

Improving term extraction with terminological resources

Sophie Aubin and Thierry Hamon

LIPN – UMR CNRS 7030
99 av. J.B. Clément, F-93430 Villetaneuse
Tél. : 33 1 49 40 40 82, Fax. : 33 1 48 26 07 12
firstname.lastname@lipn.univ-paris13.fr,
WWW home page: www-lipn.univ-paris13.fr/~lastname

Abstract. Studies of different term extractors on a corpus of the biomedical domain revealed decreasing performances when applied to highly technical texts. Facing the difficulty or impossibility to customize existing tools, we developed a tunable term extractor. It exploits linguistic-based rules in combination with the reuse of existing terminologies, *i.e.* exogenous disambiguation. Experiments reported here show that the combination of the two strategies allows the extraction of a greater number of term candidates with a higher level of reliability. We further describe the extraction process involving both endogenous and exogenous disambiguation implemented in the term extractor \LaTeX .

1 Introduction

Identifying and extracting terms from texts is now a well-known and widely explored step in the terminology building process. Different strategies can be applied: term extraction based on lexico-syntactic markers [1], chunking based syntactic frontiers and endogenous parsing [2], and distributional analysis [3]. Those different techniques show satisfying extraction results regarding the recall [4]. However, studying the outputs of three term extractors applied to an English biomedical corpus, we found that they are not adequate for highly technical texts [5]. The results of the extraction are generally noisy for different reasons. First, some errors result from the tagging of the corpus. The second limitation of such tools is their difficulty to distinguish terms or variants from nominal phrases that are not terms. Finally, they lack portability to new domains as it is difficult to define parsing patterns large enough with a good precision.

Extracting terms consists not only in identifying specific nominal phrases but also in providing a reliable syntactic analysis. The latter is commonly used to organise terminologies through a syntactic network and to compute hierarchical relationships using lexical inclusion. Manually written rules based on linguistic clues are insufficient for this task and must be combined with statistical methods.

Several strategies have been used and sometimes associated to finally extract the term candidates: statistical filtering [1], manual filtering through the tool

interface [2] or the exploitation of external resources. We propose a combination of the three methods.

The terminology extractor we implemented uses techniques comparable to state-of-the-art tools, among which chunking based on morpho-syntactic frontiers and production of the syntactic analysis of the terms extracted. We further propose new solutions for chunking and parsing by using external resources. In addition, we chose to perform positive filtering in the parsing step through the mechanism of islands of reliability (see Section 3.1). In comparison, other tools produce all parsing solutions and filter out non valid ones *a posteriori*.

We first discuss the limitations of matching existing terminologies on corpora and of automatic extraction tools. As an answer to this, we propose a combination of terminology extraction with the exploitation of testified resources. We describe the extraction process of Y_AT_EA that implements the method we propose. We finally present the results of experiments run on a biomedical corpus to characterise the effects of recycling existing terminologies in a term extractor.

2 Which approach to identify terms?

Terms can be identified in corpora regarding two approaches : matching terms issued from terminological resources, or designing automatically term extraction methods.

Using terminological resources to identify terms in texts addresses the question of the usability of resources on working corpora, namely their coverage and their adequacy. This leads to evaluate how terms issued from resources, i.e. testified terms, match in the working corpus. As terminological resources are widely available in the biomedical field, many experiments have been done on recycling terminologies to identify terms in medical and biological corpora. Coverage is generally mitigated. The coverage of well-known classifications as ICD-9, ICD-10 or SNOMED III have been observed on a 14,247 word corpus of clinical texts [6]. The evaluation leads to conclude that no classification covers sufficiently the corpus, although SNOMED has the better content coverage. Similar observations have been noted regarding the evaluation of the usability of Gene Ontology for NLP [7]. 37% of the GO terms are found in a 400,000 Medline citation corpus. Results vary depending on the GO categories from 28% to 53 % in the Medline corpus. [7] consider that this low content coverage could be due to the size of the working corpus or its narrow scope. Still, content coverage is even worse on a set of 3 million randomly selected noun phrases among 14 million terms extracted from the Medline corpus [8]: most of them are not present in UMLS. In [9], we showed that, in the context of the indexation of specialized texts, even if the combination of resources is useful to identify numerous testified terms or variants, the indexation varies greatly according to the documents.

