

# An Ontological Model for Representing Computational Lexicons

## A Componential Based Approach

Maha AL-YAHYA<sup>1</sup>, Henda AL-KHALIFA<sup>2</sup>

Information Technology  
Department,

King Saud University, Riyadh,  
Saudi Arabia

<sup>1</sup>malyahya, <sup>2</sup>hendk@ksu.edu.sa

Alia BAHANSHAL<sup>3</sup>, Iman AL-  
ODAH<sup>4</sup>

Research Institute, King  
Abdulaaziz City for Science and  
Technology  
Riyadh, Saudi Arabia

<sup>3</sup>abahanshal,  
<sup>4</sup>ialoudah@kacst.edu.sa

Nawal AL-HELWAH

Arabic Language Department  
Princess Norah University  
Riyadh, Saudi Arabia

**Abstract**— In the last decades the computational linguistics community has developed important and widely used lexical resources. Although they are very popular among the Natural Language Processing (NLP) community, they do not address two important characteristics of language. The first is that the meaning of a word in a language is a collective effort defined by the people who use the language. The second is that language is a dynamic entity (some words change their meaning, others become obsolete, new words are born). A computational model which aims to represent this real world entity should be structured in a way that allows for expansion, facilitates collaboration, and provides transparent meaning representation.

This paper addresses these two issues and provides a solution based on Semantic Web technologies. The solution is based on an ontological model for representing computational lexicons using the field theory of semantics and componential analysis. The model has been implemented on the “Time” semantic field vocabulary of the Arabic language and the results of a preliminary evaluation are presented.

**Keywords**- *Ontology, Lexicon, Arabic Language, Semantic Web*

## 1. INTRODUCTION

Language is a social and dynamic entity. Based on Saussure’s view on linguistics [1], natural languages are seen as social products, based on consensus among language users. A language lexicon is a major linguistic resource for language users, which records this consensus among language users. However, compiling a lexicon is a difficult and cumbersome task, which is expensive and involves intensive human effort. It requires collaboration and consensus among experts in the field, who are usually dispersed. It is clear that no individual or small group is able to take on this task, therefore collaboration and cooperation is a vital element in lexicon development.

Language is dynamic; its change is derived by the culture and the needs of the society in which it is spoken. New and emerging technologies, industries, products and experiences require new words to be invented and added, old words might no longer be used and therefore die and some words may endure meaning change. Traditionally, the only way for recording these changes was the periodic publishing of a new edition of a dictionary or lexicon. Advancement in technologies has resulted in the development of a digital form of these dictionaries as well as computationally based lexicons. However such

computational lexicons lack the necessary infrastructure and flexibility to accommodate the two basic elements in lexicon design and construction which stem from the nature of natural language.

As a result of the aforementioned properties of language, a language’s lexicon is alive, open and constantly expanding entity. Any computational model which aims to represent this real world entity should facilitate collaboration, allow for expansion, and provide a transparent meaning representation.

In this paper we present an ontological model which provides the foundation for a dynamic and collaborative computational lexicon. The ontological structure represents word semantics using the atomic components (features) of words. To enable the shared and open access to such a resource, we use the recent W3C standard for representing ontologies, Web Ontology Language (OWL) [2]. Moreover, using OWL for lexicon development and construction facilitates the development of semantic web services and agents that permit the community of linguists to collaborate on lexical knowledge acquisition.

This paper is organized as follows; the next section presents the background with an overview on the definition of a lexicon, and the theoretical background on the approach to semantics which forms the foundation of our proposed model. Section 3 provides a review of relevant literature on computational lexicons. Section 4 presents the design and evaluation of our ontological knowledge representation model for a computational lexicon. Finally, section 5 concludes the paper with a summary of the work done in this project and possible directions for future work.

## 2. BACKGROUND

A lexicon is a list of words in a language, the vocabulary, along with a lexical entry which gives more detail about the word [3]. The contents of each lexical entry depend on the purpose of the lexicon. A lexical entry may include any of the word properties such as meaning, relationship with other words, phonetics (sound/pronunciations), morphology, and syntax (grammatical behavior). The lexical entries that will be the focus of our computational model are those related to the semantics of the word (meaning).

Within the field of linguistics, the definition of meaning is considered one of the most ambiguous and most controversial terms in the theory of language. In general, however, there are two main schools of thought when defining meaning: the analytical (referential)

approach, and the operational approach. The analytical approach defines meaning by analyzing componential features of words, and the operational approach studies the words in usage [4].

The field theory of semantics (conceptual spheres), which forms the theoretical foundation of our computational model, follows an analytical approach. It was first introduced by Professor Jost Trier [5]. According to this theory, the meaning of a word is considered within a given view of the world. It is dependent on its relation to other words in the same semantic field (conceptual area) [6]. It assumes that the lexicon is structured into semantic fields according to a set of primitive features. Word meaning is established by the position within the field, and the relationship it has with other words in its field. A central element to this theory is componential analysis (semantic analysis) [7]. Using componential analysis, a word meaning can be defined in terms of a number of specific atomic components and decompositions which represent the distinctive features for a given word [8] [9]. These features form the basis for structuring a specific semantic field.

