

From Logical Forms to SPARQL Query with GETARUN

Rocco Tripodi, Rodolfo Delmonte

Ca' Bembo, Dorsoduro 1075 Università "Ca Foscari" 30123 – Venice, Italy
{rocco, delmont}@unive.it

Abstract. We present a system for Question Answering which computes a prospective answer from Logical Forms produced by a full-fledged NLP for text understanding, and then maps the result onto schemata in SPARQL to be used for accessing the Semantic Web. As an intermediate step, and whenever there are complex concepts to be mapped, the system looks for a corresponding amalgam in YAGO classes. It is just by the internal structure of the Logical Form that we are able to produce a suitable and meaningful context for concept disambiguation. Logical Forms are the final output of a complex system for text understanding - GETARUNS - which can deal with different levels of syntactic and semantic ambiguity in the generation of a final structure, by accessing computational lexical equipped with sub-categorization frames and appropriate selectional restrictions applied to the attachment of complements and adjuncts. The system also produces pronominal binding and instantiates the implicit arguments, if needed, in order to complete the required Predicate Argument structure which is licensed by the semantic component.

1 Introduction

Nowadays, the need of the automatic processing of information on the web has become more and more relevant in order to develop applications able to cope with unstructured information.

Semantic Web (hence SW) is the project aiming at implementing a smarter web and has its fundament in a Tim Berners-Lee paper [1]. The article describes an Artificial Intelligence task applied to the web. The idea at the heart of the project is referencing things in the real world. The referencing procedure developed over the years is based on metadata and ontology. The metadata provide a computer-readable concept specification and the ontology provides a conceptual knowledge structure, which organize concepts.

According to Wilchs [2] we could consider the SW to have an Information Extraction task at its heart. The SW task consists in relating entities to specific categories (e.g. Person, Place, Event, etc.). The formalism used to add facts in the SW is RDF¹. RDF is used in the SW to express facts by means of simple Predicate-Argument Structures (hence PAS) with subject-predicate-object structure. For instance, to express that Madonna is an artist we may use the triple below:

¹ Resource Description Framework: <http://www.w3.org/RDF/>

Subject	: ² Madonna (entertainer)
Predicate	rdfs:type
Object	dbo: ³ Artist

The example proposed has been extracted from DBpedia⁴ [3], a dataset organized over an ontology. DBpedia contains millions of facts extracted from the *Wikipedia infoboxes*⁵ and expressed in RDF triples. DBpedia is a Knowledge Base (hence KB) and also the de-facto core of the Linked Open Data⁶ (hence LOD) project [4]. This project has the ambition to become the foundational basis of the SW; it consists of several KBs linked together, in which it is possible to find the reference of millions of entities, and the facts that characterize each entity. We could, then, consider the LOD datasets as *encyclopedia*⁷ [5], where we could find information about entities, and the reference could be considered as an attribution of meaning.

The W3C⁸ standard way to access a KB on the SW is SPARQL⁹. SPARQL is used to express queries across data sources, whether the data is stored or viewed as RDF. In the Semantic Web, ontologies supply a machine-interpretable knowledge infrastructure. The real challenge does not only lie in constructing ontologies and keeping them up to date, but chiefly in linking them with the natural language [6]. In order to link automatically reference in text and entities in knowledge bases, a series of tools and heuristics are used for what can be called the semantic disambiguation task, i.e. discover the exact concept or the exact entities referenced in the text.

This is what happens in case the query to be constructed has [president,'United States'] as its goal, and the amalgam search will produce the complex concept [PresidentOfTheUnitedStates] as in example (1) below. In case no class has been recovered, as for instance in the query related to the complex structure [5th,president,'United States'] the system knows that the cardinal figure '5th' behaves like a quantifier restricting the class of [PresidentOfTheUnitedStates].

GETARUNS – the system for text understanding – produces Logical Forms (hence LFs) which are organized with a restricted ontology made up of 8 types: FOCus, PREDicate, ARGument, MODifier, ADJunct, QUANTifier, INTensifier, CARDinal PLURal. In addition, every argument has a Semantic Role to tell Subject from Object and Referential from non-Referential predicates. Another important step in the computation of the final LF, is the translation of the interrogative pronoun into a corresponding semantic class word taken from general nouns, in our case the highest concepts of WordNet hierarchy.

² PREFIX : <<http://dbpedia.org/resource/>>

³ PREFIX dbo: <<http://dbpedia.org/ontology/>>

⁴ <http://dbpedia.org/About>

⁵ <http://en.wikipedia.org/wiki/Help:Infobox>

⁶ <http://linkeddata.org/>

⁷ The encyclopedia is understood within the meaning given by Umberto Eco [6], as a network of interconnected cultural units

⁸ <http://www.w3.org/>

⁹ SPARQL Query Language for RDF: <http://www.w3.org/TR/rdf-sparql-query/>

The result is mapped into classes, properties, and restrictions (filters). as for instance in the question:

(1) Who was the wife of President Lincoln ?

which becomes the final LF:

```
be-[focus-person, arg-[wife/theme_bound], arg-['Lincoln'/theme-[mod-[pred-['President']]]]]]
```

and is then turned into the SPARQL expression,
.`?x dbpedia-owl:spouse :Abraham_Lincoln`

where "dbpedia-owl:spouse" is produced by searching the DBpedia properties and in case of failure looking into the synset associated to the concept as WIFE. In particular then, the concept "Abraham_Lincoln" is derived from DBpedia by the association of a property and an entity name, "President" and "Lincoln", which contextualizes the reference of the name to the appropriate referent in the world.

