# Developing Ontology for Arabic Blogs Retrieval

### Lilac Al-Safadi
Department of Information Technology,

College of Computer and Information Sciences,

King Saud University

### Mai Al-Badrani
Department of Information Technology,

College of Computer and Information Sciences,

Saud University

### Meshael Al-Junidey
Department of Information Technology,

College of Computer and Information Sciences,

King Saud University

## ABSTRACT

The ability to search for content on the Internet has proven to be essential for many. However, available search engines supporting Arabic language are typically limited to keyword searches and do not take in consideration the underlying semantics of the content. Semantic search engines provide searching and retrieving resources conceptually related to the user informational need. However, the presented technologies support mostly languages using Latin scripts. Arabic is still not well supported. The few works on Arabic Semantic Web applications are derived from the rule governing the traditional Arabic language and not the content available on the Web. In this research, we propose a model for representing Arabic knowledge in the Computer Technology domain using Ontologies. The model starts by elicitation users' informational needs. Ontologies will play a major role in supporting information search and retrieval processes of Arabic blogs on the Web.

## General Terms

Semantic Web, Ontology Engineering

## Keywords

Ontology, Semantic Search Engine, Knowledge Sharing, Ontology Engineering, Arabic Blog Search Engine.

## 1. INTRODUCTION

The development of the Semantic Web initiative is rapidly increasing the number of publicly available Ontologies. Ontology is the backbone of the Semantic Web; therefore, the success of any Semantic Web application depends on the design and development of Ontology. In terms of computer and information science, Ontology is a formal representation of the knowledge by a set of concepts within a domain and the relationships between those concepts (Wikipedia.org). The use of Ontologies is necessary for sharing common understanding of the structure of information among people or software agents [17, 20, 21]. Contemporary Ontologies share many structural similarities, regardless of the language in which they separate domain knowledge from operational knowledge [18]. Yet, each language has its own unshared features that need to be taken into consideration while Ontology modeling.

One of the main application areas of Ontology technology are semantic search engines and information retrieval (IR). In search engines, software agents are used to extract information to identify target Websites. Semantic search engines provide searching and retrieving resources conceptually related to the user informational needs. It performs content-based search of documents on the Web focusing on the semantic structure of the content rather than the syntactic. Swoogle (swoogle.umbc.edu), Hakia (www.hakia.com), SenseBot (www.SenseBot.com) and DeepDyve (www.deepdyve) are among the top semantic search engines which provide access to the "Deep Web" in which the digital content are not directly accessible by generic search engines. These semantic search engines have weak to no support of Arabic language. Different languages have contained the specific linguistic environment and the cultural context, which has caused the need to develop different Ontology for different information language [15]. Since most of the Ontologies publicly available are in English language, there is a strong need for Arabic language Ontologies used as the basis of Arabic Semantic applications.

Ontologies capture a domain knowledge in a formal way. The semantic field theory is based on an analytical approach, which considers the meaning of a word within a given view of the world [2]. Therefore, the process of knowledge retrieval is mainly domain dependent. Computer subjects are becoming a major concern to companies, organizations, communities, and even people in different countries. In 2010, Computer subject related blogs have top rank exceeded 35% of the total number of Arabic blogs on the Internet (www.openarab.net). Three tools were used to search for Ontologies available on the Internet, these are Swoogle (swoogle.umbc.edu), Twiki (www.twiki.org), and DAML (www.DAML.org). Our findings were that there are more than a thousand available Ontologies in the computer domain. Yet, the existing computer Ontologies proved to be unsatisfactory for the Computer Technology specific domain since they are mainly concerned with research area and academic programs. In the Computer Technology domain, there are over hundred of English Ontologies exist on the Internet, and no Arabic Ontology.

In this paper we introduce a domain-dependent Ontology for searching Arabic blogs in the Computer Technology domain. We propose a model for designing the Ontology which is based on structuring the Arabic language into a set of equivalent classes, properties and relationships. We address one of the main challenges of Ontology design which is the representation of classes. The developed Ontology can also be used as a basis for other applications such as content aggregators of Computer Technology Websites or news. Our study shows that designing and developing Arabic Ontologies need more than what is provided by keyword-based search engines and more than analyzing the morphology and grammar of the traditional Arabic language which support our hypothesis that searching Arabic blogs on the Web do need extra support.

