



# TSE-NER: An Iterative Approach for Long-Tail Entity Extraction in Scientific Publications

Sepideh Mesbah<sup>(✉)</sup>, Christoph Lofi, Manuel Valle Torre, Alessandro Bozzon,  
and Geert-Jan Houben

Delft University of Technology, Mekelweg 4, 2628 CD Delft, The Netherlands  
{s.mesbah,c.lofi,m.valletorre,a.bozzon,g.j.p.m.houben}@tudelft.nl

**Abstract.** Named Entity Recognition and Typing (NER/NET) is a challenging task, especially with long-tail entities such as the ones found in scientific publications. These entities (e.g. “WebKB”, “StatSnowball”) are rare, often relevant only in specific knowledge domains, yet important for retrieval and exploration purposes. State-of-the-art NER approaches employ supervised machine learning models, trained on expensive type-labeled data laboriously produced by human annotators. A common workaround is the generation of labeled training data from knowledge bases; this approach is not suitable for long-tail entity types that are, by definition, scarcely represented in KBs. This paper presents an iterative approach for training NER and NET classifiers in scientific publications that relies on minimal human input, namely a small seed set of instances for the targeted entity type. We introduce different strategies for training data extraction, semantic expansion, and result entity filtering. We evaluate our approach on scientific publications, focusing on the long-tail entities types *Datasets*, *Methods* in computer science publications, and *Proteins* in biomedical publications.

## 1 Introduction

The growth of domain-specific knowledge available as digital text demands more effective methods for querying, accessing, and exploring document collections. Scientific publications are a compelling example: online digital libraries (e.g. IEEE Xplore) contain hundreds of thousands documents; yet, the available retrieval functionality is often limited to keyword/faceted search on *shallow* meta-data (e.g. title, terms in abstract). A query like *retrieve the publications that used a social media dataset for food recipe recommendation* is bound to return unsatisfactory results<sup>1</sup>.

*Named entities*, obtained through an analysis of a document’s content, are an effective way to achieve better retrieval and exploration capabilities. Automatic *Named Entity* Recognition and Typing (NER/NET) is essential to unlock and

<sup>1</sup> <https://scholar.google.de/scholar?q=publications+using++social+media+datssets+for+food+recipes+recommendation>.

mine the knowledge contained in digital libraries, as most smaller domains lack the resources for manual annotation work.

To perform well, state-of-the-art NER/NET methods [3, 4, 11] either require comprehensive domain knowledge (e.g. to specify matching rules), or rely on a large amount of human-labeled training data for machine learning. Both solutions are expensive and time-consuming.

A cheaper alternative is to generate labeled training data by obtaining existing instances of the targeted entity type from Knowledge Bases (KBs) [3]. This of course requires that the desired entity type is well-covered in the KB.

**Problem Statement.** While achieving impressive performance with high-recall named entities (e.g. locations and age) [11], generic NER/NETs show their limits with domain-specific and long-tail entity types. Consider the following sentence: “We evaluated the performance of *SimFusion+* on the *WebKB* dataset”. Despite *WebKB*<sup>2</sup> being a popular dataset in the Web research community, generic NERs (e.g. *Textrazor*<sup>3</sup>) mistype it as an **Organization** instead of the domain-specific entity type **Dataset**. The entity *SimFusion+* of type **Software** is missed completely.

Literature [20, 26, 27] shows that training of *domain-specific* NER/NETs is still an open challenge for two main reasons: (1) the *long-tail* nature of such entity types, both in existing knowledge bases *and* in the targeted document collections [22]; and (2) the high cost associated with the creation of hand-crafted rules, or human-labeled training datasets for supervised machine learning techniques. Few approaches addressed these problems by relying on bootstrapping [27] or Entity Expansion [3, 11] techniques, achieving promising performance. However, how to train high-performance *long-tail* Entity Extraction and Typing with minimal human supervision remains an open research question.

**Original Contribution.** We contribute TSE-NER, an iterative approach for training NER/NET classifiers for long-tail entity types that exploits Term and Sentence Expansion, extensively expanding on [16]. TSE-NER relies on minimal human input – a seed set of instances of the targeted entity type. We introduce different strategies for training data extraction, semantic expansion, and result entity filtering. Different combinations of these strategies allow to tune the technique for either higher recall or higher precision scenarios.

We performed extensive evaluations comparing to state-of-the-art methods, and assess several sentence expansion and term filtering strategies. As our core use case, we focus on 15,994 data science publications from 10 conference series with the *Dataset* (e.g. *Imagenet*) and data processing *Methods* (e.g. *LSTM*) long-tail entity types. We show that our approach is able to consistently outperform state-of-the-art low-cost supervision methods, even with small amount of training information: with a seed set of 100 entities, our approach can achieve precision up to 0.91 when tuned for precision, and recall up to 0.41 when tuned for recall, or 0.77 and 0.30 for a balanced setting. When applied in an iterative

<sup>2</sup> <http://www.cs.cmu.edu/~WebKB/>.

<sup>3</sup> <https://www.textrazor.com/>.

fashion, our approach can achieve comparable performance with an initial seed set of only 5 entities. We show that sentence expansion and filtering strategies can provide a spectrum of performance profiles, suitable for different retrieval applications such as search (high precision) and exploration (high recall).

To study the performance of TSE-NER across scientific domains, we processed 4,525 biomedical publications focusing on *Protein* (e.g. Myoglobin) entity type. Evaluation on the Craft corpus [2] shows that TSE-NER can achieve performance comparable to existing dictionary-based systems, and obtain precision up to 0.40 and recall up to 0.28 with just 25 seed terms. TSE-NER is implemented in the *SmartPub* platform [17]; its source code is available on the companion Website [18], and its application shown in the video screencast at the following address: <https://youtu.be/zLLMwOT5sZc>.

**Outline.** The remainder of the paper is organized as follows. In Sect. 2 we cover related work. Section 3 presents our approach, and describes alternative data expansion and entity filtering strategies. The experimental setup and results are presented in Sect. 4. Section 5 concludes.