This paper is divided into two parts. In the section below we focus on providing access to the SW through Natural Language. We discuss the problems we encountered and the solutions and strategies we adopt. The second part concerns Question Answering over Linked Data. We explain our question analysis approach and give examples of how our algorithm works.

1.1 Accessing the LOD Cloud through Natural Language: problems

On the LOD Cloud the information come from different ontologies, lacking a semantic mapping among them, and many ontologies describe similar domains with different terminologies [7].

Such problems sketch two main points that we would like to address by means of semantic disambiguation technique and mapping process. Without NLP technique, access to the SW through Natural Language is allowed only using a short lexicon, which is made up of non homogeneous KBs labeling systems. This is due to the fact that a large KB has to handle with homonymy and synonymy problems.

Liu [8] noted that DBpedia contains a great number of disambiguation nodes. A disambiguation node is used to resolve conflicts when a single term can be the title of more than one article: for example the word "Mercury" can refer to several different things, including an element, a planet, an automobile brand, a record label, a NASA manned-spaceflight project, a plant, and a Roman god. Liu [8] explains that things linked by a disambiguation node are only related through rough homonymy. So when we look up a word in DBpedia we get a long list of possible candidates. Such problems are due to both word ambiguity and to the labeling system used. However, as Buitelaar claimed [9], the RDFS and OWL standards are not sufficient for the purpose of associating linguistic information with ontologies.

Besides the problem of homonymy, there is also the problem of synonymy. In DBpedia such problems are partially handled by the "redirect" property. A "redirect"

property is a property (in the RDF formalism) that links a node A to a node B, where the node B is the preferred concept for A. That property is used in DBpedia to manage misspellings, alternative spellings, tenses, punctuation, capitalizations, etc. or to redirect sub-topic in broader context [8]. In that prospective we can see that the semantic content of the synonymy is not treated, and the access to the KB through natural language is limited to a short vocabulary. To avoid such a problem Buitelard et al., [9], have proposed a solution within the ontology markup standards. The idea behind this approach is to enter linguistic information inside the ontologies. Our approach, as will be explained below, does not attempt to modify the SW standards but tries to manage them by means of NLP and IE techniques.

By now we have discussed problems related to concept names, but a KB also contains names for classes and properties. The class names are common names which specify the collocation of a concept. The property links a concept with another concept, a class or a literal.

As noted by Fu et al. [10] some relations in DBpedia have anomalous names that are hard to understand and therefore are difficult to use. Another problem concerns the fact that many relations share the same meaning, for example “dateOfBirth”, “birthDate” and “datebirth” are three variant of the same concept. So when we want to retrieve all the entities with a particular property we have to collect all the various forms of the property.

Similarly, DBpedia classes were extracted from different sources such as YAGO¹⁰, UMBEL¹¹ and Wikipedia¹². Only 170 were manually created for the project and are consistent with the DBpedia ontology [4]. Many extracted classes have the same problems of properties; besides, many classes express complex concepts with *n*-ary relations [9] such as:

- AncientGreekPhilosophers
- OlympicTennisPlayersOfTheUnitedStates
- CommandersOfTheOrderOfTheBritishEmpire

Classes of that kind have a complex semantics that is hard to use without a preprocessing phase. The first thing we do to handle these names is to split them into tokens. Then we proceed with an NLP-based analysis¹³. In particular, we analyzed them with a syntactic constituency parser and obtained the output below, where F3 is the label for fragments, SN stands for NounPhrase, SP for PrepositionalPhrase:

- f3-[sn-[Philosophers-n-sn, (mod)-[ancient_Greek-n-sn]]]
- f3-[sn-[olympic-ag-sn,tennis_players-n-sn,(mod)-[of-p-sp,the-art-sn, United_States-n-sn]]]
- f3-[sn-[Commanders-n-sn, (mod)-[of-p-sp,the-art-sn,Order-n-sn]],sp-[of-p-sp,the-art-sp,British_Empire-n-sp]]

¹⁰ Yet Another Great Ontology <http://www.mpi-inf.mpg.de/yago-naga/yago/>

¹¹ Upper Mapping and Binding Exchange Layer <http://umbel.org/>

¹² <http://www.wikipedia.org/>

¹³ The system used for the analysis is VENSES [16, 17]

The analysis identifies the head and the modifiers of the head which is the governing name. At this point the heads must be disambiguated in order to be compared with the words in text. So we can use this information with contextual information. Modifiers are used to apply consistency checks.

