e. gold standard corpora. It is tempting to apply tools and models developed for scientific articles on patents, but patent texts are quite different from scientific articles. They are typically much longer, yet have a lower word density [17]. Their writing is more difficult to understand as the protection of broad claims and a mild obfuscation of procedures are established means to increase patent value and decrease the likelihood of being reproduced [5]. Therefore, it is rather unclear whether tools developed for scientific articles perform equally well on patent data.

Since 2012, two gold standard full-text patent corpora have been published: the chapati corpus [18] and the corpus from the BioSemantics research group [19]. The field was further boosted by a recent international competition, the CEMP task at BioCreative V. For this task, two large corpora for training and development were prepared and used by 21 teams to develop patent-specific solutions [20–22], achieving the F-measure values of up to 89% (87% precision and 91% recall) using an ensemble approach. However, both corpora consist only of patent titles and abstracts, while commercially interesting applications critically depend on analyzing full texts, as a significant number of entities is not even mentioned in an abstract [17]. Furthermore, international challenges are important to make different approaches comparable and also provide a strong incentive for groups to enter a field [11], yet their performance results are difficult to extrapolate due to the relatively high homogeneity of training and test data within the competition. In contrast, real applications typically have to perform information extraction on diverse text collections without having accordingly diverse training data. To mimic such situations, cross-corpus evaluations can be used, where the performance of a tool trained on one corpus is measured on another corpus following different annotation guidelines [23].

In this paper, we take this idea one step further and perform a cross-text-genre evaluation by assessing the performance of chemical NER tools trained on scientific articles—a problem much better researched—on patent corpora. We choose tmChem [24] and ChemSpot [13], two state-of-the-art tools for chemical NER from scientific articles, and evaluate their performance (without retraining) on all four freely available gold standard

patent corpora with annotations of chemical mentions. We put emphasis on the differences between evaluations on abstracts versus full texts, showing that the latter is a considerably harder task for current tools. We also compare results on the instance level, showing that, despite having similar performance numbers, tmChem and ChemSpot actually return quite different results. This makes the creation of ensembles attractive, two of which we evaluate on all four corpora. We also contrast our cross-text-genre results to those obtained after retraining a chemical NER tool on patent corpora, showing that taking away the text-genre difference significantly boosts results, i.e., that the different characteristics of patent versus scientific texts strongly impact chemical NER performance. Overall, our results emphasize the common wisdom that using the same text genre (patents vs. scientific articles) and text type (abstracts vs. full texts) for training and application is a pre-requisite for achieving high-quality text mining results.

## Methods

In this section, we describe the characteristics of the four freely available gold standard patent corpora and present the two chemical NER systems utilized in our study. We also explain two ensemble approaches, our evaluation metrics, and the text preprocessing techniques employed.

### Patent corpora

We use all four currently existing gold standard patent corpora with annotations for chemicals. Two of them contain only the title and the abstract of patents, while the other two use complete patent documents. The two abstract corpora, denoted as CEMP\_T and CEMP\_D, contain 7000 patents each. They were originally developed for training and development purposes within the CEMP (chemical entity mention in patents) task [11] of the BioCreative<sup>3</sup> V challenge. The chapati corpus [18] is the result of a collaboration between the European Patent Office<sup>4</sup> and the CHEBI<sup>5</sup> team. Chemical entities were manually annotated for 40 complete patent documents and normalized to CHEBI identifiers. The fourth corpus, noted here BioS, was prepared by the BioSemantics<sup>6</sup> research group and covers 200 full patent documents [19]. For this corpus, patents were automatically pre-annotated and then manually curated by at least one annotator group consisting of two to ten annotators. Table 1 shows statistics on these corpora like corpus size

<sup>2</sup> <http://www.ncbi.nlm.nih.gov/pmc>.

<sup>3</sup> See <http://www.biocreative.org>.

<sup>4</sup> See <https://www.epo.org/>.

<sup>5</sup> See <https://www.ebi.ac.uk/chebi/>.

<sup>6</sup> See <http://www.biosemantics.org/>.

**Table 1 The details of the gold standard patent corpora containing the annotations for chemicals**

Corpus	Number of patents	Annotated entities	Number of annotations
CEMP training set (CEMP_T) [11, 25] ≈660 thousand token	7000 patents (title and abstract)	ABBREVIATION, FAMILY, FORMULA, TRMIAL, MULTIPLE, SYSTEMATIC, IDENTIFIERS	33543 (without normalization)
CEMP development set (CEMP_D) [11, 25] ≈650 thousand token	7000 patents (title and abstract)	ABBREVIATION, FAMILY, FORMULA, TRMIAL, MULTIPLE, SYSTEMATIC, IDENTIFIERS	32142 (without normalization)
CHEBI patent corpus (chapat) [18] ≈265 thousand token	40 full patents (title, abstract, claims, description)	CLASS, CHEMICAL, ONT, FORMULA, LIGAND, CM	18746 (normalized to CHEBI identifiers)
BioSemantic patent corpus (BioS) [19] 1,500 pages and ≈4.2 million token	200 full patents (title, abstract, claims, description)	IUPAC, SMILES, InChI, ABBREVIATION, MOA, DISEASE, FORMULA, REGISTRY NUMBER, GENERIC, TRADE-MARK, CAS NUMBER, TARGET	400125 (without normalization)

(the number of tokens separated by space), number of documents, classes of annotated chemicals and number of annotated entities.

Note that these corpora were annotated using different annotation guidelines. The corpora from the CEMP task and BioSemantics group were annotated using two specific annotation guidelines, while the chapati curators considered all entities that could be automatically mapped to a CHEBI identifier. The annotation guidelines vary in several aspects. For instance, the IUPAC name “water” should not be annotated as a chemical in the CEMP corpora but it should be in the BioS corpus [19]. Additionally, simple chemical elements are annotated in the CEMP corpora but not in the BioS corpus.

Additionally, we compare our results on the patent corpora with those achieved on two corpora consisting of scientific articles: CHEMDNER and CRAFT. The CHEMDNER corpus [25] was developed for the CHEMDNER task at the BioCreative IV challenge. The corpus consists of scientific abstracts that were annotated using the same annotation guideline used for the CEMP task. In this work, we only used the test set, containing 3000 abstracts, since the training set and the development set were used for training tmChem. The CRAFT corpus [26] consists of 97 scientific full texts, yet only 67 of these have been publicly released to date. The chemical annotations of the CRAFT corpus are limited to terms from the CHEBI database.

For illustration (see “Text-genre statistics” section), we briefly analyzed differences between patents and scientific articles in terms of the average number of words per sentence (sentence length), the average number of words per document (document length), the average number of unique/non-unique TLAs<sup>7</sup> in a document, and the average number of figures and tables per document. To this end, we used a collection of randomly selected full patent documents from European Patent Office<sup>8</sup>, and a set of randomly selected full journal articles from PubMedCentral. All texts were from year 2015; patents were selected after classification to ensure biomedical topics. We calculated the number of tables and figures by counting the number of their tags in xml format.

### Chemical NER systems

Over the last years, many tools have been presented for chemical NER, including tmChem [24], ChER [10], ChemSpot [13], becas [27], OSCAR [12] or ChemXSeer-tagger [28]. We chose two of them based on their good overall performance in a number of evaluations: (1) tmChem developed by Leaman et al. [24], as the best

system at the CHEMDNER challenge in BioCreative IV<sup>9</sup> [16], (2) ChemSpot<sup>10</sup> which was introduced in the year 2012 and outperformed all other tools for many years. Table 2 gives an overview of the two tools.