## 2 Related Work

A considerable amount of literature published in recent years addressed the *deep* analysis of text. Common approaches for *deep* analysis of publications rely on techniques such as bootstrapping [27], word-frequency analysis [25], probabilistic methods like Latent Dirichlet Allocation [8], etc. In contrast to current research [25] which limits the analysis of a publication’s content to its title, abstract, references, and authors, we extract entity instances from the much richer full text. In addition, our method does not rely on existing knowledge bases [20, 23] and it is not based on selecting the most frequent keywords [25]. More recent research [26] used both corpus-level statistics and local syntactic patterns of scientific publications to identify entities of interest. Our method uses only a small set of seed names (i.e 5–100), and automatically trained distributed word representations to train a NER in iterative steps (i.e. 2–3).

**Entity Instances Extraction.** Named Entity Recognition (NER) has been applied to identify both entity types of general interest (e.g. Person, Location, Cell, Brand, etc.) as well as for specific domains (e.g., medicine or other domain where resources for training a NER are easily available). NERs rely on different approaches such as dictionary-based, rule-based, machine-learning [26] or hybrid (combination of rule based and machine learning) [29] techniques. Despite its high accuracy, a major drawback of dictionary-based approaches is that they require an exhaustive dictionary of domain terms, which are expensive to create and many smaller domains lack the resources to do so. The same holds for rule-based techniques, which rely on formal languages to express rules and require comprehensive domain knowledge and time to create.

**Bootstrapping and Entity Set Expansion.** Most current NERs are based on Machine Learning techniques, which require a large corpus of labeled training

text [9]. Again, the high costs of data annotation is one of the main challenges in adopting specialized NER for rare entity types in specialized domains [26]. In recent years, many attempts have been made to reduce annotation costs. Active learning techniques have been proposed, asking users to annotate a small part of a text for machine learning methods [7].

Transfer learning techniques [21] use the knowledge gained from one domain and apply it to a different but related named entity type. Co-training [1] starts with a small amount of manually annotated supervised training data and attempt to increase the amount of annotated data. In contrast to previous work, we are not dependent on manually annotated supervised training data [1]; we do not require a large training corpus [21] for transfer learning; also, our approach differs from works on high-recall entity extractors (e.g. with regular expression extractors) for detecting entity types such as location and age [11].

Entity Set Expansion is a technique finding similar entities to a given small set of seed entities [3, 6, 11]. Bootstrapping [27] is another approach similar to our method that uses seed terms and extracts features such as unigrams, bigrams, left unigram, closest verb, etc. These are used to annotate more concept mentions which leads to extracting new features. This step operates in an iterative fashion until no new features are detected. Our approach is inspired by Entity Set Expansion and bootstrapping, but relies on different expansion strategies and does not require concepts already being available in knowledge bases [3].

### 3 Approach

The TSE-NER (Term and Sentence Expansion) approach for domain-specific long-tail entity recognition is organized in five steps, as shown in Fig. 1.

① An initial set of seed terms is used to identify a set of sentences used as initial *training data* (Sect. 3.1). ② *Expansion* strategies can be used to expand the set of initial seed terms, and the *training data* sentences (Sect. 3.2). ③ The *Training Data Annotation* step annotates the training data using the (possibly expanded) seed terms set (Sect. 3.3). ④ A new Named Entity Recognizer (NER) is *trained* using the annotated training data, and the newly trained NER is applied on the corpus to detect a candidate set of entities (Sect. 3.4). ⑤ The *Filtering* step refines the set candidate entities set, to improve the quality of outputted *Verified Terms* set (Sect. 3.5).

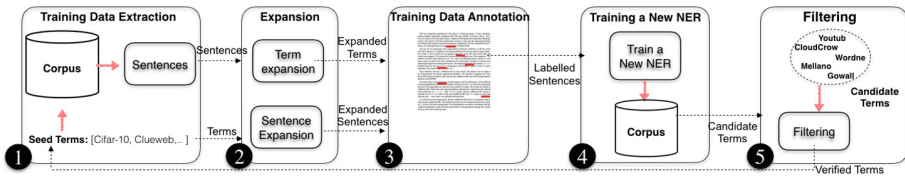


Fig. 1. Overview of the domain-specific long-tail named entities recognition approach.

TSE-NER operates under the hypothesis that there are recurring patterns in the mentions of domain-specific named entities, and that they appear in similar contexts. If this hypothesis holds, by training a classifier on the texts containing the entities, we are able to extract the instances of the entity type of interest. The process can be iterated, by repeating the first step using the newly detected terms as seeds to generate new training data. We rely on the following concepts (some are only relevant for the evaluation, and could be omitted in setups where evaluation is not necessary). The companion website [18] provides a complete unified algorithm covering the TSE-based NER training workflow.

**Known Entity Terms**  $T_{all} := T_{seed} \cup T_{test}$ : This represents a manually created set of instances of the entity type for which a NER classifier is to be trained. In this work, we split this set into a set of seed terms  $T_{seed}$  used for training, and test terms  $T_{test}$  used for evaluation purposes. In a real-life scenario not requiring a formal evaluation, of course only the seed terms would be necessary.  $T_{seed}$  may be small. In this work we consider seed sets  $5 \leq |T_{seed}| \leq 100$ . Creating  $T_{seed}$  is the only manual input required for NER training in our approach.

**Document Corpus**  $D_{all} := \{d_1, \dots, d_{|D|}\}$ : This is the complete document corpus available to our system. Parts of it can potentially be used for training, others for testing. Each document is considered to be a sequence of sentences.

**All Sentences**  $S_{all} := \{s | s \in d \wedge d \in D_{all}\}$ : This represents all sentences of the whole document corpus. Each sentence is considered to be a sequence of terms.

**Test Sentences**  $S_{test} := \bigcup_{t \in T_{test}} \{s | s \in S_{all} \wedge t \in s\}$ : These are all sentences containing any term from the test set, and they need to be excluded from any training in order to ensure the validity of our later evaluations, resulting in the set of **Development Sentences**  $S := S_{all} \setminus S_{test}$ .