Another step is done mapping the heads with synsets in WordNet, in order to expand the KB lexicon, for instance, the word “actress”, in the question:

Is Natalie Portman an actress?

matches the class:

<http://dbpedia.org/ontology/Actor>

because “actress” share the same synset of “actor”, as shown in the following term,

`dbp('Actor',[actor-n],[109765278,109767197]).`

where we associated WordNet synset labels and DBpedia classes. In particular, DBP is a Prolog compound term, where the first element corresponds to a DBpedia label, the second element adds a POS tag to the label and the last element is a list with all synset labels. WordNet mapping allows us to use hyponymy relation, for instance, the word “wife” in the question:

Who was the wife of President Lincoln?

matches the property,

<http://dbpedia.org/ontology/spouse>,

because there is an hyponymy relation between “wife” and “spouse”.

1.2 Accessing the Semantic Web through Natural Language: technique

Sowa [12] asserts that, each ontology, for practical application, must have a mapping, direct or indirect, related to and deriving from natural languages, because human knowledge is developed around human language. So an useful ontology must support a systematic mapping to and from natural languages, because such a bridge could break the static nature of a KB and make it flexible. The lack of this bridge has by now failed to achieve the hoped results in Artificial Intelligence and Knowledge Management [12].

What we have in mind is the assumption that an ontology reflects the background knowledge used in writing, reading and thinking [6]. In fact a text tells the reader which ontology to use to understand it [6]. The background knowledge, taken for granted by the author, is useful because can be used by a NLP application in order to decide a particular word sense.

Word Sense Disambiguation (hence WSD) techniques use the notion of *context* in order to decide a particular word sense. A context could differ widely across WSD methods. One may consider a whole text, a word window, a sentence or some specific words [13].

Such techniques are necessary to access a static KB because the concepts are static objects; however knowledge can then be used and developed by reasoning. This approach comes from the *dynamic construal of meaning* (DCM) [14] approach, that we follow. The fundamental assumption of DCM is that the meaning of a word changes as it is used in different contexts or language games [12].

According to Chierchia [15] we consider the computation of meaning as a set of rules that determine the reference of words. We consider common names as classes, determiners as restrictions on classes, entities as referents and verbs as relations between entities and classes. This scheme is compatible with the RDF structure and can also serve as a bridge between natural language and KBs. Our approach is also related to Wittgenstein's language games [16], in that we assume we need to use patterns of words, to access an ontology. The RDF triples are atomic facts with a simple semantic. The meaning of each fact is the result of the meaning of three components:

- **Classes:** a class could be represented by a common name. When we talk about presidents, trees, cars, or carpenters, we are talking about classes of entities.
- **Entities:** we intend an entity as his reference. To access an entity we use his label and the disambiguation is done by one or more classes to which the entity belongs.
- **Properties:** are simple or complex relations between entities, classes and literal. We need to disambiguate a property and get contextual information from it.

With our approach, we want to extract information about the meaning of text. Particularly we want to understand what specific entities are mentioned in the text. To do this we use IE techniques to identify the named entities. We can use their names as labels to access a KB in order to get all the information regarding the entities. But as we noted above the same label could refer to several entities. The solution is to use contextual information. For instance, in the following example taken from the RTE5 challenge dataset:

Proper Name + Definite Expression
Rte5 – DevSet - TH Pair 10

(CNN) -- Malawians are rallying behind **Madonna** as she awaits a ruling Friday on whether she can adopt a girl from the southern African nation. **The pop star**, who has three children, adopted a son from Malawi in 2006. She is seeking to adopt Chifundo "Mercy" James, 4. "Ninety-nine percent of the people calling in are saying, let her take the baby," said Marilyn Segula, a presenter at Capital FM, which broadcasts in at least five cities, including the capital, Lilongwe.

when we find an ambiguous entity (the pop star) we look for information that could disambiguate it. In this case, the singular definite expression "the pop star" is used to specify the entity Madonna. The definite expression consists of a determiner and a common noun that in our approach correspond to a class. At this point we have to establish which class could be associated with the noun found. This step corresponds

to a WSD procedure, which serve as a bridge between natural language and KB. This approach is particularly useful in coreference resolution task where we have an identical name but different properties as in the following example:

President Clinton leaves for Northern Ireland shortly on a diplomatic mission aimed at bolstering the 1998 power sharing accord. He's scheduled to arrive in Belfast late Tuesday after a stopover in Dublin for talks with Irish President Mary McAleese and Prime Minister Bertie Ahern. En route to Belfast, Mr. Clinton will also visit the Irish border town of Dundalk.

