

Creation and Enrichment of a Terminological Knowledge Graph in the Legal Domain^{*}

Patricia Martín-Chozas¹[0000-0002-8922-7521]

Ontology Engineering Group, Universidad Politécnica de Madrid
pmchozas@fi.upm.es
<http://oeg-upm.net>

Abstract. Domain-specific terminologies are of great use in a number of contexts, such as information retrieval from text documents or supporting humans in translation tasks. However, automated terminology extraction tools usually render plain lists with no additional information (hierarchical relations, definitions or examples of use, amongst others). The output of these tools is very often offered in non-open formats, hampering their reuse and interoperability. Moreover, terminology management tools demand a lot of manual work to curate and enrich the resources and they do not support the representation of terminological relations beyond broader/narrower. The contributions of this Thesis mitigate these problems by automating the creation of rich terminologies from plain text documents, by establishing links to external resources, and by adopting the W3C standards for the Semantic Web. The proposed method comprises six tasks: refinement, disambiguation, enrichment, relation validation, RDF conversion and relation extraction. We have applied this methodology to two different legal corpora, i.e., contracts and collective agreements. The result of this methodology will be a *Terminological Knowledge Graph* that can be exploited by different Natural Language Processing applications.

Keywords: Terminology Management · Linguistic Linked Data · Knowledge Graphs · Semantic Web.

1 Motivation

Language Resources are a remarkably valuable asset in our current multicultural and multilingual society. They are a building block in the majority of the digital media we use in our daily routines: social media, online news, audiovisual content and online shopping, to mention but a few. These activities are possible thanks to Natural Language Processing tasks such as Machine Translation, Text Annotation, Document Classification, for instance, that demand sound language resources to produce optimal results. We can find those resources all over

^{*} This work has been partially funded by the H2020 project Prêt-à-LLoD, under grant agreement No 825182, the H2020 project Lynx, under grant agreement No 780602 and a Predoctoral grant by the Consejo de Educación, Juventud y Deporte de la Comunidad de Madrid

the web, from *dictionaries* of general language to *terminologies* specialised in different domains: industry, medicine, environment, amongst others.

Why, then, the most well-known terminological resources in the legal domain are still published in physical and non machine readable formats?

These *terminological resources* lose enormous value when isolated: physical glossaries, terminologies in PDF, etc. (some examples are mentioned in Section 2). Therefore, we propose a methodology to automate terminology management processes and automatically generate interoperable resources, relying on open Semantic Web formats that allow to publish resources as Linked Data [1].

When published according to the Linked Data principles, the resources can be interlinked as machine-readable data in non-proprietary formats, giving birth to Knowledge Graphs [2] that can be used to induce information by diverse applications. Some efforts have already been devoted to this task (Section 2). However, within the legal field, we can hardly find resources online, and it is even more difficult to find them as part of the Semantic Web.

Thus, with this methodology we want to fill in the gap of linguistic legal knowledge on the web by producing sound domain-specific language resources and reusing available resources in the *Linguistic Linked Open Data* cloud¹. Throughout this document, the output of this workflow will be referred as a *Terminological Knowledge Graph* composed of *Linked Terminologies*.

The rest of the paper is structured as follows: Section 2 describes the current State of the Art regarding related tools and resources, Section 3 lists the Research Questions and Expected Contributions, Section 4 explains the proposed Methodology and Section 5 contains the initial Evaluation Plan and Conclusions.

2 State of the Art

In this section we explore current Terminology Management Approaches (2.1), Traditional Terminological Resources (2.2) and Linked Terminological Resources (2.3), that refer to terminologies in Semantic Web formats.

2.1 Terminology Management Approaches

Originally, terminology extraction has been manually performed by translators. Even with the help of Computer Assisted Translation (CAT) tools the process is not automatic: translators need to select the specific terms to be stored. For instance, the most famous CAT tool, SDL Trados Studio [3], provides a terminology management extension, MultiTerm², that allows the easy reuse, sharing and update of terminologies. However, it is a proprietary application that applies its own format (MTF.XML), which hinders the reusability by other applications. Other tools, such as GesTerm³ or the Tilde Terminology platform [4], can handle

¹ <http://linguistic-lod.org/lod-cloud>

² <https://www.sdl.com/es/software-and-services/translation-software/terminology-management/sdl-multiterm/>

³ <https://www.termcat.cat/es/gestores-terminologia>

many types of file formats and even offer collaborative options. The main drawback here is the great amount of manual work that the terminology management requires, specially in huge volumes of data. On the other hand, SketchEngine [5], that can work with large corpora and identify most frequent terminology and keywords. Still, the output is a plain list of terms with no additional information nor terminological relations.

