

A methodology to create legal ontologies in a logic programming based web information retrieval system

José Saias and Paulo Quaresma

*Departamento de Informática,
Universidade de Évora,
7000 Évora, Portugal
jsaias|pq@di.uevora.pt*

May 21, 2004

Abstract. Web legal information retrieval systems need the capability to reason with the knowledge modeled by legal ontologies. Using this knowledge it is possible to represent and to make inferences about the semantic content of legal documents.

In this paper a methodology for applying NLP techniques to automatically create a legal ontology is proposed. The ontology is defined in the OWL semantic web language and it is used in a logic programming framework, EVOLP+ISCO, to allow users to query the semantic content of the documents. ISCO allows an easy and efficient integration of declarative, object-oriented and constraint-based programming techniques with the capability to create connections with external databases. EVOLP is a dynamic logic programming framework allowing the definition of rules for actions and events.

An application of the proposed methodology to the legal web information retrieval system of the Portuguese Attorney General's Office is described.

Keywords: Ontologies, OWL, natural language processing, logic programming

1. Introduction

Modern web legal information retrieval systems need the capability to represent and to reason with the knowledge modeled by legal ontologies. In fact, the creation of ontologies allow the definition of class hierarchies, object properties, and relation rules, such as, transitivity or functionality. Using this knowledge it is possible to represent semantic objects, to associate them with legal documents, and to make inferences about them.

OWL (Ontology Web Language) is a language proposed by the W3C consortium (<http://www.w3.org>) to be used in the "semantic-web" environment for the representation of ontologies. This language is based on the previous DAML+OIL (Darpa Agent Markup Language - (W3C, 2000)) language and it is defined using RDF (Resource Description Framework - (Lassila and Swick, 1999)).



© 2004 Kluwer Academic Publishers. Printed in the Netherlands.

In this paper a methodology to automatically create an OWL ontology from a set of legal documents is proposed. The methodology is based on the following steps:

- Definition of an initial top-level ontology;
- Identification of concepts referred in the legal documents and extraction of its properties;
- Identification of relations between the identified concepts;
- Creation of an ontology using the identified concepts and relations;
- Merge of the created ontology with the initial ontology;

In the first step, an already existent top-level legal ontology was chosen. At present we are using the legal ontology from the Portuguese Attorney General's Office, consisting of around 6,000 classes and having around 10,000 relations (Quaresma and Rodrigues, 2002). However, other top-level ontologies could be used, such as, the DOLCE proposal (Guangemi et al., 2002) or the FOLaw and LRI-core proposal (Breuker and Winkels, 2003) used in the context of the IST programs E-POWER and E-COURT (Boer et al., 2002). In the future, we expect to use the results of the e-Content project LOIS – Lexical Ontologies for legal Information Sharing, which aims to create an european-wide top-level legal ontology.

In the second step, identification of concepts and its properties, several natural language processing techniques are used, namely, a syntactical parser, and a semantic analyzer able to obtain a partial interpretation of the documents. As it will be described in detail, the semantic representation allows the identification of the set of concepts that are referred in the documents and the extraction of some of their properties.

In the third step, identification of relations between concepts, an unsupervised method for acquiring word classes and relations is used (Gamallo et al., 2002b; Gamallo et al., 2002a). This method, which has some similarities with the work of (Lame, 2003), allows the identification of related and more specific concepts (subclasses). Starting from parsed documents, a subcategorisation analysis is performed and, for each word, subcategorisation patterns are extracted. Finally, a statistical analysis is performed identifying clusters of words with similar subcategorisation patterns.

In the fourth step, the results of the previous two steps are integrated in an ontology: concepts with their properties are new classes; class hierarchies and relations are created accordingly with the statistical

analysis of subcategorization patterns (several examples will be shown in the following sections).

Finally, in the fifth step, the initial top-level ontology is merged with the new one. The proposed strategy is to search for common concepts in the two ontologies and to merge the ontologies via these concepts. New classes are inserted into the top-level ontology using the information from the semantic analyser (animal, human, action, ...).

At this stage it is important to point out that the proposed methodology is based on a bottom-up approach for the definition of the legal low-level concepts: it allows the identification of the concepts and some of the relations, but additional work will be necessary to fully integrate these concepts with the upper legal ontology. We do not intend to propose any kind of standard for legal ontologies; our aim is to define a methodology to automatically create a base ontology from a specific set of legal documents.

As referred, this work has some relations with the proposal of Lame aiming to identify components of legal ontologies from the analysis of legal texts (Lame, 2003). However, we believe our proposal has a more ambitious goal: the identified components are used to create an ontology and the initial documents are enriched with instances of this ontology. This process allows the definition of semantic web agents able to query the semantic content of these documents.

As stated before, after the creation of legal ontologies expressed in OWL, documents are enriched with instances of legal classes.

