# OntoLT Version 1.0: Middleware for Ontology Extraction from Text

Paul Buitelaar[1], Michael Sintek[2]

[1] DFKI GmbH, Language Technology, Stuhlsatzenhausweg 3,
66123 Saarbrücken, Germany
[2] DFKI GmbH, Knowledge Management, Erwin-Schrödinger-Straße,
67608 Kaiserslautern, Germany
{paulb,sintek}@dfki.de

## *Introduction*

As human language is a primary mode of knowledge transfer, ontology learning from relevant text collections seems a viable option A typical approach in ontology learning from text first involves term extraction from a domain-specific corpus through a statistical process that determines their relevance for the domain corpus at hand. These are then clustered into groups with the purpose of identifying a taxonomy of potential classes. Other relations can be identified subsequently by computing a statistical measure of 'connectedness' between identified clusters.

The OntoLT approach described here follows a similar procedure, but we aim also at more directly connecting ontology engineering with linguistic analysis. Through the use of mapping rules between linguistic structure and ontological knowledge, linguistic knowledge (context words, morphological and syntactic structure, etc.) remains associated with the constructed ontology and may be used subsequently in its application and maintenance, e.g. in knowledge markup, ontology mapping and ontology evolution.

## *The  OntoLT Plug-In*

We describe OntoLT, a plug-in for Protégé[1] with which concepts (Protégé classes) and relations (Protégé slots) can be extracted automatically from linguistically annotated text collections – OntoLT was introduced in (Buitelaar et al., 2003; 2004). OntoLT provides mapping rules, defined by use of a precondition language that allow for a mapping between linguistic entities in text and class/slot candidates in Protégé.

Preconditions are implemented as Xpath-expressions over XML-based linguistic annotation. Predefined preconditions select for instance the predicate of a sentence, its linguistic subject or direct object. Preconditions can also be used to check certain conditions on these linguistic entities, for instance if the subject in a sentence corresponds to a particular lemma (the morphological stem of a word).

Selected linguistic entities may be used in constructing or extending an ontology. For this purpose, OntoLT provides operators to create classes or slots. According to the preconditions that were satisfied, corresponding operators will be activated to create a set of candidate classes and slots. Validated candidates can be integrated into a new or existing ontology.

OntoLT works with the following linguistic information: part-of-speech, morphological analysis, phrases (including head-modifier analysis) and predicate-argument structure (including grammatical functions like subject, direct object). The annotation allows for the automatic extraction of linguistic entities according to predefined preconditions: heads of nominal phrases, predicates and their subjects, etc. Linguistic objects that can be extracted in this way can be used to construct a shallow ontology of domain relevant concepts, sub-concepts and relations between concepts.

---

[1] http://protege.stanford.edu/

When installed[2], the OntoLT plug-in will appear as the right-most tab in the Protégé interface – as shown in Figure 1 below. OntoLT consists again of tabs with the following functionality:

- *Mappings*          To define a new mapping rule or refine an existing mapping rule
- *XPaths*            To (re)define Xpath-expressions (i.e. preconditions)
- *Corpora*           To upload one or more annotated document collections
- *CandidateView*     To extract concepts and/or relations from the annotated documents

Figure 1 shows the *Mappings* tab of OntoLT with two rules for mapping information from the linguistic annotation to potential Protégé classes and slots:

