

Aligning Controlled Vocabularies using a Facet based Approach for Facilitating the Linked Open Data

Ahsan Morshed
Intelligent Sensing and
System Laboratory (ICT),
CSIRO, Hobart, 7001
TAS, Australia

Md. Saddam Hossain
Mukta
Department of CSE
Islamic University of Technology
(IUT)
Dhaka, Bangladesh

Bayzid Ashik Hossain
Department of CSE
Islamic University of
Technology (IUT)
Dhaka, Bangladesh

ABSTRACT

A vocabulary stores words, synonyms, word sense definitions (i.e. glosses), relations between word senses and concepts; such a vocabulary is generally referred as the Controlled Vocabulary (CV) if choice or selections of terms are done by domain specialists. A facet is a distinct and dimensional feature of a concept or a term that allows taxonomy, ontology or controlled vocabulary to be viewed or ordered in multiple ways, rather than in a distinct way. The facet is also clearly defined, mutually exclusive, and composed by collectively exhaustive aspects, properties or characteristics of a domain. For example, a collection of rice might be represented using a name facet, place facet etc. In this paper an approach has been proposed where a facet has built for each concept by considering more general concepts (broader terms), less general concepts (narrow terms) or related concepts (related terms) that is to be called concept facet (CF). These CF's have been used for mapping two controlled vocabularies. This methodology is based on hidden semantic matching which is different from the orthodox view of matching. Finally the aim is to contribute in alignment of controlled vocabulary to enrich the Linked Data version of AGROVOC with appropriate links to other thesauri.

General Terms

Semantic Web, Linked Open Data, ontology.

Keywords

Facet, thesauri, controlled vocabulary, AGROVOC, CABI.

1. INTRODUCTION

The Semantic Web (which has gained widespread fame recently), where the underlying idea is that web contents should be expressed not only in natural language but also in a language that can be unambiguously understood, interpreted and used by software agents, thus permitting them to find, share and integrate information more easily. The central notation of the Semantic Web's idea is the ability to uniquely identify resources (with URIs) and languages (e.g. RDF/S, OWL) to formally represent knowledge (i.e. ontologies, which can simplistically be considered the taxonomies of classes representing objects, and of their inter-relationships) [12, 3].

The problem of *matching* or *aligning* (Noy, 2004) information resources [21] such as XML schemas, database schema, ontologies etc. has received much attention as a pre-requisite to data exchange. Since 2004, the Ontology Alignment Evaluation Initiative is the international event to compare on a common benchmark the state of the art matching systems. These taxonomies contain domain knowledge; the domain is

represented by a set of words and phrases used to describe concepts. A vocabulary is said to be controlled if it stores domain-specific chosen words, synonyms, word sense definitions (i.e. glosses) and relations between word senses and concepts [20]. In Controlled Vocabulary (CV), the words have denoted as “blocks from which sentences are made”, a synonym as “a word or phrase that refers to the same concept”, a sense as “a meaning of a concept” and a concept as “an abstract idea inferred or derived from specific instances”. The importance of CVs can hardly be underestimated; generally, each company or research group has its own information source e.g. databases, schemas and structures. Each of these sources has their respective set of individual CVs, creating a high level of heterogeneity. On one hand this is desirable, as it allows the involved parties to structure knowledge in a way which best fits their needs, e.g., for specific inter-office applications. On the other hand, individuals or companies also sometimes need a unified knowledge base (made up of different information sources) in order to satisfy their goals. This source of integration process requires a mapping between different CVs. Mapping between two CVs is generally a critical challenge for semantic interoperability. These CVs are used a lots as background knowledge for this data integration [7, 5]. What is more, classification are matched using CVs are lightweight ontologies, also called Formal classification (FC). In FC, lexical labels are translated to logical labels that remove ambiguities of natural language. For interested reader, we can refer to [9, 6]. In this paper, the proposed approach intended to the correspondence between concepts from two CVs, e.g., concept-to-concept mapping which includes word-to-word mapping, or synonym-to-synonym mapping. This mapping cannot be accomplished solely by a lexical comparison of two concepts using element level matcher [10, 13] that is included in SMOA Distance, Hamming Distance, Jaro Measure, SubString Distance, N-gram, JaroWinKler Measure, and Lavestein Distance; the existing semantics also need to be considered. In light of the discussion stated above, the objective of this work is to determine a fully-automated mapping between two CVs and this work may be useful for navigating vocabularies, information extraction and linking information. This paper presents work which is an extension of our previous paper in details [1]. This is major paper for AGROVOC Linked Open Data.