Then, a logic programming based framework is used to support inferences over the ontology. The logic programming framework is based on ISCO (Abreu, 2001) and EVOLP (Alferes et al., 2002). ISCO is a new declarative language implemented over GNU Prolog with object-oriented predicates, constraints and allowing simple connections with external databases. EVOLP is a dynamic logic programming language that is able to describe actions and events, allowing the system to make inferences about events, user intentions and beliefs and to be able to have cooperative interactions.

This logic programming framework seems to be quite adequate to represent and to make inferences over OWL ontologies. In fact, recent advances in the semantic web technology support this claim: some partners in the RuleML workgroup (<http://www.ruleml.org>) are adopting logic programming as its inference engine and there already exists a translator from RDFS to Prolog (Damásio, 2003). Moreover, in the scope of this work, a translator of a subset of OWL to Prolog was also built (correct accordingly with the OWL formal semantic description).

However, other inference engines could be chosen and used to answer queries about the legal knowledge conveyed by the documents. One

possible option might be to use the results of the Mandarax project (<http://mandarax.sourceforge.net>), which already supports RuleML.

Section 2 describes the proposed architecture. Section 3 describes the methodology for the creation of the ontology, namely, the natural language processing techniques used to process the documents. Section 4 describes the NLP techniques used to create the OWL instances associated with each document. Section 5 describes ISCO, the basic logic programming framework. Section 6 describes EVOLP, the dynamic logic programming framework defined over ISCO and Prolog. Section 7 describes the interaction manager and section 8 provides a simple example. Finally, in section 9 some conclusions and future work are pointed out.

2. Architecture

The system's architecture is based on several independent and modular processes. Figure 1 shows graphically these processes and their relations.

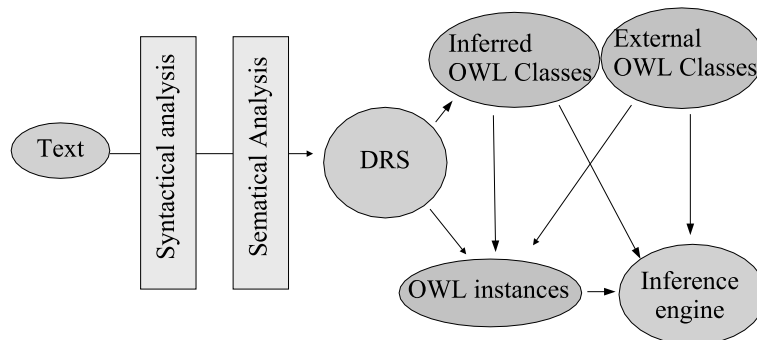


Figure 1. System's architecture

The architecture may be divided in three major modules:

- Inference of an adequate OWL ontology of classes;
- Inference of OWL instances and document enrichment;
- Inference engine.

The first module, inference of an adequate OWL ontology of classes, receives as input a top-level ontology and a set of legal documents. After a syntactical and semantical analysis, it obtains a partial semantic representation (a DRS – Discourse Representation structure (Kamp and Reyle, 1993)) for each sentence. From the DRS of each sentence,

noun expressions and verbs are extracted, and they are used to define legal classes. These classes will be clustered, classified, and merged with the initial top-level ontology (section 3 describes this step in detail).

The second module, inference of OWL instances, receives as input the DRS of each sentence and the inferred OWL ontology from the first module. With this information, and using an abductive inference mechanism, OWL instances are inferred. This step is usually called pragmatic interpretation of natural language sentences (section 4 describes these processes in more detail).

Finally, OWL classes and OWL instances are used by an inference engine, based in a logic programming framework, in order to answer queries about the semantic content of documents (sections 5, 6 and 7 describe the logic programming framework).

3. OWL ontology creation

In order to be able to deal with documents from different domains, a methodology to automatically create basic ontologies of concepts is proposed. This methodology allows the definition of a base ontology with the relevant concepts with some inferred relations. After having defined this ontology, it may be necessary to develop manual work by human experts in order to fully organize the set of extracted concepts.

The methodology is based on the following steps:

- Definition of an initial top-level ontology;
- Identification of concepts referred in the legal documents and extraction of its properties;
- Identification of relations between the identified concepts;
- Creation of an ontology using the identified concepts and relations;
- Merge of the created ontology with the initial ontology;

3.1. TOP-LEVEL ONTOLOGY

As referred in section 1, an already existent top-level legal ontology was chosen: the legal ontology from the Portuguese Attorney General's Office, consisting of around 6,000 classes and having around 10,000 relations (Quaresma and Rodrigues, 2002). As an example of some concepts in this ontology we have:

- Tribunal *Court*; properties: name, address, ...

- Tribunal Militar *Military Court* – subclass of Tribunal
- Supremo Tribunal *Supreme Court* – subclass of Tribunal

This legal ontology was merged with a general top-level ontology of concepts defined by Eckhard Bick in the VISL project ¹ (Bick, 2000), which has around 150 top concepts: animal, human, place, vehicle, concrete object, abstract object, food,

3.2. IDENTIFICATION OF CONCEPTS AND PROPERTIES

The methodology to automatically identify the concepts and the properties referred in the documents is based on the output of natural language processing tools:

- Text syntactical parsing. The documents are analyzed by the syntactical parser PALAVRAS developed by E. Bick. This parser is available for 21 different languages, including Portuguese.
- Partial semantic analysis.
- Entities extraction. From the semantic analysis output, entities and properties are extracted and represented by ontology classes.

