

# TermitUp: Generation and Enrichment of Linked Terminologies

Patricia Martín-Chozas<sup>a,\*</sup>, Karen Vázquez-Flores<sup>a</sup>, Pablo Calleja<sup>a</sup>, Elena Montiel-Ponsoda<sup>a</sup> and Víctor Rodríguez-Doncel<sup>a</sup>

<sup>a</sup> *Ontology Engineering Group, Universidad Politécnica de Madrid, Madrid, Spain*

*E-mails: pmchozas@fi.upm.es, kvazquez@delicias.dia.fi.upm.es, pcalleja@fi.upm.es, emontiel@fi.upm.es, vrodriguez@fi.upm.es*

**Editors:** Julia Bosque-Gil, University of Zaragoza, Spain; Milan Dojchinovski, Czech Technical University in Prague, Czech Republic; Philipp Cimiano, Bielefeld University, Germany

**Solicited reviews:** John McCrae, NUI Galway, Ireland; Andon Tchechmedjiev, IMT Mines Alès, France; Four Anonymous Reviewers

**Abstract.** Domain-specific terminologies play a central role in many language technology solutions. Substantial manual effort is still involved in the creation of such resources, and many of them are published in proprietary formats that cannot be easily reused in other applications. Automatic term extraction tools help alleviate this cumbersome task. However, their results are usually in the form of plain lists of terms or as unstructured data with limited linguistic information. Initiatives such as the *Linguistic Linked Open Data cloud (LLOD)* foster the publication of language resources in open structured formats, specifically RDF, and their linking to other resources on the Web of Data. In order to leverage the wealth of linguistic data in the *LLOD* and speed up the creation of linked terminological resources, we propose TermitUp, a service that generates enriched domain specific terminologies directly from corpora, and publishes them in open and structured formats. TermitUp is composed of five modules performing terminology extraction, terminology post-processing, terminology enrichment, term relation validation and RDF publication. As part of the pipeline implemented by this service, existing resources in the *LLOD* are linked with the resulting terminologies, contributing in this way to the population of the *LLOD* cloud. TermitUp has been used in the framework of European projects tackling different fields, such as the legal domain, with promising results. Different alternatives on how to model enriched terminologies are considered and good practices illustrated with examples are proposed.

**Keywords:** Terminology Generation, Terminology Enrichment, Linguistic Linked Data, Multilingualism

## 1. Introduction

International institutions have become major producers of *multilingual terminology databases*, understood as resources that account for the specialised words used in a particular field in multiple languages. Since its foundation, the European Union has maintained initiatives to cater for the collection, maintenance and creation of terminologies, thesauri or vocabularies, to cover their internal communication needs and to support translators. Some of the best known

resources are available from TermCoord<sup>1</sup> (*Terminology Coordination Unit of the European Parliament*), in charge of the interinstitutional terminology database IATE<sup>2</sup> (*InterActive Terminology for Europe*) since 2004, or the EU Vocabularies site<sup>3</sup>, maintained by the Publications Office, that is also in charge of the upkeep of the multilingual thesaurus EuroVoc<sup>4</sup>.

The creation and curation of such vocabularies has not only supported translators, documentalists and le-

---

\* Corresponding author. E-mail: pmchozas@fi.upm.es.

<sup>1</sup><https://termcoord.eu/>

<sup>2</sup><https://iate.europa.eu/>

<sup>3</sup><https://op.europa.eu/en/web/eu-vocabularies>

<sup>4</sup><http://eurovoc.europa.eu/>

gal drafters at EU institutions, but has also become a reference for translators and language professionals outside the EU. Nowadays, curated language resources have proven to be more relevant than ever in light of natural language processing (NLP) tasks that rely on sound linguistic data. For example, query expansion using WordNet<sup>5</sup>, the well-known English lexicon [1], disambiguation based on BabelNet<sup>6</sup>, a multilingual encyclopedic dictionary [2] and text classification applying DBpedia<sup>7</sup>, the semantically structured version of the Wikipedia [3], to mention but a few.

Initiatives such as the *Linguistic Linked Open Data cloud*<sup>8</sup> (henceforward *LLOD*) are focused on collecting and publishing language resources in Semantic Web formats according to the Linked Data principles [4]. When developing NLP services, one of the main challenges is finding language resources on a certain subject area with acceptable quality and ready to be reused, as revealed, for example, in previous experiments on summarisation or machine translation enhanced with terminological resources [5] [6] [7]. Consequently, our motivating scenario is focused on assisting users with different backgrounds and expertise levels to face language and related needs (see Figure 1).

In addition, with the surge in technology solutions for the legal domain, in what is called LegalTech or RegTech, such challenges have become even bigger, since resources of this sort tend to be scarce, private to companies, published in unstructured formats, or no longer available (e.g. the legal multilingual WordNets built in the LOIS project [8], the LexALP term bank on spatial planning and sustainable development [9], or the European legal taxonomy syllabus on consumer protection law [10]). From those resources that have open licenses, such as EuroVoc, most have a wider scope and do not exhaustively cover a specific area of law, or, on the contrary, may only cover a particular sub-area of law (such as the International Labour Organisation Thesaurus<sup>9</sup>); and others are only available in one language or language pair (see abundant examples of terminologies in EuroTermBank<sup>10</sup> project, now eTranslation TermBank and the Wolters Kluwer Thesaurus of Labour Law in German<sup>11</sup>). Therefore, though

highly valuable, these resources share some common drawbacks: they usually fall short of covering the specific terminological needs of a certain project or company, are not in the languages of interest, cannot be easily reused or integrated in a new application, and are sometimes only available under request.

With the aim of palliating the need for multilingual terminological resources of a specific domain or project, leveraging resources already available in the *LLOD*, we have devised a method to automatically cover the whole life cycle of the terminology creation process. Our contribution, *TermitUp*, puts together pieces of language technology previously isolated, and improves them to build a pipeline that, taking as input a domain specific corpus in one language, generates as output a multilingual terminology semantically enriched with data from the *LLOD* and published in open formats. The specific subprocesses of the method proposed include terminology extraction, terminology postprocessing, terminology enrichment, relation validation and RDF publication.

Henceforth, the paper is structured as follows: Section 2 presents relevant previous work; Section 3 exposes the linguistic foundations supporting the development of TermitUp; Section 4 lists the application requirements; Section 5 describes every component of TermitUp architecture; Section 6 exposes its current and potential impact; Section 7 contains the discussions that have arisen throughout the development and Section 8 summarises the conclusions and future work.

## 2. Related Work

This section covers previous work related to the different processes mobilised in our system, namely, automatic terminology extraction, modern terminology management tools and semi-automatic terminology enrichment approaches (2.1). We also review existing language resources in RDF and the modelling approaches they follow (2.2).

### 2.1. Terminology-related technology

There is a wide variety of ready-to-use terminology extraction tools, both proprietary (such as SDL MultiTerm Extract<sup>12</sup>, TesauroVai<sup>13</sup> and SketchEngine<sup>14</sup>) and

<sup>5</sup><https://en-word.net/>

<sup>6</sup><https://babelnet.org/>

<sup>7</sup><https://dbpedia.org/>

<sup>8</sup><https://linguistic-lod.org/>

<sup>9</sup><https://metadata.ilo.org/thesaurus.html>

<sup>10</sup><https://www.eurotermbank.com/>

<sup>11</sup><https://joinup.ec.europa.eu/solution/wkd-thesaurus-labour-law>

<sup>12</sup><https://www.sdl.com/software-and-services/translation-software/terminology-management/sdl-multiterm/>

<sup>13</sup><https://www.dail.es/shop/en/>

<sup>14</sup><https://www.sketchengine.eu/>

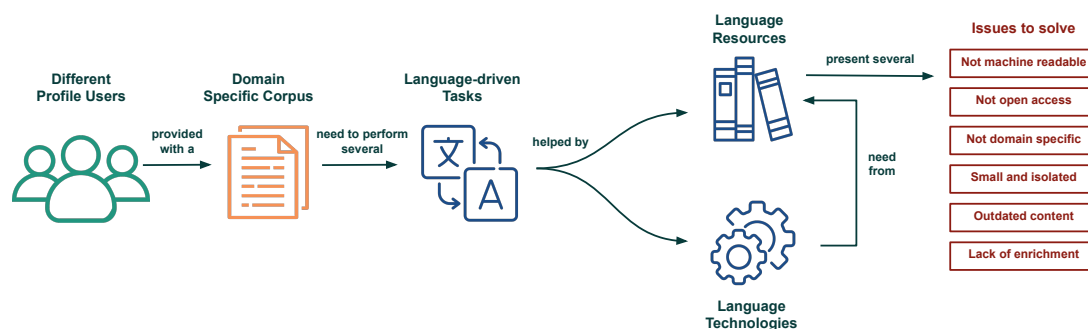


Fig. 1. Motivating scenario for the development of TermitUp.

open source (such as TermSuite<sup>15</sup>, TermoStat Web<sup>16</sup> and FiveFilters<sup>17</sup>). There are also implementations of state-of-the-art extraction algorithms, over corpora and over individual documents, such as RAKE [11], JATE [12] or TBXTools [13]. Usually, the main purpose of these tools is to generate plain lists of terms with information about their frequency in the corpus, but no additional linguistic data. Recent approaches are also trying to extract multilingual terminology across domain using transformers, which is a great step forward within the area [14].

More comprehensive terminology management tools integrate monolingual and multilingual term extraction as a starting point feature, and offer additional functionalities to enrich the extracted terms. For example, in Tilde's Terminology platform<sup>18</sup> [15], the extracted terms can be enriched with candidate translations obtained from external resources; SketchEngine<sup>19</sup> [16] identifies collocates for the extracted terms from the source corpus; PoolParty<sup>20</sup> [17] allows the manual creation of hierarchies and the manual linking to resources such as DBpedia<sup>21</sup>; Saffron<sup>22</sup> [18] suggests hierarchical relations between terms, to be afterwards supervised, and VocBench<sup>23</sup> [19] [20] allows the collaborative manual edition of vocabularies.

