

# Domain Ontology Construction by Partial Import for Document Annotation

Tayybah Kiren, Muhammad Shoaib, and Sang-Jo Lee

**Abstract**—Semantic annotation is a useful technique to understand the underlying meanings of the document. Domain specific knowledge provided by the domain ontologies is needed to semantically annotate the document. But the main problem is the availability of an appropriate ontology for the domain. Ontology construction from scratch is protracted and labor intensive job. Therefore, it is advantageous to construct the ontologies by reusing the existing ontologies. This paper proposes a technique for constructing domain ontology for semantic annotation of a document by partially importing the existing ontologies.

**Index Terms**—Semantic annotation, partial reuse, domain ontology.

## I. INTRODUCTION

Finding the right information from the huge corpus of documents available on the World Wide Web is very difficult. Documents contain the valuable knowledge for particular domain. But this knowledge cannot be efficiently exploited by machines for automation purpose because of the unstructured nature of the documents. Thus the creation of semantic metadata related to document content seems to be a way to exploit this knowledge and extract implicit and explicit information. The process of creating semantic metadata is called “Semantic annotation”.

Ontologies are defined as formal explicit specifications of shared conceptualizations. In Semantic web, domain ontologies are considered to be the major source of semantic metadata creation. But unfortunately constructing domain ontologies from scratch is very time consuming task which makes the semantic annotation process very slow. Hence there is a need of constructing the domain ontologies by partially reusing the existing ontologies on the web. This paper proposes an approach for constructing domain ontology for semantic annotation by partially reusing the existing relevant ontologies on the web. Section II presents the existing approaches for domain ontology construction and for semantic annotation, proposed approach is described in Section III and evaluated in Section IV. Section V presents the conclusion.

## II. RELATED WORK

Constructing domain ontologies from text document is a non-trivial task. Most of the existing approaches

construct ontology from scratch. Sleeman et al proposed an approach for extracting ontological knowledge from the existing knowledge bases [1]. In [2] text corpus is used as background knowledge for constructing ontologies. Problem with these approaches that the knowledge sources used for ontology building do not express the knowledge in explicit way. This can limit the constructed ontologies. Moreover the construction of ontologies from scratch is quite time consuming task. A solution to these problems is reusing the existing relevant ontologies. There are some approaches that construct the ontologies by importing the existing relevant ontologies. The approach presented in [3] partition the large ontologies into small portions according to their class hierarchy structure. Difficulty with this approach is that it is not suitable for small size ontologies. Grau *et al.* [4] introduced a module extraction approach based on logic. This approach extracts locality based modules. There are many approaches for ontology based semantic annotation of text documents. Most of these approaches use existing ontologies for semantic annotation. Ontoglass [5] is an ontology based annotation system that performs morpheme level annotation. Amaya [6] is performs manual semantic annotation using RDF data but it just provides information like author, creator or title of document. OntoAnnotate [7] is another annotation tool, which keeps a copy of document to be annotated along with the annotation in the annotation database, thus making the annotation database heavy. All these semantic annotation approaches use some existing ontology for making annotations. Thus it is evident from the literature that there are many approaches for ontology construction from scratch and ontology partitioning and modularization. But there is no complete methodology for constructing domain ontologies by partially reusing the existing ontologies and using them to semantically annotating the unstructured document.

## III. PROPOSED APPROACH

Architecture of the proposed approach for domain ontology construction for semantic annotation is shown in Fig. 1 Preprocessing of document includes stemming and removing the common words as they are not good candidate to be the domain concepts. The domain concepts for which no relevant ontologies exist are defined specifically according to the domain ontology and others are retrieved in the form of relevant modules from existing ontologies. Stepwise description of the proposed methodology is given in the following text.

### A. Step 1: Finding Domain Concepts

To construct the domain ontologies first we will find out the domain concepts from given document. Keywords are

Manuscript received July 9, 2013; revised September 16, 2013.

Tayybah Kiren is with University of Engineering and Technology Lahore, Pakistan (e-mail: Tayybah@uet.edu.pk).

born concepts in the domain to which they belong [8]. Hence we first find out the keywords from the input document. Keywords are extracted on the basis of format (whether it is bold or highlighted) and position in the document. Position of a word in the document also shows its significance and relatedness to the domain of document [9]. For example if a word appears both in introduction and conclusion it generally carries more information. Thus domain concepts are extracted by using Position weight formula as presented in equation 1 which extracts keyword based on linguistic features.

$$PW(t,d) = \sum_{i=1}^n pw(t_i) \quad (1)$$

where  $pw(t_i)$  is calculated by taking into consideration the weights assigned to the paragraph

$$Rank(O) = \sum_{i=1}^n Match(O, k_i) \quad (2)$$

where  $Match(O, k_i) =$

- 1 if there is exact match between  $O$  and  $K_i$
- 0 if there is no match between  $O$  and  $K_i$
- 0.5 if there is match between  $O$  and  $K_i$  through WordNet synonym

and sentence to which the word  $t_i$  belongs, and  $d$  is the document.