Even tools, such as the PoolParty Semantic Suite [6], that is specially designed handle language resources in Semantic Web formats and allows the creation of hierarchies involves a lot of manual efforts: terms and relations amongst them need to be individually selected by the user.

2.2 Traditional Terminological Resources

One of the most important resources of this kind, at European level, is *IATE*, the terminological database of the European Union, originally built in TBX (TermBase eXchange format). The terms contained belong to several domains and languages, covering the activities of the European Union (agriculture, politics, sociology, medicine, etc.). At a national level, *Terminesp* is also a great effort developed by the Spanish Association of Terminology⁴. It contains multilingual terms related UNE Spanish Standards that can be searched through an online portal. A more specific resource are the glossaries from the *Terminología Oberta* service developed by the Catalan Terminological Centre (*TERMCAT*)⁵, that also cover very different domains, but mainly at a regional level.

Traditional Terminological Resources in the Legal Domain. As mentioned before, some of the most valuable terminological assets in the legal domain nowadays are still published in physical formats. This is the case of *Black's Law Dictionary*, a monolingual legal dictionary widely used by translators [7].

However, the great part are published in online portals, such as the *Dudario jurídico de la ONU*⁶, developed by the Translation department of the United Nations, that gives information about the correct usage of a term in different contexts. Similarly, the *United Nation Terminology Database (UNTERM)*⁷ provides terminology and nomenclatures used in the work of the United Nations in eight different languages.

2.3 Linked Terminological Resources

A fundamental remark at this point is the distinction between “RDF Resource” and “Linked Resource”. An “RDF Resource” can be isolated, but a “Linked Resource” is published in RDF and interconnected with other resources. Thus, here we will analyse Linked Terminologies as they are the output of this work.

⁴ <http://www.aeter.org/>

⁵ <http://www.termcat.cat/en>

⁶ <https://onutraduccion.wordpress.com/pref/dudario-juridico/>

⁷ <https://unterm.un.org/UNTERM/portal/welcome>

Some of the resources mentioned in Section 2.2 are exposed as online portals but they have also been published as Linked Data:

- IATE was converted into RDF, following the *lemon* model⁸ and linked with the European Migration Network glossary [8].
- Terminesp and TERMCAT glossaries were transformed and linked generating the *TerminotecaRDF* platform [9] [10].

Linked Terminological Resources in the Legal Domain. With the aim of enriching the legal knowledge gap in the Semantic Web, some experiments have already converted and linked legal language resources. For instance, the linking of *IATE*, *Creative Common licenses*, documents from the *World Intellectual Property Organization (WIPO)* and other relevant resources [11].

Another significant effort is the publication of *EuroVoc* as Linked Data, following the SKOS vocabulary [12]. This thesaurus is maintained by the Publications Office of the European Commission and it contains a great number of terms from the legal domain. It has been linked with resources such as the *UNESCO* and the *GEMET* thesauri, amongst others. EuroVoc is also available through a *SPARQL endpoint*⁹, supported by PoolParty¹⁰.

3 Research Questions and Expected Contributions

Based on our motivation and the needs raised from the state of the art we can formulate the following research questions:

1. *How can terminology management processes be enhanced by the use of Semantic Web technologies?*
2. *Is it possible to guarantee the quality and specificity to the legal domain of the resulting terminological knowledge graph?*
3. *Which applications can benefit from terminological resources in Semantic Web formats?*

Consequently, our main expected contribution is summarised as the *Creation and Enrichment of a Terminological Knowledge Graph in the Legal Domain*. Due to the lack of legal terminological resources in the web in general and in the Semantic Web in particular, we have applied this approach to the legal field, but we propose a domain independent methodology that can be applied to other areas of knowledge. It is comprised by the following sub-contributions:

- Refinement of automatically extracted terms.
- Enrichment of such term lists with disambiguated data from external Knowledge Bases
- Identification of new terminological relations and validation of existing ones