3.2.1. *Syntactical analysis*

The parser developed by E. Bick is based on the Constraint Grammars (Karlsson, 1990) formalism and covers a major portion of the Portuguese sentences. However, its output is in a non-standard format and it was necessary to transform it into a structured form, like XML and Prolog terms. A translation tool from the VISL output into XML and Prolog terms was developed and it is available to the VISL users (a detailed description of this tool was presented in (Gasperin et al., 2003)).

As an example, suppose the following sentence:

O bombeiro Manuel salvou a criança. *The fireman Manuel saved the child.*

This sentence has the VISL output:

```
STA:fcl
SUBJ:np
=>N:art('o' M S)      0
=H:n('bombeiro' M S)  bombeiro
```

¹ <http://visl.hum.sdu.dk/visl>

```

==N<:prop('Manuel' M S) Manuel
P:v-fin('salvar' PS 3S IND)      salvou
ACC:np
=>N:art('a' F S)                a
=H:n('criança' F S)             criança

```

As it can be seen, the subject, predicate and direct object were correctly parsed. From this output, the XML translator produces three files:

1. The first file links each word with a *word* tag with a specific *id*.

```

<!DOCTYPE words SYSTEM "words.dtd">
<words>
<word id="word_1">0</word>
<word id="word_2">bombeiro</word>
<word id="word_3">Manuel</word>
<word id="word_4">salvou</word>
<word id="word_5">a</word>
<word id="word_6">criança</word>
<word id="word_7">.</word>
</words>

```

2. The second file associates each *word* with its part-of-speech information.

```

<!DOCTYPE words SYSTEM "wordsPOS.dtd">
<words>
<word id="word_1">
<art canon="o" gender="M" number="S"/>
</word>
<word id="word_2">
<n canon="bombeiro" gender="M" number="S"/>
</word>
<word id="word_3">
<prop canon="Manuel" gender="M" number="S"/>
</word>
<word id="word_4">
<v canon="salvar">
<fin tense="PS" person="3S" mode="IND"/>
</v>
</word>
<word id="word_5">
<art canon="a" gender="F" number="S"/>
</word>
<word id="word_6">

```

```
<n canon="criança" gender="F" number="S"/>
</word>
</words>
```

3. The third file has the parsing structure.

```
<!DOCTYPE text SYSTEM "text_ext.dtd">
<text>
<paragraph id="paragraph_1">
<sentence id="sentence_1" span="word_1..word_7">
<chunk id="chunk_1" ext="subj" form="np"
      span="word_1..word_3">
<chunk id="chunk_2" ext="n" form="art" span="word_1">
</chunk>
<chunk id="chunk_3" ext="h" form="n" span="word_2">
</chunk>
<chunk id="chunk_4" ext="n" form="prop" span="word_3">
</chunk>
</chunk>
<chunk id="chunk_5" ext="p" form="v_fin" span="word_4">
</chunk>
<chunk id="chunk_6" ext="acc" form="np"
      span="word_5..word_6">
<chunk id="chunk_7" ext="n" form="art" span="word_5">
</chunk>
<chunk id="chunk_8" ext="h" form="n" span="word_6">
</chunk>
</chunk>
</sentence>
</paragraph>
</text>
```

3.2.2. *Semantic analysis*

Each syntactical structure is translated into a First-Order Logic expression. The technique used for this analysis is based on DRS's (Discourse Representation Structures (Kamp and Reyle, 1993)) and it was described in more detail in (Quaresma and Rodrigues, 2003). The partial semantic representation of a sentence is a DRS built with two lists, one with the rewritten sentence and the other with the sentence discourse referents.

At present, we are only dealing with a restricted semantic analysis and we are not able to handle every aspect of the semantics: our focus is on the representation on concepts (nouns and verbs) and the correct extraction of its properties (modifiers, agents, objects).

From the XML structure, using XSL transformations, it is possible to obtain the semantic representation of each sentence.

The semantic representation of the example presented in the previous sub-section is:

```
sentence(doc1, [fireman(A), name(A,'Manuel'), child(B), save(A,B)], [ref(A), ref(B)]).
```

This structure represents an instance of a fireman *A*, named 'Manuel', and an instance of a child *B* which are related by the action *to save*.

A general tool able to obtain similar semantic partial representations for every sentence was developed and it was applied to the full set of legal documents of the Portuguese Attorney General's Office (7000 documents).

3.2.3. *Entities extraction*

From the sentence semantic representation, entities are extracted and they are the basis for the creation of an ontology of concepts. In fact, for each new concept, a new class, subclass of a top-class 'Entity', is created.