In the following, we introduce the iterative version of our approach, representing the current iteration number as  $i$  whereas initially  $i = 0$ . Each iteration  $i$  uses its own term list  $T_i$ , which initially is  $T_0 \subseteq T_{seed}$  (the size of the subset of  $T_{seed}$  depends on the desired use case, as discussed in Sect. 4.3).

### 3.1 Training Data Extraction

As a first step, a set of training data sentences  $S_i$  for the current iteration is created by extracting suitable sentences from  $S$ . At this stage, this is realized by selecting all sentences containing any of the seed terms. Therefore,  $S_i$  provides examples of the positive classification class as they are guaranteed to contain a desired entity instance. To better capture the usage context of the seed entity, we also extract surrounding sentences in the text:  $S_i := \bigcup_{t \in T_i} \{s | s \in S \wedge (t \in s \vee t \in \text{successor}(s) \vee t \in \text{predecessor}(s))\}$ .

### 3.2 Expansion

The small size of the seed term set  $T_{seed}$  has two obvious shortcomings that can greatly hinder the accuracy and recall of the trained NERs: (1) the amount

of training data sentences  $S_i$  is limited; and (2) there are only few examples of mentions of the entity instances of the given type. In addition, the *generalization* capability of the NER for identifying new named entities can also be affected: an insufficient amount of positive examples can lead to entities of the targeted type being labeled negatively; while the extraction of sentences in the training data that are related to seed terms will cause a shortage of negative examples. To account for these issues, we designed two *expansion* strategies.

**Term Expansion (TE).** Term Expansion is designed to increase the number of known instances of the desired entity type before training the NER. An expanded set of entities will provide more positive examples in the training data, thus ideally improving the precision of the NER. In scientific documents, it is common for domain-specific named entities to be in close proximity, e.g. to enumerate alternative solutions, or list technical artifacts. The *Term Expansion* (TE) strategy is therefore designed to test and exploit this hypothesis.

We introduce the interface  $expandTerms(terms_s)$ , with  $terms_s \subseteq terms_i$ . While many different implementations for this interface are possible, in this work we use *semantic similarity*: terms which are semantically similar to terms in the seed list should be included in the expansion. For example, given the dataset seed terms `Clueweb` and `cim-10`, the expansion should add similar terms like `trec-2005`.

We exploit the distributional hypothesis [10] stating that terms frequently occurring in similar context are semantically related, using the popular *word2vec* implementation of skip-n-gram word embeddings [19]. In essence, *word2vec* embeds each term of a large document corpus into low-dimensional vector space (100 dimensions in our case), and the cosine distance between two vectors has been shown to be a high-quality approximation of semantic relatedness [14]. In our implementation, we trained the *word2vec* model on the whole development sentence collection  $S$ , as described in [19], learning all uni- and bigram word vectors of all terms in the corpus. Then, in its most basic version, we select all terms from all sentences, and cluster them with respect to their embedding vectors using K-means clustering. Silhouette analysis is used to find the optimal number  $k$  of clusters. Finally, clusters that contain at least one of the seed terms are considered to (only) contain entities the same type (e.g. *Dataset*).

Initial experiments have shown that this naive approach is slow, and that it can potentially introduce many false positives due to (1) the large number of considered terms, and (2) the sometimes faulty assumption that all terms in cluster are indeed similar as *word2vec* relatedness is not always reliable for similarity measurements [14].

---

**Algorithm 1.** TE using Semantic Relatedness

---

```

function EXPANDTERMS( $terms_s$ )
   $T_{entity} := \{t \in S \mid s \in S \wedge isEntity(t)\}$ 
   $\triangleright$  All entities in  $S$ 
   $clusters := cluster(word2vec(T_{entity}))$ 
   $\triangleright$  Cluster the embeddings
   $clusters_{correct} := \{c \in clusters \mid t \in terms_s$ 
     $\wedge t \in c\}$ 
   $\triangleright$  Select clusters containing any initial term
  return  $\bigcup_{c \in clusters_{correct}}$ 
end function

```

---

To improve, in the following we only consider terms which are likely to be named entities by using NLTK entity detection to obtain a list of all entities  $E_{all}$  contained in  $S^4$ . This results in the Algorithm 1.

**Sentence Expansion (SE).** A second (optional) measure to increase the size and variety of the training set is the *Sentence Expansion* (SE) strategy. It addresses the problem of the over-representation of positive examples resulting from selecting only sentences with instances of the desired type (see Sect. 3.1). The goal is to include negatives sentences not containing instances of the desired type, but are still very similar in semantics and vocabulary.

We rely on *doc2vec* document embeddings [13], a variant of *word2vec*, to learn vector representations of all sentences. For each sentence in  $S$ , we use the cosine distance to discover the most similar sentences filtered to those not containing any known instance of the targeted type. As such sentences might contain an unknown instance of that type, we always combine SE with term expansion to minimize the risk of accidentally mislabeling them as negative examples.

### 3.3 Training Data Annotation

The annotation of training data from the (expanded) seed terms is performed automatically, with no human intervention. After obtaining an (expanded) set of instances  $T_i$  (the current term list) and training sentences  $S_i$ , we annotate each term  $A_{T_i} := \text{annotate}_{T_i}(S_i)$  in all training sentences if they are a positive instance of the targeted entity type, i.e. if the term  $\in T_i$ . Using  $A_{T_i}$ , any state-of-the-art supervised NER can be trained.

### 3.4 NER Training

For training a new  $NER_i$ , we used the Stanford NER tagger<sup>5</sup> to train a Conditional Random Field (CRF) model. As the focus of this paper is the process of training data generation, we do not consider additional algorithms. CRF has shown to be an effective technique on different NER tasks [12]; the goal of CRF is to learn the hidden structure of an input sequence. This is done by defining a set of feature functions (e.g. word features, current position of the word labels of the nearby word), assigning them weights and transforming them to a probability to detect the output label of a given entity. The features used in the training of the model are listed in the companion website. After a NER for the current iteration  $N_i$  is trained, it is used to annotate the whole development corpus  $S$ , i.e.  $A_{NER_i} := \text{annotate}_{NER_i}(S)$ . All positively annotated terms are considered newly discovered instances of our desired type.