- **HeadNounToClass_ModToSubClass** maps a head-noun to a class and in combination with its modifier(s) to one or more sub-class(es)
- **SubjToClass_PredToSlot_DObjToRange** maps a linguistic subject to a class, its predicate to a corresponding slot for this class and the direct object to the "range" of this slot.
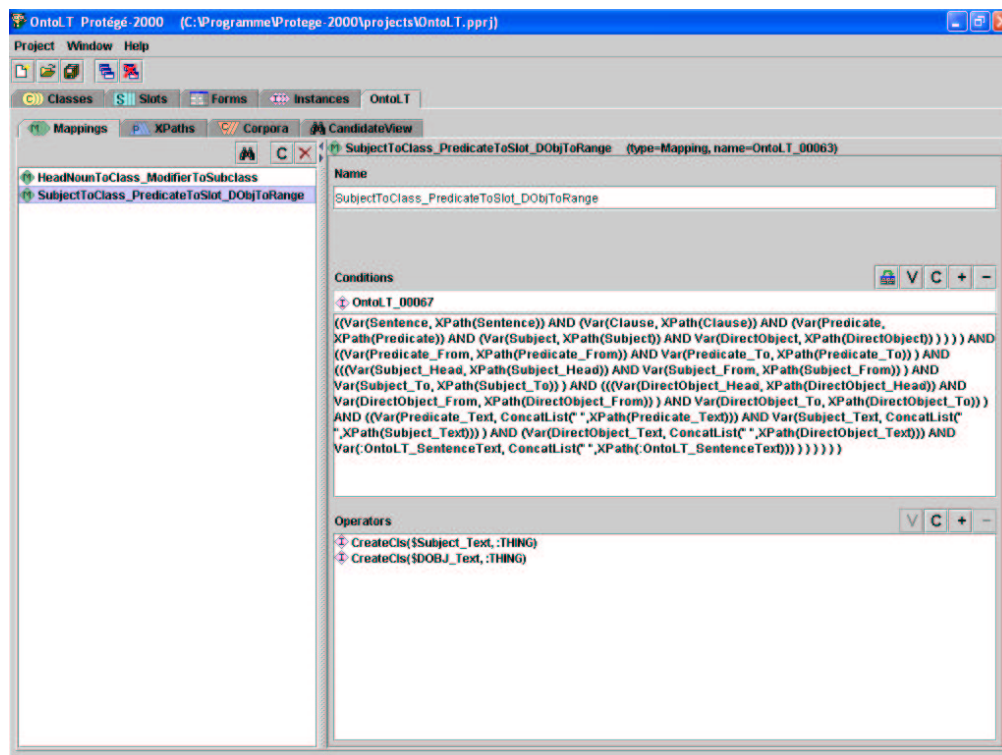


**Figure 1: The *Mappings* Tab in the OntoLT Plug-In**

## *Extract Candidates*

To extract concept and/or relation candidates, the user uploads the domain-specific text collection that they want to use for the ontology extraction task. For this purpose, a linguistically annotated version of this document collection should be uploaded in the *Corpora* tab[3].

---

[2] OntoLT is available as a plug-in for Protégé version 1.8 (http://olp.dfki.de/OntoLT/OntoLT.htm)

By running the mapping rules as defined in the *Mappings* tab over the annotated document collection, corresponding candidates will be extracted and displayed in the *CandidateView* tab – as shown in Figure 2 below. Here, the user can inspect extracted candidates for inclusion in the ontology. For this purpose, the user is provided with the linguistic context (sentence) from which the candidate was extracted, the mapping rule that was applied and the results of this. For further ease of use, extracted candidates can be organized according to frequency or in alphabetical order.
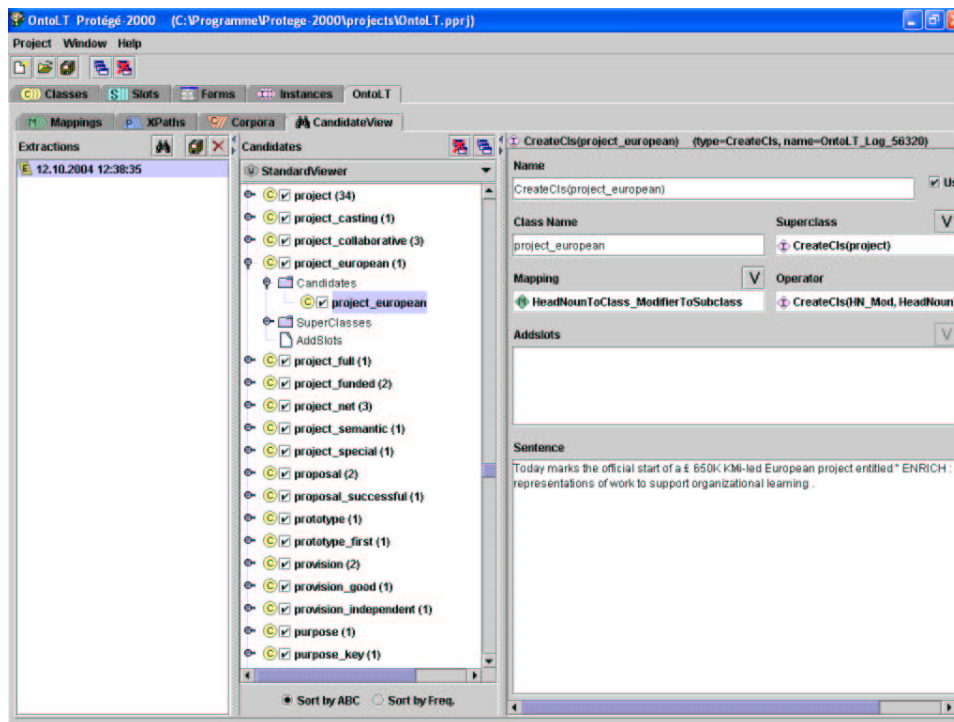


**Figure 2: Extracted Candidates in the *CandidateView* tab**

## Statistical Relevance

In ontology development for particular tasks and/or domains there is a need to focus only on those concepts and relations that are (highly) relevant for the task or domain in question. For this purpose, OntoLT includes a functionality for the computation of a statistical relevance score for extracted words by comparison of their frequencies in a (domain-specific) corpus with frequencies in a reference corpus. In this way, word use in a particular domain is contrasted with that of more general word use.