Alternatives, based on the automatic extraction of terms, have been widely proposed since the 90's. [4] give an overview of the proposed term extractors. These term identification methods generally exploit linguistic information like boundaries or, more often, patterns. Such approaches are difficult to evaluate

without a golden standard and evaluations vary according to the methods. However, the recall is generally good ([2] estimates the silence to 5%), while the precision is rather low ([2] rejects 50% of the extracted term candidates, the system discussed in [10] has an error rate of 20%).

Pure term extraction methods rarely use terminological resources. Such domain information is rather exploited at the filtering step [10]. However, the usefulness of terminological resources in a term extraction process is demonstrated in FASTR [11]. Results of this term variant extraction system are rather good as term variation acquisition increases the terminological resource coverage. The limitation of this approach is the acquisition of terms unrelated to testified ones.

Regarding the works discussed above, it seems obvious that terminological resources provide precious information that must be used in a term identification task. However, exploiting terminological resources requires their availability and adequacy on the targeted corpus. On the opposite, automatic term extraction approaches suffer from a necessary human validation step. In that respect, we aim at combining both approaches by developing a term extraction method that exploits terminological resources when available.

3 Strategy of term extraction

The software $\text{Y}_{\text{A}}\text{T}_{\text{E}}\text{A}$, developed in the context of the ALVIS¹ project, aims at extracting noun phrases that look like terms from a corpus. It provides their syntactic analysis in a head-modifier format. As an input, the term extractor requires a corpus which has been segmented into words and sentences, lemmatized and tagged with part-of-speech (POS) information. The implementation of this term extractor allows to process large corpora. It is not dependent on a specific language in the sense that all linguistic features can be modified or created for a new language, sub-language or tagset. In the experiments described here, we used the GENIA tagger² [12] which is specifically designed for biomedical corpora and uses the Penn TreeBank tagset.

The main strategy of analysis of the term candidates is based on the exploitation of simple parsing patterns and endogenous disambiguation. Exogenous disambiguation is also made possible for the identification and the analysis of term candidates by the use of external resources, *i.e.* lists of testified terms.

This section includes the presentation of both endogenous and exogenous disambiguation strategies. We also describe the whole extraction process implemented in $\text{Y}_{\text{A}}\text{T}_{\text{E}}\text{A}$.

3.1 Endogenous and exogenous disambiguation

Endogenous disambiguation consists in the exploitation of intermediate extraction results for the parsing of a given Maximal Noun Phrase (MNP).

¹ European Project STREP IST-1-002068-STP, <http://www.alvis.info/alvis/>

² <http://www-tsujii.is.s.u-tokyo.ac.jp/GENIA/tagger/>

All the MNPs corresponding to parsing patterns are parsed first. In a second step, remaining unparsed MNPs are processed using the MNPs parsed during the first step as *islands of reliability*. An *island of reliability* is a subsequence (contiguous or not) of a MNP that corresponds to a shorter term candidate in either its inflected or lemmatized form. It is used as an anchor as follows: the subsequence covered by the island is reduced to the word found to be the syntactic head of the island. Parsing patterns are then applied to the simplified MNP.

This feature allows the parse of complex noun phrases using a limited number of simple parsing patterns (80 patterns containing a maximum of 3 content words were defined for the experiments described below). In addition, islands increase the degree of reliability of the parse as shown in Figure 1.

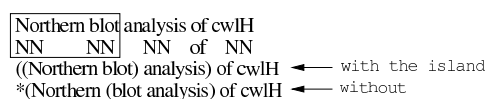


Fig. 1. Effect of an island on parsing

YATEA allows exogenous disambiguation, *i.e.* the exploitation of existing (testified) terminologies to assist the chunking, parsing and extraction steps.