With regard to semi-automatic terminology enrichment, we also find several approaches in the literature. In [21], the enrichment consists of adding terms

to a source thesaurus by exploiting parallel corpora. In [22], WordNet is used to establish *hierarchical* relations between the source terms. Oliveira and Gomes [23] propose a method to automatically enrich a Portuguese thesaurus with synonyms extracted from dictionaries. Some efforts have also been devoted to further specialise the *related to* relation that is common in thesauri with specific semantic relations, as in [24]. In the reviewed works, the scope of the proposed solutions has been limited to one aspect of the terminological resource (synonyms), one external resource (WordNet), or one specific language or language pair. In any case, these efforts deal with one specific feature of the resource or for certain languages, that cannot always be easily extrapolated to other domains or languages.

## 2.2. Language resources in the Semantic Web

Concerning existing language resources published in RDF, general domain resources are the most valuable assets in the *LLOD* cloud. WordNet<sup>24</sup>, for instance, is a well known general lexicon of the English language that has been converted into RDF following the *lemon* model [25] and linked with many other resources within the cloud. BabelNet is one of the resources that exploits the linked version of WordNet. Combining Wikipedia and other resources, BabelNet constitutes a multilingual semantic network of encyclopedic and language content that covers several domains [26]. The *lemon* model was also used in the conversion of the Apertium bilingual dictionaries into RDF, a smaller but very relevant work in this area [27].

Apart from the general resources mentioned above, the *LLOD* cloud also gathers some domain specific resources. One of the most important contributions of

<sup>15</sup><http://termsuite.github.io/>

<sup>16</sup><http://termostat.ling.umontreal.ca/>

<sup>17</sup><https://www.fivefilters.org/term-extraction/>

<sup>18</sup><https://term.tilde.com/>

<sup>19</sup><https://www.sketchengine.eu/>

<sup>20</sup><https://www.poolparty.biz/>

<sup>21</sup><https://wiki.dbpedia.org/>

<sup>22</sup><https://saffron.insight-centre.org/>

<sup>23</sup><http://vocbench.uniroma2.it/>

<sup>24</sup><https://wordnet.princeton.edu/>

1 this kind is the RDF dump of IATE, an effort described  
 2 in [28]. The complete resource is available through a  
 3 Search API, but not structured in RDF. There have also  
 4 been efforts to automatically enrich these data [29]  
 5 with machine translated definitions. IATE offers trans-  
 6 lations, synonyms and definitions for terms in various  
 7 domains, but it lacks relations amongst terms.

8 Some type of term relations are, however, present in  
 9 *EuroVoc*<sup>25</sup>, which gathers data from 21 different do-  
 10 mains, with half being closely related to legal activi-  
 11 ties. Several scientific works are devoted to the con-  
 12 version of EuroVoc into RDF [30–32] and it is now  
 13 publicly available through a SPARQL Endpoint hosted  
 14 by the Publications Office. Although it is not officially  
 15 part of the *LLOD*, there are several mapping efforts  
 16 with resources in the cloud. Yet, from the point of  
 17 view of resources that can be used for NLP tasks, Eu-  
 18 roVoc is highly valuable as it contains translations,  
 19 synonyms and term relations, but lacks other types of  
 20 linguistic descriptions such as morphosyntactic infor-  
 21 mation or definitions. Also, for domain-specific NLP  
 22 tasks, frequently, the terms contained are too general,  
 23 for instance, to process specialised legal documents.  
 24 Similar issues can be encountered in related resources  
 25 such as the TheSoz Thesaurus for Social Sciences  
 26 [33] and the STW Thesaurus for Economics<sup>26</sup>, both  
 27 of them modelled according to SKOS<sup>27</sup>. Unlike Eu-  
 28 roVoc, their content is focused on one specific domain,  
 29 and can be of great help when processing legal docu-  
 30 mentation. They have, however, an additional limita-  
 31 tion: while EuroVoc contains terms in 22 languages,  
 32 TheSoz is only available for English, French, German  
 33 and Russian, and STW is bilingual (English-German).  
 34 The same issue concerns the UNESCO Thesaurus<sup>28</sup>,  
 35 which provides multidisciplinary terminology in En-  
 36 glish, French, Spanish and Russian. Finally, the In-  
 37 ternational Labour Organisation Thesaurus<sup>29</sup> collects  
 38 specific terminology for the labour law domain. Unfor-  
 39 tunately, terms are only published in English, French  
 40 and Spanish, synonyms and definitions are scarce, and  
 41 data is only available under request.

42 In summary, to ease the creation of terminological  
 43 resources, we can make use of state-of-the-art termi-  
 44 nology extraction tools, although only a few of them  
 45 provide additional linguistic or semantic data to fur-  
 46

1 ther describe the terms. To relieve this situation, there  
 2 have been some approaches pursuing automatic termi-  
 3 nology enrichment, yet, they are targeted at one spe-  
 4 cific type of information, and most of them involve  
 5 manual efforts. In this paper, we present TermitUp, an  
 6 automatic approach to generate Multilingual Seman-  
 7 tically Enriched Legal Terminologies from corpora in  
 8 Semantic Web formats. With TermitUp, terms are au-  
 9 tomatically enriched with translations, term variants  
 10 or synonyms, definitions, examples of use, informa-  
 11 tion about frequency and hierarchical relations, and are  
 12 linked with other resources in the *LLOD* cloud.

### 13 3. Theoretical Underpinnings 14

15 The pipeline implemented by TermitUp is in line  
 16 with the terminographical methods proposed by well-  
 17 established Terminology theories for the compilation  
 18 of terminological resources (communicative theory of  
 19 terminology [34], socioterminology [35], sociocogni-  
 20 tive theory of terminology [36] or frame-based theory  
 21 [37]). In the most common scenario, the starting point  
 22 in a terminological work is a corpus of specialised  
 23 texts. The more care taken in constructing the corpus,  
 24 the better. According to Barrière [38], texts should be  
 25 domain relevant and contain *knowledge-rich contexts*  
 26 (a notion defined by Meyer as "sentences that are of in-  
 27 terest to terminologists because they contain important  
 28 terms and *knowledge patterns*", i.e., expressions of se-  
 29 mantic relationships [39]). In our approach, the corpus  
 30 construction task is a manual task assigned to users,  
 31 who may not be so interested in the knowledge-rich  
 32 value of texts, but on the relevance of the documents  
 33 for a certain endeavour.

34 The next step consists in identifying terminological  
 35 units in those documents. These can correspond to dif-  
 36 ferent part-of-speech (noun, verb, adjective, adverb),  
 37 and participate in multi-word expressions or phraseo-  
 38 logical units. Deciding if a lexical unit has a termino-  
 39 logical status is not devoid of difficulties. To assist ter-  
 40 minologists in this step, several authors propose guide-  
 41 lines in the form of criteria that lexical units have to  
 42 satisfy to be considered terms [34] [40]. The meaning  
 43 of a unit is to be discovered in text and constructed  
 44 through relations to other terminological units. This al-  
 45 lows terminologists to derive the conceptual structure  
 46 underlying those designations, which enables transla-  
 47 tors or any other language professionals (documental-  
 48 ists, technical writers, subject specialists, etc.) to un-  
 49 derstand an area of knowledge. Such a structure can  
 50  
 51

25 <https://publications.europa.eu/en/web/eu-vocabularies>

26 <https://zbw.eu/stw>

27 <https://www.w3.org/TR/skos-reference/>

28 <http://vocabularies.unesco.org/browser/thesaurus>

29 <https://metadata.ilo.org/thesaurus/>

1 take the form of an ontology, as suggested in [37], and  
2 is the approach taken by the so called *terminological*  
3 *knowledge bases*, as dubbed in [41], in which a knowl-  
4 edge base component is enriched with a linguistic (termi-  
5 nological) component. Some well-known examples  
6 of terminological knowledge bases in different areas  
7 are GENOMA-KB [42], OncoTerm [43] or EcoLexi-  
8 con [44].

9 These theories also propound that terms are to be  
10 analysed as used in real communication by experts in  
11 the domain, and that this may result in identifying var-  
12 ious forms of designations (synonyms or term vari-  
13 ants). Variants are to be accounted for in terminolog-  
14 ical resources, as well as the causes for that varia-  
15 tion [45]. Depending on the purpose of the resource at  
16 hand, additional linguistic descriptions are also com-  
17 mon in terminological resources, namely, source of the  
18 term, morphosyntactic information, definition, refer-  
19 ences to other terms (which can be of different na-  
20 ture, e.g. synonyms, hyponyms, antonyms), usage con-  
21 texts (that show how the term behaves in real texts), us-  
22 age notes, or phraseology. Terms are usually assigned  
23 to a domain, and all sources from which information  
24 has been obtained are referenced, together with other  
25 metadata (author, date, reliability degree, etc.).

26 When considering the multilingual perspective, best  
27 practices in terminology work recommend that equiv-  
28 alents in other languages are also collected from  
29 domain-specific corpus in the languages of interest, as  
30 well as the rest of linguistic descriptions [34]. An exact  
31 equivalence relation is assumed when terms in mul-  
32 tiple languages are related to a source term, although  
33 language and culture differences may be captured in  
34 the form of notes. However, previous works on mul-  
35 tilingual terminological knowledge bases in the legal  
36 domain show how important it is to define culture-  
37 specific knowledge as intermediate representations as-  
38 sociated with a common shared ontology [46].