2. Facet Controlled Vocabulary

2.1 Facet

A facet is like a diamond that is consisting of different faces. Its distinct features allow thesauri, classifications or taxonomies to be organized in different ways, rather than in a single way. The facet is also clearly defined, mutually

exclusive, and composed by collectively exhaustive aspects of properties or characteristics of a domain. For example, a collection of rice might be classified using cultural and seasonal facets.

A Facet is constructed by following two steps [7]:

2.1.1 *Domain analysis*: First analysis of the term by consulting domain experts. This process is called idea plane, the language independent conceptual level, where simple concepts are identified. Each identified concept is expressed in the verbal plane of a given language. For example in English, trying to articulate the idea coextensively, namely identifying a term which exactly and unambiguously expresses the concept.

2.1.2 *Term collections and organization*: Secondly, collect terms and make an order of homogenous terms according to their characteristics, and order them (in hierarchies) in a meaningful sequence. The set of homogenous terms form a facet. For example, cow and milk form a facet called Dairy System (these entities are part-of relation with Dairy System).

Above steps construct a faceted knowledge organization system and corresponding to background knowledge, namely the a-priori knowledge which must exist in order to preserve effective semantics. Notice that the grouping of terms of step 2 have real world semantics, namely they are ontologies, classification and thesauri which are formed using partOf, isA, isSubclassOf and instanceOf relationships.

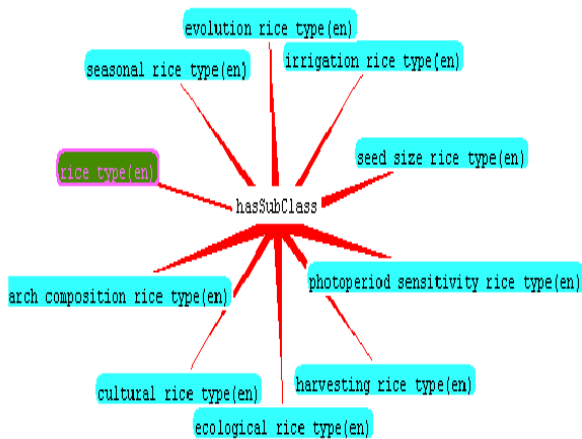


Fig 1: Types of Rice.

To properly consider a facet we need to consider the following elements.

-To show an independent modular domain can describe specific characteristics of a topic which can be seen as independent modularization of domain. For instance, dairy product can be seen in Nutrition.

S.R. Ranganathan [14, 15] was the first to present the notion “facet” in library and information science (LIS). He proposed five different aspects to consider for building facet, PMEST: Personality (P), Matter (M), Energy (E), Space (S) and Time (T). However, his student Bhattacharyya [2] proposed a refinement which consists of four main categories, called DEPA: Discipline (D) (what we call now a domain), Entity (E), Property (P), and Action (A).

In details DEPA can be visualized in the following way:

Discipline (Domain): it includes established field of studies (e.g., Library Science, Mathematics and Physics), applications of traditional pure disciplines (e.g., Engineering and Agriculture), any aggregates of such fields (e.g., Physical Science and Social Sciences), or also more modern terms, fields like music, sports, computer science, and so on.

Entity: the elementary category entity is manifested in conceptual existence. Basically the concept represents the core idea of a domain treated as under this element category. For example: Rice is an entity or concept in Agriculture domain.

Property: it includes characteristic denoting quantities or qualitative characteristics. For example, quality, quantity, Measure, Weight, Taste, etc.

Action: every concept should be considered with the notion of “doing”. It includes processes and steps of doing. An action can manifest as “Self-action” or “External action” which is an action done by some agent (explicit or implicit) on or by itself. For example Imagination, Interaction, Reaction, Reasoning, Thinking and so on are implicit action. An external action is an action done by some agent (explicit or implicit) in a concept of any of the elementary categories described above. For example, Organization, Cooperation, Classification, Cataloging, Calculation, Design and so on are explicit action.