The user can use this functionality in the *Mappings* tab by selecting: the specific mapping they want to refine; the XPath element for which frequencies should be compared; the corpus to be examined; the reference corpus to be considered. Results of the statistical relevance measure are displayed in a choice box with the topmost 1% of the words pre-selected. The user can further refine this list by (de-) selection of words. After validation, selected words are included in the mapping rule as a further precondition.

---

[3] Linguistic annotation is currently provided by SCHUG, a rule-based system for German and English analysis (Declerck, 2002). SCHUG is not an integrated part of OntoLT, but could be accessed through a web service.

## *Ontology Extraction*

After refinement of the mapping rules the user can run the extraction again, extract candidate classes and slots, validate them (select or deselect them) and generate ontology fragments accordingly by pushing the right-most button in the middle pane. The extracted classes and slots will be displayed in the *Classes* tab of Protégé – as shown in Figure 3 below. Using standard Protégé functionalities, the extracted ontology fragments can be further edited and exported in one or more formats.
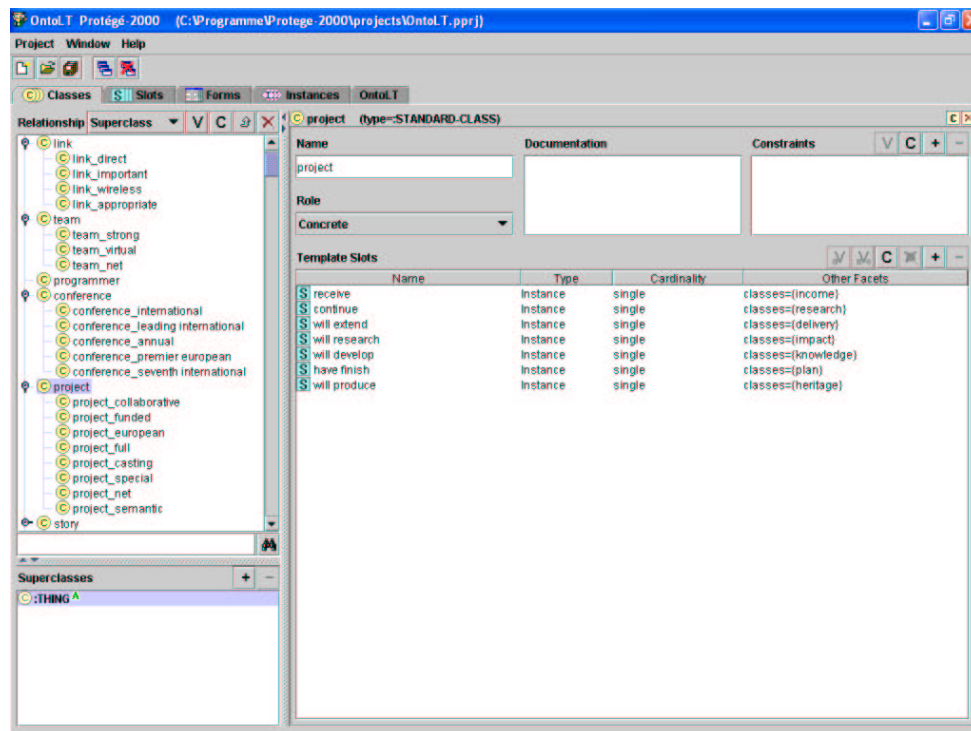


**Figure 3: Extracted Ontology Fragment**

## *References*

Buitelaar P., Olejnik D. and Sintek M. *OntoLT: A Protégé Plug-In for Ontology Extraction from Text* In: Proceedings of the Demo Session of ISWC-2003, Sanibel Island, Florida, October 2003.

Buitelaar P., Olejnik D. and Sintek M. *A Protégé Plug-In for Ontology Extraction from Text Based on Linguistic Analysis* In: Proceedings of the 1st European Semantic Web Symposium (ESWS), Heraklion, Greece, May 2004.

Declerck Th. *A set of tools for integrating linguistic and non-linguistic information.* Proceedings of the SAAKM workshop at ECAI, Lyon, 2002.

## *Acknowledgements*