On the other hand, from the output of the semantic analyser it is possible to identify some potential class properties:

- Modifiers, such as adjectives, are candidates to be properties of nouns;
- Direct objects of transitive verbs are candidates to be properties of the associated verbs;

For instance, for the expression *the black cat* it is possible to identify *color* as a property of *cat*, because it is known that *black* is an instance of a color (from the correspondent semantic tag in the dictionary).

In the referred example it would be possible to extract the following entities:

- bombeiro *fireman*, with a property: 'name'
- salvar *to save*
- criança *child*

3.3. IDENTIFICATION OF RELATIONS BETWEEN CLASSES

As it was shown in the previous sub-section, the identification of concepts does not allow the creation of relations, hierarchical or others, between them.

Our approach is to use an unsupervised method for acquiring word classes and relations (Gamallo et al., 2002b; Gamallo et al., 2002a). The goal is to learn, for each word, what kind of modifiers and what kind of heads it subcategorises. For instance, the word *republic* may appear as an head of a noun phrase, such as *republic of Ireland*, *republic of Portugal*, or as a modifier, like *president of the republic*, *government of the republic*. The obtained subcategorisation patterns are clustered into classes and relations are extracted.

Using this approach it is possible to identify hierarchical relations, such as the existent between *republic* and *republic of Portugal* and also to identify other semantic relations, such as the ones between *lei – law* and *norma – norm*. The strategy is to use statistical analysis to identify clusters of words with similar subcategorisation patterns (words which have similar modifiers and heads).

As methodology, we start from the parsed documents and, for each word, subcategorisation patterns are extracted and clusters and relations are identified. A detailed description of the methodology is described in (Gamallo et al., 2002a).

Note that this approach has some limitations and it is not able to identify correctly what kind of relations exist between two concepts. For instance, two related concepts may be synonyms or the opposite. A more deep knowledge-aware approach is needed to handle these kind of problems.

The inferred relations are used to create an hierarchy of classes in the ontology and to link them via a *related* relation.

3.4. CREATION OF THE ONTOLOGY

In the fourth step of the methodology, the results of the previous two sub-sections are integrated into a new ontology:

- Concepts with their properties are new classes;
- Class hierarchies and relations are created accordingly with the results of the previous sub-section;

For instance, using this approach to the previous example, *republic of Portugal* will be a sub-class of *republic*.

3.5. MERGE OF THE ONTOLOGIES

Finally, in the fifth step, the initial top-level ontology is merged with the new one.

In this process new classes are inserted into the top-level ontology using their names and information from the semantic analyser:

- If a class exists with an equal name in the top-level ontology, then the two classes are merged;
- Otherwise, a search is made in the top-level ontology for a class with semantically compatible information and the new class is created as a sub-class of the existent one.

For instance, if a new class named *fireman* is classified to be a *human* concept by the NLP analysers, then the new class will be a sub-class of the *human* top-level concept.

The overall strategy is to search for common concepts in the two ontologies and to merge the ontologies via these concepts.

4. OWL instances creation

After having defined an ontology of classes, it is necessary to extract and to represent instances of those classes and to associate them with documents.

The proposed methodology tries to infer instances of those ontologies using the following three steps:

- Translation of the OWL ontologies into a logic programming form;
- Definition of logic programming rules allowing the inference of instances;
- Generation of OWL instances.

4.0.1. OWL translation

The first step, translation of OWL ontologies into Prolog, was implemented in Java and it creates a Prolog term for each OWL class, subclass, or property. The translation of this subset on OWL is correct accordingly with the OWL formal semantic description (Saias, 2003).

For instance, suppose there exists a definition in OWL for a class *citizen* and for a sub-class *military*. After the translation, we'll have:

```
class(citizen, 'external.owl#citizen').
class(military, 'external.owl#military').

subclass('external.owl#military', 'external.owl#citizen').
```

Moreover, suppose class *military* has a property of having a *rank*, which can have one of several possible values: general, colonel, ...

In this situation, we'll have the following Prolog terms:

```
property(rank, 'external.owl#rank',
          'external.owl#military').

hasPossibleValue('external.owl#rank', general).
hasPossibleValue('external.owl#rank', colonel).
```

4.0.2. Prolog rules

In the second step of this methodology, logic programming rules are defined allowing the inference of instances from the DRS representation of each sentence and the Prolog representation of the ontology.

One of these rules allows the inference of class and properties from values:

```
infer(Value, Class, Property) :-
    hasPossibleValue(PropertyURI, Value),
    property(Property, PropertyURI, ClassURI),
    class(Class, ClassURI).
```

In this LP rule, *Value* is the name of an entity (input) and *Class* and *Property* are identifiers of classes and properties that may have this value (output).