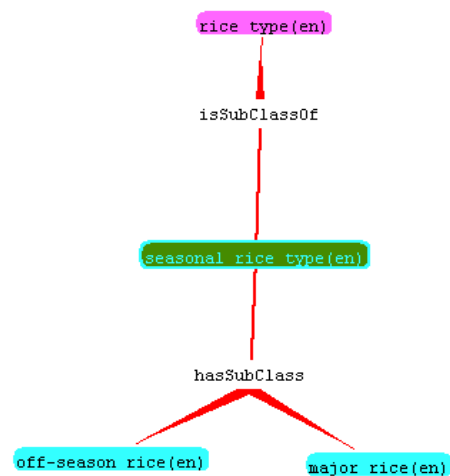


Fig 2: Types of Seasonal Rice.

To build a concept facet, discipline can be taken and so can be entity from DEPA model. Other properties will not be considered in this case. This process can be called semantic factoring. For example, in this experiment the chosen domain or discipline is Agriculture science. In this domain rice is entity or concept. Different kinds of rice exist in this world. Figure 1 [4] shows a distinct module of rice type which is lying in seasonal rice type, cultural rice type, seed size rice type and so on. These types depend on cultural, size, seasonal and others factors. Each of which can be considered as different facet.

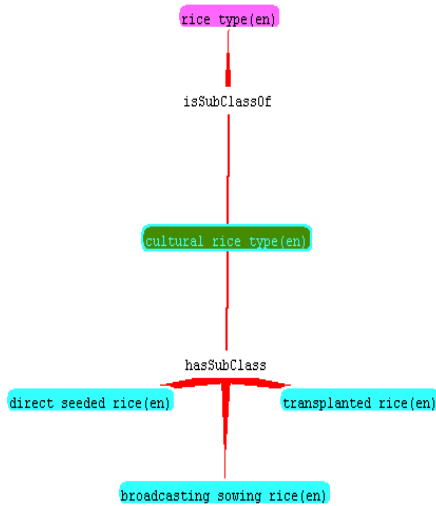


Fig 3: Types of Cultural Rice.

Figure 2 shows one module of rice type which is seasonal types of rice. These types of rice are mostly cultivated in Asian countries like India, Bangladesh, Nepal and Pakistan. These kinds of rice are cultivated during rainy season and it comes out after two or three months. It is totally depended on time factor. Figure 3 shows Cultural rice type. This class of rice mostly cultivated in Thailand. Seeds are cultivated one time in the one place of land. After that it comes out from seeds directly; this type of rice is called direct seed rice. On the other hand, some seeds are cultivated two times. One place is for growing a certain levels of seeds and then another place is for growing fully and then it becomes paddy; these kinds of rice are called transplant rice.

3. Controlled Vocabulary Matching

The problem revolves around the concept of CV matching based on the semantic matching idea described in [8]. The key intuition behind matching controlled vocabularies is the determination of mapping by computing syntactic and semantic relations which hold between the entities of any given two CVs [8, 18]. Let us consider matching 4-tuples $\langle ID_{ij}, c_i, d_j, R \rangle$, $i = 1, \dots, N_C$; $j = 1, \dots, N_D$ where ID_{ij} is a unique identifier of the given mapped element; c_i is the i -th node of the CV1, N_C is number of nodes in the CV1, d_j is the j -th node of the CV2, N_D is the number of nodes in the CV2 and R specify a semantic relation which may hold between the concepts at nodes c_i and d_j . Therefore, the CV matching is defined with the following problem in the light of above discussion,: given two CV T_C and T_D compute the $N_C \times N_D$ mapped element ID_{ij}, c_i, d_j, R with $c_i \in T_C, i = 1, \dots, N_C, d_j \in T_D, j = 1, \dots, N_D$ and R is the strongest semantic relation holding between concepts at node c_i, d_j . Since we look for the $N_C \times N_D$ correspondence, the cardinality of mapping between elements can be determined to be $1 : N$. If necessary, these can also be decomposed straightforwardly into mapping elements with the 1:1 cardinality. For example: We can find out the relationship between cereal and food if we have a mapped vocabulary.