The paper is organized as follows; section 2 discusses works in Arabic Semantic Web. Section 3 analyzes the Arabic language used on the Web and the content need to be presented in the Arabic Ontology. The method for designing the Ontology is illustrated in section 4. The Ontology implementation is described in section 5. Section 6 presents experiment and test results of the proposed Ontology. Section 7 discusses those results. Finally, Section 8 concludes our paper.

## 2. WORKS IN ARABIC SEMANTIC WEB

Due to the increasing number of Arabic content on the Web, an application is needed to exploit the large amount of information. Certain considerations need to be taken when designing Arabic-based applications. Most of the works in Arabic Semantic Web are driven by the traditional Arabic Language structure and rules governing the formation of its vocabulary [3, 8, 9, 10, 11, 14]. Although these works are promising, no experiments were provided on real data over the Web to decide if these designs cater to the thousands of the blogs users on the Web.

Arabic Ontology is the foundation of the creation of semantic-based applications that supports content in Arabic language. Belkredim [8, 9] focuses on developing Ontology using verbs and roots. Verbs are classified according to derivation rules of the Arabic language. The root is a set of three or four letters فعل and يفعل. 85% of Arabic words are derived from triliteral roots [19]. Belkredim works describe a theoretical model but no implementation is provided to support this hypothesis.

Basing an Ontology on roots is imprecise because some terminologies with different meaning have the same root. For instance, words like {معلم، تعليم، علم} (what is written, Science, Education and Teacher) are three totally different concepts in a knowledge domain with different relations, all share the same root and derived from علم (to know). Another example in the Computer Technology domain is illustrated in the below figure. It shows that the root خدم (to serve) is shared by a number of different concepts {خدمة، خادم، مستخدم} (what is written, User, Server, and Service). In addition, some words have no roots for instance دولاب (what is written, wheel or trundle).
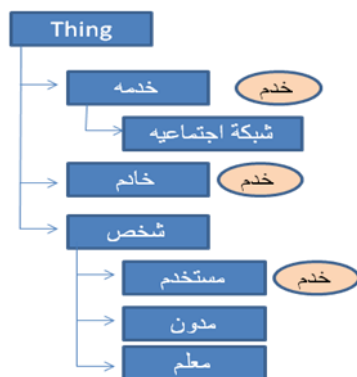


**Fig 1: Classification of concepts which share the same root**

## 3. UNDERSTANDING THE NEEDS OF ARABIC BLOG USERS

Often, knowledge is elicited from the subject matter experts in a particular domain. One of the most important things to remember

is that it is not enough to check the domain knowledge, but also what the Ontology are going to be used for, the types of questions the information in the Ontology should provide answers, and the users of the Ontology [12]. We focus on a user-centric approach where we emphasize on eliciting the need of the user before an expert is consulted in order to know what kind of knowledge is required to satisfy the user's needs.

We have conducted experiments on a number of randomly selected real-world Arabic blogs in the Computer Technology domain. A sample of the set of Arabic blogs is listed in table 1 below.

**Table 1. Sample of analyzed Arabic blogs in the Computer Technology Domain**

| Blog name | URL |
|---|---|
| أنا سيسكو | http://iam.cisco.com.sa |
| فص ولزق | http://www.qaswlasq.com |
| اتكلم | http://itkallem.com |
| تكنولوجيا العرب | http://tech-arabic.blogspot.com |
| Web2.0 | http://khaledalhourani.com/blog |

We preprocessed the collected data set by removing stop words, and numbers and focusing on terminologies that appear more than three times in the corpus. We found that 34.38% of data set are English words, 10.63% of data set are words expressed by users (modern slang made from an English word) such as { ويب سايت}, which has equivalence in tradition Arabic language to {موقع ويب and موقع}, 30.62% of data set are expressed in tradition Arabic language such as {جهاز and محرك بحث}, which is equivalent to search engine and hardware, and 24.37% are slang words with no equivalence in Arabic language such as { قوقل and غوغل}, which is equivalent to Google. The tables below show a sample of our collected data set.