39 Finally, we briefly refer to the theoretical stud-  
40 ies (and practical applications) made by terminolo-  
41 gists about terminological or conceptual relationships  
42 between terms. Conscious of the importance of ac-  
43 counting for such relationships in termbanks, termino-  
44 graphers have characterised them, studied them in  
45 particular domains, and created methods for iden-  
46 tifying them in corpora. The most important rela-  
47 tions in this regard are the so-called hierarchical re-  
48 lationships (hyperonymy-hyponymy and meronymy).  
49 However, several non-hierarchical relationships have  
50 been intensively studied in some particular domains  
51 (cause-effect, entity-function), and others have also

1 been considered for inclusion in terminological re-  
2 sources (antonymy, synonymy, derivative relationships,  
3 co-occurrences and collocations). For a nice overview  
4 we refer the interested reader to [47].

## 4. Requirements

9 The development of the first version of TermitUp  
10 was guided by a set of requirements derived from  
11 the study of existing language technologies, specifi-  
12 cally those that deal with terminology, and the obser-  
13 vation of their results, as well as from numerous dis-  
14 cussions between the linguists, computer scientists and  
15 researchers involved in this project.

16 **R1. Enrichment.** When confronted with domain  
17 specific data, there is a need for identifying the specific  
18 terms used in texts, as well as their meaning. Plain lists  
19 of terms tend not to suffice if they are to be used for  
20 annotation, classification or disambiguation and other  
21 complex NLP tasks. Definitions, morpho-syntactic in-  
22 formation, term variants and explicit relations amongst  
23 terms can significantly contribute to improving perfor-  
24 mance of subsequent text processing tasks.

25 **R2. Multilingualism.** As already mentioned, inter-  
26 national institutions have catered for the creation of  
27 multilingual terminologies or thesauri to meet their  
28 needs. However, these do not necessarily cover the  
29 needs of a company or project in terms of languages,  
30 or the purposes of the system being developed. This re-  
31 sults in the need for systems that assist in the creation  
32 of ad-hoc terminologies for certain language combina-  
33 tions. There have been some initial attempts to devel-  
34 oping terminology extractors that work on multilingual  
35 corpora, but results are still preliminary.

36 **R3. Disambiguation.** Although traditional theories  
37 to terminology and language planning have backed the  
38 approach that the terms in a domain are unambiguous,  
39 unique and semantically precise, corpus-based termi-  
40 nology studies have demonstrated that term variation  
41 or synonymy is common also in domain specific ar-  
42 eas, and that texts may also vary in the degree of speci-  
43 ficity. Additionally, external language resources (see  
44 Requirement 4) may contain different senses of a term,  
45 since they are usually of a general character rather than  
46 domain specific. This translates to a necessity for a  
47 disambiguation step when matching corpus-extracted  
48 terms with external ones.

49 **R4. Reusability and Standardisation.** Knowledge  
50 reuse is the cornerstone of Linked Open Data [4] and  
51 the main goal of TermitUp. To meet this objective, this

1 service extracts knowledge from existing resources in  
2 the *LLOD* cloud and publishes the resulting terminologies  
3 in a structured and open-licensed manner, agreed  
4 by the community, so they can be freely reused.

5 **R5. Data provenance.** When working with texts  
6 from a specific domain, it is of utmost importance to  
7 guarantee the univocity of the terms managed. There-  
8 fore, knowing the source from which each term has  
9 been extracted is equally essential, since by knowing  
10 these sources, the final user of the terminology has the  
11 freedom to choose which term to use depending on the  
12 confidence level of such sources. Taking into account  
13 that we are dealing with terminologies enriched with  
14 heterogeneous external resources, we must maintain  
15 traceability not only of the terms themselves, but of  
16 each piece of information associated with them: syn-  
17 onyms, translations, definitions, usage examples, etc.

18 **R6. Open source and easy access.** Following the  
19 philosophy of Linked Open Data, we highlight open  
20 source as one of the requirements for this service. All  
21 the code used will be openly exposed in a Github  
22 repository to allow collaboration between users and  
23 developers. In addition, TermitUp will be published as  
24 a web service for easy integration with other processes.

25 Throughout this paper, we describe TermitUp func-  
26 tionalities and expose how their specific features com-  
27 ply with each of the requirements above mentioned.

## 30 5. TermitUp Architecture

31  
32 With the aim of satisfying the requirements spelled  
33 out in the previous section, we present TermitUp, a ser-  
34 vice to generate domain specific terminologies directly  
35 from corpus, enriched with disambiguated terminolog-  
36 ical data from existing language resources in the *LLOD*  
37 cloud. This section describes the five interdependent  
38 modules that compose TermitUp architecture.

### 39 5.1. Module 1: Terminology Extraction

40  
41 This module allows to obtain a list of the most rep-  
42 resentative terms from a given corpus. After analysing  
43 and testing several open source automatic term ex-  
44 traction (ATE) tools, and also proprietary software, as  
45 mentioned in Section 2, we chose to implement the  
46 TBXTools service<sup>30</sup> [48]. TBXTools is a *fast and flexi-*  
47 *ble* tool that offers statistical and linguistic approaches  
48 to term extraction. In addition, it is published as a

51 <sup>30</sup><https://sourceforge.net/projects/tbxtools/files/>

1 Python library that we could easily implement and  
2 modify to satisfy our specific needs (i.e. language and  
3 maximum number of tokens per term). The part-of-  
4 speech tagging in the linguistic approach is supported  
5 by Freeling<sup>31</sup>. However, the performance of the tag-  
6 ger in a preliminary testing phase was not satisfactory  
7 compared to other state-of-the-art part-of-speech tag-  
8 gers for Spanish: the application is developed in C++  
9 and its implementation is very time-consuming. More-  
10 over, the results obtained by the statistical method were  
11 of good quality, and we decided to rely on the statisti-  
12 cal method only.

13 Originally, TBXTools is intended to process English  
14 texts but we fine-tuned the tool to work with Spanish  
15 texts (a need arose from our use cases, Requirement  
16 2). We added lists of Spanish stop words and a set of  
17 exclusion regular expressions to avoid noisy construc-  
18 tions, which can be consulted in the repository<sup>32</sup>.

### 19 5.2. Module 2: Terminology Post-processing

20  
21 Regardless of previously mentioned improvements,  
22 we manually reviewed the automatically raw extracted  
23 terms and noticed recurrent patterns in Spanish that did  
24 not correspond to any multi-word term. For this pur-  
25 pose, we relied on some works that have studied the  
26 most common structure of terms in English and Span-  
27 ish, specifically in the legal domain [34] [49] [50].

28 Traditionally, nouns were considered the main parts  
29 of speech to be included in terminological resources  
30 [51], since their main purpose was to label concepts.  
31 However, linguistic approaches to terminology argue  
32 that terms can belong to different parts of speech  
33 (nouns, verbs, adjectives, and sometimes adverbs), of-  
34 ten with closely related meanings (for instance, the  
35 verb *to contract* and the noun *contract*) [40].

36 With the objective of filtering common term patters  
37 from invalid structures, we designed a post-processing  
38 stage in which a *terminology filtering algorithm* relies  
39 on part-of-speech annotations to remove structures that  
40 do not correspond to valid terms in Spanish. In this  
41 regard, a set of 42 linguistic patterns were compiled  
42 to detect what we call *non-terminological* structures.  
43 Examples of such patterns can be found in Table 1.

44 Additionally, we also implemented Añotador<sup>33</sup> [52],  
45 a service to identify dates and temporal expressions, so

46 <sup>31</sup><http://nlp.lsi.upc.edu/freeling/>

47 <sup>32</sup><https://github.com/Pret-a-LLOD/termitup/tree/master/data>

48 <sup>33</sup><https://annotador.oeg.fi.upm.es/>

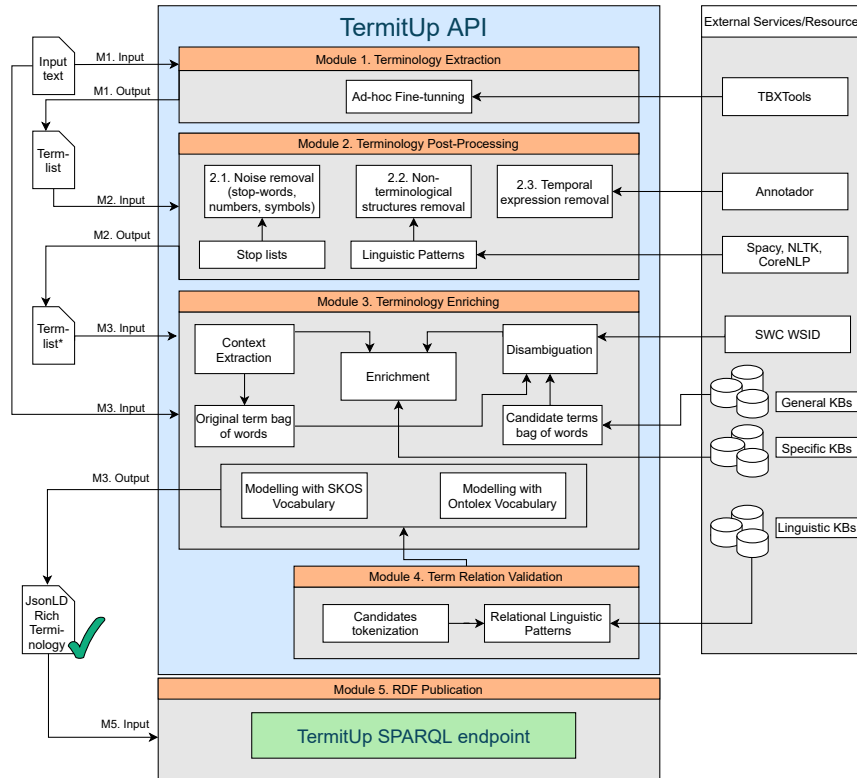


Fig. 2. TermitUp Architecture

that we could also remove them, together with some additional noisy elements.