Here, in fact, through the class name derived from the word President we get the right reference for Bill Clinton. Otherwise, if we had used only the label we could not clearly identify the entity, because there are a lot of entities in DBpedia with the same label. In this way, coreference resolution is performed in parallel with entity identification. Consider another example below, with a text taken from the same RTE5 dataset:

Definite Expression + Proper Name
Rte5 – TestSet - TH Pair 269

The eruption happened at around 1:30 PM local time, the United States Geological Survey reported. The volcano had erupted four times on Friday, billowing ash up to 51,000 feet up into the air. These are the latest in a series of eruptions from Mount Redoubt, which started on March 22. The volcano had not erupted since a four-month period in 1989-90. The Alaska Volcano Observatory set its alert level at red, the highest possible level, meaning that an eruption is imminent, and that it would send a "significant emission of volcanic ash into the atmosphere."

In this example the name "Mount Redoubt" could refer to different entities:

- Mount Redoubt (Alaska) in Alaska, United States
- Mount Redoubt (Washington) in Washington, United States
- Redoubt Mountain in Banff National Park, Canada

but the characteristic of being a volcano belongs only to one entity:

http://dbpedia.org/resource/Mount_Redoubt

We use abduction to guess a new hypothesis that explains some fact. More on this in the following sections.

2 Question Answering over Linked Data

We start from the assumption that, any system for Information Extraction, or Question Answering, working under the hypothesis of open domain, unlimited vocabulary, and unstructured text, needs access to world knowledge. The encyclopedic knowledge we are referring to is the one that could be represented by web KB and in particular by the LOD project. Accessing KBs is done with the RDF triples structure in mind, which would correspond strictly to a PAS; and the disambiguation task is done using background information derived from the text.

2.1 Question analysis

As said above, question analysis is performed using GETARUNS¹⁴ [17, 18], the system for text understanding developed at the University of Venice, which is organized as a pipeline that includes two versions of the system: what we call the Partial and the Deep GETARUNS. At first we will present the Deep version, which is equipped with three main modules: a lower module for parsing, where sentence strategies are implemented; a middle module for semantic interpretation and discourse model construction which is cast into Situation Semantics; there is also a higher module where reasoning and generation takes place which we use to answer questions from text.

The system is based on LFG (Lexical Functional Grammar) theoretical framework and has a highly interconnected modular structure. The Closed Domain version of the system is a top-down depth-first DCG (Definite Clause Grammars) based parser written in Prolog Horn Clauses, which uses a strong deterministic policy by means of a lookahead mechanism with a WFST (Weighted Finite State Transducer) to help recovery when failure is unavoidable due to strong attachment ambiguity.

It is divided up into a pipeline of sequential but independent modules which realize the subdivision of a parsing scheme as proposed in LFG theory where a c-structure is built before the f-structure can be projected by unification into a DAG (Direct Acyclic Graph). In this sense we try to apply in a given sequence phrase-structure rules as they are ordered in the grammar: whenever a syntactic constituent is successfully built, it is checked for semantic consistency. In case the governing predicate expects obligatory arguments to be lexically realized they will be searched and checked for uniqueness and coherence as LFG grammaticality principles require.

Syntactic and semantic information is accessed and used as soon as possible: in particular, both categorial and sub-categorization information attached to predicates in the lexicon is extracted as soon as the main predicate is processed, be it adjective, noun or verb, and is used to subsequently restrict the number of possible structures to

¹⁴ The system has been tested in STEP competition, and can be downloaded in two separate places. The partial system called VENSES in its stand-alone version is available at http://www.aclweb.org/aclwiki/index.php?title=Textual_Entailment_Resource_Pool. The complete deep system is available at <http://project.cgm.unive.it/html/sharedtask/>.

be built. Adjuncts are computed by semantic compatibility tests on the basis of selectional restrictions of main predicates and adjunct heads.

The output of grammatical modules is fed then onto the Binding Module which activates an algorithm for anaphoric binding. Antecedents for pronouns are ranked according to grammatical function, semantic role, inherent features and their position at f-structure. Eventually, this information is added into the original f-structure graph and then passed on to the Discourse Module (hence DM).

The grammar is equipped with a core lexicon containing most frequent 5000 fully specified inflected word forms where each entry is followed by its lemma and a list of morphological features, organised in the form of attribute-value pairs. However, a morphological analyser for English is also available with a root dictionaries (25,000 for English) which can only provide guesses for syntactic sub-categorization, though. In addition to that, there are all lexical form provided by a fully revised version of WordNet and COMLEX¹⁵, and in order to take into account phrasal and adverbial verbal compound forms, we also use lexical entries made available by UPenn and TAG encoding. Their grammatical verbal syntactic codes have then been adapted to our formalism and are used to generate a sub-categorization scheme with an aspectual and semantic class associated to it – however no selectional restrictions can reasonably be formulated on arguments of predicates. Semantic inherent features for Out of Vocabulary Words, be they nouns, verbs, adjectives or adverbs, are provided by a fully revised version of WordNet - plus EuroWordnet¹⁶, with a number of additions coming from computer, economics, and advertising semantic fields - in which we used 75 semantic classes similar to those provided by CoreLex¹⁷.