<sup>4</sup> NLTK entity detection is based on grammatical context. It does not perform any typing, and due to its simplicity, has high recall values.

<sup>5</sup> <https://github.com/dat/stanford-ner>.

### 3.5 Filtering

After applying the NER to the development corpus, we obtain a list of new candidate terms. As our process relied on several steps which might have introduced noise and false positives (like the expansion steps, but also the NER itself), the goal of this last (optional) step is to filter out candidate terms that are unlikely of the targeted type using a set of external heuristics with different assumptions:

**Wordnet + Stopwords (WS) Filtering.** In the domain-specific language of scientific documents, it is common for named entities to be “proper” of that domain (like *Simlex-999*), or to be expressed as acronyms (like *Clueweb*, *SVM*, *RCV*). In this strategy, named entities are assumed to be not relevant if they are part of the “common” English language, either as proper nouns (e.g. *software*, *database*, *figure*), or a Stopwords (e.g. *on*, *at*). This is achieved by performing lookup operations in WordNet<sup>6</sup> and in common lists of stopwords<sup>7</sup>. As both sources focus on general English language, only domain-specific terms should be preserved.

**Similar Terms (ST) Filtering.** In order to distinguish between different entity types that pertain to a given domain (e.g. *SVM* is of type *Method*, while *Clueweb* is of type *Dataset*), this filtering strategy employs an approach similar to the one used in the *Term Expansion* (TE) strategy. The idea is to cluster entities based on their embedding feature using K-means clustering, and keep all the entities that appear in the cluster that contains a seed term.

**Pointwise Mutual Information (PMI) Filtering.** This filtering strategy adopts a semantic similarity measure derived from the number of times two given keywords appear together in a *sentence* in our corpus. The heuristic behind this filter is vaguely inspired by Hearst Patterns [24], as we manually compile a list of context terms/patterns *CX* which likely indicate the presence of an instance of our desired class (e.g., “we evaluate on x” typically indicates a dataset). Unlike the other filters, it does increase the manual resource costs for training.

Given a set of candidate entities  $CT_i$  and the context term set *CX*, we measure the PMI between them using  $\log \frac{N(ct, cx)}{N(ct)N(cx)}$  with  $ct \in CT_i \wedge cx \in CX$ , and  $N(ct, cx)$  being the number of sentences in which both a candidate entity (*ct*) and a given keyword (*t*) occur (analogously,  $N(ct)$  counts the number of occurrences of *ct*). Finally, candidate terms are filtered and excluded if their PMI value is below a given threshold value.

**Knowledge Base Lookup (KBL) Filtering.** Our target are long-tail domain-specific entities, i.e. entities that are not part of existing knowledge bases. Named entities that could be linked to a knowledge base could be assumed incorrect, and therefore amenable to exclusion from the final named entity set. In the KBL approach we exclude the entities that have a reference in the DBpedia.

<sup>6</sup> <http://wordnet.princeton.edu/>.

<sup>7</sup> <http://www.nltk.org/book/ch02.html>.



**Ensemble (EN) Filtering.** Different filtering strategies are likely to remove different named entities. To reduce the likelihood of misclassification, the *Ensemble* (EN) filtering strategy combines the judgment of multiple filtering strategies, to preserve candidate entities that are considered correct by one or more strategy. Intuitively, if each strategy makes different errors, then a combination of the filters’ judgment can reduce the total error. We preserve the entities that are passed through two out of three selected filtering strategies.

## 4 Evaluation

This section reports on an empirical evaluation to assess the performance of the approach (and its variants) described in Sect. 3, and the ability to utilize it for long-tail named entity recognition. Sect. 4.1 describes the experimental set-up, followed by the results (Sect. 4.2), and their discussion (Sect. 4.3).

### 4.1 Experimental Setup

**Corpora.** Our main evaluation, shown in the following sections, is performed on the data science (15,994 papers from 10 conference series) domain. To assess the performance of TSE-NER in other scientific domains, at the end of the section we describe an experiment over 4,525 publications from 10 biomedical journals. The full description of the corpora is described in the companion Web site [18]. Publications are processed using GROBID [15], to extract a structured full-text representation of their content.

**Long Tail Entity Types Selection.** Scientific publications contain a large quantity of long-tail named entities. Focusing on the data science domain, we address the entity types *Dataset* (i.e. dataset presented or used in a publication), and *Methods* (i.e. algorithms – novel or pre-existing – used to create/enrich/analyze a dataset). Both entities types are scarcely represented in existing knowledge bases<sup>8</sup>. To evaluate the performance of our approach, we create a set of 150 seed instances  $T_{all}$  for each targeted type, collected public from public websites<sup>9</sup>.

For each type, 50 of those are selected as test terms for that type  $T_{test}$ , while 100 are used as seed terms  $T_{seed}$ .

**Evaluation Dataset.** As discussed in Sect. 3, in the training process all test sentences  $S_{test}$  (i.e. sentences mentioning terms in  $T_{test}$ ) in the corpus  $D_{all}$  are removed. For evaluation, we manually created a type-annotated test set: for

<sup>8</sup> In DBPedia, the type `dbo:database` features 989 instances, but mostly related to biology, economy, and history. The type `dbo:software` contain names of several algorithms, but the list is clearly incomplete.

<sup>9</sup> For instance: <https://github.com/caesar0301/awesome-public-datasets>. The full list of seed entity instances, as well as the list of sources are available on the companion Website.

each test term, we select all sentences in which they are contained including any adjacent sentence, forming the set of annotated sentences  $S_{annotated} := \cup_{t \in T_{test}} \{s | s \in S_{test} \wedge (t \in s \vee t \in successor(s) \vee t \in predecessor(s))\}$ . An expert annotator labeled each term as an instance of the target type to create the test annotation set used for evaluation  $A_{test} := annotate_{expert}(S_{annotated})$ .

Details of statistics on sentences used for training and testing can be found in the companion Web site. For training, depending on the seed set size between 5 and 100, we used between 198 and 2863 sentences for the *dataset* entity type and 617 to 18545 sentences for the *Method* entity type.