For instance, the sentence

The colonel saved the child.

has the following DRS form:

```
sentence(d1, [colonel(X), child(Y), save(X,Y)],
          [ref(X),ref(Y)]).
```

From this DRS form and, using the Prolog rules, it is possible to infer the following new form (because *colonel* is a possible value for the *rank* property of the *military* class):

```
sentence(d1, [military(X), rank(X,colonel), child(Y),
             save(X,Y)], [ref(X),ref(Y)]).
```

This process is usually called, in the natural language processing field, pragmatic interpretation of sentences and it can be seen as an abductive process where properties (antecedents) are inferred from values (consequents) (Hobbs et al., 1990).

Similar approaches can be applied to capture different natural language sentences characteristics.

For instance it is possible to relate adjectives or prepositional phrases with the head nouns:

The colonel from the army saved the child.

This sentence has the following DRS form:

```
sentence(d1, [colonel(X),
              army(Z), rel(X,Z),
              child(Y), save(X,Y)],
          [ref(X),ref(Z),ref(Y)]).
```

From this DRS form and, using Prolog rules, it is possible to infer the following new form:

```
sentence(d1, [military(X),rank(X,colonel),belong(X,army),
              children(Y), save(X,Y)], [ref(X),ref(Y)]).
```

4.0.3. OWL generation

In the third step, the results of the pragmatic interpretation of each sentence are transformed into correspondent OWL instances. For instance, for the last example of the previous sub-section, the following OWL instances would be created:

```
<pgr:Military rdf:ID="m11">
  <pgr:rank rdf:resource="external.owl#Colonel"/>
  <pgr:belong rdf:resource="external.owl#Army"/>
</pgr:Military>

<pgr:Child rdf:ID="c2">
</pgr:Child>

<pgr:ToSave rdf:ID="s5">
  <pgr:subject rdf:resource="#m11"/>
  <pgr:object rdf:resource="#c2"/>
</pgr:ToSave>
```

These instances define and relate a military (colonel and from the army), a child, through the instance of the action *to save*.

As a final result of this step, every document is enriched with the OWL instances obtained from the pragmatic interpretation of its sentences.

5. ISCO

In this section, the logic programming framework that is going to be used as the inference engine for answering queries about the semantic content of documents (OWL instances) is briefly described.

ISCO (Abreu, 2001) is a logic based development language implemented in GNU Prolog that gives the developer several distinct possibilities:

- It supports Object-Oriented features: classes, hierarchies, inheritance.
- It supports Constraint Logic Programming. Specifically, it supports finite domain constraints in ISCO queries.
- it gives a simple access to external relational databases through ODBC. It has a back-end for PostgreSQL and Oracle.
- It allows the access to external relational databases as a part of a declarative/deductive object-oriented (with inheritance) database. Among other things, the system maps relational tables to classes – which may be used as Prolog predicates.
- It gives a simple database structure description language that can help in database schema analysis. Tools are available to create an ISCO database description from an existing relational database schema and also the opposite action.

The proposed system uses ISCO's capability to establish connections from Prolog to relational databases in an efficient and simple way. For example, the following SQL table:

```
CREATE TABLE "document" (
  "number" int4 NOT NULL,
  "title" text,
  Constraint "number_pkey"
    Primary Key ("number")
);
```

Maps into the following ISCO class definition (and vice-versa):

```
external(pgr,document) class document.
    number: int. key.
    title: text.
```

Taking this ISCO feature into account, a translator from OWL into ISCO class definitions was developed. This translator was applied to every OWL class described in the previous section and, as a consequence, correspondent SQL tables and ISCO classes definitions were obtained. Moreover, each OWL class instance was transformed into an SQL table row and an ISCO logic programming fact. As an example, the *toSave* presented previously is translated into the following fact:

```
toSave(ID=s5, subject='#m11', object='#c2').
```

For each defined class a set of Prolog predicates implementing the four basic operations are created: query, insert, update and delete.

Variables occurring in queries are mapped to SQL and may carry CLP(FD) constraints, which will be expressed in SQL, whenever possible. For example, suppose variable *X* is an FD variable whose domain is (1..1000), the query

$$\text{document}(\text{number} = X, \text{title} = Y) \quad (1)$$

will return all pairs (*X*, *Y*) where *X* is a document number and *Y* is the document's title. *X* is subject to the constraints that were valid upon execution of the query, ie. in the range 1 to 1000.

ISCO class declarations feature inheritance, simple domain integrity constraints and global integrity constraints.

6. EVOLP

As it was described in the previous section, ISCO allows a declarative representation of ontologies and object instances. However, there is also a need to represent actions and to model the evolution of the knowledge.

In (Alferes et al., 1999) it was introduced a declarative, high-level language for knowledge updates called *LUPS* (Language of UPdateS) that describes transitions between consecutive knowledge states. Recently, a new language, *EVOLP* (Alferes et al., 2002), was proposed having a simpler and more general formulation of logic program updates. In this section a brief description of the *EVOLP* language will be given. A detailed description of the language and of its formalization is presented at the cited article.