When each sentence is parsed, tense aspect and temporal adjuncts are accessed to build the basic temporal interpretation to be used by the temporal reasoner. Eventually, two important modules are fired: Quantifier Raising and Pronominal Binding. QR is computed on f-structure which is represented internally as a DAG. It may introduce a pair of functional components: an operator where the quantifier can be raised, and a pool containing the associated variable where the quantifier is actually placed in the f-structure representation. This information may then be used by the following higher system to inspect quantifier scope. Pronominal binding is carried out at first at sentence internal level. DAGs will be searched for binding domains and antecedents matched to the pronouns if any to produce a list of possible bindings. Best candidates will then be chosen.

GETARUNS, has a linguistically based semantic module which is used to build up the Discourse Model (hence DM). Semantic processing is strongly modularized and distributed amongst a number of different submodules which take care of Spatio-Temporal Reasoning, Discourse Level Anaphora Resolution, and other subsidiary processes like Topic Hierarchy (Centering) which cooperate to find the most probable antecedent of coreferring and cospecifying referential expressions when creating

¹⁵ <http://nlp.cs.nyu.edu/comlex>

¹⁶ <http://www.illc.uva.nl/EuroWordNet/>

¹⁷ CoreLex:-<http://www.cs.brandeis.edu/~paulb/CoreLex/corelex.html>

semantic individuals. These are then asserted in the DM, which becomes a useful knowledge representation used to solve nominal coreference. The system uses two resolution submodules which work in sequence: they constitute independent modules and allow no backtracking. The first one is fired whenever a free sentence external pronoun is spotted; the second one takes the results of the first submodule and checks for nominal anaphora. They have access to all data structures contemporarily and pass the resolved pair, anaphor-antecedent to the following modules. Semantic Mapping is performed in two steps: at first a Logical Form is produced which is a structural mapping from DAGs onto unscoped well-formed formulas. These are then turned into situational semantics informational units, infons which may become facts or sits. Each unit has a relation, a list of arguments which in our case receive their semantic roles from lower processing – a polarity, a temporal and a spatial location index.

Logical Forms derived from DAGs or f-structure sentence level representations are simplified in order to be useful for the question answering task. In particular, we come up with a non-recursive linear representation at propositional level where we introduce prefixes for each semantic head which are very close to DRS-conditions:

– PRED, QUANT, CARD, PLUR, ARG, MOD, ADJ, FOC

where Foc contains the question type derived from a mapping of the Wh- word, its possible nominal or adjectival head and a restricted set of semantic general classes, like MEASURE, MANNER, QUANTITY, REASON etc.

2.2 From Logical Form to SPARQL query

Our system produces a LF of natural language questions by means of GETARUNS. From LF, the system extracts the semantic elements needed to produce a SPARQL query that is then used to address LOD endpoint.

LFs produced by GETARUNS are all expressed as complex Prolog terms, and can be decomposed into three subparts: there is a Pred - the main verb predicate of the question -, a Focus - this is the question head expressed in the question which may correspond to an interrogative pronouns or may have a nominal head -, and then Arguments - this slot contains argument head and its internal modifiers and attributes like Quantifier, Cardinality, Plural. This slot may also contain other Arguments or entities and so on recursively. For instance, consider the following example:

```
Quest Which are the presidents of the United States of America?  
Pred [be]  
Focus [person]  
Arg [president/theme_bound-[['United_States_of_America']]]
```

As can be gather from the example, the Question is decomposed into three subelements, these are then used to build the SPARQL query. Predicate [be] can be regarded as the fact "belonging to a class". Focus [person] tells us that the reply foreseen by the question must be of Type Person, important feature which is easily expressed in SPARQL. We then look for the elements in Arg inside the two ontologies, DBpedia and YAGO and we obtain the class:

yago: PresidentOfTheUnitedStates

At this point, we can start building the query according to the schema:

```
?x a [Focus]
?x a [Class]
```

As explained above, there is no unique way of expressing the relation between properties and classes, and Person may belong to a number of different classes that have the same meaning. In order to cover the all of them in the KB we need to address them all in the query and consequently we come up with a multiple recursive query of the kind that we show below, where triples are conjoined by the clause UNION.

```
select distinct ?x ?string WHERE {
  {?x a dbpedia-owl:Person . ?x a yago:PresidentsOfTheUnitedStates} union
  {?x a foaf:person . ?x a yago:PresidentsOfTheUnitedStates} union
  {?x a yago:Person100007846 . ?x a yago:PresidentsOfTheUnitedStates}
OPTIONAL {?x rdfs:label ?string . FILTER (lang(?string) = "en")}}
```

This query has a Precision of 1.0 and a Recall of 0,93 which leads to believe that missing results are due mainly to the difficulty of addressing all the different ways of expressing the same concept in DBPedia (as for instance, "being the president of the United States" or "being a Person").