For testing 50 seed terms were used for both dataset (i.e. 3149 sentences) and method (i.e. 1097 sentences) entity type. The evaluation protocol is described in Algorithm 2, where the *seed\_size* values can be initialized with different values. Our analysis was not limited to the 50 test seed terms, we further evaluated 200 entities recognized by TSE-NER via a pooling technique.

---

**Algorithm 2.** Evaluation Protocol

---

```

function EVALUATE(seed_size)
   $T \subseteq_{seed\_size} T_{seed}$ 
   $NER_{final} := longtailTrain(T, S_{all})$ 
   $A_{final} := annotate_{NER_{final}}(S_{annotated})$ 
   $result := analyze(A_{final}, A_{test})$ 
end function

```

---

## 4.2 Results

For a given entity type (*Dataset* and *Method*), we test the performance with differently sized seed sets and expansion strategies to create the training data for generating the NER model, and different filtering strategies. We report the performance of the basic WS, PMI, and EN strategies, plus a combination of the WS, ST, and KBL strategies, as listed in Table 1. The complete evaluation results for all the seed set size and the filtering techniques can be found in the companion Web site. We investigate iterative performance, and results on the manually annotated test from the previous section.

Tables 1 and 2 summarize the performance achieved for *Dataset* and *Method* entity types. In Table 2, the *No Expansion* and *Term Expansion* figures for the *Method* type are omitted for brevity's sake. Our approach is able to achieve excellent precision [89% – 91%] with both entity types, and good recall (up to 41%) with the *Dataset* type. The lower recall obtained with the *Method* type can be explained with the greater diversity (in terms of n-grams and use of acronyms) of method names.

The expansion strategies lead to an average +200% (SE – *Dataset*) and +300% (TE – *Dataset*) increase in recall, thus demonstrating their effectiveness for generalization. On average, filtering decrease recall, but with precision improvements up to +20% (PM – *Method*). These are promising figures, considering the minimal human supervision involved in the training of the NERs. We can also show the different trade-offs our approach can strike: different configurations of filtering and expansion lead to different results with respect to precision and recall values, allowing for example a high-precision slightly-lower

recall setup for a digital library, and a higher recall lower precision setup for a Web retrieval system.

**Expansion Strategies.** Expansion strategies increase the size and variety of training datasets, thus improving the precision and recall. Both strategies achieve the expected results, although with different performance increase: compared to *NE* strategy, both *TE* and *SE* achieve a considerable performance boost ( $\mu = +190\%$ ) for recall, but at cost of lower precision ( $\mu = -8.7\%$ ). We account the better recall performance of *TE* to the contextual similarity (and proximity) of named entities of the same type in technical documents (e.g. *Gov2*, *Robust04*, *ClueWeb* and *Wt10g*). The precision decrease in *TE* can be accounted to treating some terms incorrectly as positive instances due to their presence in the same embedding clusters as the seed terms (see also Sect. 3.2). The *SE* strategy shows lower recall ( $\mu = +210\%$  over *NE*), but with less precision loss ( $\mu = -5.2\%$  than *NE*). We account this positive behaviour to the presence of more quality negative examples, helping to maintain the generalization capabilities of the NER, while refining the quality of its recognition.

**Filtering Strategies.** We observe no significant improvement in precision with the *WS* filtering approach. Manual inspection of results reveal that most of the false positives are already domain-specific terms (e.g. *Pagerank*, *Overcite* for *Dataset*, and *NDCG* for *Method*) which are not included in Wordnet, but that are of the wrong type. *SS* slightly increases the precision by keeping only the entities that appear in the same cluster as the seed names; however, this comes at a cost,

**Table 1.** *Dataset* entity type: precision/recall/F-score on evaluation dataset. Legend: *NE* – No Expansion; *TE* – Term Expansion; *SE* – Sentence Expansion; *NF* – No Filtering; *WS* – Wordnet + StopWords; *SS* – Similar Terms + WS; *KS* – Knowledge Base Lookup + SS; *PM* – Point-wise Mutual Information; *EN* – Ensemble.

<i>Strategy</i>	<i>#S</i>	<i>NF</i>	<i>WS</i>	<i>SS</i>	<i>KS</i>	<i>PM</i>	<i>EN</i>
<i>NE</i>	5	.83/.05/.10	.84/.04/.08	.86/.03/.07	.75/.01/.01	.90/.04/.09	.86/.04/.08
	25	.84/.08/.16	.83/.07/.13	.86/.07/.13	.78/.01/.03	.91/.08/.15	.85/.07/.13
	100	.85/.15/.26	.85/.13/.22	.87/.12/.22	.82/.03/.07	.91/.13/.24	.86/.12/.22
<i>TE</i>	5	.76/.14/.25	.78/.13/.22	.79/.11/.20	.74/.04/.09	.83/.13/.23	.80/.13/.22
	25	.72/.29/.42	.73/.28/.40	.75/.27/.40	.73/.17/.28	.77/.27/.40	.75/.27/.40
	100	.69/.41/.51	.70/.39/.50	.71/.38/.50	.71/.28/.40	.74/.38/.50	.72/.38/.50
<i>SE</i>	5	.83/.07/.14	.84/.06/.12	.86/.05/.10	.82/.01/.02	.91/.07/.13	.86/.06/.11
	25	.81/.22/.35	.80/.18/.29	.83/.17/.29	.77/.04/.08	.89/.20/.33	.82/.18/.29
	100	.77/.30/.43	.77/.24/.37	.80/.23/.36	.78/.07/.13	.86/.26/.40	.79/.24/.37