EVOLP allows the specification of a program's evolution, through the existence of rules which indicate assertions to the program. *EVOLP*

programs are sets of generalized logic program rules defined over an extended propositional language L_{assert} , defined over any propositional language L in the following way (Alferes et al., 2002):

- All propositional atoms in L are propositional atoms in L_{assert}
- If each of L_0, \dots, L_n is a literal in L_{assert} , then $L_0 \leftarrow L_1, \dots, L_n$ is a generalized logic program rule over L_{assert} .
- If R is a rule over L_{assert} then $assert(R)$ is a propositional atom of L_{assert} .
- Nothing else is a propositional atom in L_{assert} .

The formal definition of the semantics of EVOLP is presented at the referred article, but the general idea is the following: whenever the atom $assert(R)$ belongs to an interpretation, i.e. belongs to a model according to the stable model semantics of the current program, then R must belong to the program in the next state. For instance, the following rule form:

$$assert(b \leftarrow a) \leftarrow c \quad (2)$$

means that if c is true in a state, then the next state must have rule $b \leftarrow a$.

EVOLP has also the notion of external events, i.e. assertions that do not persist by inertia. This notion is fundamental to model interaction between agents and to represent actions. For instance, it is important to be able to represent actions and its effects and pre-conditions:

$$assert(Effect) \leftarrow Action, PreConditions \quad (3)$$

If, in a specific state, there is the event *Action* and if *PreConditions* hold, then the next state will have *Effect*.

7. Interaction Management

The interaction manager is built on the ISCO+EVOLP logic programming framework.

As final goal, we aim to handle the following kind of questions:

- Situations where action A is performed
- Situations where action A is performed having subject S
- Situations where S is the subject of an action

Note that the inference engine needs to be able to deal with the ontology relations. For instance, the question "situations where action A is performed having subject S" means "situations where action A (or any of its sub-classes) is performed having subject S (or any of its sub-classes)".

The interaction manager is composed by the following main tasks:

- Query management
- Interaction management

7.1. QUERY MANAGEMENT

The analysis of a natural language query is split in three subprocesses: Syntax, Semantics, and Pragmatics.

7.1.1. *Syntax*

As syntactic analyser we are using the analyzer developed by E. Bick and referred previously (Bick, 2000). The VISL output is translated into Prolog facts by the same translator referred in section 3. This translation can be handled by the same translator because there is a direct relation between the XML structure and the Prolog term structure.

As an example, the following query:

```
Quem salvou crian\c{c}as?
‘‘Who saved children?’’
```

Has the following syntactical structure:

```
sentence(syn(que(fcl,
  subj(pron_indp('quem', 'M/F', 'S', '<interr>'), 'Quem'),
  p(v_fin('salvar', 'PS', '3S', 'IND'), 'salvou'),
  acc(n('criança', 'F', 'P', '<H>'), 'crianças', '?')))).
```

7.1.2. *Semantics*

As referred in section 3, each syntactical structure is translated into a First-Order Logic expression (DRS). The semantic representation of a sentence is a DRS built with two lists, one with the rewritten sentence and the other with the sentence discourse referents. For instance, the semantic representation of the sentence above is the following expression:

```
child(B), toSave(A,B),
```

and the following discourse referents list:

A : [ref(A),ref(B)]

These structures represent instances of children *B* related with instances of the *toSave* action.

Note that, at present, we are not able to deal with general unrestricted queries and to translate them from a syntactical into a semantic structure. In fact this a quite complex NLP problem and we have decided to deal only with specific subsets of the Portuguese language, namely, with interrogatives about specific domains.

7.1.3. Pragmatic Interpretation

The pragmatic module receives the semantic query representation and tries to interpret it in the context of the database information, which was constructed from the translation of the OWL instances into ISCO facts (as described previously in section 5).

In order to achieve this behavior the system tries to find the best explanations for the sentence logic form to be true in the knowledge base. As already referred, this strategy for interpretation is known as “interpretation as abduction” (Hobbs et al., 1990) and this process was described in detail in (Quintano et al., 2001).

From the description of the OWL (and ISCO) classes it is possible to obtain the correspondent ISCO query:

```
child(id=B),
toSave(id=C, subject=A, object=B),
```

This query was obtained using additional logic programming rules for the interpretation of actions in the context of the ontology class descriptions:

```
assert action(id=C,subject=B,object=C) <-
    action(A,B), entity(A), entity(B).
```

Note that the ontology hierarchy was used to infer that *children* are entities and *to save* is an action.

The interpretation of the ISCO predicates is done by accessing the knowledge base in order to collect (and constraint) all entities identifiers:

```
- $A=_\#(104..109:156..157)$ -- A constrained to all
entities with the desired properties
```

The above expression contains the possible interpretations of the query in the context of the knowledge base.

7.2. INTERACTION MANAGEMENT

The interaction manager has to represent the actions associated with the queries (*informs* or *request*), and to model the user attitudes (intentions and beliefs).