ChemSpot employs a hybrid approach, in which the results of a CRF model trained to recognize IUPAC entities are combined with dictionary matching to find other chemical names. TmChem uses ensembles of two CRF models, called Model1 and Model2, with different setups and configurations. As the implementation of the ensembles were not freely available, we performed our experiments using the individual models. We limited our study to the results from Model1, noted as tmChem, as its performance was always very close to or higher than that of the Model2 [22, 24]. Both tmChem and ChemSpot build on BANNER as CRF implementation [30], but use different feature sets, tokenization methods, and training sets. Both tools were trained on scientific abstracts, but the training corpora comprise different articles and were annotated based on different annotation guidelines [31]. Note that the annotation guideline used to annotate the training set of tmChem is very similar to the one used to annotate the two CEMP patent corpora. Both tools normalize extracted entities. TmChem maps the entities to identifiers from CHEBI and MESH, whereas ChemSpot maps them to further databases, like InChI and DrugBank.

### Ensemble NER systems

A comparison of the concrete set of entities returned by tmChem and ChemSpot (see “Comparison at instance level” section), respectively, showed significant divergence. Since ensembles of NER tools often outperform individual tools [28], we also measure the performance of two ensembles produced by merging the results of tmChem and ChemSpot. One ensemble system, called Ensemble-I, accepts a mention as a chemical name when both tmChem and ChemSpot recognize it as such. The second ensemble, noted Ensemble-U, considers a span as a chemical name when it is recognized by at least one of the two systems.

### Evaluation metrics

Performance values were computed in terms of precision, recall, F-measure, and true positive (TP), false positive (FP), and false negative (FN) counts; in all cases, only exact span matches were considered. Precision measures the ratio of correctly predicted chemical entities to all predicted entities; recall is defined as the ratio

<sup>7</sup> TLA is defined as any three-letter word with letters all in uppercase form.

<sup>8</sup> See <https://www.epo.org/>.

<sup>9</sup> Although the performance of ChER was very close to that of tmChem in the BioCreative IV challenge, we were not able to include it in our study as it is not freely available.

<sup>10</sup> Note that ChemSpot was not evaluated in the BioCreative IV challenge.

**Table 2 Details on the chemical NER tools in terms of training sets, databases to which the entities are normalized, classes of chemicals addressed, and tokenization methods**

NER tool	Training set	Databases	Classes	Tokenization method
tmChem [24]	CHEMDNER corpus at BioCreative IV (training and development sets)	CHEBI MESH	SYSTEMATIC FORMULA FAMILY TRIVIAL IDENTIFIER MULTIPLE ABBREVIATION	Tokenization at every non-letter and non-digit characters, number-letter changes and lower case letter followed by an uppercase letter
ChemSpot [13]	A subset of SCAI Corpus [29] containing only IUPAC	ChemIDplus CHEBI CAS NUMBER PubChem InChI DrugBank KEGG Human Metabolome MESH	SYSTEMATIC FORMULA FAMILY TRIVIAL IDENTIFIER MULTIPLE ABBREVIATION	Tokenization at every non-letter and non-digit characters and number-letter changes

of correctly predicted entities to all annotated entities within a corpus. F-measure is the harmonic mean of the precision and the recall values.

We measured performance values using the conllval<sup>11</sup> script run over the prediction and reference annotation files in IOB format. We also compared the different methods with respect to the execution time on patents and scientific articles.

### Text preprocessing

The different gold standard corpora were available in different text formats which we homogenized before running the NER tools. In a first step, each document was converted to plain text format and stored in a single file. Then we transformed these files to the input format defined by each NER tool. For evaluation purposes, we tokenized each prediction and gold standard annotation files as suggested by Klinger et al. [32]<sup>12</sup>, and represented each token in IOB format. We used the Stanford parser<sup>13</sup> [33] to split the text into sentences and to parse a sentence.

### Results

We first provide evaluation scores for the models trained on the abstract of scientific articles when applied to patents. Then we present an analysis of entities frequently recognized incorrectly as chemicals or non-chemicals by the two systems. Afterwards, we describe results of the two ensemble systems. Finally we present the results of cross-corpus and intra-corpus evaluations to study the impact of the use of different text genres and text types as training and test sets in patent mining.

### Cross-genre evaluation

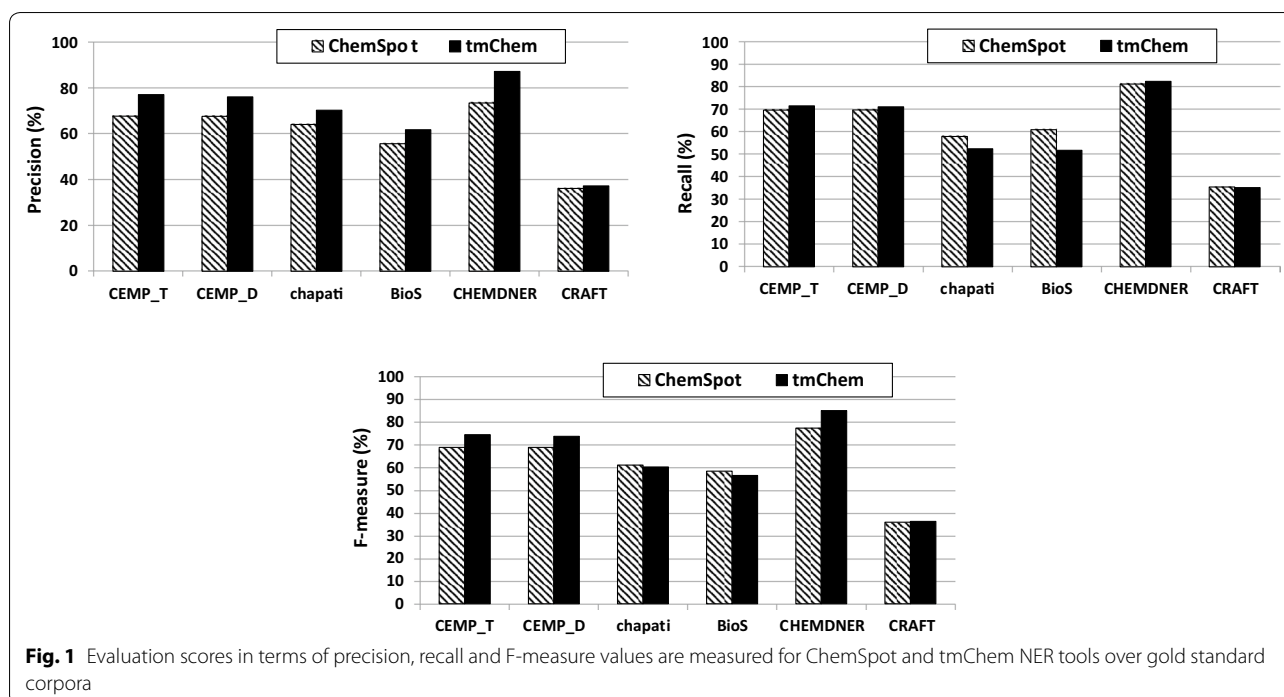
We compared the performance of tmChem and ChemSpot on four patent corpora and two corpora consisting of scientific articles in terms of precision, recall, and F-measure. The results are shown in Fig. 1. The following ranking of corpora regarding their assessability by the tools can be inferred: CRAFT<BioS<chapati<CEMPs<CHEMDNER.

We first observe that the performances of both tools are much lower on the CRAFT corpus than on the other corpora. The reason for this discrepancy seems to be that the version of the CHEBI database used for annotating the CRAFT corpus is quite different from that used for training both NER tools [25]. On the CHEMDNER corpus, both tools, despite the use of different training data, have higher performance values than on the other corpora, indicating that they perform best on scientific abstracts—the type of texts they also were trained on.

<sup>11</sup> See <http://www.cnts.ua.ac.be/conll2000/chunking>.

<sup>12</sup> In this format, every non-letter and non-digit character, and all number-letter changes are split.

<sup>13</sup> Version (stanford-parser 3.5) is available in <http://nlp.stanford.edu/software/lex-parser.shtml>.