In some cases, no useful class can be derived from Args produced by the LF. In that case, we need to introduce what can be regarded as FILTERS, which we derive from quantifiers and other restrictions to predicates, in order to narrow down the search, as for instance in the question:

```
Quest Who has been the 5th president of the United states of America?
Pred [be]
Focus [person]
Arg [[president],card-'5th',['United_States']]
```

where we understand that the element individuated by Card, "5th", behaves like a restriction that operated on the class yago:PresidentsOfTheUnitedStates. Since there is no way to express such a restriction in SPARQL, we create a FILTER that looks into short literals for the specific word "5th", "president", "United States". This FILTER will be added to the previous query, like this:

```
?x ?prop ?lbl .
FILTER
  (?prop != dbpedia-owl:abstract &&
  ?prop != rdfs:comment %26%26
  regex(?lbl, "(^)president(|$)","i") &&
  regex(?lbl, "(^)5th(|$)","i") &&
```

```
regex(?lbl, "(^ )United States( |$)","i" ).
```

2.2.1 YES/NO QUESTIONS

In case the LF does not produce a Focus element, the system understands that the question type is yes/no. In this case, the system will create a query of type ASK, which is meant to verify the existence of one or more RDF triples. Suppose the question is the following,

```
Quest Is Christian Bale starring in Batman Begins?
Pred [be]
Focus []
Arg ['Christian_Bale'/theme_bound-[mod-[pred-[star], ['Begins'/theme-[mod-[pred-
['Batman']]]]]]]]
```

by analyzing the Arg element we realize that there are two entities and one property. In the organization of the final query, we proceed by looking for entities first: this we do because we find it important to verify the existence of a given concept before proceeding to submit the actual query containing it. In this preliminary phase, we search for the concepts related to the entities "Christian Bale" and "Batman Begins" in order to contextualize them. Then we also look for the predicate "star" in a special mapping we built where DBpedia properties are linked to WordNet verb synsets. When building this mapping, we found out that in many cases there was no possible correspondance between the information present in WordNet and the amalgamated labels of DBpedia. So we had to proceed manually.

The ASK query we produce for the above example is based on the simple scheme:

```
Ent Prop Obj
```

which produces the following query

```
ASK {
  { :Christian_Bale dbpedia-owl:starring :Batman_Begins. } Union
  { :Batman_Begins dbpedia-owl:starring :Christian_Bale . } }
```

as can be seen, we reverse the order of the two arguments of the predicate STAR, because we do not know whether it is being used in the active or the passive form.

In other questions we proceed by disambiguating a property contained in the LF before proceeding to build the corresponding query. This is the case of the example below,

```
Quest Who was the wife of President Lincoln?
Pred [be]
Focus [person]
Arg [wife/theme_bound-[['Lincoln'/theme-[(mod)-[pred-['President']]]]]]]]
```

Here, the system finds at first one property "being wife" (which is not expressed as a

class but as a DBpedia property) and another element which consists of a label [Lincoln] and a class [President]. This latter property helps us to disambiguate the entity expressed by the question, because it contextualizes the reference, and it allows us to recover the actual intended entity, i.e. Abraham_Lincoln, by means of the procedure previously indicated. In this query, the scheme is the following one:

```
?x Prop Ent
```

and it allows us to build the following query:

```
select distinct ?x ?string WHERE {
  { ?x dbpedia-owl:spouse :Abraham_Lincoln . } Union
  { :Abraham_Lincoln dbpedia-owl:spouse ?x . }
  OPTIONAL { ?x rdfs:label ?string . FILTER (lang(?string) = "en") }
```

Also in this case we use the reversed version of the query, which counts as the logically derivable statement "President Lincoln has a wife x".

2.2.2 FILTERS: GRADABLE ADJECTIVES AND QUANTIFIERS

There are other special cases of queries which require some filtering of the results, as shown in questions where the relevant property is expressed by a comparative or superlative adjective as in,

Quest What is the highest mountain?

```
LF [[be],focus-[mountain],[mountain/theme_bound-[(mod)-[pred-[highest]]]]]
```

Quest Which mountains are higher than the Nanga Parbat?

```
LF [[be],focus-[mountain],[higher/theme_bound-['Parbat'/theme_bound-[ (mod)-[pred-['Nanga']]]]]]
```

In both cases we have a superlative which is mapped through a specific filter: in the first case we have a scheme like,

```
?x a Class.
?x prop ?value
ORDER BY DESC(?value) LIMIT 1
```

which is transformed in the following query:

```
select distinct ?x ?string WHERE {
  {?x a dbpedia:Mountain. ?x dbpedia-owl:elevation ?value. } Union
  {?x a dbpedia:Mountain. ?x dbpedia2:elevationM ?value. }
  OPTIONAL { ?x rdfs:label ?string . FILTER (lang(?string) = "en") }
  ORDER BY DESC(?value) LIMIT 1
```

In the second case the presence of a superlative induces a slightly different scheme:

```
?x a Class
ent prop ?valueE
?x prop ?valueX
  FILTER (?valueX > ?valueE) .
```

which is transformed in the following query:

```
select distinct ?x ?string WHERE {
  {?x a :Mountain . dbpedia:Nanga_Parbat dbpedia2:elevationM ?y1. ?x
  dbpedia2:elevationM ?y2.}
  {?x a :Mountain . dbpedia:Nanga_Parbat dbpedia-owl:elevation ?y1. ?x dbpedia-
  owl:elevation ?y2.}
  FILTER (?y2 > ?y1) .
  OPTIONAL {?x rdfs:label ?string . FILTER (lang(?string) = "en")}}
```