⁸ <https://lemon-model.net/>

⁹ <https://lynx.poolparty.biz/PoolParty/sparql/Eurovoc4.3>

¹⁰ <https://www.poolparty.biz/>

4 Research Methodology and Approach

This work is the continuation of a Master Thesis, aimed at building a Linguistic Linked Open Data cloud on the Legal domain (see [13]), that served as the foundation of the current work. This work proposed a semi-automatic approach to create Legal Linked Terminologies relying on proprietary software such as SketchEngine¹¹ and also open-source applications such as OpenRefine¹². A remarkable contribution of this work was the exhaustive collection of existing legal language resources performed, that can be found here¹³. The huge amount of manual work involved in managing the datasets found with the above mentioned tools was the definite impulse to research on an automatic workflow.

The suggested approach is composed of six subtasks, as illustrated in Figure 1; some of them are ongoing and others are still pending. The base input is a corpus of documents that needs to be processed through a Terminology Extraction step. This task is out of the scope of the contribution, since it is not the goal of our research: there are already very good terminology extraction algorithms with a high performance (such as TTF-IDF or CValue). We have, however, worked on linguistic patterns to adapt an open source extraction software to the legal terminology [14]. Consequently, the input of our workflow is a raw list of terms previously extracted.

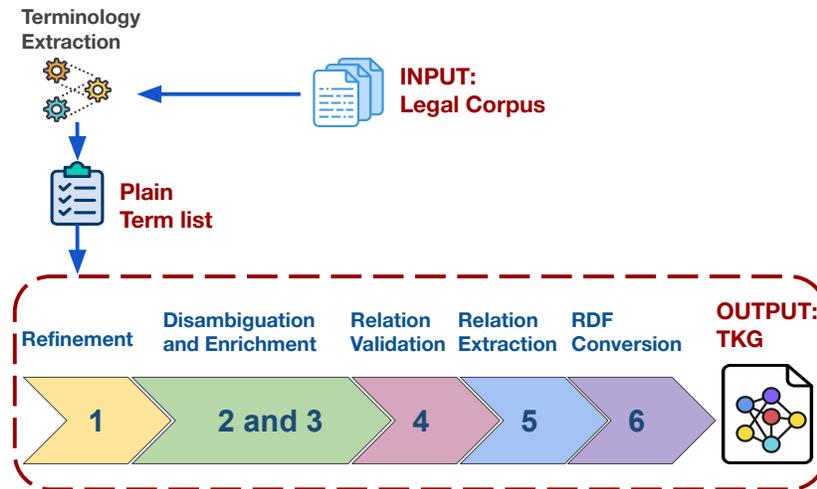


Fig. 1. Terminological Knowledge Graph Generation Workflow

¹¹ <https://www.sketchengine.eu/>

¹² <http://openrefine.org/>

¹³ <http://data.lynx-project.eu/dataset?organization=oeg>

Task 1: Refinement. After analysing the raw output of several terminology extraction tools [13], we noticed that they tend to include noisy terms that need to be filtered. Consequently, we propose a series of automatic refinement suggestions that include lemmatization, removing non terminological structures, removing duplicates, unifying caps, creating top concepts (such as “business” as the top concept of “business partner”, “business unit”, etc.) and removing Named Entities (such as “Ms Robertson”).

Tasks 2 and 3: Disambiguation and Enrichment. Once the terms lists are filtered, they can be enriched with additional information by querying external knowledge bases (IATE, Wikidata, EuroVoc, for instance). However, we first need to make sure that the source term in our terminology and the target term in the queried knowledge base refer to the same lexical sense or concept; this is, terms need to be *disambiguated*.

For this task, we are researching on disambiguation techniques based on sense embeddings, such as BERT [15]. The idea is to generate sense embeddings from the source and the target terms and compare them: if both vectors are similar, then we assume they refer to the same sense, link both terms and extract relevant information such as translations, synonyms, related terms, etc.

Task 4: Relation Validation. In previous enrichment experiments using Wikidata¹⁴ as the external knowledge base, we noticed many issues concerning the data collected under the *also known as* property. The data gathered under this property should be aliases¹⁵ (spelling variants, scientific names and nicknames) and should be categorised as *synonyms* of the source term. However, in many occasions we found broader, narrower and related terms contained under this property, so we have developed a series of axioms to verify each type (Table 1). In this step, we also need to query a second knowledge based specialised in linguistic data, such as ConceptNet¹⁶, BabelNet¹⁷ or WordNet¹⁸ (see Figure 2).