**Table 2. Sample data set – Words and Arabic root**

| Word | Arabic Root |
|---|---|
| قاعدة بيانات | قعد / بين |
| مدونة | دون |
| تطوير التطبيقات | طور / طبق |
| متصفح | صفح |
| المجلدات | جلد |
| حاسب محمول | حسب / حمل |
| لوحة مفاتيح | لوح / فتح |
| تقنيين | تقن |
| تصميم | صمم |

**Table 3. Sample data set – Modern Arabic words with equivalence in Traditional Arabic language**

| Modern Arabic words |
|---|
| بلاك بيري |
| بلاي ستيشن |
| ديفيانت ارت |
| الانترنت |
| توينتر |
| جوجل ريدر |
| آبل |
| فايلاتي |
| كلاسات |

**Table 4. Sample data set – Modern Arabic words with no equivalence in Traditional Arabic language**

| Modern Arabic words |
|---|
| سريدبروو |
| الويب |
| لسكبي اجمي |
| لينكسس |
| وي |
| أجاكس |
| جافا سكريبت |
| فوؤل |
| ابفون |

The high percentage of slang words in our resulted data set is an impact of the Internet and social networking on the Arabic language. Our experiment conform that users need additional information while searching Arabic blogs. The identification of users' needs will make a future built system to be more acceptable and more likely to be used [16]. Existing works Arabic Ontologies proved to be unsatisfactory for the specific purpose that centers on Computer Technology. They were deemed inadequate because they focus on traditional Arabic language other than been developed for their own specific purpose.

Since the Ontology is a core element for Semantic Web applications, this paper was dedicated to discuss the modeling approach that was adopted to cover all possible aspects needed in creating the Arabic Ontology that is based on this model. The Ontology contains all needed concepts and logical rules and requirements that form the basis of the application. Based in the results of the conducted experiments, the approach used in this paper is developing Ontology using nouns in singular form, and combining slang words as classes in the Ontology with *equivalent* type of association.

# 4. ONTOLOGY ENGINEERING

Ontological engineering approach was followed in developing the proposed Arabic Ontology [7]. The ontology was designed based on WSMO [5] framework for modeling semantic web services. A model driven architecture is used [4] and forward engineering approach was adopted where we started by modeling the Ontology first and then using this ontology as a domain model to form the basis of the generation of the Semantic search engine. The process of Ontology engineering encompasses many phases. It starts with eliciting the domain knowledge to be represented by the Ontology.

A domain knowledge is represented by basic categorization of terminologies in the domain. The interrelationship between one terminology and another that relates to its meaning can also result to the presentation of knowledge. Usually Ontology can be built using domain experts or learned from information available in a corpus of the domain. The goal of Ontology learning is to automatically extract relevant concepts and relations from the given corpus or other kinds of data sets to form Ontology (Wikipedia, 2011). A user-centric approach was used in developing our proposed Arabic Ontology, which takes into
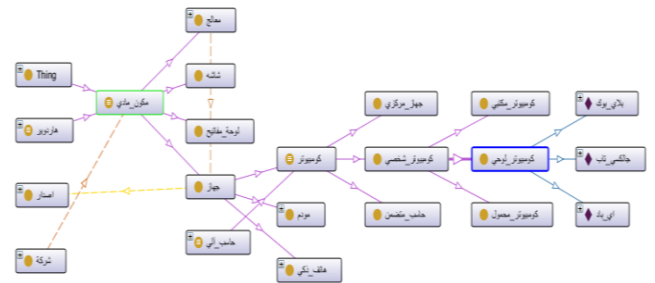
consideration the user's of the Ontology and the questions the Ontology is expected to answer.

The life cycle of Ontology development can also be subdivided into the following stages: extracting terms, discovering synonyms, obtaining concepts, extracting concept hierarchies, defining relations among concepts, deducing rules or axioms [15]. The result is a formal explicit description of *concepts* in a domain of discourse, *properties* of each concept describing various features and attributes of the concept, and *relationships* between classes. Ontology together with a set of *individuals* of classes constitutes a knowledge base.