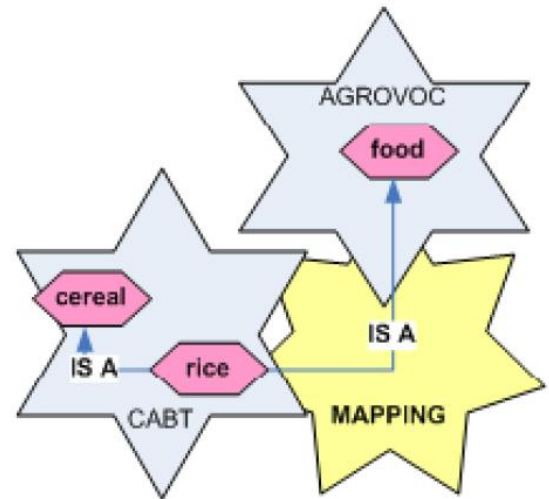


Fig 4: Controlled Vocabulary Mapping and Matching.

4. Concept Facet Matcher (CF-Matcher)

A Concept Facet (CF) contains distinct features for each concept: it includes combined relations, $CF = \langle lg, mg, R \rangle$, where lg identifies less general concepts (one or more), mg identifies more general concepts (one or more) and R identifies related concepts (one or more). In order to realize a matching between two vocabularies (CV1, CV2), we consider the CF from all given CVs's concepts: for every CF of CV1, we check the matching with all CFs of CV2. These concept facets are stored in tables for matching purpose. The methodology of the matching algorithm applied to every concept, can be represented as figure 5..

The matching between two concept facets follows the top-down approach and used several lexical comparison algorithms (SMOA Distance, Hamming Distance, Jaro Measure, SubString Distance, N-gram, JaroWinKler Measure, and Lavestein Distance)[10, 13]. Firstly, comparing started with the more general concepts; if they match (they have same lexicalizations or they are synonyms) assumption has been taken that the concepts under investigation belongs to same concept (they match). Secondly (either we got match or not), Comparing started with the less general concepts based on the results of two mentioned matching, we may obtain exact match (in case more general and less general concepts match), partial match (in case of only one match), or not match. Related concepts of CFs are considered to validate the previous results.

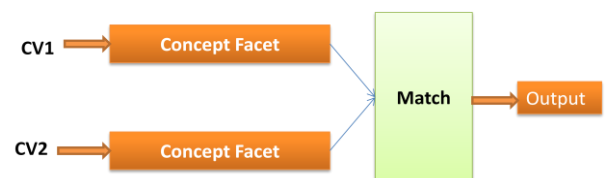


Fig 5: CV Matcher.

In short, it can express CF-matcher the in following algorithm :

Algorithm 1 *buildCFacet(CV)*

```

for i = 0 to CV do
store cF ← (Mg,Lg;R)
end for
return cF

```

In algorithm 1, each controlled vocabulary is taken and stored each concept information in cF. cF is containing more general concepts (BT), less general concepts (NT) and related concepts (RT).

Algorithm 2 *MatchingFacet(CV1,CV2)*

```

cF1=BuildCFacet(CV1)
cF2=BuildCFacet(CV2)
for i = 0 to cF 1 do
for j = 0 to cF 2 do
cfmatcher=elementLevelMatcher(cF1,cF2)
end for j
end for i

```

In algorithm 2, two concept facets are compared using element level matchers and store all matching information in cfmatcher.

5. Results and Evaluation: the AGROVOC and CABI

In our experiments, we used the AGROVOC and the CABI thesaurus because there is no complete mapping between them. The results of the mapping will be published online so that users can use them for better indexing, searching and information retrieval [11, 19].

5.1 AGROVOC

AGROVOC is a multilingual controlled vocabulary designed to cover the terminology of all subject fields in agriculture, forestry, fisheries, food and related domains (e.g. the environment). The AGROVOC Thesaurus was developed by FAO and the Commission of the European Communities in the early 1980s. Since then it has been updated continuously by FAO and local institutions in member countries. It is mainly used for indexing and retrieval data in agriculture information systems both inside and outside FAO. It has approximately 20,000 concepts and four types of relations derived from the ISO standard. Among the available format, it is used the XML version for our task [16].