This observation indicates that models trained on scientific abstracts are not quite as capable of recognizing chemical entities from patents. Drawing a conclusion regarding the difference between scientific abstracts and full texts, however, is difficult due to the quite different scope of chemical annotations in the CRAFT corpus [25].

The results, obtained from both tools, also report a ~10% higher performance on patent abstracts compared to full patents, indicating that there is more similarity between patent abstracts and scientific abstracts than between scientific abstracts and full patents.

TmChem has higher F-measure values than ChemSpot (at least 5%) on the CHEMDNER and the CEMP's corpora, both of which follow an annotation guideline similar to that of the tmChem training set. In contrast, ChemSpot was trained on a corpus with a different annotation guideline. However, the F-measure values of both tmChem and ChemSpot were very close on full patent corpora, annotated using guidelines which differ from those of the systems' training sets. This means that the improvement obtained by tmChem on patent abstracts is likely due to the similarity of the annotation guidelines and not to the superiority of the method. Thus, we cannot conclude which tool is better suited for chemical NER on patents.

#### Comparison at instance level

We studied the top 10 entities most frequently incorrectly recognized (FPs and FNs) by tmChem and ChemSpot,

respectively. These entities along with their FP and FN counts are shown in Tables 3 and 4. Mistakes made frequently by both tools are in italic font.

Counting the number of common errors, we find around 50% overlap between the top-10 entities with highest FP counts and nearly 70% overlap between entities with highest FN values between the tools. However, the individual error frequencies are very different. For example, the number of times that the entity "alkyl" is incorrectly recognized as a chemical entity by ChemSpot is around 6 times higher than that of tmChem (in patent abstracts). There are several similar cases in the corpora containing full patents. Similarly, there are also a number of common entities with diverging FN values. By excluding common entities with highly different frequencies, the overlaps between these two tools for patent corpora are reduced to around 40 (for FP counts) and 50% (for FN counts), which indicates the two tools perform rather differently.

#### Error distribution

We observed many common entities on the lists of errors arisen from tmChem and ChemSpot in Tables 3 and 4, but having highly different frequencies. This observation motivated us to compare patent full texts and abstracts in terms of the distributions of the FP and FN counts measured for unique entities. We depicted the distributions of these values which were sorted from high to low, and covered 25% of error cases, in Figs. 2 and 3.

**Table 3** The top 10 entities with highest FP for each chemical NER tool on the four different corpora

CEMP_T		CEMP_D	
ChemSpot	tmChem	ChemSpot	tmChem
Water 951	<i>Sodium</i> 128	Water 842	<i>Sodium</i> 117
<i>Alkyl</i> 260	Sugar 66	<i>Alkyl</i> 194	<i>Nucleotide</i> 74
<i>Sodium</i> 186	CH <sub>2</sub> 56	<i>Sodium</i> 194	<i>Ester</i> 49
DEG 155	<i>Sulfate</i> 43	Peptide 153	<i>Calcium</i> 49
Peptide 107	NO 42	Chitosan 130	O 46
Chitosan 91	Solvate 40	DEG 108	NO 45
Starch 81	<i>Alkyl</i> 39	Parkinson 80	N 44
<i>Calcium</i> 74	Hydrogen 38	<i>Calcium</i> 76	<i>Alkyl</i> 37
<i>Sulfate</i> 66	<i>Calcium</i> 35	<i>Nucleotide</i> 72	Sulfate 37
Parkinson 60	Beta-cyclodextrin 34	<i>Ester</i> 67	Beta-cyclodextrin 36
Chapati		BioS	
ChemSpot	tmChem	ChemSpot	tmChem
Factor H 121	CO 127	<i>Hydrogen</i> 6246	<i>Hydrogen</i> 6179
<i>Serine</i> 108	<i>Serine</i> 108	<i>1H</i> 5034	<i>Carbon</i> 5518
<i>Alkyl</i> 81	N 88	<i>Carbon</i> 5004	H 3091
<i>Amino acid</i> 66	NH-SO <sub>2</sub> 64	<i>3H</i> 3928	<i>1H</i> 2785
SO <sub>2</sub> -NR<21>R<22 62	<i>NH-CO-R&lt;21 63</i>	<i>Alkyl</i> 3812	<i>3H</i> 2643
CO-R<23 60	<i>Amino acid</i> 61	<i>2H</i> 2946	<i>Nitrogen</i> 2619
<i>NH-CO-R&lt;21 55</i>	Carbon 57	<i>Nitrogen</i> 2878	<i>Silica</i> 1466
Ci-I0 54	<i>Nitroxide</i> 52	<i>Silica</i> 2011	CDCl <sub>3</sub> 1320
CO-NR<21>R<22 53	C 51	DMSO-d <sub>6</sub> 1652	<i>2H</i> 1259
<i>Nitroxide</i> 52	H 46	<i>Oxygen</i> 1203	<i>Oxygen</i> 1110

Common mistakes are shown in italic

The distributions over patent abstracts showed that 25% of FP counts are produced by around 90 unique entities for tmChem, and only 30 unique ones for ChemSpot. Moreover, the shapes of the distributions were quite different. Similarly, by comparing the distributions of FP counts over full patents, we observed that the number of unique entities leading to 25% of FPs for tmChem is around 20 and for ChemSpot, it is nearly 30. The shapes of tmChem and ChemSpot distributions were very similar over full patents, but they were different from the ones obtained for ChemSpot over patent abstracts. These results confirm that the distributions of FP counts are different over full patents and patent abstracts.

On the contrary, the shapes of the distributions drawn for FN counts were very similar for both systems over all corpora, but the number of unique entities leading to 25% of FNs over full patents is around 25 while it is 75 for abstracts.

#### Impact of simple chemical elements

Interestingly, there are quite a number of simple chemical elements in these lists (e.g. H, N, S). The appearance of simple chemical elements in both lists of errors indicates

that these entities are generally ambiguous and difficult to be correctly predicted, although they are rather irrelevant for many applications in areas such as cell biology or omics studies. This observation encouraged us to study the impact of simple chemical elements on the performance values of different types of patent texts.

First, we computed the FP and FN counts for simple chemical elements normalized with the FP and FN counts of all entities extracted by each NER tool from each corpus as shown in Fig. 4. The results obtained by the two tools demonstrate that the normalized FP values of simple elements are higher than those of FN values for full patents, while they are approximately analogous for patent abstracts. It implies that simple chemical elements are frequently recognized incorrectly as chemicals on full patents.

Following this observation, we recalculated the performance values, i.e., precision, recall, and F-measure, by excluding the annotations of all simple chemical elements<sup>14</sup> from the gold standard corpora, and also filtering the simple elements predicted by ChemSpot and

<sup>14</sup> See <http://www.chemicalelements.com>.