### 5.3. Module 3: Terminology Enrichment

The next step in this approach is to take full advantage of the information in the *LLOD* relative to the previously filtered terms. Since most of the available resources have a wider scope, either covering several legal areas or general encyclopedic knowledge, a disambiguation process becomes necessary. To this end, we implemented an available word sense disambiguation (WSD) algorithm<sup>34</sup> based on BERT<sup>35</sup>.

At this point, we introduce the concept of *sense indicator*, that refers to any word in the surroundings of a term that can be used to disambiguate its meaning.

The algorithm receives as input a *source sense indicator* and several candidate *target sense indicators*

from the queried external resources. In TermitUp, the source sense indicator is built by the term  $t$  and its surrounding context (up to 100 tokens) from the input corpus  $Ct$ . For each term we retrieve up to five contexts ( $Ct_1...Ct_5$ ). The candidate target sense indicators ( $s_1...s_n$ ) consist of any information items related to target terms, such as definitions, synonyms, broader, narrower or related terms, etc.

At first, we assumed that good target sense indicators could be definitions, since definitions contain other relevant words or terms in the domain. For instance: a *training contract* is a particular type of *employment contract* drawn up between an *employer*, a *training organisation* and an *apprentice*. However, we observed that not all the accessed resources contained definitions, so we decided to take every other possible piece of information that could indicate the sense of a term: broader terms, term variants (synonyms) and domain descriptors (see Figure 3). We intentionally avoided using narrower and related terms since often they included terms from neighbouring domains that

<sup>34</sup>[https://github.com/semantic-web-company/ptlm\\_wsdl](https://github.com/semantic-web-company/ptlm_wsdl)

<sup>35</sup><https://github.com/google-research/bert>

Table 1

Examples of Spanish Non-terminological Patterns and Temporal Expressions, and their approximate translation into English for the sake of understanding.

Exclusion patterns	Examples in Spanish	Temporal Expressions in Spanish
[ADV]	simultáneamente (simultaneously)	12 de febrero (February 12th)
[ADV] + [ADJ]	inmediatamente posteriores (immediately after)	diez semanas (ten weeks)
[ADJ] + [ADV]	ininterrumpida inmediatamente (uninterrupted immediately)	quince días (fifteen days)
[NOUN] + [AUX]	partes deberán (parts shall)	nueve meses (nine months)
[NOUN] + [VERB]	consultas corresponderá (enquiries will correspond)	febrero de 2012 (February 2012)
[VERB] + [ADJ]	quedar constituida (be established)	meses siguientes (following months)
[VERB] + [NOUN]	produzcan cambios (produce changes)	
[ADJ] + [ADV] + [ADJ]	objetivas debidamente motivadas (objective duly motivated)	
[ADJ] + [SCONJ] + [ADV]	negociadora si bien (negotiating as well)	
[NOUN] + [ADV] + [ADJ]	discriminación tanto directa (discrimination both direct)	
[NOUN] + [ADV] + [SCONJ]	trabajadores siempre que (workers as long as)	
[NOUN] + [AUX] + [ADJ]	negociadora estará integrada (negotiating is integrated)	
[NOUN] + [AUX] + [VERB]	partes deberán negociar (partners should negotiate)	
[NOUN] + [VERB] + [VERB]	trabajadores podrán acordar (workers could agree)	
[VERB] + [NOUN] + [ADJ]	concurren causas económicas (economic causes concur)	

misled the algorithm. For instance, for the term *promoter*, in the sense of *a person who supports the development of a company*, we get as narrower term *DNA promoter*, as *part of the DNA that starts transcription*.

Table 2 shows an example of the five contexts for the term *hearing* obtained from the input corpus, three sense indicators built with domain descriptors from the queried resource and the resulting *weights*, returned by the WSD implementation. From these weights, the highest refers to the sense that is closest to our domain of interest. From the terms that refer to the sense in question, we can therefore establish a link and enrich our terminology with all the related information available in the queried resources, namely, definitions, translations, synonyms, broader, narrower and related terms. Through this approach, we satisfy Requirement 1: Enrichment; Requirement 2: Multilingualism; and Requirement 3: Disambiguation.

Table 3 lists the *LLOD* language resources exploited and the type of data retrieved from each of them.

#### 5.4. Module 4: Term Relation Validation

Some of the resources accessed were originally created and curated by experts. Others, however, were the result of collaborative efforts by users with different levels of expertise. This is why some of the data contained in these resources are not always correct, as it is the case of synonyms and hierarchical relations ob-

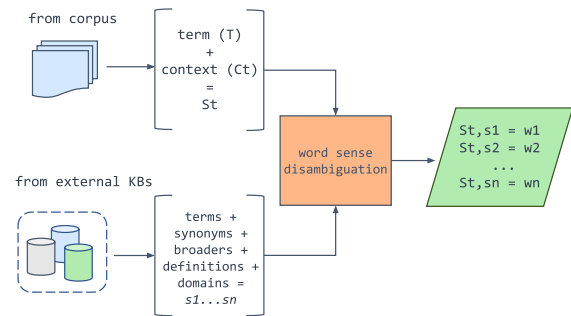


Fig. 3. Representation of the word sense disambiguation workflow

tained, for instance, from Wikidata<sup>36</sup>. The aim of this module is, thus, to check if such relations are correct. Prospective experiments to this module were already published in *LREC 2020* conference [53], where authors describe the terminology theories that support this work, approach and evaluation of the results.

This approach is inspired by the X-bar theory, that states that the formation of multi-word terms follows a hierarchical structure [34]. The approach suggests a comparison amongst the tokens of terms  $t_1$  and  $t_2$ , and the token synonyms  $s1t_1...snt_1$  and  $s1t_2...snt_2$  that are retrieved from a linguistic knowledge base. If a synonymy relation is found amongst tokens of two terms, these terms present a terminological relation. The synonyms in this approach were retrieved from

<sup>36</sup><https://www.wikidata.org/>



Table 2

WSD example for the term *hearing*, with five different contexts representing the sense of the term, and three candidate sense indicators from the queried knowledge base (IATE in this case). The results show that *s2* is the closest sense and *Ct4* the context that better represents it.

Context		Results			
<b>Ct1</b>	the difficulty of retaining the hearing date arising from the practical difficulties for the witness	<b>s1</b>	<b>s2</b>	<b>s3</b>	
<b>Ct2</b>	after consideration on the papers by Her Honour Judge Stacey, the ET hearing has since been postponed	<b>Ct1</b>	4.45	<b>6.10</b>	5.58
<b>Ct3</b>	it seems that there had been an early case management hearing on 10 April 2017	<b>Ct2</b>	7.44	<b>7.46</b>	7.02
<b>Ct4</b>	the Tribunal may order any person in Great Britain to attend a hearing to give evidence	<b>Ct3</b>	6.22	<b>7.79</b>	6.88
<b>Ct5</b>	an application for a witness order may be made at a hearing or by an application in writing to the Tribunal	<b>Ct4</b>	<b>7.17</b>	<b>7.94</b>	<b>7.82</b>
<b>Ct5</b>		<b>Ct5</b>	6.48	7.53	<b>7.73</b>
Senses					
<b>s1</b>	[hearing, parliamentary procedure]				
<b>s2</b>	[hearing, European Union law]				
<b>s3</b>	[hearing, audition, medical science]				

Table 3

List of resources exploited in the legal use case of TermitUp, and the type of information extracted from each of them. All of them are modelled in SKOS and accessed through SPARQL endpoints, except for IATE, whose RDF version is limited and outdated, and its JSON API offers more complete and up-to-date data.

Resource Name	Type of information available
<i>IATE*</i>	Translations, Synonyms, Definitions, Language Notes and Related Terms
<i>Eurovoc</i>	Translations, Synonyms, Hierarchical Relations and Related Terms
<i>UNESCO Thesaurus</i>	Translations, Synonyms, Hierarchical Relations and Related Terms
<i>International Labour Organisation Thesaurus</i>	Translations, Synonyms, Definitions, Hierarchical Relations and Related Terms
<i>STW Thesaurus</i>	Translations, Synonyms, Definitions, Hierarchical Relations and Related Terms
<i>Thesoz Thesaurus</i>	Translations, Synonyms, Definitions, Hierarchical Relations and Related Terms
<i>Wikidata</i>	Translations, Synonyms, Definitions, Hierarchical Relations and Related Terms

ConceptNet<sup>37</sup>, a large multilingual knowledge graph that brings together data from many open-domain lexical sources (DBpedia, Wiktionary and Open Multilingual WordNet, amongst others). This module can also be used to discover terminological relations amongst the initial term list (see Figure 4).

Additionally, we have implemented a set of rules based on POS-tagging and stemming to generate relations between word forms belonging to the same word family, also known as derivatives. This allows us to group word forms that belong to the same family and gather them under the same concept. Thus, every time we find two terms that follow the patterns *noun-noun*, *noun-adj*, *noun-verb*, *adj-adj*, *noun-verb* that share the same stem, we generate a *related* relation.