**Table 2.** *Method* entity type: precision/recall/F-score. Legend as in Table 1.

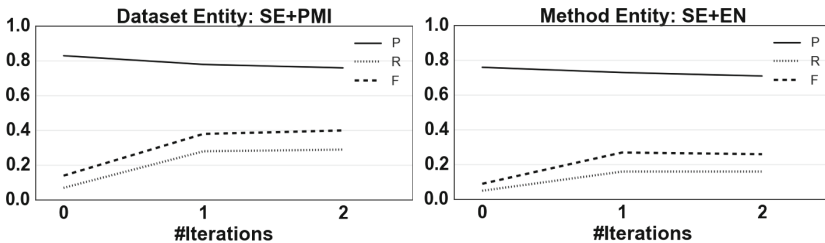
<i>Strategy</i>	<i>#S</i>	<i>NF</i>	<i>WS</i>	<i>SS</i>	<i>KS</i>	<i>PM</i>	<i>EN</i>
<i>SE</i>	5	.76/.04/.08	.77/.03/.07	.77/.01/.01	.84/.01/.01	.86/.01/.03	.84/.03/.05
	25	.77/.14/.24	.77/.12/.21	.79/.09/.16	.87/.05/.09	.86/.05/.09	.85/.09/.17
	100	.68/.15/.25	.67/.14/.23	.65/.12/.20	.84/.07/.13	.85/.05/.10	.83/.10/.19

as the recall is also penalized by the exclusion of entities of interest that are in other clusters. KB excludes popular entities that are contained in the knowledge base (e.g. *Wordnet*, *Dailymed*), but also some rare entities that are mistyped.

For instance, the *Dataset* entities *Ratebeer*<sup>10</sup> or *Jester* can be retrieved from DBpedia using the lookup search, although the result points to another entity. This is a clear limitation with the adopted lookup technique, which could be avoided with a more precise implementation of the lookup function. PMI usually gets the highest precision; the strategy proved effective in removing false positives, but penalizes recall by excluding entities that do not appear with the words in the context list *CX*. For instance, *Unigene (Dataset)* often appears in with the term *data source*, which is not in our context list and thus filtered out. The EN strategy keeps only the entities that are preserved by two out of three (WS, KB and PMI) filtering strategies. While reducing the number of false positives, this proves to be too restrictive; for instance *Dataset* names such as *Yelp*, *Twitter*, *Foursquare* and *Nasdaq* are removed by both the WS and KB strategies.

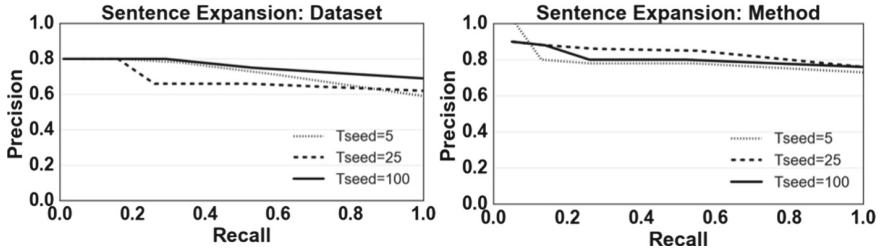
**Seed Set Size.** We randomly initialize  $T \subseteq T_{seed}$  with  $|T| = 5, 10, 25, 50, 100$  (see Algorithm 2). We execute the evaluation cycle 10 times for each size of  $T$ , and again vary expansion and filtering strategies. The recall performance sharply increase with the number of seeds term ( $\mu = +340\%$  from 5 to 100 seeds): this is due to the increase in the number of sentences available for NER training, and is an expected behaviour. The decrease in precision is an average of  $-6\%$  from 5 to 100 seeds, with an average value of  $-5.1\%$  for *Dataset* and  $-6.9\%$  for *Methods*. Noteworthy are the good performance with as little as 5 seed entities (*Datasets*: 0.25 F-score with TE strategy and no filtering).

**Iterative NER Training.** Figure 2 shows the result of the iterative NER training using Sentence Expansion with 5 seeds. We report the results with the PMI (*Dataset*) and EN (*Methods*) filtering, as they are the ones offering the most balanced performance in both precision and recall. Despite the small initial seed seed, it is possible to achieve precision and recall comparable to the ones obtained with an initial set of 100 seeds in only 2 iterations.



**Fig. 2.** Dataset (L) and Method (R) entity: iterative NER training using 5 initial seeds.

<sup>10</sup> <http://lookup.dbpedia.org/api/search/KeywordSearch?QueryClass=&QueryString=ratebeer>.



**Fig. 3.** *Dataset* (L) and *Method* (R): precision and recall for ranked top 10, 25, 50, 100 and 200 entities, varying seeds sizes.

**Analysis of Recognized Entities.** To widen the scope of our evaluation, we extended our result analysis beyond the 150 named entities in  $T_{all}$ . We manually investigated up-to-now unknown named entities which have been recognized by the NER after training. We applied a method inspired by the pooling technique typically used in information retrieval research: given a list of seed terms  $T_{seed}$  of a given type, and a list of recognized potential filtered terms  $FT$  of an yet unknown type, the idea is to rank the items in the list of candidate terms  $FT$  according to their embedding similarity to the items in the seed set  $T_{seed}$  and collect the top  $K$ . As a result, the obtained precision and recall measurements are only approximate values. The similarity is measured based on the cosine similarity between the *word2vec* embedding vectors. Each entity in the lists has been manually checked by an expert. Figure 3 shows the precision and recall of the top  $K = 10, 25, 50, 100,$  and 200 retrieved entities using the SE approach. As in the previous experiment, we used the PMI and EN filtering strategies respectively for *Dataset* and *Method* types. Precision performance are consistently high at all level of recall. Note that we randomly selected  $T \subseteq T_{seed}$  with  $|T|=5, 25, 100$  seed terms and used them to train the NER using the SE strategy. Variations in precision performance in Fig. 3 are therefore accountable on the initial seed term used in each configuration (seed terms might bring in more false positives).

The *Dataset* entities *mslr-web10* (a benchmark collection for learning to rank method) and *ace2004* (ACE 2004 Multilingual Training Corpus); and *Method* entities such as *TimedTextTank* and *StatSnowball* are a sample of extracted entities. More examples can be found in the companion website. Some examples of incorrect detected entities are due to ambiguous nature of the sentence. Consider the following sentence: “*The implementation of scikitlearn toolkit was adopted for these methods*”, since it is similar to a sentence that contains a method entity, the entity *scikitlearn* was detected as a method although its a software library. In another sentence: “*The Research Support Libraries Programme (RSLP) Collection Description Project developed a model.*”, *RSLP* (a project) was detected as a dataset due to its surrounding words (e.g. *collection, libraries*).