## 4.1 Classes

The developed Ontology consists of 110 classes, 80.9% are specific to the Computer Technology domain. It is impossible to cover all the terminologies of Computer Technology domain in one Ontology. Instead, we try to provide an initial Ontology specialized in Arabic Computer Technology domain. Enabling reuse of domain knowledge was one of the driving forces behind recent surge in Ontology research [14]. Others can simply reuse the developed Ontology and extend it to describe their domain or serve their application of interest.

Figure 2 shows a sub graph of the developed Ontology with classes and their relationship. Rectangles represent classes, while lines ending with an arrow indicate an association between the two classes. Class names are in Arabic. Some represents objects in traditional Arabic language and others in modern Arabic language associated with *equivalent* relationship.



**Fig 2: Sub-graph of the proposed Arabic Ontology in the Computer Technology Domain**

The Ontology relevant terms were gathered from domain users. The terminologies were expanded using sources like Computer Terms Dictionary [1], Computer Technology Ontologies such as ittags.owl[1] and 07Jun10.owl[2], and Google sets utility (labs.google.com/sets). Google English-Arabic translation tool (translate.google.com) was used for automatic translation of the content of the English Ontology, which translated the English classes into 15 Arabic classes. The rest of translations were made manually.

We also looked at domain specific articles on the Web like wikipedia.org, buzzle.com, webopedia.com and articlesbase.com. As an example, the word انترنت (Internet) was expanded to شبكة

---

[1] www.ittags.com/ontology/ittags.owl

[2] what.csc.villanova.edu/twiki/pub/Main/OWLFileInformation/07Jun10.owl

عنكبوتية (Web) in the domain specific articles from the Web, and شبكة اتصال (Network) was expanded to شبكة عالمية (Global Network) in the computer terms dictionary, and مدونة (Blog) was expanded to انترنت (Internet) using Google sets utility. These tools proved quite useful because they not only sorted most of the words already found under specific categories, but they provided even more concepts for the Arabic Computer Technology Ontology. Nevertheless, the terms and definitions in the initial terminology did not represent the final state of the terms and definitions that were included in the domain Ontology. They were rather a first-draft or gloss for the sake of getting the relevant information organized and assembled in a single place.

The next two steps are developing the class hierarchy and defining properties of concepts.

## 4.2  Individuals and Properties

Individuals in the Ontology are instances of predefined classes. For instance فيسبوك (Facebook) is an individual of شبكة اجتماعية (Social Network). The developed Ontology consists of 78 individuals. Individuals must have features or describes the object itself. For instance, اسم شركة (Company name) is an object property of object شركة (Company).
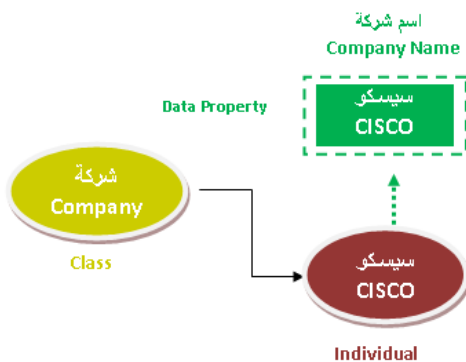


**Fig 3: The representation of class شركة (Company), individual سيسكو (CISCO), and property اسم شركة (Company Name)**

## 4.3  Relationships

Ontology combines the "things" (classes and instance) that might appear together under a certain category. In our case, the category is the Computer Technology field. After analyzing the classes that determine the semantic dimensions, we organized the terminologies into a hierarchy associated with components via Ontological relations. Sample of the relationships which exist between instances in the Ontology are تنتجه شركة, لها شعار and يستعمل (what is written, has logo, produced by, and use). 48 object relations were defined in the proposed Ontology of different types;

- The simplest relation which appears in any Ontology is inheritance. It relates concepts in generalization and specialization relation. The inheritance relation usually organizes concepts in a tree form where child nodes are connected to parent nodes by means of unidirectional "is-a" relationships (inheritance). For instance, جهاز (Device) is a sub class of مكون مادي (Hardware).
- Equivalent classes, the relationship between these nodes is a bidirectional "is-a" object property. This relationship is used in

our proposed Ontology to associate concept representation in Modern Arabic and their equivalence in Traditional Arabic. For instance, هاردوير is an equivalent class of مكون مادي (Hardware).