**Table 4** The top 10 entities with highest FN for each chemical NER tool on the four different corpora

CEMP_T		CEMP_D	
ChemSpot	tmChem	ChemSpot	tmChem
<i>H</i> 227	<i>Alkyl</i> 226	<i>H</i> 233	<i>Alkyl</i> 246
<i>Aryl</i> 170	<i>Aryl</i> 179	<i>Aryl</i> 174	<i>Aryl</i> 183
<i>C1-6 alkyl</i> 115	<i>H</i> 173	<i>Heterocyclic</i> 133	<i>H</i> 179
<i>Heteroaryl</i> 82	<i>C1-6 alkyl</i> 121	<i>Heteroaryl</i> 87	<i>Heterocyclic</i> 135
<i>Alkyl</i> 74	<i>S</i> 86	<i>N</i> 76	<i>S</i> 86
<i>N</i> 71	<i>Cyano</i> 85	<i>C1-6 alkyl</i> 69	<i>C1-6 alkyl</i> 76
<i>Alkoxy</i> 67	<i>Heterocyclic</i> 62	<i>Alkoxy</i> 63	<i>Cyano</i> 71
<i>Cyano</i> 62	<i>Halo</i> 55	<i>Alkyl</i> 59	<i>N</i> 56
<i>Heterocyclic</i> 61	<i>Oligonucleotides</i> 50	<i>Aromatic</i> 51	<i>Halo</i> 52
<i>Halo</i> 51	<i>Opioid</i> 50	<i>Cyano</i> 45	<i>Aromatic</i> 51
Chapati		BioS	
ChemSpot	tmChem	ChemSpot	tmChem
<i>Drug</i> 234	<i>Water</i> 264	<i>Alkyl</i> 5295	<i>Alkyl</i> 8145
<i>Ci-10 alkyl</i> 160	<i>Drug</i> 234	<i>Aryl</i> 4698	<i>Water</i> 7142
<i>NR</i> 139	<i>Peptide</i> 205	<i>DMSO</i> 3184	<i>Aryl</i> 5426
<i>Insulin</i> 107	<i>Ci-10 alkyl</i> 160	<i>Heteroaryl</i> 2435	<i>Ph</i> 1995
<i>Aptamer</i> 92	<i>NR</i> 139	<i>Alkoxy</i> 1993	<i>H</i> 1921
<i>Polypeptide</i> 88	<i>Insulin</i> 107	<i>H</i> 1869	<i>Brine</i> 1822
<i>SO2R</i> 65	<i>CN</i> 94	<i>Brine</i> 1777	<i>DMSO</i> 1490
<i>NH-CO-R</i> 63	<i>Aptamer</i> 92	<i>Inhibitors</i> 1468	<i>Ethyl acetate</i> 1473
<i>SO2-NR</i> 63	<i>Polypeptide</i> 89	<i>Substituted</i> 1447	<i>Inhibitors</i> 1472
<i>NH-SO2-R</i> 62	<i>SO2R</i> 65	<i>Lower alkyl</i> 1422	<i>Substituted</i> 1447

Common mistakes are shown in italic

tmChem. The results in Fig. 5 show that the performance values of the corpora with the same text type converge into highly close values after removing simple chemical elements. However the impact the removal of simple elements has on the precision values is insignificant, except for the BioS corpus, on which the precision improves by up to 5%. Recall and F-measure values are not strongly affected by ignoring simple chemical elements.

#### Ensemble performance

As we found substantial differences in the concrete results computed by tmChem and ChemSpot, we decided to measure the performance of the two ensemble systems obtained by (a) intersecting (Ensemble-I) and (b) unifying (Ensemble-U) the results of the two systems (see “Ensemble NER systems” section).

We provide precision, recall, and F-measure values calculated for tmChem, ChemSpot, and for the intersection and the union of their outputs in Fig. 6. As expected, on all corpora, the highest precision is obtained by intersecting the results of the two tools, while the highest recall is provided by unifying the results of the two systems. The results also show that the Ensemble-U provides the

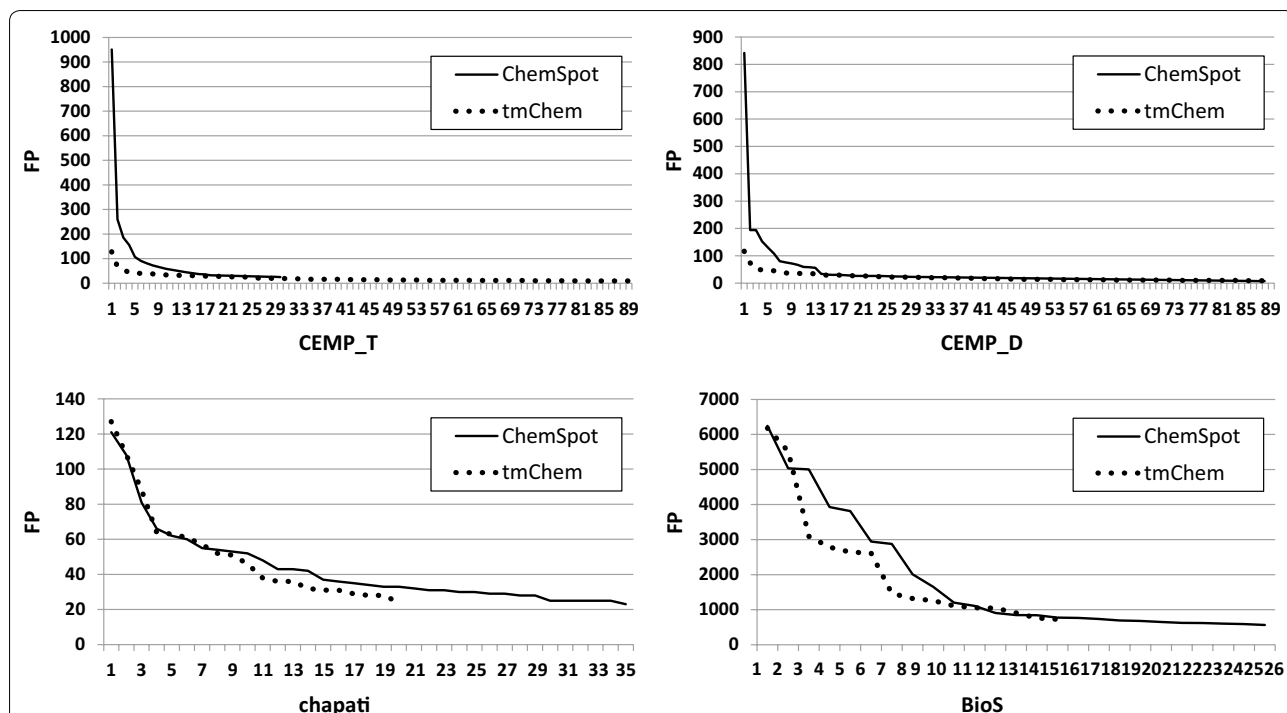
highest F-measure value on full patents, while tmChem has the highest F-measure scores on patent abstracts. This can probably be attributed to the use of similar annotation guidelines for training and test sets.

#### Cross-text-genre to cross-corpus evaluation

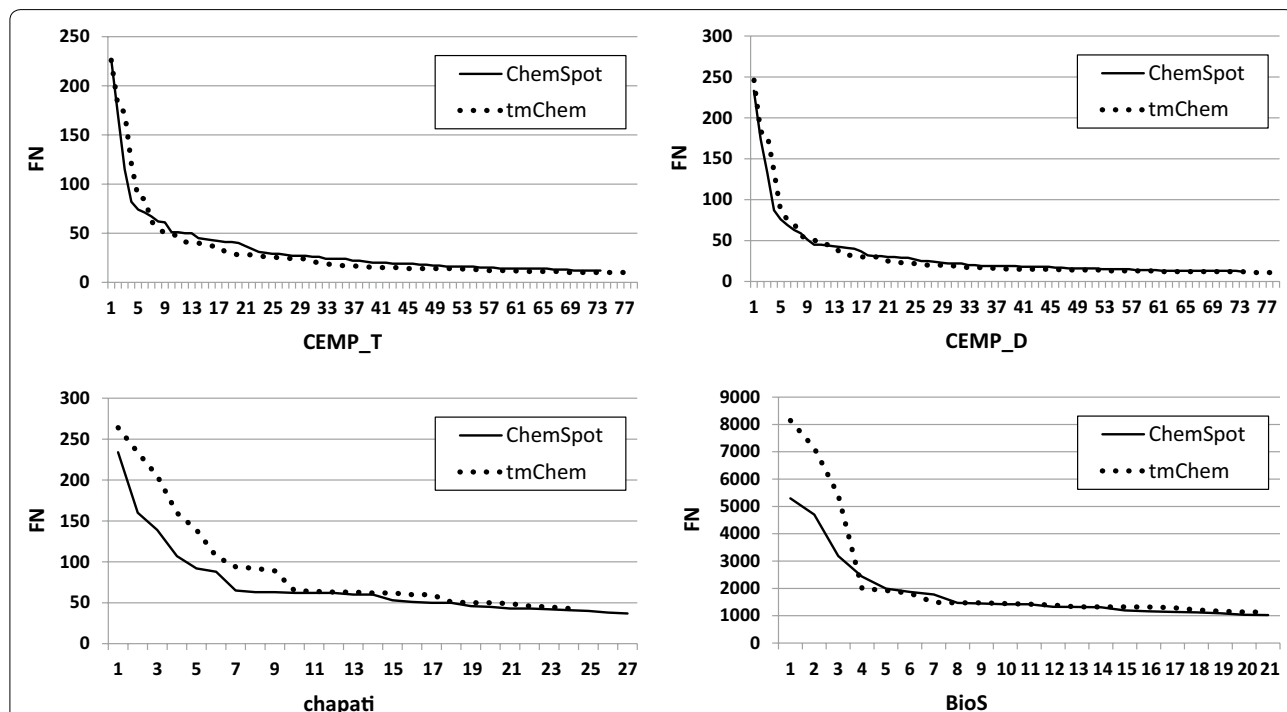
We measured the performances of different models obtained by retraining tmChem using patent corpora. We retrained only tmChem because of the well documented process in its public API. We first performed a cross-corpus evaluation, where tmChem was trained on one corpus and evaluated on other corpora. In addition to cross-corpus evaluation, we performed intra-corpus evaluation by assessing the performances of the models using fourfold cross validation. The performance values of the models trained on patent corpora and the tmChem default model are depicted in Fig. 7.