During chunking, sequences of words corresponding to testified terms are identified. They cannot be further split or deleted. Their POS tags and lemmas can be corrected according to those associated to the testified term. If an MNP corresponds to a testified term for which a parse exists (provided by the user or computed using parsing patterns), it is recorded as a term candidate with the highest score of reliability. Similarly to endogenous disambiguation, subsequences of MNPs corresponding to testified terms are used as islands of reliability in order to augment the number and quality of parsed MNPs.

3.2 Term candidate extraction process

A noun phrase is extracted from the corpus and considered a term candidate if at least one parse is found for it. This is performed in three main steps, (1) *chunking*, *i.e.* construction of a list of Maximal Noun Phrases from the corpus, (2) *parsing*, *i.e.* attempts to find at least one syntactic parse for each MNP and, (3) *extraction* of term candidates. The result of the term extraction process is two lists of noun phrases: one contains parsed MNPs, called *term candidates*, the other contains MNPs for which no parse was found. Both lists are proposed to the user through a validation interface (ongoing development).

1. **Chunking:** the corpus is chunked into Maximal Noun Phrases.

The POS tags associated to the words of the corpus are used to delimit the MNPs according to the resources provided by the user: chunking frontiers and exceptions, forbidden structures and potentially, testified terms.

Chunking frontiers are tags or words that are not allowed to appear in MNPs, e.g. verbs (VBG) or prepositions (IN). *Chunking exceptions* are used to refine frontiers. For instance, "of" is a frontier exception to prepositions, "many" and "several" being exceptions to adjectives. *Forbidden structures* are exceptions for more complex structures and are used to prevent from extracting sequences that look like terms (syntactically valid) but are known not to be terms or parts of terms like "of course". MNPs and subparts of MNPs corresponding to testified terms (when available) are protected and cannot be modified using the chunking data. For instance, the tag FW is *a priori* not allowed in MNPs. However, if an MNP is equal to or contains the testified term "in/IN vitro/FW", it will be kept as such.

2. **Parsing:** for each identified MNP type, except monolexical MNPs, different parsing methods are applied in decreasing order of reliability. Once a method succeeds in parsing the MNP, the parsing process comes to an end. Still, one method can compute several parses for the same MNP, making the parsing non-deterministic if desired. We consider 3 different parsing methods:
 - TT-COVERED: the MNP inflected or lemmatized form corresponds to one or several combined testified terms (TT);
 - PATTERN-COVERED: the POS sequence of the (possibly simplified) MNP corresponds to a parsing pattern provided by user;
 - PROGRESSIVE: the MNP is progressively reduced at its left and right ends by the application of parsing patterns. Islands of reliability from term candidates or testified terms are also used to reduce the MNP sequence of the MNP to allow the application of parsing patterns.
3. **Extraction** of term candidates: MNPs that received a parse in the previous processing step are considered term candidates. Statistical measures will further be implemented to order MNPs according to their likelihood to be a term in order to facilitate their validation by the user.

4 Experiments

To characterise the effects of resources on term extraction, we compare the results provided by YATEA using or not existing terminologies on a biomedical corpus. We present and comment the effects on chunking, parsing and extraction of the term candidates.

4.1 Materials

Working corpus We carry out an experiment on a corpus of 16,600 sentences (438,513 words) describing genomic interaction of the model organism "*Bacillus subtilis*". The corpus was tagged and lemmatized using the GENIA tagger [12].

Terminological resources To study the reuse of terminologies in the term extractor, we tested two types of resources: terms from two public databases and a list of terms extracted from the working corpus. We first selected and

merged two specialized resources covering genomic vocabulary: Gene Ontology [13] and MeSH [14], both issued from the december 2005 release of UMLS [15]. The Gene Ontology resource³ (henceforth GO) aims at proposing a controlled vocabulary related to the genomic description of any organism, prokaryotes as well as eukaryotes [16]. GO proposes a list of 24,803 terms. The Medical Subject Headings thesaurus⁴ (henceforth MeSH) is dedicated to the indexation of the Medline database. The UMLS version of the MeSH offers 390,489 terms used in the medical domain [17].