### 5.5. Module 5: RDF Publication

The publication in RDF of the resulting terminological data does not constitute a module of the API itself, but is part of the enrichment module (Module 3), that directly returns a list of files in JSON-LD for-

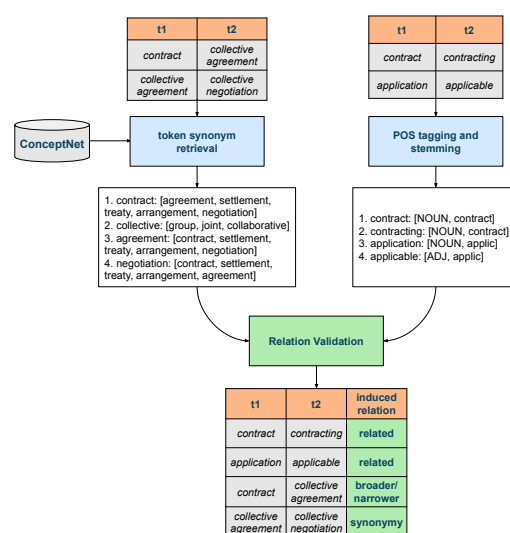


Fig. 4. Relation validation process

mat for each of the terms processed. Users can choose the vocabulary to represent such files: either SKOS or Ontolex. We consider this choice a fundamental piece of the application, because depending on the future

<sup>37</sup><http://conceptnet.io/>

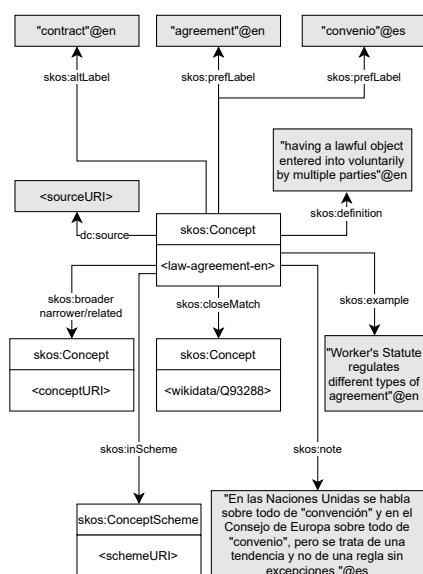


Fig. 5. Example of term modelled with SKOS

application of the terminologies, one model will be more suitable than the other. For example, if the user wants to use this terminology with a tool designed to specifically manage taxonomies, such as PoolParty or VocBench, it is necessary to represent the terminology with SKOS. If, on the contrary, the user intends to enrich the terms with morphological information, then the Ontolex model<sup>38</sup> [54] will be the most appropriate. Figures 5 and 6 exemplify the representation models followed, in which grey boxes represent literals and white boxes represent classes. Some of the white boxes are divided into two parts, where the upper part shows the name of the class and the lower contains some of the properties attached to that class.

Once the user has chosen their preferred RDF vocabulary, the publication module (Module 5) enables the publishing of the results in a Virtuoso Query SPARQL Editor<sup>39</sup> that can be subsequently accessed and queried by the user. The publication is, of course, optional, as the user may want to review the terminology before its publication. The modular architecture of TermitUp allows the human intervention at any point of the pipeline, meaning that the result of each process could be reviewed before starting the next one. In fact, in the future, we would like to developing a terminology editing platform connected to the TermitUp triple

<sup>38</sup><https://www.w3.org/2016/05/ontolex/>

<sup>39</sup><https://termitup.oeg.fi.upm.es/sparql>

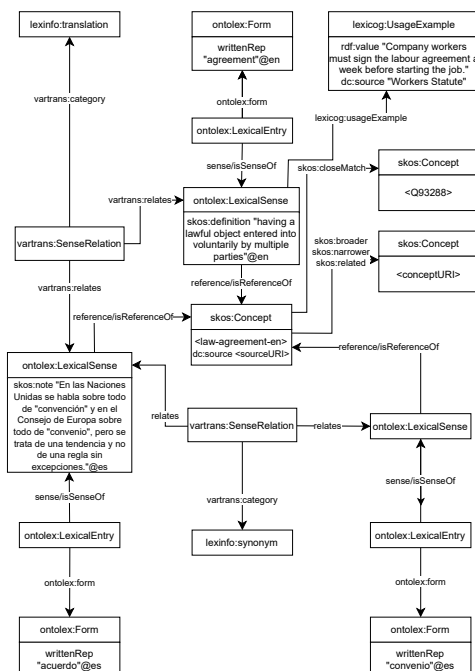


Fig. 6. Example of term modelled with Ontolex

store, that allows accessing the terminologies through a user interface, so that users can update them whenever necessary.

The combination of the exploitation of LLOD resources and publication of results in JSON-LD of Module 3, and the publication service represented by Module 5 completely satisfy Requirement 4: Reusability and Standardisation.

## 6. Impact

TermitUp has been developed in the framework of the H2020 Prêt-à-LLOD<sup>40</sup> project, whose objective is to promote the generation and adoption of linguistic technologies that reuse Linked Data. TermitUp contributes to achieving this goal by reducing the human effort necessary to create high quality, rich multilingual terminologies as linked data. In this project, with three pilots of disparate nature, spanning radically different domains such as the pharmaceutical and the e-government ones, TermitUp could be used in a number of different contexts.

<sup>40</sup><https://pret-a-llood.eu/>

The main use case of TermitUp has been in the framework of the H2020 Lynx project<sup>41</sup>, to produce a labor law terminology, with the intention of improving legal information retrieval tasks –synonyms and hyponyms being of the highest importance. This multilingual terminology (Dutch, English, German and Spanish), after a manual curation made by the domain experts, has been thus validated and published as a SKOS concept scheme. The results are accessible either through the Lynx Terminology platform<sup>42</sup> or downloadable as a static bulk dataset in Zenodo<sup>43</sup>. The main purpose is to contribute to the query expansion process implemented in the Question and Answering Module (SEAR and/or QADocservices), and for navigation purposes amongst documents in different languages. More information about the role of this terminology in the cross-lingual search pilot of Lynx can be found in this deliverable [55].

To evaluate TermitUp's enrichment we have compared this labor law terminology with a gold standard generated from the same corpus (see Table 4). In this gold standard, terms have been manually extracted, semi-automatically enriched and manually reviewed by two Spanish terminology experts. Afterwards, an expert in knowledge management from an international law firm has reviewed and validated the quality of the work. In the context of the project supported by Grupo CPonet<sup>44</sup>, TermitUp is also being used to generate a terminology on crime, where one single punishable event is referred with a surprisingly high number of forms in different language registers.

But the impact of TermitUp goes beyond these domain-specific applications. Its use as a streamlined component in composite workflows suggests a wider range of applications. TermitUp might be used to create user-specific terminologies, contribute to the linguistic analysis of a community, or create more precise vector models, with new features corresponding to the links discovered by TermitUp. In its latest application within the SmarTerp project<sup>45</sup>, TermitUp-craft terminologies will support interpreting professionals by providing them extra information on the discourse. The idea of applying TermitUp in this project is based on the analysis of interpreters' needs in terms of domain knowledge presented in [56]. This manual contains a

chapter called *Ad hoc Knowledge Acquisition in interpreting*, which explains the documentation phase prior to the interpretation process, highlighting the importance of corpora and terminologies. For this reason, TermitUp fits perfectly as a supporting tool in this documentation phase providing the interpreter with translations, synonyms and related terms to enhance their performance. TermitUp also serves as a means to improve the access and conservation of the glossaries created during the interpretation, helping solve the problem mentioned by Gile: "*Often, because of time pressure, interpreters just write down entries as they encounter them in documents or during the conference, sometimes on sheets of paper they happen to have on hand. They generally do not sort entries manually because of the time this would take. [...] most interpreters either threw away a large proportion of their glossaries prepared for specific conferences or collected them in a disorganized way and lost access to much of the information*". The application of TermitUp in SmarTerp is still preliminary, with a number of challenges related to efficiency pending to be solved, since the project just started.

TermitUp is available in a public GitHub repository<sup>46</sup>, as a Python project licensed under Apache License 2.0 terms. The functionality is also available through a HTTP REST API, thus satisfying Requirement 6. These web services are described using OpenAPI<sup>47</sup>, and they are running in web servers supported by the Prêt-à-LLOD project. A stable release of the software has also been published in the Zenodo platform<sup>48</sup>, favoring the preservation and reproducibility of the research work.

## 7. Discussion

The main limitation found during the development of this service is related to the publication of enriched terminologies in RDF, i.e., to Requirement 5. The objective of this requirement is to maintain the traceability of the data, since the provenance of the information is an essential indicator of its quality. Thus, TermitUp endeavours to store all sources of the collected data.

In the following, we analyse the different type of data collected by the service and the representation possibilities that SKOS and Ontolex offer:

<sup>41</sup><https://lynx-project.eu/>

<sup>42</sup><http://lkg.lynx-project.eu/kos>

<sup>43</sup><https://zenodo.org/communities/lynx/?page=1&size=20>

<sup>44</sup><https://www.grupocponet.com/>

<sup>45</sup><https://kunveno.digital/our-proyect/>

<sup>46</sup><https://github.com/Pret-a-LLOD/termitup>

<sup>47</sup><http://termitup.oeg.fi.upm.es/swagger>

<sup>48</sup><https://doi.org/10.5281/zenodo.4461806>

Table 4

Comparison of the enrichment numbers of the semi-automatically generated gold standard and the Labor Law terminology automatically generated with TermitUp. We are comparing five types of enrichment and the approximate generation time.

	prefLabels	altLabels	definitions	broader/narrower	related	Estimated Time
<b>Gold standard</b>	723	1229	308	398	162	~ 160 hours
<b>Labor Law Terminology</b>	710	943	264	475	272	~ 11 hours
<i>Accuracy</i>	0.982	0.767	0.857	1.193	1.679	

- *Terms, synonyms and translations*: In SKOS, they are treated as literals, represented with the properties `skos:prefLabel` and `skos:altLabel`, that do not allow to attach any additional information to them. SKOS-XL<sup>49</sup>, on the other hand, extends the model to treat these properties as classes, being able to preserve the source. In Ontolex, terms, synonyms and translations are represented as classes, allowing the representation of its source.
- *Context*: the context of a term is treated as an example of how it is used within a text. Therefore, the most suitable property to represent it in SKOS is `skos:example` (subproperty of `skos:note`<sup>50</sup>), that allows representing text strings but no additional information. In Ontolex, on the other hand, the Lexicography module [57] considers this need and introduces the `lexicog:UsageExample` class, that on the contrary, allows representing more information beyond the text itself.
- *Term note*: this is a key element of traditional terminology cards that provides additional information, such as usage recommendations and domain data. Some of the modern language resources do not use term notes anymore, but others still keep them, thus, we consider them valuable pieces of knowledge for language professionals that need to be preserved. In SKOS, term notes can be modeled with `skos:note` and in Ontolex with `ontolex:usage`, both object properties pointing to literals. This implies that if we collect term notes from different language resources, we would not be able to model their provenance.
- *Definitions*: the same occurs with definitions, since SKOS vocabulary applies `skos:definition`, that is also a subproperty of `skos:note`, therefore an object property that points to a literal. Ontolex does not propose any class for definitions either, and also employs `skos:definition`. We therefore have the same issue to model its provenance.