The F-measure values of all models evaluated on chapati are nearly identical. Additionally, the F-measure scores of the model trained using the chapati corpus are very close on other corpora, perhaps for its small number of instances. Thus, we limit our analysis to the remaining corpora.

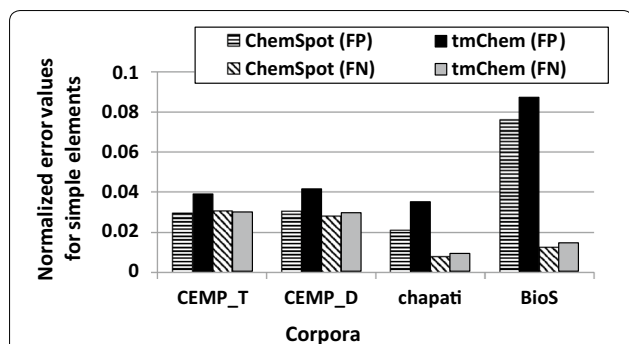




**Fig. 2** Distributions of FP counts from high to low, for unique entities covering 25% of cases, obtained by tmChem and ChemSpot over all corpora. The x-axis represents the number of unique entities. The distributions are notably different for full patents compared to patent abstracts



**Fig. 3** Distributions of FN counts from high to low, for unique entities covering 25% of cases, obtained by tmChem and ChemSpot over all corpora. The x-axis represents the number of unique entities. The distributions are very similar for full patents and patent abstracts



**Fig. 4** The FP and FN counts of simple chemical elements normalized by the FP and FN counts obtained for the entire entities by tmChem and ChemSpot over all corpora

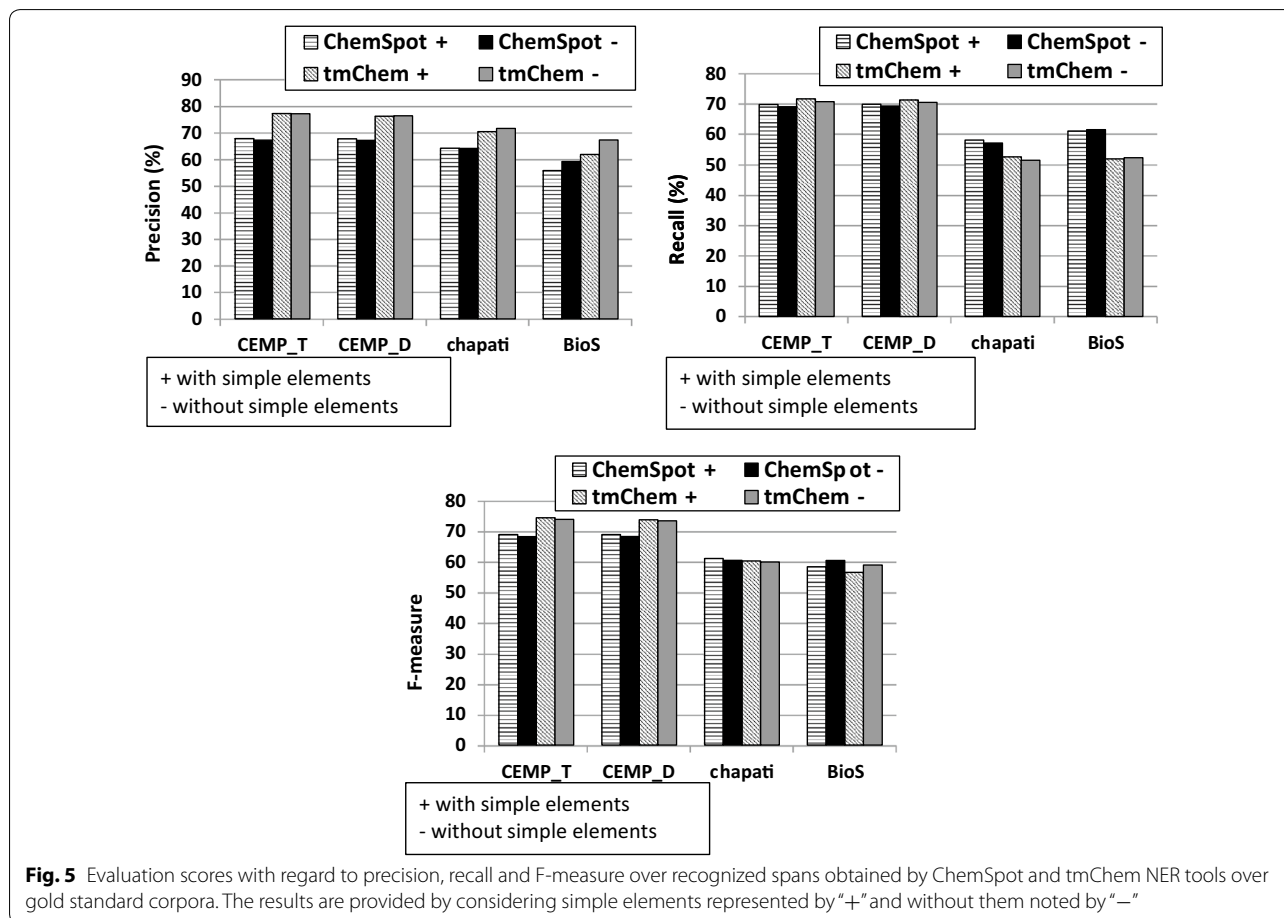
The results show that the performance values on the BioS corpus, containing full patents, are demoted when using models trained on CEMP\_T or CEMP\_D, both of which contain patent abstracts and were annotated using a different annotation guideline. Likewise,

the performance values on CEMP\_T and CEMP\_D are higher when using models trained on CEMPs corpora, and not on BioS. However, we cannot conclude that the models trained on abstracts are not suitable to identify entities from full texts and vice versa, as the corpora are annotated using different annotation guidelines.