The TAC (Terms Acquired in Corpus) resource is a list of 515 terms extracted from our working corpus using three term extractors [5]. The 515 terms occur at least 20 times in the corpus and were validated by a biologist.

4.2 Results

We present and comment the results of Y_{ATE} using no resource, the combination of GO and MeSH (GO+MeSH) and finally the TAC resource.

Chunking is affected by resources in several ways. As shown in Table 1, they allow the identification of new MNPs that were originally rejected due to their POS tag(s). In addition, the MNPs tend to be longer and monolexical terms less numerous. As MNPs are more complex, the number of types of POS sequences to be parsed is augmented. However, this increase in diversity is expected to be compensated by the parsing mechanism related to islands of reliability.

Table 1. Effects of resources on chunking

Version	MNPs		Monolexical		Words/ complex MNP	POS sequences types
	types	occ	types	occ		
no resource	45,716	84,810	6,989	30,815	3.61	2,965
GO+MeSH	46,079	85,004	6,949	30,272	3.63	3,256
TAC	46,315	84,918	6,934	29,695	3.65	3,500

Parsing MNPs is also affected by the use of resources that increase the reliability of parses since testified terms are used as islands of reliability. The contribution of each parsing method is presented in Figure 2 regarding the total types and occurrences of MNPs. Interestingly, the TT-COVERED method discharges the PROGRESSIVE method which is the least reliable. The increase in the contribution of the PATTERN-COVERED method is explained by the extraction of new short terms like species names, e.g. "*Escherichia/FW coli/FW*", the expansion of monolexical terms like "*DNA/NNP*" to "*DNA/NNP binding/VBG*"

³ <http://www.geneontology.org/>

⁴ <http://www.nlm.nih.gov/mesh>

that results from tag correction (VBG replaced by NN) and the simplification of MNPs using islands before the application of the parsing patterns.

The comparison of the diagrams on types and occurrences shows that both resources cover frequent terms. Still, GO+MeSH unsurprisingly contributes little (1777 terms out of 415,292 are used) compared to the corpus-tuned resource (TAC).

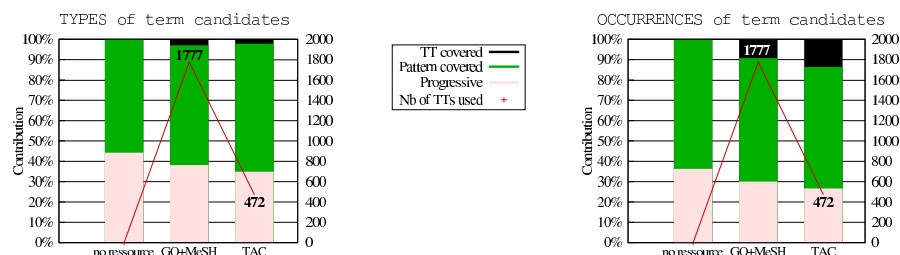


Fig. 2. Contribution of parsing methods

Extraction of term candidates is dependent on both preceding steps as an MNP found during chunking is considered a term candidate if at least one parse is found for it. Statistical filtering methods, that will be further implemented, are expected to provide qualitative information on term candidates and to allow the extraction of monolexical terms. On a quantitative point of view, using existing terminologies results in the extraction of a greater number of term candidates.

5 Conclusion and future works

Term extractors on the one hand and terminology matching techniques on the other hand show limitations in term acquisition and term exploitation respectively. To both reduce noisy results of the extraction and augment the coverage of existing terminologies, we proposed to combine both techniques in a term extractor. With a first experiment on a biomedical corpus, we showed that the exploitation of existing terminologies in a term extractor positively influences the identification of maximal noun phrases, their parsing and finally the extraction of lists of term candidates. The result of the extraction is a corpus-tuned list of term candidates. It is composed of a subset of the external resource(s) augmented with term candidates acquired in the corpus in conformity with the former. As future works, we intend to add statistical features to assist the endogenous and exogenous disambiguation. The handling of coordinations is also about to be integrated. Finally, a precise evaluation of the outputs of \LaTeX through a validation interface is planned.