**Comparison with State-of-the-Art.** We compared our method with: (1) the BootStrapping (BS) based concept extraction approach [27], a commonly used state-of-the-art technique in scientific literature; the experiments where executed

with the code and the parameters  $(k, n, t)$  to  $(2000, 200, 2)$  provided in [27], and with 100 seeds. And, (2) improved and expanded Hearst Pattern (HP) [24] for automatically building or extending knowledge bases extracting type-instance relations e.g., X such as Y as in “we used datasets such as twitter”. Intuitively, the performance of BS decreases with less number of seed terms. For the HP we kept type-instance pairs related to dataset or method (i.e. the context words in  $CX$ ). Experiments on our evaluation dataset shown that TSE-NER achieved better performance in terms of precision/recall/f-score for the *dataset* entity type (0.77/0.30/0.43) compared to *BS* (0.08/0.13/0.10) and *HP* (0.92/0.15/0.27) as well as for the *method* entity (*TSE-NER*: 0.68/0.15/0.25, *BS*: 0.11/0.32/0.16, *HP*: 0.64/0.04/0.07). The high precision and low recall in *HP* is explained by the limited set of *HP* patterns. We infer that different expansion strategies augmented the performance of our technique compared to the *BS* which just relies on features such as unigrams, bigrams, closest verb, etc. Finally we also evaluated the performance of traditional supervised annotation. The supervised approach can achieve precision/recall/f-score of 0.82/0.35/0.49 for *dataset* entity type and 0.70/0.17/0.28 for *method* entity type using training data from 100 seeds.

**Biomedical Domain.** To test the performance of TSE-NER on another domain, we processed 4,525 biomedical publications from 10 journals focusing on the *Protein* entity type. The seed terms were selected from the protein ontology.<sup>11</sup> We excluded test terms appearing in the Craft corpus [2] (a manually annotated corpus containing 67 full-text biomedical journals) and kept only those with references in the publications (see companion site). We randomly initialized  $T \subseteq T_{seed}$  with  $|T| = 5, 25, 100$  and employed the SE strategy and a simple WS filtering. The evaluation cycle has been executed 10 times for each size of  $T$ , and results are averaged. TSE-NER can achieve precision/recall/f-score of 0.57/0.08/0.14 using 5 seeds, 0.40/0.28/0.32 using 25 seeds, and 0.38/0.46/0.41 with 100 seeds. The latter results are comparable to extensive dictionary-based systems [28] (0.44/0.43/0.43) [5] (0.57/0.57/0.57) where existing ontologies in the biomedical domain are used for matching *Protein* entities of the text.

### 4.3 Discussion

The design goal of the TSE-NER approach was minimizing the training costs in scenarios where the targeted entity types are rare, and little to no resources (for manual annotations) are available. In these cases, relying on dictionaries or knowledge-bases is not feasible, and common techniques like supervised learning cannot be applied. We believe to have successfully reached that goal, as we could show that even with small seed lists  $T_{seed}$  with little as 5 or 25 terms, high-precision NERs could be trained.

Nonetheless, this ease-of-training comes at a price: recall values are low, and are unlikely to be able to compete with known much more elaborately trained NERs for popular types. However, by selecting different configurations for filtering and expansion, recall can be moderately improved at the cost of precision.

<sup>11</sup> <http://obofoundry.org/ontology/pr.html>.

Also, the effectiveness of such changes of configurations seems to slightly differ between the *Dataset* and *Method* entity types. As a result, we cannot identify one clear best configuration as TSE-NER seems to benefit from some entity type-specific tuning. However, this also provides some flexibility to tune with respect to different quality and application requirements.

Furthermore, some of our underlying assumptions, heuristics and implementation choices, are designed as a simplistic prove-of-concepts, and deserve further discussion and refinement. As an example, consider WS WordNet filtering: we assumed domain-specific named entities would not be part of common English language. While this is true for many relevant domain-specific entities, several datasets (for instance) do indeed carry common names like the **census dataset**. For a production system, more complex implementations and tailored crafting is necessary for reaching better performance values. Another restriction is related to the core heuristics found in the term and sentence expansion, where we assume that similar types of entities occur in similar contexts – which is not necessarily always the case.

**Threats To Validity.** Our evaluation has been performed on an extensive document corpus, covering two distinctively different domains. However, we focused only on a limited set of entity types. The hypothesis described in Sect. 3 hold for *Datasets*, *Methods*, and *Proteins*, but further experiments are needed for other entity types in the same domains (e.g. *Software*) or in other domains. Despite the good performance achieved, it could already be noted that even between those three types, no single TSE-NER configuration is clearly the best. In order to obtain a complete understanding of the full capabilities, limitations, and trade-offs of our approach, more studies addressing additional domains and entity types are necessary.

## 5 Conclusion

We presented a novel approach for the extraction of domain-specific long-tail entities from scientific publications. A limiting factor in this scenario is the lack of resources and/or available explicit knowledge to allow for established NER training techniques. We explored techniques able to limit the reliance on human supervision, resulting in an iterative approach that requires only a small set of seed terms of the targeted type. Our core contributions, in addition to the overall approach, are a set of expansion strategies exploiting semantic relatedness between terms to increase the size and labelling quality of the generated training dataset, as well as several filtering techniques to control the noise.

In our evaluation, we could show that we can reach a precision of up to 0.91, or a recall of up to 0.41 – a good result considering the very cheap training costs. Furthermore, we could show that recall can be traded for more precision to a moderate extend by changing the configuration of our NER training process.

For future work, additional evaluation addressing more domains and entity types is of importance to better understand the range of applicability of our approach. Also, many of our currently still simplistic heuristics and implementation choices can benefit from (domain-specific) improvement and optimization.