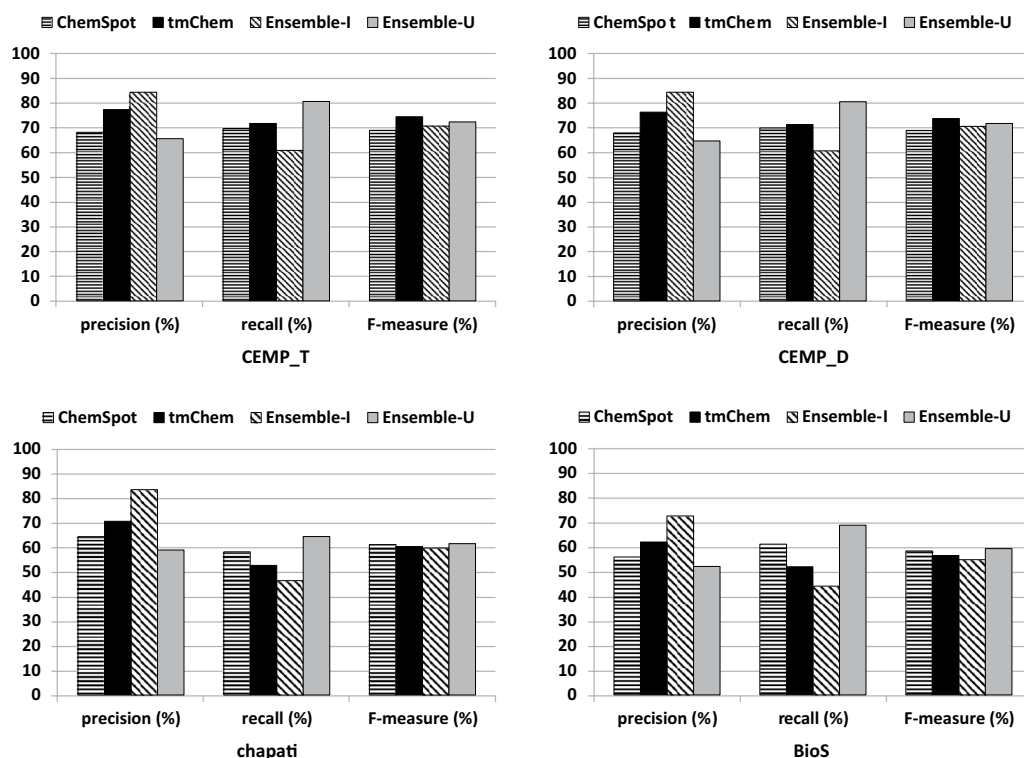
The precision, recall and F-measure values achieved in intra-corpora evaluation, shown in Fig. 7, indicate that the precision values on full patents are at least 10% lower than on patent abstracts. Although the recall value on BioS is close to the ones obtained on patent abstracts, its F-measure value is still lower than those of patent abstracts. These results imply that identifying chemical names from full patents is more difficult compared to that of patent abstracts.

## Discussion

We have empirically shown that significant differences exist between the results of chemical NER on patents and scientific articles and even between different types of



**Fig. 5** Evaluation scores with regard to precision, recall and F-measure over recognized spans obtained by ChemSpot and tmChem NER tools over gold standard corpora. The results are provided by considering simple elements represented by "+" and without them noted by "-"



**Fig. 6** Evaluation scores with regard to precision, recall, and F-measure values over recognized spans obtained by ChemSpot, tmChem, the area of their intersection and union over gold standard corpora

patent texts. Our study has demonstrated that identifying chemical entities from patent full texts is more complex than from patent abstracts or scientific abstracts. In the following, we assess the complexity of this task on patents, especially on patent full texts.

#### Difference between scientific articles and patents

The performance values attained in cross-text-genre evaluations show that the F-measure values of the models trained on the abstracts of scientific articles decrease by around 10% when tested on patent abstracts and by nearly 18% when applied to patent full texts.

The lower F-measure scores obtained by tmChem on patent abstracts compared to that of scientific abstracts, while both have annotation guidelines very similar to those of the tmChem training set, show that there are several chemical entities in patent abstracts that cannot be recognized by the models trained using scientific articles. This finding emphasizes the difficulty of the chemical NER task on patents.

The F-measure scores of ChemSpot trained on scientific abstracts annotated using a guideline different from the ones used for the patent corpora, indicate that these models are not adequate to recognize entities from

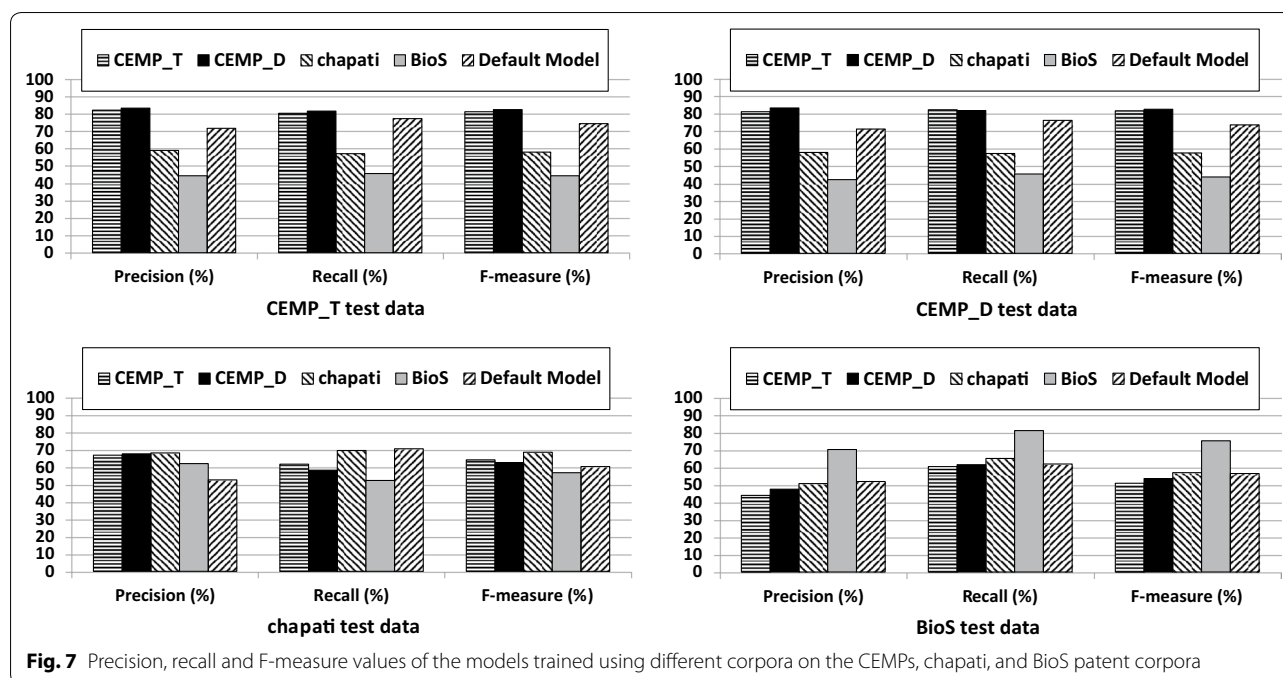
patents, and accentuate the need for more annotated patent corpora for chemical NER.

#### Execution time analysis

We compared the execution time of tmChem and ChemSpot over 10 complete patents and 10 journal articles randomly selected from European Patent Office and PubMedCentral, respectively (see “Patent corpora” section). The tools were run single-threaded on a machine with 1TB RAM and Linux operating system. The execution time values, in seconds, are reported in Table 5.

ChemSpot required 149 s to complete the task on scientific articles, around four times faster than the time required for patents. Similarly, tmChem needed approximately 50% more time to finish the task on patents compared with scientific articles. The main reason is the difference in their lengths (see next section).

Then we estimated the execution time of the two systems that one would have to expect on 10 million patents and 10 million full scientific articles, assuming 8 parallel threads by extrapolating the above values. The results show that tmChem would take around 3 months over patents and 2 months over journals while ChemSpot would take approximately 2 years for patents and



7 months for journals. We conclude that large parallel systems are required for patent chemical NER.