References

1. Daille, B.: Conceptual structuring through term variations. In Bond, F., Kohonen, A., Carthy, D.M., Villaciencio, A., eds.: *Proceedings of the ACL'2003 Workshop on Multiword Expressions: Analysis, Acquisition, and Treatment*. (2003) 9–16
2. Bourigault, D.: An endogeneous corpus-based method for structural noun phrase disambiguation. In: *Proceedings of the EACL'93, Utrecht, The Netherlands (1993)* 81–86
3. Bourigault, D., Fabre, C.: Approche linguistique pour l'analyse syntaxique de corpus. *Cahiers de Grammaire* (25) (2000) 131–151
4. Cabré, M.T., Estopà, R., Vivaldi, J.: Automatic term detection: a review of current systems. In: *Recent Advances in Computational Terminology*. John Benjamins, Amsterdam, Philadelphia (2001)
5. Aubin, S.: Recommandations sur l'utilisation des outils terminologiques. Technical report, *Projet ExtraPloDocs (2003)* <http://www-lipn.univ-paris13.fr/~poibeau/Extra/D31b.pdf>.
6. Chute, C.G., Cohn, S.P., Campbell, K.E., Olivier, D.E., Campbell, J.R.: The content coverage of clinical classifications. *Journal of American Medical Informatics Association* **3** (1996) 224–233
7. McCray, A.T., Browne, A.C., Bodenreider, O.: The lexical properties of the gene ontology (GO). In: *Proceedings of the AMIA 2002 Annual Symposium*. (2002) 504–508
8. Bodenreider, O., Rindfleisch, T.C., Burgun, A.: Unsupervised, corpus-based method for extending a biomedical terminology. In: *Workshop on Natural Language Processing in the Biomedical Domain (ACL2002)*. (2002) 53–60
9. Hamon, T.: Indexer les documents spécialisés : les ressources terminologiques contrôlées sont-elles suffisantes ? In: *6^{ème} rencontres Terminologie et Intelligence Artificielle*, Rouen, France (2005) 71–82
10. Enguehard, C., Malvache, P., Trigano, P.: Indexation de textes : l'apprentissage des concepts. In: *Proceedings of COLING'92, Nantes, France (1992)* 1197–1202
11. Jacquemin, C., Klavans, J.L., Tzoukermann, E.: Expansion of multi-word terms for indexing and retrieval using morphology and syntax. In: *Proceedings of the ACL'97/EACL'97, Barcelona, Spain (1997)* 24–31
12. Tsuruoka, Y., Tateishi, Y., Kim, J.D., Ohta, T., McNaught, J., Ananiadou, S., Tsujii, J.: Developing a robust part-of-speech tagger for biomedical text. In: *Proceedings of Advances in Informatics - 10th Panhellenic Conference on Informatics*. LNCS 3746 (2005) 382–392
13. Consortium, T.G.O.: Gene ontology: tool for the unification of biology. *Nature genetics* **25** (2000) 25–29
14. MeSH: Medical subject headings. Library of Medicine, Bethesda, Maryland, WWW page <http://www.nlm.nih.gov/mesh/meshhome.html>, (1998)
15. National Library of Medicine, ed.: *UMLS Knowledge Source*. 13th edn. (2003)
16. Consortium, T.G.O.: Creating the Gene Ontology Resource: Design and Implementation. *Genome Res.* **11**(8) (2001) 1425–1433
17. Côté, R.A.: *Répertoire d'anatomopathologie de la SNOMED internationale, v3.4*. Université de Sherbrooke, Sherbrooke, Québec. (1996)