Table 1. Axioms for inducing semantic relations between alternative labels (A) of a term (T) using term synonyms (S)

Axiom	Induction
$ T = A \wedge [\forall t_j \in T, \exists! a_i \in A, t_j = a_i \vee a_i \in S_{t_j}]$	T and A are synonyms
$ T < A \wedge \forall t_j \in T, \exists a_i \in A, t_j = a_i \vee a_i \in S_{t_j}$	A is a narrower term of T
$ T > A \wedge \forall a_i \in A, \exists t_j \in T, a_i = t_j \vee a_i \in S_{t_j}$	A is a broader term of T
$\exists t_j \in T, a_i \in A, t_j = a_i \vee S_{t_j} \in A$	T and A are related

¹⁴ https://www.wikidata.org/wiki/Wikidata:Main_Page

¹⁵ <https://www.wikidata.org/wiki/Help:Aliases>

¹⁶ <http://conceptnet.io/>

¹⁷ <https://babelnet.org/>

¹⁸ <https://en-word.net/>

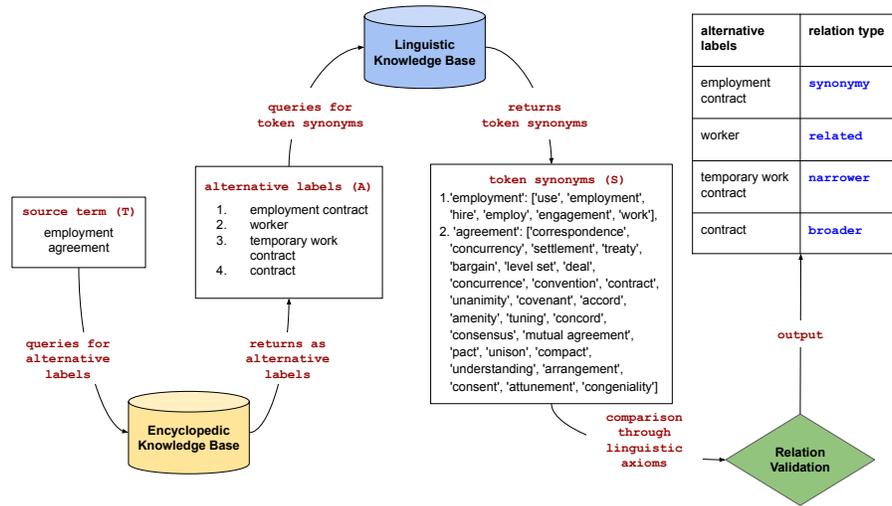


Fig. 2. Relation Validation Example

Task 5: Relation Extraction. In this task, our aim is to discover which terminological relations can be found under the *related* property assigned in the previous stage. As an example, in Table 2 we have identified terminological relations for some possible "related" terms of *employment agreement*.

Table 2. Example of Legal Terminological Relations.

Term 1	Legal Relation	Term 2
Employment Agreement	signed by	Worker
Employment Agreement	negotiated with	Company
Company	provides	Service
Worker	earns	Salary
Worker	works	Overtime

Task 6: RDF conversion. The terminologies are being represented using the SKOS vocabulary¹⁹ since it is an intuitive model whose properties can be used to represent most of the term attributes (*skos:concept*, *skos:prefLabel*, *skos:altLabel*, *skos:description*, *skos:broader*, *skos:narrower* and *skos:related*). However, we still need to research on additional RDF models to represent the properties to be extracted in Task 6 (see Figure ??).

The resulting Terminological Knowledge Graph will be serialised as JSON-LD²⁰, since it is an easy format both for human and machines to interoperate.

¹⁹ <https://www.w3.org/TR/swbp-skos-core-spec/>

²⁰ <https://json-ld.org/>

Our first conversion experiments were done by applying an ad-hoc script; however, to avoid scalability issues, we are researching on mapping language tools that interpret the RDF Mapping Language (RML [16]) and have already been successfully used to transform semi-structured data into Knowledge Graphs.