Besides the difficulties stated above, we face another modelling decision, since we find different types of sources at different levels. This is, the language resources with which the terms are enriched (i.e. IATE) can be understood as *intermediate sources*, that could be represented with the `schema:provider` property<sup>51</sup>. Intermediate sources are different from *original sources*, that could be either a corpus (i.e. European Legislation), an organisation (i.e. European Commission), an application (i.e. Definition Extractor) or an individual (i.e. John Doe, terminologist). For their representation, we consider properties such as `dc:source` and `dct:BibliographicResource` from DublinCore<sup>52</sup> and the classes `prov:Entity`, `prov:Agent`, `prov:Person` and `prov:Organization` from PROV ontology<sup>53</sup>.

Another discussion that arose from the modelling stage debates was whether the `skos:definition` (and related documentation properties) should be attached either to the `skos:Concept` or to the `ontolex:LexicalSense`. The SKOS specification remains vague in this point, and both approaches are at least syntactically sound – neither `skos:definition` nor its superproperty `skos:note` declare a `rdfs:domain`. This freedom suggests a flexible use which might be suitable to capture some subtleties.

First, when terminological data is transformed from different sources, definitions sometimes seem attached to concepts (e.g. data imported from Wikidata qualifies concepts), sometimes lexical senses (e.g. data imported from WordNet). We suggest the application of `skos:definition` in a flexible manner, being its subject a `skos:Concept` or a `ontolex:LexicalSense` at discretion.

Second, this loosen specification brings about the opportunity to distinguish reference and sense, in *fregean* terms. In his famous essay *Über Sinn und Bedeutung* (1892), Gottlob Frege told apart the reference and the sense of expressions [58]. In this writing, Frege uses the example of Venus: both "the morning star" and "the evening star" refer to the same object, Venus, but

<sup>51</sup><https://schema.org/provider>

<sup>52</sup><https://www.dublincore.org/specifications/dublin-core/dcmi-terms/>

<sup>53</sup><https://www.w3.org/TR/prov-0/>

<sup>49</sup><https://www.w3.org/TR/skos-reference/#x1>

<sup>50</sup><https://www.w3.org/TR/skos-reference#notes>

1 the thought they express is rather different. The sense  
2 is a mode of presentation, illuminating only a single  
3 aspect of the referent. We wonder whether computers  
4 can capture these nuances. We can certainly make such  
5 an effort, reserving the objective information about  
6 Venus for its skos:Concept (e.g. radius = 3000 km),  
7 but administer the different subjective perceptions the  
8 different components of the synset. Perhaps we want  
9 to attribute the ontolex:LexicalSense "Venus" a rela-  
10 tively neutral subjective value related to celestial bod-  
11 ies, and give the ontolex:LexicalSense "morning star"  
12 a hotter affective valence, possibly related to a more  
13 poetic context. These definitions and affective valences  
14 will be necessarily stereotypes, not reflecting subjec-  
15 tive values (which are different for each mind), but in-  
16 tersubjective, namely, reflecting common perceptions  
17 and images (we refer the reader to [59] for more infor-  
18 mation about emotional words).

19 We wonder whether personalized lemonizations will  
20 ever be possible, describing the linguistic realities of  
21 specific individuals, perhaps inferred from personal  
22 big data such as personal email inboxes or alike.  
23 But this endeavour is well beyond the scope of this  
24 paper; we only stress the opportunity of attributing  
25 skos:definition (and other triples) to skos:Concept or  
26 ontolex:LexicalSense in the most beneficial manner; in  
27 this sense, the *ontoterminology* theory may be a nice  
28 point of discussion [60].

29 We have therefore gathered such ongoing discus-  
30 sions on modelling issues in a proposal for good prac-  
31 tices to model terminological resources, published as a  
32 *Terminology* draft in the wiki of the Ontology-lexicon  
33 Community Group in the W3C<sup>54</sup>, where we expose  
34 background, motivation and use cases, and suggest  
35 complementary elements to the existing models. Such  
36 modelling modifications, naturally, need to be agreed  
37 by the community before applying them.

## 38 8. Conclusions and Future Work

39  
40  
41  
42  
43 The automation in the generation of language re-  
44 sources (specifically, terminological resources) is a  
45 challenging task still unresolved. Automatic terminol-  
46 ogy extraction and terminology management tools pro-  
47 vide a good starting point and excellent assistance both  
48 for terminology experts and language professionals,  
49 but substantial manual effort is still required.

1 This contribution intends to lighten such manual ef-  
2 forts, firstly by automating the post-processing step  
3 that terminologists usually need to perform over auto-  
4 matically extracted terms, and secondly, by exploiting  
5 the wealth of linguistic and terminological knowledge  
6 available in the *Linguistic Linked Open Data* cloud.  
7 The fact that such resources are published according to  
8 Semantic Web standards and open licences contributes  
9 to their simple and immediate integration in language  
10 technology solutions. However, the majority of them  
11 are too general, and do not contain domain-specific  
12 terms nor rich linguistic descriptions.

13 TermitUp helps covering those gaps by extracting  
14 and post-processing terms from domain specific cor-  
15 pora, and enriching them with translations, synonyms,  
16 definitions, usage notes and terminological relations.  
17 Consequently, this application establishes links to the  
18 resources exploited, contributing to the population of  
19 the *LLOD* with domain expert knowledge. Addition-  
20 ally, the tool offers a module that helps validate the  
21 terminological relations retrieved, that sometimes may  
22 be imprecise. Finally, the tool structures the resulting  
23 enriched terminologies, either following SKOS or On-  
24 tolex model; and stores them in a Virtuoso SPARQL  
25 Editor so that they can be freely accessed.

26 If we make an overall comparison with the terminology-  
27 related technology presented in Section 2.1, we can  
28 notice that TermitUp tackles some issues that they do  
29 not observe, which makes TermitUp not a competitor  
30 but a complementary application:

- 31 – Tilde’s Terminology platform extracts terms from  
32 corpus and it is able to look for translations in  
33 other resources. It, however, does not enrich with  
34 definitions, synonyms, usage contexts or rela-  
35 tions, and it returns unstructured data.
- 36 – SketchEngine is a tool specialised in corpus man-  
37 agement. It is also very well known for its termi-  
38 nology extraction algorithm. Although it gives in-  
39 formation about term co-occurrences and contex-  
40 tual information, the tool does not perform termi-  
41 nology enrichment nor semantic representation.
- 42 – PoolParty is a powerful thesaurus management  
43 tool that allows creating hierarchical relations  
44 amongst terms, representing resources in SKOS  
45 and linking them to existing ones such as DBpe-  
46 dia. Still, all the work needs to be manually per-  
47 formed through a user interface. In this case, Ter-  
48 mitUp could be used to speed up the terminology  
49 generation process and PoolParty would enable  
50 the manual revision by experts.

51 <sup>54</sup><https://www.w3.org/community/ontolex/wiki/Terminology>

- 1 – Saffron was originally a tool for taxonomy extraction. Recent improvements on the tool include terminology extraction, linking to DBpedia and knowledge graph generation. Saffron features are similar to those of TermitUp; it is however intended to work over scientific publications, and the added value is not terminology enrichment as in TermitUp, but "author and content" oriented.
- 2 – VocBench is a tool for collaborative management of ontologies and thesauri. It does not generate terminological assets, but helps curate them. As PoolParty, VocBench seems a complementary tool to manage resulting terminologies from the TermitUp workflow.

Furthermore, a remarkable technical benefit of TermitUp is that its development is open source and the community can improve, contribute to or adapt it to their specific uses cases. Also, as it is based on a REST API architecture, TermitUp can be easily integrated with other state-of-the-art technologies or tools.

On the other hand, throughout the development of the service, we have faced several modelling challenges, concretely those related to the provenance of each type of data. With the current vocabularies to model linguistic linked data, not every piece of linguistic information is represented by a class, specifically notes and definitions. This means that no additional information can be added to them, such as the resource from which they have been retrieved. As a consequence, we have discussed and proposed an improvement of the existing models and good practices to accurately structure terminological resources built from heterogeneous data sources to the W3C Ontology-Lexicon Community Group.

During this development, we have also noticed that there is room for improvement in the quality of open (language) knowledge bases available in the *LLOD* - a fact that affects the performance of services relying on them. This is due to the fact that some of the biggest resources, such as Wikidata and ConceptNet, have been semi-automatically built, and their data have not been curated. On the contrary, those manually reviewed, such as KDictionaries' RDF version [61], can only be accessed under permission. We therefore continue pursuing the publishing of high-quality language data in open formats, such as the complete version of IATE RDF, and encourage data owners to do it as well.