- Others, the relationship between these nodes is different between each pairs according to their features.
  • Hypernym, the relationship between nodes is a unidirectional "type-of". For instance, كومبيوتر (Computer) is kind-of جهاز (Device).
  • Hyponymy shares a "type-of" relationship with its hypernym. For instance, لوحي (Computer Tablet) is type-of كومبيوتر (Computer) and كومبيوتر محمول (Laptop) is type-of كومبيوتر (Computer).
  • Coordinate terms, the relationship between nodes is a bidirectional "X & Y shares a hypernym". For instance, كومبيوتر محمول (Computer Tablet) and كومبيوتر لوحي (Laptop).
  • Holonym, the relationship between nodes is "part-of". For instance, انترنت (Internet) is part-of ويب (Web).
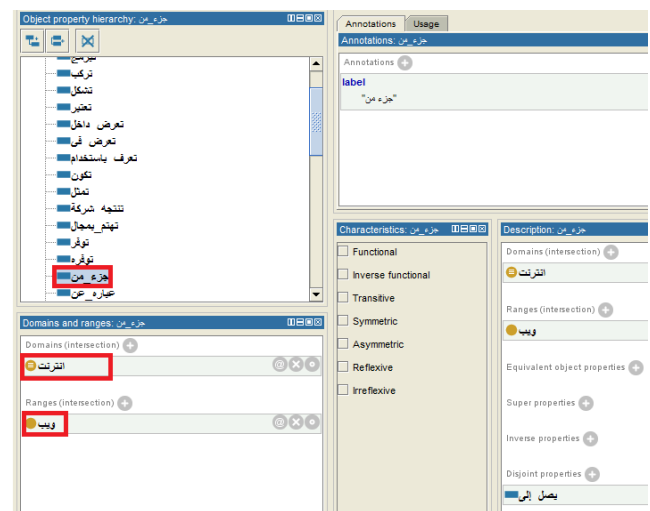


**Fig 4: The representation of the triple ويب (Web), جزء من (part-of), انترنت (Internet) in protégé 4.1**

## 5.  ONTOLOGY IMPLEMENTATION

To enable interoperability, sharing and reuse of this valuable resource, the developed Ontology is implemented using Web Ontology Language (OWL) [6]. OWL is the most recent development in standard Ontology languages, endorsed by the World Wide Web Consortium (W3C) to promote the Semantic Web vision.  We build the Ontology using Protégé-4.1 as Ontology-editing environments.  Protégé-4.1 (protege.stanford.edu) is a free open-source platform used to describe Ontologies declaratively, stating explicitly what the class hierarchy is and to which classes individuals belong. Protégé-4.1 was capable of displaying Arabic script, yet it does not include a SPARQL Query Panel. We had to use protégé 3.4.4 to test our sample queries using SPARQL query.

One of the main challenges faced us in using Semantic Web tools was there support of Arabic script. Jena is not compatible with Arabic language. In NetBeans the Arabic text does not appeared in the proper Arabic text as shown in figure 5.

ResultSetFormatter Jena's method was used to solve part of the support as well as setting the system locale, see figure 6.
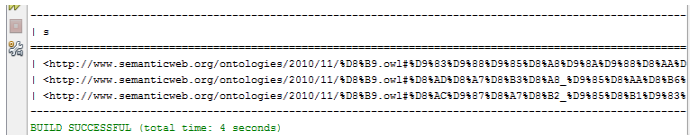


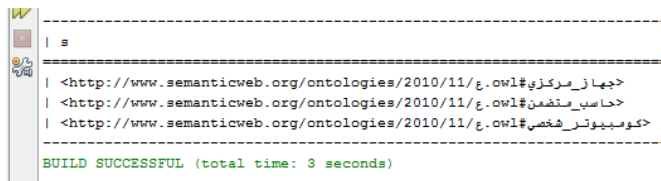**Fig 5: The representation of Arabic script in Netbeans**.