## References

1. Agerri, R., Rigau, G.: Robust multilingual named entity recognition with shallow semi-supervised features. *Artif. Intell.* **238**, 63–82 (2016)
2. Bada, M., et al.: Concept annotation in the craft corpus. *BMC bioinf.* **13**(1), 161 (2012)
3. Brambilla, M., Ceri, S., Della Valle, E., Volonterio, R., Acero Salazar, F.X.: Extracting emerging knowledge from social media. In: *International Conference on World Wide Web*, pp. 795–804 (2017)
4. Derczynski, L., Nichols, E., van Erp, M., Limsopatham, N.: Results of the WNUT2017 shared task on novel and emerging entity recognition. In: *Proceedings of the 3rd Workshop on Noisy User-Generated Text*, pp. 140–147 (2017)
5. Funk, C., et al.: Large-scale biomedical concept recognition: an evaluation of current automatic annotators and their parameters. *BMC bioinf.* **15**(1), 59 (2014)
6. García-Pablos, A., Cuadros, M., Rigau, G.: W2VLDA: almost unsupervised system for aspect based sentiment analysis. *Expert Syst. Appl.* **91**, 127–137 (2018)
7. Goldberg, S., Wang, D.Z., Grant, C.: A probabilistically integrated system for crowd-assisted text labeling and extraction. *J. Data Inf. Qual. (JDIQ)* **8**(2), 10 (2017)
8. Griffiths, T.L., Steyvers, M.: Finding scientific topics. *Proc. Nat. Acad. Sci.* **101**(suppl 1), 5228–5235 (2004)
9. Habibi, M., Weber, L., Neves, M., Wiegandt, D.L., Leser, U.: Deep learning with word embeddings improves biomedical named entity recognition. *Bioinformatics* **33**(14), i37–i48 (2017)
10. Harris, Z.: Distributional structure. *Word* **10**, 146–162 (1954)
11. Kejriwal, M., Szekeley, P.: Information extraction in illicit web domains. In: *International Conference on World Wide Web*, pp. 997–1006 (2017)
12. Lafferty, J., McCallum, A., Pereira, F.C.: Conditional random fields: probabilistic models for segmenting and labeling sequence data. In: *International Conference on Machine Learning*, vol. 951, pp. 282–289 (2001)
13. Le, Q., Mikolov, T.: Distributed representations of sentences and documents. In: *International Conference on Machine Learning (ICML-14)*, pp. 1188–1196 (2014)
14. Lofi, C.: Measuring semantic similarity and relatedness with distributional and knowledge-based approaches. *Inf. Media Tech.* **10**(3), 493–501 (2015)
15. Lopez, P.: GROBID: combining automatic bibliographic data recognition and term extraction for scholarship publications. In: Agosti, M., Borbinha, J., Kapidakis, S., Papatheodorou, C., Tsakonas, G. (eds.) *ECDL 2009. LNCS*, vol. 5714, pp. 473–474. Springer, Heidelberg (2009). [https://doi.org/10.1007/978-3-642-04346-8\\_62](https://doi.org/10.1007/978-3-642-04346-8_62)
16. Mesbah, S., Fragkeskos, K., Lofi, C., Bozzon, A., Houben, G.-J.: Semantic annotation of data processing pipelines in scientific publications. In: Blomqvist, E., Maynard, D., Gangemi, A., Hoekstra, R., Hitzler, P., Hartig, O. (eds.) *ESWC 2017. LNCS*, vol. 10249, pp. 321–336. Springer, Cham (2017). [https://doi.org/10.1007/978-3-319-58068-5\\_20](https://doi.org/10.1007/978-3-319-58068-5_20)
17. Mesbah, S., Lofi, C., Bozzon, A., Houben, G.-J.: SmartPub: a platform for long-tail entity extraction from scientific publications. In: *The Web Conference* (2018)
18. Mesbah, S., Lofi, C., Bozzon, A., Houben, G.-J.: TSE-NER companion page (2018). <https://sites.google.com/view/iswc2018/>
19. Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., Dean, J.: Distributed representations of words and phrases and their compositionality. In: *Advances in Neural Information Processing Systems*, pp. 3111–3119 (2013)



20. Osborne, F., de Ribaupierre, H., Motta, E.: TechMiner: extracting technologies from academic publications. In: Blomqvist, E., Ciancarini, P., Poggi, F., Vitali, F. (eds.) EKAW 2016. LNCS (LNAI), vol. 10024, pp. 463–479. Springer, Cham (2016). [https://doi.org/10.1007/978-3-319-49004-5\\_30](https://doi.org/10.1007/978-3-319-49004-5_30)
21. Qu, L., Ferraro, G., Zhou, L., Hou, W., Baldwin, T.: Named entity recognition for novel types by transfer learning. In: EMNLP (2016)
22. Reinanda, R., Meij, E., de Rijke, M.: Document filtering for long-tail entities. In: Proceedings of the 25th ACM International on Conference on Information and Knowledge Management, pp. 771–780. ACM (2016)
23. Sateli, B., Witte, R.: What’s in this paper?: Combining rhetorical entities with linked open data for semantic literature querying. In: International Conference on World Wide Web, pp. 1023–1028 (2015)
24. Seitner, J., et al.: A large database of hypernymy relations extracted from the web. In: LREC (2016)
25. Shubankar, K., Singh, A., Pudi, V.: A frequent keyword-set based algorithm for topic modeling and clustering of research papers. In: 2011 3rd Conference on Data Mining and Optimization (DMO), pp. 96–102. IEEE (2011)
26. Siddiqui, T., Ren, X., Parameswaran, A., Han, J.: FacetGist: collective extraction of document facets in large technical corpora. In: International Conference on Information and Knowledge Management, pp. 871–880. ACM (2016)
27. Tsai, C.-T., Kundu, G., Roth, D.: Concept-based analysis of scientific literature. In: International Conference on Information Knowledge Management. ACM (2013)
28. Tseytlin, E., Mitchell, K., Legowski, E., Corrigan, J., Chavan, G., Jacobson, R.S.: Noble-flexible concept recognition for large-scale biomedical natural language processing. *BMC bioinf.* **17**(1), 32 (2016)
29. Tuarob, S., Bhatia, S., Mitra, P., Giles, C.L.: Algorithmseer: a system for extracting and searching for algorithms in scholarly big data. *IEEE Trans. Big Data* **2**(1), 3–17 (2016)