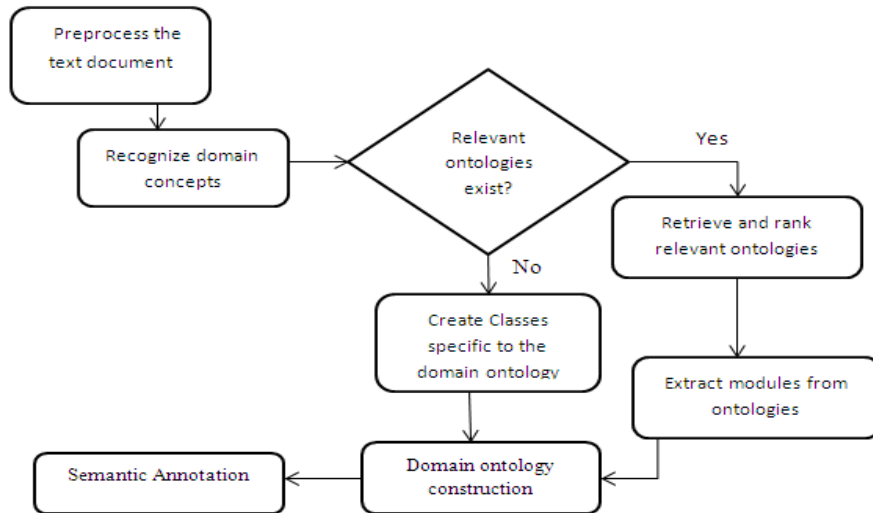


Fig. 1. Proposed system architecture.

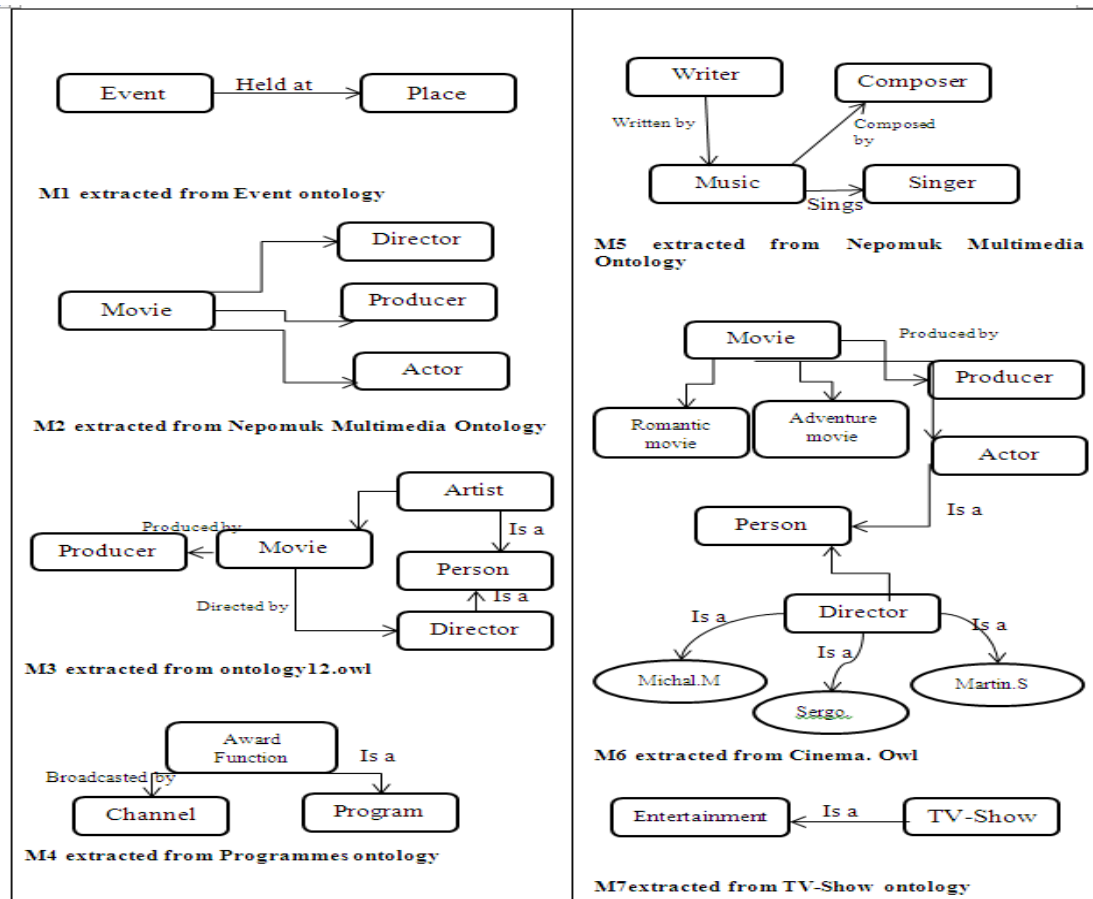


Fig. 2. Modules extracted from relevant ontologies.