Regarding the publishing of the results, an immediate step is to resume the work started in the Terminoteca RDF project [62], whose objective is the

creation of a repository of multilingual terminologies. That is, to link different terminologies in a single graph so that they can be queried from a single entry point. Therefore, it seems logical that, since the objective of TermitUp is to generate rich multilingual linked terminologies, the next step would be to publish them in Terminoteca RDF, that would also allow to browse the terms through a graphic interface.

On the other hand, we observed that traditional terminological resources (such as TERMIUM and IATE) do not make explicit the relations that may exist between terms, that are to be inferred by the user from the information contained in definitions or usage notes. Terminological knowledge bases or thesauri, which follow a more conceptual approach, intend to classify concepts in a conceptual structure and include hierarchical relations (broader-narrower term relations), as well as an unidentified type of relation that flags that two terms are somehow related (see "related to" in EuroVoc or Agrovoc). Frame-semantics and other Lexicon driven approaches to terminology (see Section 3) agree on the interest of capturing terminological relations, including *domain-specific relations*, that describe how two terms interact with each other in a given area of expertise. The most generic relations include cause-effect and object-function, for instance.

Consequently, the next version of TermitUp is thought to contain an additional module that allows performing automatic domain-specific relation extraction amongst the terms in the terminology, based on the study of their behaviour in the corpus.

Finally, challenging the current domain-specific application of the tool, we have already two potential projects of very different domains, in which TermitUp will take part: 1) Authors have recently worked jointly with the DFKI research center<sup>55</sup>, on the conversion of the terminological base from the Deutsche Bahn (main railway German Company)<sup>56</sup> into Semantic Web formats. This resource lacks Spanish terminological data, and TermitUp will be used to enrich it with Spanish data on the domain. 2) Authors are also involved in a project to transform the main database of Spanish emotional words, Emofinder [63], into a knowledge graph for better access, update and conservation. In this context, TermitUp will not handle *terms* but *words* from the general domain, and it will enrich the resource mainly with translations. In the first case, TermitUp

<sup>55</sup><https://www.dfki.de/>

<sup>56</sup><https://www.bahn.de/>

will query terminological resources from the transport domain, while in the latter, it will access general purpose resources, adding some important ones such as DBpedia and BabelNet.

## Acknowledgements

This work has been supported by the European Union's Horizon 2020 research and innovation programme through the Prêt-à-LLOD<sup>57</sup> project, with grant agreement No. 825182.

## References

- [1] D. Pal, M. Mitra and K. Datta, Improving query expansion using WordNet, *Journal of the Association for Information Science and Technology* **65**(12) (2014), 2469–2478, <https://doi.org/10.1002/asi.23143>.
- [2] R. Navigli and S.P. Ponzetto, Multilingual WSD with Just a Few Lines of Code: the BabelNet API, in: *Proceedings of the ACL 2012 System Demonstrations*, Association for Computational Linguistics, Jeju Island, Korea, 2012, pp. 67–72. <https://aclanthology.org/P12-3012>.
- [3] J. Flisar and V. Podgorelec, Improving short text classification using information from DBpedia ontology, *Fundamenta Informaticae* **172**(3) (2020), 261–297, <https://doi.org/10.3233/FI-2020-1905>.
- [4] C. Bizer, T. Heath and T. Berners-Lee, Linked data: The story so far, in: *Semantic services, interoperability and web applications: emerging concepts*, IGI global, 2011, pp. 205–227, <https://doi.org/10.4018/jswis.2009081901>.
- [5] J. Vivaldi, I. Da Cunha, J.-M. Torres-Moreno and P. Velázquez-Morales, Automatic Summarization Using Terminological and Semantic Resources., in: *Proceedings of the Eleventh International Conference on Language Resources and Evaluation, LREC, 2010*. ISBN 2-9517408-6-7.
- [6] F. Schmedding, P. Klügl, D. Baehrens, C. Simon, K. Simon and K. Tomanek, EuroVoc-Based Summarization of European Case Law, in: *AI Approaches to the Complexity of Legal Systems*, Springer, 2015, pp. 205–219, [https://doi.org/10.1007/978-3-030-00178-0\\_13](https://doi.org/10.1007/978-3-030-00178-0_13).
- [7] M. Arcan, M. Turchi, S. Tonelli and P. Buitelaar, Leveraging bilingual terminology to improve machine translation in a CAT environment, *Natural Language Engineering* **23**(5) (2017), 763–788, <https://doi.org/10.1017/S1351324917000195>.
- [8] D. Tiscornia, The Lois Project: Lexical Ontologies for Legal Information Sharing, in: *Proceedings of of the V Legislative XML Workshop, European Press Academic Publishing, 2007*, pp. 189–204. ISBN 9788883980466.
- [9] V. Lyding, E. Chiochetti, G. Sérasset and F. Brunet-Manquat, The LexALP information system: Term bank and corpus for multilingual legal terminology consolidated, in: *Proceedings of the workshop on multilingual language resources and interoperability*, Association for Computational Linguistics, 2006, pp. 25–31. ISBN 9781932432824.
- [10] G. Ajani, G. Boella, L. Di Caro, L. Robaldo, L. Humphreys, S. Praduroux, P. Rossi and A. Violato, The European legal taxonomy syllabus: a multi-lingual, multi-level ontology framework to untangle the web of European legal terminology, *Applied Ontology* **11**(4) (2016), 325–375, <https://doi.org/10.3233/AO-170174>.
- [11] S. Rose, D. Engel, N. Cramer and W. Cowley, Automatic keyword extraction from individual documents, *Text mining: applications and theory* **1** (2010), 1–20, <https://doi.org/10.1002/9780470689646.ch1>.
- [12] Z. Zhang, J. Gao and F. Ciravegna, JATE2.0: Java Automatic Term Extraction with Apache Solr, in: *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, European Language Resources Association (ELRA), 2016. ISBN 978-2-9517408-9-1.
- [13] M. Vázquez and A. Oliver, Improving term candidates selection using terminological tokens, *Terminology. International Journal of Theoretical and Applied Issues in Specialized Communication* **24**(1) (2018), 122–147, <https://doi.org/10.1075/term.00016.vaz>.
- [14] C. Lang, L. Wachowiak, B. Heinisch and D. Gromann, Transforming Term Extraction: Transformer-Based Approaches to Multilingual Term Extraction Across Domains, in: *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, 2021, pp. 3607–3620, <https://doi.org/10.18653/v1%2F2021.findings-acl.316>.
- [15] T. Gornostay, Terminology management in real use, in: *Proceedings of the 5th International Conference Applied Linguistics in Science and Education*, 2010, pp. 25–26.
- [16] A. Kilgariff, V. Baisa, J. Bušta, M. Jakubiček, V. Kovář, J. Michelfeit, P. Rychlý and V. Suchomel, The Sketch Engine: ten years on, *Lexicography* **1**(1) (2014), 7–36, <https://doi.org/10.1007/s40607-014-0009-9>.
- [17] T. Schandl and A. Blumauer, PoolParty: SKOS thesaurus management utilizing linked data, in: *Extended Semantic Web Conference*, Springer, 2010, pp. 421–425, [http://dx.doi.org/10.1007/978-3-642-13489-0\\_36](http://dx.doi.org/10.1007/978-3-642-13489-0_36).
- [18] G. Bordea, P. Buitelaar and T. Polajnar, Domain-independent term extraction through domain modelling, in: *The 10th international conference on terminology and artificial intelligence (TIA 2013), Paris, France*, 10th International Conference on Terminology and Artificial Intelligence, 2013. ISBN 978-2-9174-9025-9.
- [19] A. Stellato, S. Rajbhandari, A. Turbati, M. Fiorelli, C. Caracciolo, T. Lorenzetti, J. Keizer and M.T. Paziienza, VocBench: a web application for collaborative development of multilingual thesauri, in: *European semantic web conference*, Springer, 2015, pp. 38–53, [https://doi.org/10.1007/978-3-319-18818-8\\_3](https://doi.org/10.1007/978-3-319-18818-8_3).
- [20] A. Stellato, M. Fiorelli, A. Turbati, T. Lorenzetti, W. Van Gemert, D. Dechandon, C. Laaboudi-Spoiden, A. Gerencsér, A. Waniart, E. Costetchi et al., VocBench 3: A collaborative Semantic Web editor for ontologies, thesauri and lexicons, *Semantic Web* **11**(5) (2020), 855–881, <https://doi.org/10.3233/sw-200370>.
- [21] H. Déjean, E. Gaussier, J.-M. Renders and F. Sadat, Automatic processing of multilingual medical terminology: applications to thesaurus enrichment and cross-language informa-