This task is also achieved through the use of the EVOLP language (see (Quaresma and Rodrigues, 2001; Quaresma and Lopes, 1995) for a more detailed description of these rules). For instance, the rules which describe the effect of an inform, and a request speech act are:

$$\text{assert}(\text{bel}(A, \text{bel}(B, P))) \leftarrow \text{inform}(B, A, P). \quad (4)$$

$$\text{assert}(\text{bel}(A, \text{int}(B, \text{Action}))) \leftarrow \text{request}(B, A, \text{Action}). \quad (5)$$

These rules mean that if an agent A is informed of a property P , then it will start to believe that the other agent believes in P ; additionally, if B requests A to perform an action Action , then A starts to believe that B intends Action to be performed.

In order to represent collaborative behavior it is also necessary to model the transference of information between the agents:

$$\text{assert}(\text{bel}(A, P)) \leftarrow \text{bel}(A, \text{bel}(B, P)). \quad (6)$$

$$\text{assert}(\text{int}(A, \text{Action})) \leftarrow \text{bel}(A, \text{int}(B, \text{Action})). \quad (7)$$

These two rules means that if an agent A believes another agent believes in P , then it will start to believe in P (it is a cooperative, credulous agent); moreover, it will also adopt the intentions of the other agents.

There is also the need for a rule linking the system intentions and the accesses to the databases:

$$\text{assert}(\text{inf}(A, B, P)) \leftarrow \text{int}(A, \text{inf}(A, B, P)), \text{isco}(P). \quad (8)$$

$$\text{assert}(\text{not int}(A, B, \text{inf}(A, B, P))) \leftarrow \text{inf}(A, B, P). \quad (9)$$

The first rule defines that, if the system intends to inform the user about some property, then it will access the ISCO database and it will perform an inform action. The second rule means that the inform action will end the intention to perform the inform action!

8. Examples

Considering the already presented query:

```
Quem salvou crianças?
‘‘Who saved children?’’
```

The interaction manager receives the query pragmatic interpretation:

```
Q = [child(id=B), toSave(id=C, subject=A, object=B)].
```

After having the sentence rewritten into its semantic representation form, the speech act is recognized:

```
request(user, system, inform(user, system, Q))
```

Using the *request* and the transference of intentions rules the following property is supported:

```
int(system,inform(system, user, Q))
```

Now, using the rules presented in the previous section, the system accesses the ISCO databases and it is able to obtain the final constraints to the discourse referent variables:

```
- $A=_\#(104..109:156..157)$ -- A constrained to all
    entities with the desired properties
```

Using the inferred constraints it is possible to obtain the set of solutions to the user query and to answer:

```
m11: Military - rank: colonel; belong: army.
```

If, for instance, the user asked the more general query:

```
Quem salvou quem?
‘‘Who saved somebody?’’
```

The unique difference would be the pragmatic interpretation and (probably) the system's answer:

```
Q = [ save(id=C,subject=A,object=B) ].
```

9. Conclusions and Future Work

A methodology to automatically create legal ontologies was proposed.

The methodology uses syntactical, semantical and pragmatical analysers to obtain sentence representations and to identify entities and entity relations. The obtained ontologies are merged with other externally defined top-level ontologies.

The obtained new ontology is used, with the semantic representation of sentences, to infer class instances and to enrich documents with this semantic information. The inference of the instances associated with each sentence is done via an abductive process – interpretation as abduction.

Ontologies and the inferred instances are represented in the OWL language.

On the other hand, translators from OWL into ISCO/Prolog were developed and a logic programming based interaction manager was developed. The interaction manager uses many important features from its base LP framework: objects, constraints, inheritance.

At present, the implementation of the system is in a prototype/test phase and it needs work in many areas:

- Ontology creation. The ontology was created automatically but it was not possible to identify many relations between the classes. In order to be able to define these relations we intend to extend the statistical analysis of word subcategorisation to take into account semantic information from the dictionary and existent Wordnets.
- Normalisation of concepts. The parsing process was not able to completely identify and eliminate duplicates and incorrections. A more sophisticated analysis is needed.
- OWL translation into ISCO/Prolog. A full translation of the OWL language needs to be implemented and its correction has to be proved.
- Evaluation. The system needs to be fully evaluated and to be tested by users. Moreover it should be applied to other legal documents, such as, legislation.

References

- Abreu, S.: 2001, 'ISCO: A Practical Language for Heterogeneous Information System Construction'. In: *Proceedings of INAP'01*. Tokyo, Japan, INAP.