**Fig 6: The representation of Arabic script after using Jena's method**

## 6. EXPERIMENT RESULTS

A number of experiments were carried out to validate in practice the Ontology's ability to bridge the semantic gap. Below we describe an Arabic query imposed on our proposed Ontology.

The first experiment show the results of an Arabic query using a property description "تفاحة مَضوعة" (what is written Bitten Apple), which is the description of the logo of Apple Company. The query structured as follows

```
SELECT ?products
WHERE
  { ?products  ع:تنتجه_شركة  ?s
    { ?x  rdfs:label  "تفاحة مقضومة".
      ?s ?w ?x.  } }
```

This Query structure was tested in Protégé 3.4.4 SPARQL Query Panel. The average precision rate of experiment results is 50%.
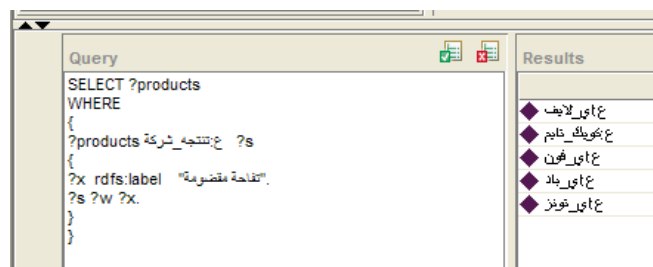


**Fig 7: The query result in protégé 3.4.4**

Follows is the conversion of the above query into Jena API format built as a pilot program to test the support for the Arabic query in Jena.

```
String CqueryString = "PREFIX dc:<http://purl.org/dc/elements/1.1/>"+
  "PREFIX
ع: <http://www.semanticweb.org/ontologies/2010/11/ع.owl#>"+
  "PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>"+
  "SELECT ?products WHERE { ?products ع:تنتجه_شركة ?s
  { ?x  rdfs:label \""+inword+"\"." + "?s ?w ?x. }}
  OFFSET 0 LIMIT 1000";
```

## 7. DISCUSSION

In the Arabic Semantic Web world, morphological analysis based on traditional Arabic would not return the appropriate result. User requirement is an essential step and need to be integrated in. A combination of traditional and modern analysis of Arabic content can enhance the accuracy of our results. In this research, we attempt to build Arabic Ontology where we integrate both traditional Arabic with modern Arabic words widely used in the domain knowledge.

Our focus in this paper is to implement the Arabic Ontology. We do recognize the importance of combining the vocabulary and the morphology of a language with its semantics. Our future prototype will focus on presenting a complete framework and that would combine the develop Arabic Ontology with NLP, and investigate how this would affect the precision and recall of Arabic search engines.

## 8. CONCLUSION

In this paper, we have described the development of an Arabic Ontology in a Computer Technology domain to serve semantic-based search and retrieval of Arabic blogs on the Web. We analyzed the Arabic language on the Web and investigated the existing Arabic support offered by Semantic Web applications and research. The analysis showed weak support for traditional Arabic language and almost no support for modern Arabic language, which is becoming today's blogs language. Thus, the need for developing domain-based Ontologies for combining traditional Arabic and modern Arabic is crucial. We listed the steps in the Ontology-development. One of the most important things to remember is that it is not enough to check the domain knowledge, but also to analyze the content in which the Ontology are going to use be used for, and the users of the Ontology. Building semantic based search engine for blogs using traditional Arabic language undoubtedly insufficient. An application is as good as the Ontology built for.

## 9. ACKNOWLEDGMENTS

## 10. REFERENCES

**Arabic Reference**

[1] الكياـني, تيسير :معجم الكياـني لمصطلحات الكومبيوتر والإنترنت. طـ2, لبنان: مكتبة لبنان ناشرون2004.

**English Reference**

[2] A. Lehrer, "Semantic Fields and Lexical Structure". New York: American Elsevier, 1974.