#### Text-genre statistics

To better analyze the evident complexity of the NER task on patents apparent from the lower performance and longer execution time of NER tools compared to those of scientific articles, we quantified the differences in their text structures by the average sentence length, document length, and the average number of unique/non-unique TLAs, figures and tables per document (see “Patent corpora” section). We measured the values using 17,000 patent documents and 17,000 journal articles which have been chosen randomly from the European Patent Office and PubMedCentral respectively. The values are provided in Table 6.

The average number of words per sentence is almost the same for both patents and journals. However, the average number of words of a patent document is approximately five times higher than that of a journal article, which is in agreement with the findings obtained by Aras et al. [34]. We also observed that the number of TLAs is four times higher in patents than in journal articles, on average. This huge number of TLAs per document makes the NER task on patents harder because of the inherent ambiguity of acronyms. Moreover, the number of tables and figures in patents are more than those in scientific articles. This also makes the extraction of entities from patent documents more difficult than from journal articles [35].

**Table 5** The execution time, in seconds, of NER tools over 10 full patent documents and 10 journal articles

Text genre	Chemical NER tool	
	ChemSpot	tmChem
10 Patent documents	562	66
10 Scientific articles	<i>149</i>	<i>42</i>

The execution time values of both systems are lower on scientific articles shown in italic compared with patents

#### Difference between patent full texts and patent abstracts

The intra-corpus evaluation scores obtained by retraining tmChem (see “Cross-text-genre to cross-corpus evaluation” section) show that the precision (F-measure) values on abstracts are at least 12% (6%) higher than those on full texts. Since both training and test sets contain documents of the same type (abstracts vs. full texts) annotated with the same annotation guideline, we can conclude that the NER task over patent full texts is more complex than that on patent abstracts.

Moreover, the comparisons at instance level indicate that the patterns of errors observed for FP counts are generally different for different types of patent texts, while they are nearly identical for FN counts. We also infer that filtering just a small number of cases correctly as non-chemicals could reduce the FP or FN values significantly. However, achieving such a filtering is difficult, as shown in the following section.

**Table 6 Statistical measurements calculated over 17,000 patent documents and 17,000 journal articles**

Text genre	Sentence length	Document length	Number of unique TLAs	Number of TLAs	Number of tables	Number of figures
Patents	21.12	<i>17,736.00</i>	<i>26.75</i>	<i>187.47</i>	<i>5.34</i>	<i>7.10</i>
Articles	<i>21.70</i>	3512.30	8.47	44.73	2.03	2.97

The largest values are represented in italic for each measurement

### Highly ambiguous entities

Results in “[Comparison at instance level](#)” section have shown that there are several entities which are frequently observed in both lists of entities with highest FP and FN counts. These are entities whose occurrences can, but need not indicate a chemical depending on their local context in patents. In Tables 7 and 8, we provided full confusion matrices for two entities with this property. The entity “alkyl” is observed on both error lists of the corpora containing patent abstracts, and the entity “H” is found in both error lists of corpora containing full patents. The results show that FP and FN counts are in close proximity for both cases which means that the recognition of the corresponding entities is rather difficult.

### Impact of different annotation guidelines

By comparing the results obtained at the instance level shown in Tables 3 and 4, we noticed that some of the errors are produced due to the differences in the annotation guidelines of NER training sets and patent test sets (see “[Patent corpora](#)” section). For example, in these tables, the word “water” is not correctly recognized as a chemical entity by tmChem from BioS corpus or is incorrectly considered chemical by ChemSpot from CEMPs corpora due to the differences in the annotation guidelines.

Moreover, there are many simple chemical elements in the list of entities with high FP counts obtained by tmChem for BioS in Table 3, because simple elements

**Table 7 The full confusion matrix for the ambiguous entity “alkyl” calculated for ChemSpot and tmChem over CEMP\_T and CEMP\_D corpora**

ChemSpot CEMP_T	Predicted “alkyl”	Predicted others	tmChem CEMP_T	Predicted “alkyl”	Predicted others
Actual	TP	FN	Actual	TP	FN
“Alkyl”	354	74	“Alkyl”	206	226
Actual	FP	TN	Actual	FP	TN
Others	260	599	Others	39	816
ChemSpot CEMP_D	Predicted “alkyl”	Predicted others	tmChem CEMP_D	Predicted “alkyl”	predicted others
Actual	TP	FN	Actual	TP	FN
“Alkyl”	372	59	“Alkyl”	187	246
Actual	FP	TN	Actual	FP	TN
Others	194	560	Others	37	715

**Table 8 The full confusion matrix for the ambiguous entity “H” calculated for ChemSpot and tmChem over chapati and BioS corpora containing complete patent documents**

ChemSpot chapati	Predicted “H”	Predicted others	tmChem chapati	Predicted “H”	Predicted others
Actual	TP	FN	Actual	TP	FN
“H”	33	37	“H”	36	34
Actual	FP	TN	Actual	FP	TN
Others	11	789	Others	46	754
ChemSpot BioS	Predicted “H”	Predicted others	tmChem BioS	Predicted “H”	Predicted others
Actual	TP	FN	Actual	TP	FN
“H”	344	1869	“H”	309	1921
Actual	FP	TN	Actual	FP	TN
Others	905	135213	Others	3091	133010

are annotated as chemicals in the training set used for tmChem while they are not labeled as chemicals in BioS corpus. The impact of different rules for annotating simple chemical elements is also observed from the improvement obtained by tmChem in the precision of the BioS corpus after excluding simple chemical elements from both reference and prediction files in “Impact of simple chemical elements” section.

## Conclusion

In this paper, we performed a cross-text-genre evaluation by measuring the tagging quality of the two NER baselines trained on the abstract of scientific articles when evaluated on patent corpora. We noticed that the results are significantly worse on patent corpora compared to scientific abstracts. Although intra-corpus evaluation has shown that training on patent corpora will improve the performance results, performance values are still below the ones achieved for scientific abstracts. Our findings clearly confirm that there are major differences in the NER task between patent and scientific abstracts, and emphasize the complexity of this task on patents.

Moreover, we compared patent abstracts and full texts and addressed the differences between them using various evaluation metrics such as intra-corpus evaluations, and comparison of errors observed at the instance level. We showed that the results on patent abstracts are not extendable to patent full texts which are more important in practice. Therefore, the preparation of more annotated patent full texts is a major requirement for further research in this area.

## Authors' contributions

The experiments are conceived and designed by MH, FS and UL. The experiments are performed by MH and DLW. The results are analyzed by MH and UL. The paper is written by MH and UL. All authors read and approved the final manuscript.

## Author details

<sup>1</sup> Knowledge Management in Bioinformatics, Humboldt-Universität zu Berlin, 12489 Berlin, Germany. <sup>2</sup> Averbis GmbH, 79106 Freiburg, Germany.

## Authors' information

Maryam Habibi received a B.Sc. degree in computer engineering in 2008, and a M.Sc. degree in computer engineering (artificial intelligence) in 2010 both from Sharif University of Technology, Tehran, Iran. She received the Ph.D. degree in electrical engineering from École Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland in 2015. She is currently a post doctoral researcher in the Laboratory of Knowledge Management in Bioinformatics (WBI), Humboldt-Universität, Berlin, Germany. Previously, she was a research assistant at the Idiap research Institute, Martigny, Switzerland. Her current research activities and interests include Natural Language Processing (NLP), especially Information Extraction.

David Luis Wiegandt is a student of computer science at Humboldt-Universität zu Berlin, Germany, where he received his bachelor's degree in June 2016. In June 2015, he joined Knowledge Management in Bioinformatics as a student research assistant.

Florian Schmedding works at Averbis GmbH and is responsible for customer applications and research projects concerning text mining and classification technologies. Before joining Averbis in 2013, he received a PhD

about his research on the SPARQL query language from the University of Freiburg and worked as software engineer in the field of automotive information systems.