### B. Step 2: Searching and Ranking Relevant Ontologies

In this step the selected domain concepts are used as query to retrieve the reusable relevant ontologies from the web. We have used semantic web search engine SWOOGLE<sup>1</sup> as well as Google for retrieving the ontologies. The result set of ontologies obtained after searching the web is ranked according to the relevancy of ontology to the domain concepts. The ranking measure is defined as in equation 2, where O is the ontology whose rank is being calculated and  $k_i$  is the set of keywords

### C. Step 3: Extracting Relevant Modules from the Ontologies

A module M from an ontology O is defined as “Minimal Portion of an ontology O that if integrated with ontology O1 yields the same effect as if the whole ontology O was integrated with O1”. We have proposed an algorithm for module extraction that matches the ontology classes with the set of domain concepts WordNet [10] is used to get the synonyms of the ontology classes so that if a class label is similar in meaning to any of the domain concept is also retrieved.

Algorithm: Extracting Modules from an input ontology against a set of Domain concepts.

Input: Ontology O in OWL format, domain Concepts= K

Output: Module M

Procedure: Module Extraction (O, K)

Check\_List: Container for all nodes in the ontology O.

Module\_List: Container for nodes included in the extracted module.

Load the ontology

Check\_List = Get all classes their relations, individuals and restrictions.

Recent= first element of Check\_List.

While (Check\_List! = Empty)

x = recent. Class

AccessWordNet for synonyms

S= getSyno(x)

If ((x = any of the K OR any member of S = any of the K) AND (x is not the member of Module\_List)

Module\_List. Class= Recent. Class

Module\_List. Relations= Recent. Relations

Module\_List. Individuals = Recent. Individuals

Module\_List. Restrictions= Recent. Restrictions

End if

Recent = Next element of Check\_List

End While

### D. Step 4: Constructing Domain Ontology

In this step the extracted modules are assembled with the classes which are defined specifically for the desired domain ontology. Extracted modules and specifically created concepts are integrated in incremental way, thus removing any inconsistencies in ontologies parallel with the ontology construction process.

### E. Step 5: Semantic Annotation

At this step the input web page is semantically annotated with the domain ontology. The domain ontology is looked up

to find matching between the ontology concepts and extracted domain concepts and relationships. If an ontology concept is not directly matched to the domain concepts then subsumption relationships are used to map them.

## IV. RESULTS AND DISCUSSION

TABLE I: EXTRACTED RELEVANT ONTOLOGIES

S .No	Ontology
1	Nepomuk Multimedia Ontology (NMM)
2	<a href="http://labotalc.loria.fr/~kasimir/downloads/owl/cinema.owl">http://labotalc.loria.fr/~kasimir/downloads/owl/cinema.owl</a>
3	<a href="http://tw.rpi.edu/wiki/Special:ExportRDF/Category:TV_Show.owl">http://tw.rpi.edu/wiki/Special:ExportRDF/Category:TV_Show.owl</a>
4	<a href="http://raimond.me.uk/c4dm/music.owl">http://raimond.me.uk/c4dm/music.owl</a>
5	<a href="http://lsdis.cs.uga.edu/proj/traks/ontologies/ontology_12.owl">http://lsdis.cs.uga.edu/proj/traks/ontologies/ontology_12.owl</a>
6	<a href="http://motools.sourceforge.net/event/event.html(event Ontology)">http://motools.sourceforge.net/event/event.html(event Ontology)</a>
7	<a href="http://www.aktors.org/ontology/portal">http://www.aktors.org/ontology/portal</a>
8	<a href="http://www.bbc.co.uk/ontologies/programmes/2009-04-17.shtml (programmes ontology)">http://www.bbc.co.uk/ontologies/programmes/2009-04-17.shtml (programmes ontology)</a>

For the narration of proposed technique using case study, we have constructed an ontology for a document that describes the event of an award show that is stored in a local directory. Domain concepts are extracted by applying the formulas given in proposed technique for domain concept extraction. These domain concepts are then used as query to the search engine for retrieving the relevant ontologies. Set of domain concepts are stored as a set K.

K: { Entertainment, TV-show , Award Show , Award distribution, Performances , Actor, Director, Producer, Movie, Music, Nominees, Award winners, Audience , Guests }.

Retrieved ontologies are then ranked according to their degree of relevancy to the set K by using the proposed Rank measure. Table I shows the list of top ranked ontologies. By applying the proposed algorithm for module extraction, modules extracted from ontologies are shown in Fig. 2. By combining the ontology modules shown in Fig. 2 with the classes defined specifically for the “Award function ontology” the construction of domain ontology is completed. Graphical representation of the domain ontology is given in Fig. 3. Note that in the constructed domain ontology of Award Show some concepts (e.g. Nominees, Award distribution etc.) are created from scratch because these were not matched to any of the existing ontology on the web. Now by tagging the words in text document with ontology concepts make the text document machine understandable. For example if some person searches for the word Filmfare or Oscar then it is known to the search engine that it is an award function and an award function is broadcasted by some channel . By this context knowledge, achieved by ontology based annotation, the search engine will retrieve the pages which do not explicitly include the word filmfare or Oscar but they are about the award achievements of a movie or director.

<sup>1</sup> <http://swoogle.umbc.edu/>