In this case, at first we recover the class to which the prospective answers belongs, by means of DBpedia ontology, and then, after we have analysed the superlative, we look for the properties it may be referred to and the kind of filter to use. Properties are recovered by means of our mapping onto DBpedia. As to the mapping of the two adjectives "higher" and "highest", they will be mapped both onto dbpedia2:elevationM and dbpedia-owl:elevation; because they are understood as belonging to the domain of :Place, which is the class right superior to :Mountain.

Information present in the Focus element allow us to build expectations and filters for a specific type of answer. In particular in case we have a question like:

```
Quest How many films did Leonardo DiCaprio star in?'
LF [[do],focus-[quantity],pred-[star],(mod)-[pred-['Leonardo_DiCaprio']],[films]]
```

The Focus [quantity] requires us to count the number of results obtained from the query. Building the query then is done by using the remaining part of the question, where we have an entity [Leonardo_diCaprio], a predicate [star], and a class name [films]. Eventually we come up with the following scheme:

```
?x a Class.
?x prop Ent
```

just because the Focus is not a class, we can use the class found in the Arg to produce the final query:

```
select count(?x) WHERE {
  ?x a dbpedia-owl:Film
  { :Leonardo_DiCaprio dbpedia-owl:starring ?x. } union
  { ?x dbpedia-owl:starring :Leonardo_DiCaprio. } union
  { ?x dbpedia-owl:starring "Leonardo DiCaprio"@en. } }
```

Here again we reverse subject and object and we add a third entry which is referred to the label associated to the name of the entity. In fact, in many cases, DBPedia refers to an entity with one of its label rather than with referring to a unique link. This fact is the reason why we lose sometimes points in the computation of recall, since literals may be missing when we impose a certain class to results of the search.

When the predicate used in the question is not a copular verb, we come up with different schemes, as for instance in:

Quest Which books were written by Danielle Steel?
LF [[write],focus-[book],[Steel'/ [(mod)-[pred-['Danielle']]]]]

or

Quest Which actors were born in Germany?
LF [[bear],focus-[actor],adj-[pred-['Germany']]]

The underlying scheme would be,

?x a [Focus]
?x Pred [Arg]

from which we build two different queries: in the first case,

```
select distinct ?x ?string WHERE {  
  ?x a dbpedia-owl:Book  
  { :Danielle_Steel dbpedia-owl:author ?x. } union  
  { ?x dbpedia-owl:author :Danielle_Steel. } union  
  { ?x dbpedia-owl:author "Danielle Steel"@en. }  
  OPTIONAL { ?x rdfs:label ?string . FILTER (lang(?string) = "en") }
```

in the second example,

```
select distinct ?x ?string WHERE {  
  ?x a dbpedia-owl:Actor  
  { ?x dbpedia-owl:birthDate:Germany. } union  
  { ?x dbpedia-owl:birthPlace :Germany. } union  
  { ?x dbpprop :birthPlace :Germany. } union  
  { ?x dbpprop:birthDate :Germany. } union  
  { ?x dbpprop:birthDate :Germany. } union  
  { ?x dbpprop:placeOfBirth :Germany. }  
  OPTIONAL { ?x rdfs:label ?string . FILTER (lang(?string) = "en") }
```

In the latter case, as in previous ones, we added recursively as many triples as there are properties linked to the Pred. Also note that in this case, subject and object are not reversed, and this is due to the nature of the complement which is computed as ADJunct or Oblique and not as Object or Xcomp(lement) or open complement for

predicative structures.

2.2.3 PROBLEMS

In our system the major problems we had have been with the ability to recover complex concepts, as for instance in the question:

Give me all female German chancellors!

where we try to decompose the meaning into three different but intertwined queries:

?x Female
?x German
?x Chancellor

But we don't get desired results and the reason is that DBPedia does not contain the male/female distinction. Probably there are amalgams which can express the complex concept to be a woman and be the head of the German government, but at the moment, our mapping strategy has not been able to find a class for the concept. On the contrary, it worked fine in the case of `yago:PresidentsOfTheUnitedStates` and in many others.

Other problems regard the use of literals in place of unique identifiers. For instance in the question:

Quest In which programming language is GIMP written?
LF [[write],focus-[programming_language],['GIMP']

we use the scheme

?x a [Focus]
?x Prop Ent

But this is not correct since the reply for this question (C and GTK+) is not expressed with two unique references but with a literal, and literals cannot belong to any class. In this case the system does not receive a result and a second scheme is used, which consists in the elimination of the Focus and the reversal of subject and object:

Ent Prop ?x

But also in this case we jump into a problem because we use as Prop [write] and this verb has a mapping which does not allow us to obtain the desired result: in fact, the property needed to obtain the correct result (C and GTK+) is `dbpprop:programmingLanguage`, and it is very difficult to derive from the Pred element [write].