Ulf Leser holds the Chair of Knowledge Management in Bioinformatics at the department for Computer Science of Humboldt-Universität zu Berlin, Germany. His main topics of research are biomedical text mining, statistical analysis of -omics data sets and scalable workflows for the analysis of next-generation sequencing data.

## Acknowledgements

The authors are grateful to the Federal Ministry for Economic Affairs and Energy (BMWi) for its financial support through the BioPatent project [KF2205219BZ4].

## Competing interests

The authors declare that they have no competing interests.

Received: 1 July 2016 Accepted: 18 October 2016

Published online: 28 October 2016

## References

- Eisenberg RS (2003) Patents, product exclusivity, and information dissemination: how law directs biopharmaceutical research and development. *Fordham Law Rev* 72(3):477
- Smith BL, Mann SO (2004) Innovation and intellectual property protection in the software industry: an emerging role for patents? *Univ Chic Law Rev* 71(1):241–264
- Granted Patents 2006-2015 per field of technology. <https://www.epo.org/about-us/annual-reports-statistics/statistics.html>. Accessed 23 May 2016
- US Patent Statistics Chart Calendar Years 1963–2015. [http://www.uspto.gov/web/offices/ac/ido/oeip/taf/us\\_stat.htm](http://www.uspto.gov/web/offices/ac/ido/oeip/taf/us_stat.htm). Accessed 23 May 2016
- Adams S (2010) The text, the full text and nothing but the text: Part 1—standards for creating textual information in patent documents and general search implications. *World Pat Inf* 32(1):22–29
- Gurulingappa H, Müller B, Klinger R, Mevissen HT, Hofmann-Apitius M, Fluck J, et al. (2009) Patent retrieval in chemistry based on semantically tagged named entities. In: The eighteenth text REtrieval conference (TREC 2009) Proceedings
- Itoh H, Mano H, Ogawa Y (2003) Term distillation in patent retrieval. In: Proceedings of the ACL-2003 workshop on Patent corpus processing, vol 20. Association for Computational Linguistics, 2003, pp 41–45
- Hansen P, Järvelin K (2005) Collaborative information retrieval in an information-intensive domain. *Inf Process Manag* 41(5):1101–1119
- Mukherjee S, Bamba B (2004) BioPatentMiner: an information retrieval system for biomedical patents. In: Proceedings of the Thirtieth international conference on Very large data bases-Volume 30. VLDB Endowment, 2004, pp 1066–1077
- Batista-Navarro R, Rak R, Ananiadou S (2015) Optimising chemical named entity recognition with pre-processing analytics, knowledge-rich features and heuristics. *J Cheminform* 7(S1):1–13
- Krallinger M, Rabal O, Lourenço A, Perez MP, Rodriguez GP, Vazquez M, et al. (2015) Overview of the CHEMDNER patents task. In: Proceedings of the fifth BioCreative challenge evaluation workshop, pp 63–75
- Jessop DM, Adams SE, Willighagen EL, Hawizy L, Murray-Rust P (2011) OSCAR4: a flexible architecture for chemical text-mining. *J Cheminform* 3(1):1
- Rocktäschel T, Weidlich M, Leser U (2012) ChemSpot: a hybrid system for chemical named entity recognition. *Bioinformatics* 28(12):1633–1640
- Krallinger M, Leitner F, Rabal O, Vazquez M, Oyarzabal J, Valencia A (2013) Overview of the chemical compound and drug name recognition (CHEMDNER) task. In: BioCreative challenge evaluation workshop, vol. 2
- Eltyeb S, Salim N (2014) Chemical named entities recognition: a review on approaches and applications. *J Cheminform* 6(1):1
- Krallinger M, Leitner F, Rabal O, Vazquez M, Oyarzabal J, Valencia A (2015) CHEMDNER: the drugs and chemical names extraction challenge. *J Cheminform* 7(S1):1

17. Müller B, Klinger R, Gurulingappa H, Mevissen HT, Hofmann-Apitius M, Fluck J et al (2010) Abstracts versus full texts and patents: a quantitative analysis of biomedical entities. *Advances in multidisciplinary retrieval*. Springer, New York, pp 152–165
18. Chapati corpus. <http://chebi.cvs.sourceforge.net/viewvc/chebi/chapati/>. Accessed 10 2015
19. Akhondi SA, Klenner AG, Tyrchan C, Manchala AK, Boppana K, Lowe D et al (2014) Annotated chemical patent corpus: a gold standard for text mining. *PLoS One* 9(9):e107477
20. Matos S, Sequeira J, Campos D, Oliveira JL (2015) Identification of chemical and gene mentions in patent texts using feature-rich conditional random fields. In: *Proceedings of the fifth BioCreative challenge evaluation workshop*, pp 76–81
21. Akhondi SA, Pons E, Zubair Afzal Hv, Mulligen JA (2015) Patent mining: combining dictionary-based and machine-learning approaches. In: *Proceedings of the fifth BioCreative challenge evaluation workshop*, pp 102–109
22. Leaman R, Wei CH, Zou C, Lu Z (2015) Mining patents with tmChem, GNormPlus and an ensemble of open systems. In: *Proceedings of the fifth BioCreative challenge evaluation workshop*, pp 140–146
23. Tikk D, Thomas P, Palaga P, Hakenberg J, Leser U (2010) A comprehensive benchmark of kernel methods to extract protein-protein interactions from literature. *PLoS Comput Biol* 6(7):e100837
24. Leaman R, Wei CH, Lu Z (2015) tmChem: a high performance approach for chemical named entity recognition and normalization. *J Cheminform* 7(S-1):1–10
25. Krallinger M, Rabal O, Leitner F, Vazquez M, Salgado D, Lu Z et al (2015) The CHEMDNER corpus of chemicals and drugs and its annotation principles. *J Cheminform* 7(S1):1–17
26. Bada M, Eckert M, Evans D, Garcia K, Shipley K, Sitnikov D et al (2012) Concept annotation in the CRAFT corpus. *BMC Bioinform* 13(1):1
27. Campos D, Matos S, Oliveira JL (2013) Chemical name recognition with harmonized feature-rich conditional random fields. In: *BioCreative challenge evaluation workshop*, vol 2
28. Khabsa M, Giles CL (2015) Chemical entity extraction using CRF and an ensemble of extractors. *J Cheminform* 7(S1):1–9
29. Bagewadi S, Bobić T, Hofmann-Apitius M, Fluck J, Klinger R (2014) Detecting miRNA mentions and relations in biomedical literature. *F1000Res* 3(205):1–33
30. Leaman R, Gonzalez G (2008) BANNER: an executable survey of advances in biomedical named entity recognition. In: *Pacific symposium on bio-computing*, vol 13, pp 652–663
31. Dieb TM, Yoshioka M (2015) Comparison of different strategies for utilizing two CHEMDNER corpora. In: *Proceedings of the fifth BioCreative challenge evaluation workshop*, pp 110–115
32. Klinger R, Kolářik C, Fluck J, Hofmann-Apitius M, Friedrich CM (2008) Detection of IUPAC and IUPAC-like chemical names. *Bioinformatics* 24(13):i268–i276
33. Chen D, Manning CD (2014) A fast and accurate dependency parser using neural networks. In: *EMNLP 2014: conference on empirical methods in natural language processing*, pp 740–750
34. Aras H, Hackl-Sommer R, Schwantner M, Sofean M (2014) Applications and challenges of text mining with patents. In: *Proceedings of the first international workshop on patent mining and its applications (IPaMin 2014)*
35. Zimmermann M (2011) Chemical structure reconstruction with chemOCR. In: *Proceedings of the twentieth text retrieval conference, TREC*

Submit your manuscript to a SpringerOpen® journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Immediate publication on acceptance
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

---

Submit your next manuscript at ► [springeropen.com](http://springeropen.com)

---