5.2 CABI

CABI is a monolingual controlled vocabulary designed to cover the terminology of all subject fields in agriculture, forestry, horticulture, soil science, entomology, mycology, parasitology, veterinary medicine, nutrition and rural studies. The CABI thesaurus was developed by CABI which is a not for profit purpose, science-based development and information organization. It has 48,000 concepts and four

types of relationship derived from the ISO standard. We obtained data as text format and converted it to XML format for experiment purposes [17].

5.3 Results and Evaluation Descriptions

We started our experiments using 1000 concepts from each controlled vocabulary. Managing all concepts was a challenge because the two vocabularies are not organized in the same structure. We converted each vocabulary to the same format in order to conduct the test. We obtained 325 exact matches, 550 partial matches and 125 not matches concepts from FALCON-AO. Also, we obtained 175 exact matches from tested CF-Matching algorithm, but we found different numbers of partial matches from eight element label matchers. SMOA Distance matcher gives more partial matches than others. Hamming distance, Jaro Measure, SubString Distance, and N-gram which do not give a satisfactory numbers of matches. JaroWinKler Mesasure and Levestein Distance produce quite similar results. However, we got 465 partial matches (average) and 360 not matches (average) concepts. Furthermore, we choose FALCON-AO (Automatic Ontology Matching tool) because it has given the best results according to mapping evaluation report. In our experiments, we considered 0.19 as our given threshold value for partial match and 1.0 for exact match. Figure 6 shows exact match between AGROVOC concept number c 635 and CABT concept number 11576. Similarly, figure 7 shows partial match between AGROVOC concept number c 3500 and CABT concept number 42585. We got these results from FALCON-AO and CFMatcher. But according to our Domain expert at FAO, figure 6 shows matching results.

Matching Results

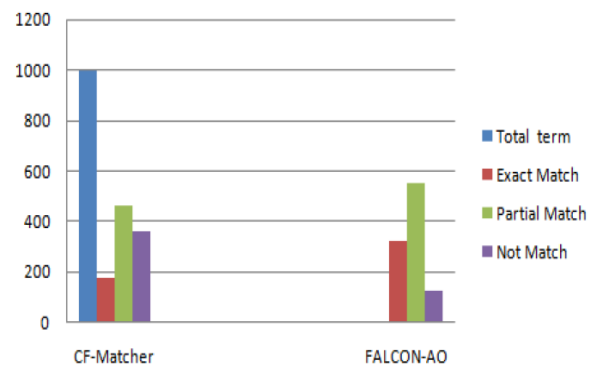


Fig 6: CV Matching results.

Figure 7 shows correct results and figure 8 partial results. Because there is no relationship between “Basella” and “Ballasts”. From Dictionary, In Figure 9, we presented a human readable prototype so that one can access concepts information from two thesauri and see mapping results. The domain experts can validate the results and this information is stored into the database. We faced lots of challenges during our experiments. Overall, lots of data were overlapping and two automatic tools gave some partial matches which were not correct according

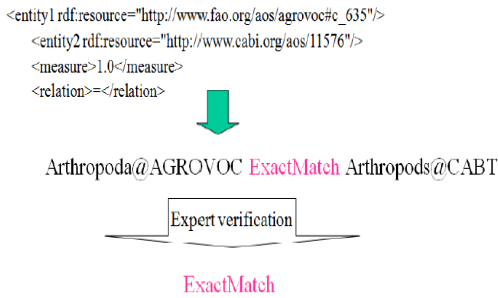


Fig 7: Exact Match.

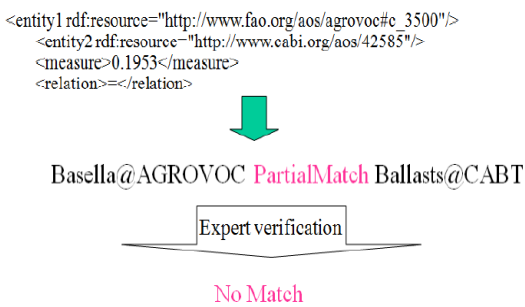


Fig 8: Partial Match.