<sup>57</sup><https://pret-a-llod.eu/>

- tion retrieval, *Artificial Intelligence in Medicine* **33**(2) (2005), 111–124, <https://doi.org/10.1016/j.artmed.2004.07.015>.
- [22] L. Hollink, V. Malaisé and G. Schreiber, Thesaurus enrichment for query expansion in audiovisual archives, *Multimedia Tools and Applications* **49**(1) (2010), 235–257, <https://doi.org/10.1007/s11042-009-0400-y>.
- [23] H.G. Oliveira and P. Gomes, Towards the automatic enrichment of a thesaurus with information in dictionaries, *Expert Systems* **30**(4) (2013), 320–332, <https://doi.org/10.1111/essy.12029>.
- [24] Y. Wu, Enriching a thesaurus as a better question-answering tool and information retrieval aid, *Journal of Information Science* **44**(4) (2018), 512–525, <http://dx.doi.org/10.1177/0165551517706219>.
- [25] J. McCrae, C. Fellbaum and P. Cimiano, Publishing and Linking WordNet using lemon and RDF, in: *Proceedings of the 3rd Workshop on Linked Data in Linguistics*, 2014.
- [26] R. Navigli and S.P. Ponzetto, BabelNet: Building a very large multilingual semantic network, in: *Proceedings of the 48th annual meeting of the association for computational linguistics*, 2010, pp. 216–225. ISBN 978-1-932432-66-4.
- [27] J. Gracia, M. Villegas, A. Gomez-Perez and N. Bel, The aperture bilingual dictionaries on the web of data, *Semantic Web* **9**(2) (2018), 231–240, <https://doi.org/10.3233/SW-170258>.
- [28] P. Cimiano, J.P. McCrae, V. Rodríguez-Doncel, T. Gornostay, A. Gómez-Pérez, B. Siemoneit and A. Lagzdins, Linked terminologies: applying linked data principles to terminological resources, in: *Proceedings of the eLex 2015 Conference*, 2015, pp. 504–517, ISBN 978-961-93594-3-3.
- [29] M. Arcan, E. Montiel-Ponsoda, J.P. McCrae and P. Buitelaar, Automatic Enrichment of Terminological Resources: the IATE RDF Example, in: *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*, 2018, <https://doi.org/10.5281/zenodo.2599942>.
- [30] L.P. Paredes, J. Rodríguez and E.R. Azcona, Promoting government controlled vocabularies for the Semantic Web: the EUROVOC thesaurus and the CPV product classification system, *Semantic Interoperability in the European Digital Library* (2008), 111.
- [31] M.M. Martínez González, B. Pérez León, M.L. Alvite Díez et al., SKOS en la integración de conocimiento en los sistemas de información jurídica, *Actas del Taller de Trabajo Zoco'09/IISBD* **3**(6) (2009), 56.
- [32] L.A. Díez, B. Pérez-León, M. Martínez-González and D.-J.V. Blanco, Propuesta de representación del tesoro Eurovoc en SKOS para su integración en sistemas de información jurídica, *Scire: representación y organización del conocimiento* (2010), <https://doi.org/10.54886/scire.v16i2.4015>.
- [33] B. Zopilko, J. Schaible, P. Mayr and B. Mathiak, The-Soz: A SKOS representation of the thesaurus for the social sciences, *Semantic Web* **4**(3) (2013), 257–263, <https://doi.org/10.3233/SW-2012-0081>.
- [34] M.T. Cabré and M.T.C. i Castellví, *La terminología: teoría, metodología, aplicaciones*, Editorial Antártida/Empúries, 1993.
- [35] F. Gaudin, *Socioterminologie*, Publication Univ Rouen Havre, 1993, <https://doi.org/10.4000/praxematique.2188>.
- [36] R. Temmerman, *Towards new ways of terminology description: The sociocognitive-approach*, Vol. 3, John Benjamins Publishing, 2000, <http://dx.doi.org/10.4314/lex.v14i1.51442>.
- [37] P. Faber, Frames as a framework for terminology, *Handbook of terminology* **1**(14) (2015), 14–33, <http://dx.doi.org/10.1075/hot.1.02fra1>.
- [38] C. Barriere and A. Agbago, TerminWeb: a software environment for term study in rich contexts, in: *International Conference on Terminology, Standardisation and Technology Transfer (TSTT 2006)*, 2006.
- [39] I. Meyer, Extracting knowledge-rich contexts for terminography, *Recent Advances in Computational Terminology* **2** (2001), 279, <https://doi.org/10.1075/nlp.2.15mey>.
- [40] M.-C. L'Homme, *Lexical semantics for terminology: an introduction*, Vol. 20, John Benjamins Publishing Company, 2020, <https://doi.org/10.1075/tilp.20>.
- [41] I. Meyer, D. Skuce, L. Bowker and K. Eck, Towards a new generation of terminological resources: an experiment in building a terminological knowledge base, in: *COLING 1992 Volume 3: The 15th International Conference on Computational Linguistics*, 1992, <https://doi.org/10.3115/992383.992410>.
- [42] M.T. Cabré, C. Bach, R. Estopà, J. Feliu, G. Martínez and J. Vivaldi, The GENOMA-KB Project: Towards the Integration of Concepts, Terms, Textual Corpora and Entities., in: *Proceedings of the Eleventh International Conference on Language Resources and Evaluation, LREC*, 2004. ISBN 2-9517408-1-6.
- [43] P. Faber et al., ONCOTERM: sistema bilingüe de información y recursos oncológicos, *La traducción científico-técnica y la terminología en la sociedad de la información* (2002), 177, <http://dx.doi.org/10.6035/EstudisTraduccio.2002.10>.
- [44] P. Faber, P. León-Araúz and A. Reimerink, Representing environmental knowledge in EcoLexicon, in: *Languages for Specific Purposes in the Digital Era*, Springer, 2014, pp. 267–301, [http://dx.doi.org/10.1007/978-3-319-02222-2\\_13](http://dx.doi.org/10.1007/978-3-319-02222-2_13).
- [45] M.T. Cabré, J. Freixa, M. Lorente and C. Tebé, La terminología hoy: replanteamiento o diversificación, *Organon* **12**(26) (1998). ISBN 9788475964058.
- [46] K. Kerremans, R. Temmerman and J. Tummers, Representing multilingual and culture-specific knowledge in a VAT regulatory ontology: Support from the termtography method, in: *OTM Confederated International Conferences "On the Move to Meaningful Internet Systems"*, Springer, 2003, [https://doi.org/10.1007/978-3-540-39962-9\\_68](https://doi.org/10.1007/978-3-540-39962-9_68).
- [47] M.-C. L'Homme and E. Marshman, Terminological Relationships and Corpus-based Methods for Discovering them: An Assessment for Terminographers, L. Bowker (Éd.), *Lexicography, Terminology, and Translation. Text-based studies in honour of Ingrid Meyer*, University of Ottawa Press, Ottawa (2006), 67–80, <http://dx.doi.org/10.2307/j.ctt1ckpgs3.8>.
- [48] A. Oliver and M. Vázquez, TBXTools: a free, fast and flexible tool for automatic terminology extraction, in: *Proceedings of the international conference recent advances in natural language processing*, 2015, pp. 473–479. ISSN 1313-8502.
- [49] E. Alcaraz and B. Hughes, El español jurídico, *Barcelona: Ariel* (2002). ISBN 978-84-344-1872-1.
- [50] E. Alcaraz and B. Hughes, *Legal translation explained*, Routledge, 2014, <https://doi.org/10.4324/9781315760346>.
- [51] M.-C. L'Homme, What can verbs and adjectives tell us about terms?, *Terminology and Knowledge Engineering, TKE 2002. Proceedings* (2002), 28–30, <http://dx.doi.org/10.13140/2.1.4075.3927>.
- [52] M. Navas-Loro and V. Rodríguez-Doncel, Annotador: a temporal tagger for Spanish, *J. Intell. Fuzzy Syst.* **39** (2020), 1979–1991, <https://doi.org/10.3233/JIFS-179865>.



- [53] P. Martín-Chozas, S. Ahmadi and E. Montiel-Ponsoda, Defying Wikidata: Validation of terminological relations in the web of data, in: *The 12th International Conference on Language Resources and Evaluation (LREC)*, 2020. ISBN 979-10-95546-34-4.
- [54] J.P. McCrae, J. Bosque-Gil, J. Gracia, P. Buitelaar and P. Cimiano, The Ontolex-Lemon model: development and applications, in: *Proceedings of eLex 2017 conference*, 2017, pp. 19–21, ISSN 2533-5626.
- [55] P. Boil, E. Gómez and P. Calleja, Lynx D5.7 Demonstrator for pilot 3, Zenodo, 2020. doi:10.5281/zenodo.4300691.
- [56] D. Gile, *Basic concepts and models for interpreter and translator training*, Vol. 8, John Benjamins Publishing, 2009. ISBN 978-9027224323.
- [57] J. Bosque-Gil, J. Gracia and E. Montiel-Ponsoda, Towards a Module for Lexicography in OntoLex., in: *LDK Workshops*, 2017, pp. 74–84.
- [58] G. Frege, Über sinn und bedeutung, *Zeitschrift für Philosophie und philosophische Kritik* **100** (1892), 25–50. ISBN 9783150195826.
- [59] M.M. Bradley and P.J. Lang, Affective norms for English words (ANEW): Instruction manual and affective ratings, Technical Report, The center for research in psychophysiology, University of Florida, 1999.
- [60] C. Roche, M. Calberg-Challot, L. Damas and P. Rouard, Ontoterminology: A new paradigm for terminology, in: *International Conference on Knowledge Engineering and Ontology Development*, 2009, pp. 321–326, <https://doi.org/10.5220/0002330803210326>.
- [61] J. Bosque-Gil, J. Gracia, E. Montiel-Ponsoda and G. Aguado-de-Cea, Modelling multilingual lexicographic resources for the Web of Data: The K Dictionaries case, in: *GLOBALEX 2016 Lexicographic Resources for Human Language Technology Workshop Programme*, 2016, p. 65.
- [62] J. Bosque-Gil, E. Montiel-Ponsoda, J. Gracia and G. Aguado-de Cea, Terminoteca RDF: a gathering point for multilingual terminologies in Spain, in: *Proceedings of TKE 2016 the 12th International conference on Terminology and Knowledge Engineering*, 2016, pp. 136–146. ISBN 9788799917907.
- [63] I. Fraga, M. Guasch, J. Haro, I. Padrón and P. Ferré, EmoFinder: The meeting point for Spanish emotional words, *Behavior Research Methods* **50**(1) (2018), 84–93, <https://doi.org/10.3758/s13428-017-1006-3>.