- Alferes, J., A. Brogi, J. Leite, and L. Pereira: 2002, 'Evolving Logic Programs'. In: S. Flesca, S. Greco, N. Leone, and G. Ianni (eds.): *JELIA'02 – Proceedings of the 8th European Conference on Logics and Artificial Intelligence*. pp. 50–61, Springer-Verlag LNCS 2424.
- Alferes, J. J., L. M. Pereira, H. Przymusinska, T. C. Przymusinski, and P. Quaresma: 1999, 'Preliminary exploration on actions as updates'. In: M. C. Meo and M. Vilares-Ferro (eds.): *Procs. of the 1999 Joint Conference on Declarative Programming (AGP'99)*. L'Aquila, Italy, pp. 259–271.
- Bick, E.: 2000, *The Parsing System "Palavras"*. *Automatic Grammatical Analysis of Portuguese in a Constraint Grammar Framework*. Aarhus University Press.
- Boer, A., R. Hoekstra, R. Winkels, T. van Engers, and F. Willaert: 2002, 'Proposal for a Dutch Legal XML Standard'. In: *EGOV2002 – Proceedings of the First International Conference on Electronic Government*.
- Breuker, J. and R. Winkels: 2003, 'Use and reuse of legal ontologies in knowledge engineering and information management'. *Journal of Artificial Intelligence and Law*. In this issue.
- Damáσιο, C.: 2003, 'W4 – Well-founded semantics for the World Wide Web'. In: H. Boley, G. Grosz, S. Tabet, and G. Wagner (eds.): *Rule Markup Techniques for the Semantic Web*. Dagstuhl, Germany.
- Gamallo, P., A. Agustini, and G. Lopes: 2002a, 'Using co-composition for acquiring syntactic and semantic subcategorisation'. In: *ACL-SIGLEX'02*. Philadelphia, USA.
- Gamallo, P., A. Agustini, P. Quaresma, and G. Lopes: 2002b, 'Using semantic word classes in text information retrieval systems'. In: S. Pinto (ed.): *SBIE'2002 – XII Simpósio Brasileiro de Informática na Educação, Workshop de Ontologias*. Porto Alegre, Brasil, pp. 593–597, Unisinos. ISBN 85-7431-133-2.
- Gasperin, C., R. Vieira, R. Goulart, and P. Quaresma: 2003, 'Extracting XML syntactic chunks from Portuguese corpora'. In: *TALN'2003 - Workshop on Natural Language Processing of Minority Languages and Small Languages of the Conference on "Traitement Automatique des Langues Naturelles"*. Batz-sur-Mer, France.
- Guangemí, A., N. Guarino, C. Masolo, A. Oltramari, and L. Schneider: 2002, 'Sweetening ontologies with DOLCE'. In: A. Gomez-Perez and V. R. Benjamins (eds.): *Proceedings of the EKAW'2002*. pp. 166–181, Springer-Verlag.
- Hobbs, J., M. Stickel, D. Appelt, and P. Martin: 1990, 'Interpretation as Abduction'. Technical Report SRI Technical Note 499, 333 Ravenswood Ave., Menlo Park, CA 94025.
- Kamp, H. and U. Reyle: 1993, *From Discourse to Logic*. Dordrecht: Kluwer.
- Karlsson, F.: 1990, 'Constraint grammar as a framework for parsing running text'. In: H. Karlgren (ed.): *13th International Conference on Computational Linguistics*, Vol. 3. Helsinki, Finland, pp. 168–173.
- Lame, G.: 2003, 'Using text analysis techniques to identify legal ontologies' components'. In: *Workshop on Legal Ontologies of the International Conference on Artificial Intelligence and Law*.
- Lassila, O. and R. Swick: 1999, 'Resource Description Framework (RDF) - Model and Syntax Specification'. W3C.
- Quaresma, P. and J. G. Lopes: 1995, 'Unified logic programming approach to the abduction of plans and intentions in information-seeking dialogues'. *Journal of Logic Programming* **54**.
- Quaresma, P. and I. Rodrigues: 2001, 'Using Logic Programming to model Multi-Agent Web Legal Systems – An Application Report'. In: *Proceedings of the*

- ICAIL'01 - International Conference on Artificial Intelligence and Law*. St. Louis, USA, ACM.
- Quaresma, P. and I. P. Rodrigues: 2002, 'PGR: Portuguese Attorney General's Office Decisions on the Web'. In: Bartenstein, Geske, Hannebauer, and Yoshie (eds.): *Web-Knowledge Management and Decision Support*. Springer-Verlag.
- Quaresma, P. and I. P. Rodrigues: 2003, 'A natural language interface for information retrieval on semantic web documents'. In: E. Menasalvas, J. Segovia, and P. Szczepaniak (eds.): *AWIC'2003 - Atlantic Web Intelligence Conference*. Madrid, Spain, pp. 142–154, Springer-Verlag.
- Quintano, L., I. Rodrigues, and S. Abreu: 2001, 'Relational Information Retrieval through Natural Language Analysis'. In: *Proceedings of INAP'01*. Tokyo, Japan, INAP.
- Saias, J.: 2003, 'Uma Metodologia para a construção automática de Ontologias e a sua aplicação em Sistemas de Recuperação de Informação – A methodology for the automatic creation of ontologies and its application in information retrieval systems'. Master's thesis, University of Évora, Portugal. In Portuguese.
- W3C: 2000, 'DAML+OIL – DARPA Agent Markup Language'. www.daml.org.