AGROVOC			CABI		
Search	<input type="text"/>	Search	Search	Land	Search
Available Terms@AGROVOC			Available Terms@CABI		
Page:1			Page:1		
Letter	Concept Name	Term Code	Letter	Concept Name	Term Code
A	Agribusiness	1001	A	Agribusiness	441
A	Agricultural Land	1002	L	Land	442
			B	Biology	4433
Navigation			Navigation		
prev 1 next			prev 1 next		
AGROVOC		Relationship	CABI		
Agricultural Land		=	Land		
Exact Match		<input checked="" type="checkbox"/> Valid	<input type="button" value="submit"/>		

Fig 9: CV Matching system prototype.

to experts. However, our correct mapping results, after verification and validation, will be used for searching purpose.

6. Application toward Linked Open Data

The development of the Web of Data, built by applying Linked Data (LD) (Berners-Lee, 2011) (Heath, 2011) principles and using Semantic Web technologies, is gaining great attention in the academic as well as the industrial world. This is the frontier[21] of data integration and sharing. In a web where each piece of data is published by means of standard technologies and data formats, and where each piece of data can be univocally named and located, data integration (understood as the possibility of programmatically accessing

data residing in different sources) is perceived to be closer now than ever before.

For the bibliographic and librarian world, Linked Data [21] offers the technology and the social attention needed to publish and interlink metadata sets. If, for example, a term in the AGROVOC thesaurus is linked with a term in the GEMET thesaurus, all documents indexed by the same term in the document repositories related to AGROVOC and GEMET are also potentially linked. Using appropriate applications, information queries can be submitted against both repositories, and data results presented (and processed) to the user in a unified way. For this reason, many thesauri are adopting the Linked Data approach to data publishing. This paper presented a work on aligning AGROVOC with one relevant thesaurus, in order to publish AGROVOC as Linked Data. However, it is possible to check the rest of the thesauri with the same logical consequences.

In this paper, it has been shown a system for automatic vocabulary matching using concept facets. The proposed system convinced that it helps for better information searching, browsing, and extraction in agriculture and related domains. There are some open research issues: the semantic heterogeneity between two controlled vocabularies in a single domain; the multi-word concepts; the possibility of automatically link non-matched concepts to external reliable resources such as public thesauri, encyclopedia or dictionaries. Now, current work is extending for semantic search for Agricultural domain as near future focus.

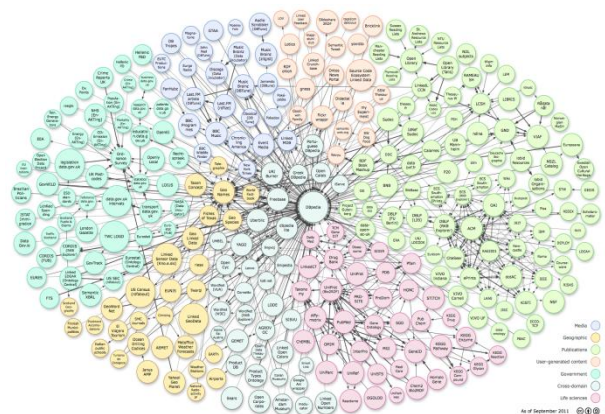


Fig 10: AGROVOC in Linked Open Data cloud.

Steps to align Linked data [21] version of AGROVOC with any other controlled vocabulary:

- All the thesauri are loaded as single local triple store in the format of SKOS-RDF.
- Then compute the algorithm-2 which is mentioned earlier in section 4. The proposed approach supports similar type of language for the time being.
- Rather than only comparing the string similarity, proposed approach does computation by eliciting the relation between the associated nodes.
- In order to combine these similarity values into a single number, the proposed approach computed the arithmetic average of all similarity values, as the simplest way to combine several values, which seemed to us appropriate for a first attempt. Finally, an empirically identified threshold was applied to

select the candidate matches to pass to human evaluation.

- This allows the system to facilitate AGROVOC together with its entire outbound links at the same time. The approach used Pubby¹ to serve as frontend of the data repository: AGROVOC is now published in the style of Linked Data publishing.

The work presented in this paper was not very user-friendly due to the constraint of existing visualization and interaction methods used in the project. Increasing usability of the interface is the future plan of the proposed approach. What we have tried to do was to demonstrate its potential for alternative use that was not anticipated as a concrete idea in the project agenda. The work was completed, tested and accepted in the project for a full-scale implementation and integration in future version, certifying its usefulness and quality.