On the other hand, we are considering different ontologies, such as the PROV-O²¹ and the Web Annotation Ontology²² in order to keep track of the provenance of the data.

5 Evaluation plan and Conclusions

Evaluation It comprises one of the main complications of this work, since the most appropriate evaluation should be user based, involving the people for whom the application is intended; in this case, translation and law professionals, students and small enterprises. The issue here is that users need to evaluate the final tool, thus, middle evaluations are more difficult to perform.

In this thesis, we can find two main objects of evaluation: on the one hand, we need to evaluate the Linked Terminology Creation Workflow proposed against other terminology management applications; and in the other hand, we need to evaluate the output, this is, the Terminological Knowledge Graph. For both of them, we will keep track of the data related to the task completeness, efficiency, effectiveness and quality of the result. Additionally, we need to evaluate the maintenance of the Knowledge Graph, this is: research on how to keep the information of the graph updated during the time. An additional task devoted to this aim should be added to the pipeline.

Conclusions On the whole, this work remarks the need of 1) great improvements in terminology management workflows, 2) language resources published in Semantic Web formats and 3) specially in the legal domain. For these reasons we aim at an automatic workflow to generate this kind of resources, since, currently, even the most famous terminology management tools involve a huge amount of manual efforts. Another major drawback is that not many of those tools are intended to manage terminologies in Semantic Web formats. We have also spotted a gap on the automatic extraction of terminological relations, since most of the related work is focused on conceptual and ontological relations. Consequently, our goal is to automatically enrich the terminologies with this kind of relations as well.

References

1. T. Berners-Lee, “Design issues: Linked data (2006),” URL <http://www.w3.org/DesignIssues/LinkedData.html>, 2011.
2. L. Ehrlinger and W. Wöß, “Towards a definition of knowledge graphs.” *SEMANTiCS (Posters, Demos, SuCCESS)*, vol. 48, 2016.

²¹ <https://www.w3.org/TR/prov-o/>

²² <https://www.w3.org/ns/oa>

3. A. Walker, *SDL Trados Studio—A Practical Guide*. Packt Publishing Ltd, 2014.
4. T. Gornostay, “Terminology management in real use,” in *Proceedings of the 5th International Conference Applied Linguistics in Science and Education*, 2010, pp. 25–26.
5. A. Kilgarriff, V. Baisa, J. Bušta, M. Jakubíček, V. Kovář, J. Michelfeit, P. Rychlý, and V. Suchomel, “The sketch engine: ten years on,” *Lexicography*, vol. 1, no. 1, pp. 7–36, 2014.
6. T. Schandl and A. Blumauer, “Poolparty: Skos thesaurus management utilizing linked data,” in *Extended Semantic Web Conference*. Springer, 2010, pp. 421–425.
7. L. Biel, “Legal terminology in translation practice: dictionaries, googling or discussion forums,” *SKASE Journal of Translation and Interpretation*, vol. 3, no. 1, pp. 22–38, 2008.
8. 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,” 2015.
9. J. Bosque-Gil, E. Montiel-Ponsoda, J. Gracia, and G. Aguado-de Cea, “Terminoteca rdf: a gathering point for multilingual terminologies in spain,” *TERM BASES AND LINGUISTIC LINKED OPEN DATA*, p. 136.
10. J. Bosque-Gil, J. Gracia, and A. Gómez-Pérez, “Linked data in lexicography,” *Kernerman Dictionary News*, vol. 24, pp. 19–24, 2016.
11. V. Rodríguez-Doncel, C. Santos, P. Casanovas, and A. Gómez-Pérez, “A linked term bank of copyright-related terms.” in *JURIX*, 2015, pp. 91–100.
12. 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*, vol. 16, no. 2, pp. 47–51, 2010.
13. P. M. Chozas, E. M. Ponsoda, and V. Rodríguez-Doncel, “Towards a linked open data cloud of language resources in the legal domain,” 2017.
14. P. Martín-Chozas and P. Calleja, “Challenges of terminology extraction from legal spanish corpora,” 2018.
15. J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
16. A. Dimou, M. Vander Sande, P. Colpaert, R. Verborgh, E. Mannens, and R. Van de Walle, “Rml: A generic language for integrated rdf mappings of heterogeneous data.” *Ldow*, vol. 1184, 2014.