3 Evaluation

We have tested our system on the training set made available by QALD-1 workshop organizers, referred to DBPedia. The training set contains 50 question expressed in natural language to submit to DBPedia. We obtained correct answers (Precision and Recall = 1) to 23 questions over 50, with a final overall Precision and Recall equal to 0.46.

We looked into the mistakes and found out that:

- a. in 13 cases, we did build up an efficient and adequate query;
- b. in 4 cases we got a Precision ranging 0.80-0.98;
- c. in 5 cases we got a Recall ranging 0.85-0.99.

In case a. we obtained partial results and the Recall ranged between 0.4-0.8. In case b. results are due to the presence of literals, which duplicate reference to the same entity with different names though: this could be avoided building up filters that eliminate multiple reference. In case c. we did not get any result. We assume that this is due to the fact that DBPedia allows to refer to the same entity or concept using different properties which however were not present in our mapping, thus preventing some elements not to be included in our results.

We include here below the table for the final results we obtained for the 50 dbpedia questions.

Number of constructed Queries	Number of correct Answers	Number of wrong Answers	Global Recall	Global Precision	Global F-Measure
50	23	27	0.46	0.46	0.46

And here below are the IDs of failed questions:

Failed questions (IDs)

7, 11, 37, 31, 3, 40, 46, 50, 6, 35, 19, 34, 15, 25, 5, 33, 28, 45, 36, 17, 44, 48, 1, 13, 4, 16, 24.

4 Conclusion

The remaining 6 cases are due to problems that our system has encountered for a lack of a strong mapping to many DBpedia properties. We have to understand the meaning that some properties have in DBpedia and then to move that information to the system, as we have already done automatically with classes and WordNet synset. Work is underway to improve on the mapping from SPARQL and with properties.

5 References

1. Berners-Lee, T., Hendler, J., Lassila, O. 2001 The Semantic Web, Scientific American.
2. Wilks, Y. 1997. Information extraction as a core language technology. In M.-T. Paziienza (Ed.), Information Extraction. Springer, Berlin.

3. Bizer, C., Lehmann, J., Kobilarov, G., Auer, S., Becker, C., Cyganiak, R., Hellmann S. 2009. DBpedia – A Crystallization Point for the Web of Data. *Journal of Web Semantics*:
4. Berners-Lee, T. 2006. Linked Data - Design Issues, <http://www.w3.org/DesignIssues/LinkedData.html>
5. Eco, U. 2007. Dall'albero al labirinto. Studi storici sul segno e l'interpretazione (From the tree to the labyrinth), Bompiani.
6. Brewster C., Ciravegna F., Wilks Y. 2003. Background and foreground knowledge in dynamic ontology construction. In *Semantic Web Workshop, SIGIR'03*, Toronto, Canada.
7. Doan, A., Madhavan, J., Dhamankar, R., Domingos, P., Halevy, A. 2003. Learning to Match Ontologies on the Semantic Web. *VLDB Journal* 12(4):303– 319.
8. Liu, O. 2009. "Relation Discovery on the DBpedia Semantic Web", TER 2009, supervised by Jérôme Euzenat.
9. Buitelaar, P., Cimiano, P., Haase, P., Sintek, M. 2009. Towards linguistically grounded ontologies. In *Procs. Of European Semantic Web Conference*.
10. Fu, L., Wang, H., Yu, Y. 2009. Towards Better Understanding and Utilizing Relations in DBpedia
11. Suchanek, F. M., Kasneci G., Weikum, G. 2007. "Yago - A Core of Semantic Knowledge", 16th international World Wide Web conference
12. Sowa, J. F. 2010. The role of logic and ontology in language and reasoning, Chapter 11 of *Theory and Applications of Ontology: Philosophical Perspectives*, edited by R. Poli & J. Seibt, Berlin: Springer, pp. 231-263.
13. Xiaobin, L., Szpakowicz, S., Matwin, S. 1995. A wordnet-based algorithm for word sense disambiguation. In *Proceedings of IJCAI-95*, pages 1368-1374, Montreal, Canada, August.
14. Cruise, D. A. 2002. Microsenses, default specificity and the semantics-pragmatics boundary, in *Axiomathes* 1, 1-20.
15. Chierchia G. 1997. *Le strutture del linguaggio. Semantica, (The structures of language. Semantics)* il Mulino, Bologna.
16. Wittgenstein, L. 1953. *Philosophical Investigations*, Basil Blackwell, Oxford.
17. Delmonte, R. 2007. *Computational Linguistic Text Processing - Logical Form, Semantic Interpretation, Discourse Relations and Question Answering*. New York: Nova Science Publishers.
18. Delmonte, R. 2008 *Computational Linguistic Text Processing - Lexicon, Grammar, Parsing and Anaphora Resolution*. New York: Nova Science Publishers.