7. Conclusion

In this paper, it has shown with the proposed system for automatic vocabulary matching using concept facets. It can be convinced that it helps for better information searching, browsing, and extraction in agriculture and related domains. There are some open research issues: the semantic heterogeneity between two controlled vocabularies in a single domain; the multi-word concepts; the possibility of automatically link non-matched concepts to external reliable resources such as public thesauri, encyclopedia or dictionaries. As future plan, we are extending our work for semantic search and semantic tagging for Agricultural domain.

Acknowledgement

Authors would like to thank Prof. Fausto Giunchiglia, Johannes Keizer, and Gudrun Johannsen for their valuable suggestions. Also, they would like to thank Shaun Hobbs of CABI for kindly providing the data files.

8. REFERENCES

- [1] A.Morshed and M.Sini. Creating and aligning controlled vocabularies. In *Advance Technology for Digital Libraries, AT4DL*, Trento, Italy, 2009.
- [2] Bhattachary.G. Popsi:its fundamentals and procedure based on a general theory of subject indexing language. In *Library Science with a slant to Documentation*, volume 16, pages 1–34, 1979.
- [3] P. Bouquet, L. Serafini, and S. Zanobini. Semantic coordination: a new approach and an application. In *Proc. of the 2nd International Semantic Web Conference (ISWO'03)*. Sanibel Islands, Florida, USA, October 2003.
- [4] D AGROVOC concept server. <http://naist.cpe.ku.ac.th/agrovoc/>.
- [5] P. Shvaiko F. Giunchiglia and M. Yatskevich. Discovering missing background knowledge in onology matching. In *17th European Conference on Artificial Intelligence (ECAI 2006)*, volume 141, pages 382–386, 2006.
- [6] M.Marchese F.Giunchiglia and I.Zaihrayeu. Encoding classifications into lightweight ontologies. *Data Semantics VIII*, pages 57–81, 2007.
- [7] F. Giunchiglia, B.Dutta, and V.Maltese. Faceted lightweight ontologies. In *LNCS*, 2009.
- [8] F. Giunchiglia, P. Shvaiko, and M. Yatskevich. S-match: An algorithm and an implementation of semantic matching. In *Proceedings of ESWS'04*, 2004.
- [9] F. Giunchiglia and I. Zaihrayeu. Lighweight ontologies. Technical report at DIT, the University of Trento, Italy, October 2007.
- [10] J.Euzenate and P.Shaviko. *Ontology Matching*. Springer, 1st edition, 2007.
- [11] Sini M. Chang C. Li S. Lu W. He C. Liang, A. and J. Keizer. The mapping schema from chinese agricultural thesaurus to agrovoc. In *Proceedings of the fifth Conference of the European Federation for Information Technology in Agriculture, Food and Environment and the thirdWorld Congress on Computers in Agriculture and Natural Resources*, 2005.
- [12] Bernardo Magnini, Luciano Serafini, and Manuela Speranza. Making explicit the semantics hidden in schema models. In: *Proceedings of the Workshop on Human Language Technology for the Semantic Web and Web Services*, held at ISWC-2003, Sanibel Island, Florida, October 2003.
- [13] Natalya F. Noy. Semantic integration: a survey of ontology-based approaches. *SIGMOD Rec.*, 33(4):65–70, 2004
- [14] S.R Ranganathan. *Element of library classification*. Asia Publishing house.
- [15] S.R Ranganathan. *Prolegomena to library classification*. Asia Publishing house.
- [16] Agrovoc thesaurus. <http://www.fao.org/agrovoc>
- [17] CABI thesaurus. <http://www.cabi.org/>.
- [18] Duncan J. Watts and Steven H. Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 393:440–442, 1998.
- [19] L.Finch H. Kolb W.Hage, M.Sini and G.Schreiber. The oaei food task:an analysis of a thesaurus alignment task.
- [20] I. Zaihrayeu, L. Sun, F. Giunchiglia, W. Pan, Q. Ju, M. Chi, and X. Huang. From web directories to ontologies: Natural language processing challenges. In *ISWC/ASWC*, 2007.
- [21] Ahsan Morshed, Caterina Caracciolo, Gudrun Johannsen, Johannes Keizer. *Thesaurus Alignment for Linked Data Publishing*. *Proc. Int'l Conf. on Dublin Core and Metadata Applications 2011*.

¹ <http://www4.wiwiss.fu-berlin.de/pubby/>