[3] B. Hammo, H. Salem, S. Lytinten, and M, Evens. "QARAB: A Question Answering System to Support the Arabic Language". In: Proc. of the workshop on Computational approaches to Semitic languages, ACL, pages 55-65, Philadelphia, 2002

[4] C. Miller, J. Mukerji, C. Burt, D. Dsouza and K. Duddy et al., 2001. Model Driven Architecture (MDA). Document number ormsc/2001-07-01, Architecture Board ORMSC, OMG

[5] D. Fensel, "Ontologies: Silver Bullet for Knowledge Management and Electronic Commerce". 1st edition. Springer-Verlag, Berlin, 2001.

[6] D. McGuinness, et al, "OWL Web Ontology Language Overview (W3C Recommendation 10 February 2004)," http://www.w3.org/TR/2004/REC-owl-features-0040210/, 2006-01-20.

[7] D. Nicola, M. Missikoff, and R. Navigli, "A Software Engineering Approach to Ontology Building", Information Systems, 34(2009), pp. 258–275.

[8] F. Belkredim and F. Meziane, "DEAR//-ONTO: A DErivational ARabic Ontology Based on Verbs", International Journal of Computer Processing of Languages, 21(3):279-291, 2008.

[9] F. Belkridem, and A. El Sebai, "An Ontology Based Formalism for the Arabic Language Using Verbs and Derivatives", Communications of the IBIMA, 11(2009), pp. 44–52.

[10] H. Aliane, Z. Alimazighi, and A. Cherif Mazari, "Al-Khalil: The Arabic Linguistic Ontology Project", in Proceedings of the Seventh Conference on International Language Resources and Evaluation (LREC'10), Valletta,Malta, 2010.

[11] H. Al-Khalifa, M. Al-Yahya, et al. "SemQ: A Proposed Framework for Representing Semantic Opposition in the Holy Quran using Semantic Web Technologies". In proceedings of CTIT09 IEEE CNF, Dec 2009.

[12] K. Darwish, "Building a Shallow Arabic Morphological Analyzer in One Day", Proceedings of the Computational Approaches to Semitic Languages, A workshop affiliated with ACL-2002, University of Pennsylvania, USA, 2002.

[13] L., Al-Safadi, N. Abdulateef, "Educational Advertising Ontology: A Domain-Dependent Ontology for Semantic Advertising Networks" Journal of Computer Science 6(10): 1041-1048, 2010.

[14] M. Al-Yahya, H. Al-Khalifa, A. Bahanshal, I. Al-Odah, N. Al-Helwah, "An Antological Mdoel for Representing Semantic Lexicons: An Application on Times nouns in the Holy Quran". The Arbaian Journal for Science and Engineering, 35(2c), 2010, 21-37

[15] M. Beseiso, A. Ahmad, R. Ismail, "A Survey of Arabic Language Support in Semantic Web". International Journal of Computer Applications. 9(1), 2010, 35-40

[16] M. Ismail, M. Yaakob, and S. Kareem, "Semantic Search Engine in Institutional Repository : An Ontological Approach", Proc. International Conference on Library and Information Science (ICOLIS 2007), 26 - 27 June 2007, Kuala Lumpur Malaysia, pp. 55-63. ISSN 978-983-43491-0-3.

[17] M. Musen, "Dimensions of knowledge sharing and reuse. Computers and Biomedical Research", 1992, 25: 435-467.

[18] M. Nieto, "An overview of Ontologies", Technical report, Center for Research in Information and Automation Technologies and Interactive and Cooperative Technologies Lab, Universidad De Las Americas Puebla, starlab.vub.ac.be/teaching/Ontologies_overview.pdf> 20 July 2007.

[19] S. Elkateb, W. Black, H. Rodriguez, M. Alkhalifa, P. Vossen, A. Pease, and C. Fellbaum, "Building a WordNet for Arabic", Proceedings of The fifth international conference on Language Resources and Evaluation, 2006

[20] T. Gruber, "A Translation Approach to Portable Ontology Specification". Knowledge Acquisition, 1993, 5: 199-220.

[21] "Why would someone want to develop an Ontology?", All1Source Technologies http://www.all1sourcetech.com/develop-Ontology.