Using this approach to meaning enables us to differentiate between different words within the same semantic field. The following example illustrates the componential analysis of the words "man", "woman", "boy", and "girl". These words all belong to the same semantic field of "human", and they are defined in terms of two semantic dimensions "adulthood" and "gender". The meaning of the individual words can be expressed as a combination of these features [10], [11]:

- Man= Human+ Adult +Male
- Woman= Human+ Adult +Female
- Boy= Human+ Child+ Male
- Girl= Human+ Child+ Female

These formulae are called componential definitions of the semantic units, and can be regarded as formalized dictionary definitions [10].

Language change takes on a variety of forms [12] ; Lexical changes (new words added to the lexicon), Semantic change (old words in new use), Phonetic change (sounds of words change and their pronunciation), Spelling changes (the spelling of the word changes), or Syntactic change (the basic grammatical structure of the sentence). In our research we focus on the first two, lexical change and semantic change, since they are both concerned with the two important views on language.

Lexical change occurs when a new word is invented and added to the lexicon. For example, the verb "google" is a relatively new word which means "to search the internet using the google search engine". Semantic change occurs when the meaning of an existing word changes or a new sense is added. For example, the word "wiki" is a Hawaiian word meaning "quick", the same word is used to refer to the "wiki" technology. It is evident that when coining a referent for a new concept, the field theory of semantics and the componential analysis approach to meaning plays an important role, since people naturally think of features and properties when describing something [11][12].

### 3. RELATED WORK

The earliest forms of computational lexicons were designed as a dictionary in a machine-readable format, as in the ACQUILEX project [13]. As technologies advanced, the computational linguistics community sought more advanced representations that are capable of satisfying the requirements for natural language processing tasks, such as syntactic features COMLEX [14], word senses WordNet [15] and EuroWordNet [16] and semantic frames such as FrameNet [17] [18], and VerbNet [19], [20].

WordNet [15] is probably one of the most common and widely used lexicon. It is one of the earliest, developed since the early 80's. WordNet was built by lexicographers on the basis of analysis of language. The design is based on psycholinguistic and computational theories of human lexical memory. WordNet uses a semantic network structure representing words and concepts as an interrelated system consistent with the way humans organize their mental lexicon. The main taxonomic structure of these resources consists in a hierarchy of hyponyms. The synset is the set of synonyms that plays the central role of lexical concept in WordNet; it acts as the root for other semantic relations. Following on the success of WordNet, word nets for other languages have been developed, such as EuroWordNet [16], which contains networks for Dutch, Italian, Spanish, English, Czech, Estonian, French, and German.

Another lexical resource developed to be used in natural language processing applications is FrameNet [17]. FrameNet is a lexicon for the English language. As the name implies, its design is based on frame semantics and corpus analysis. FrameNet provides a computational representation of semantic and syntactic combinatory possibilities of each sense of a word.

VerbNet [20] is an on-line verb lexicon for the English language. It is based on a hierarchical structure of verb classes. Verbs in the lexicon are linked to other lexical resources such as WordNet and FrameNet. Each class is described using thematic roles, restrictions on the arguments, and frames consisting of a syntactic description and semantic predicates with a temporal function. Each frame is associated with explicit semantic information, expressed as a conjunction of Boolean semantic predicates such as motion, contact, or cause. Since its release, it has been used in a number of NLP applications for characterizing verbs and verb classes.

With recent advancements in Semantic Web technologies, a new paradigm for computational lexicon development is emerging. A growing trend is linking computational lexicons to upper (general) ontologies such as SUMO (Suggested Upper Merged Ontology) [21] and CYC [22], and linguistic foundational ontologies such as DOLCE [23].

There have been a number of projects where an ontological structure has been adopted in lexicon design. For example, WordNet has been restructured according to the principles of formal ontology in the OntoWordNet project [24] and has been represented using the W3C standard Web Ontology language (OWL)[2]. WordNet has also been linked to SUMO, the suggested upper merged ontology. In addition, FrameNet has been linked to the

Suggested Upper Merged Ontology (SUMO) [25], and has been represented using OWL [26].

Another computational lexicon development project strongly influenced by Semantic web technologies is the multilingual MILE project [27]. The design of the lexicon is based on defining lexical classes for creating objects to be used in building MILE conformant lexical entries. Lexical objects include semantic and syntactic features, semantic relations, syntactic constructions, predicate and arguments, etc. Lexical Classes are organized in a hierarchy and defined using RDF schema. LexOnto [28] is a lexicon ontology developed for the lexicon engineer to map information in domain ontologies to natural language lexical frames (sub-categorization frames) which can aid in NLP applications. It aims at associating lexical structures to ontological classes and properties found in domain ontologies, in other words translate an ontological structure into a lexical frame for NLP applications. A crucial difference between our model and LexOnto is that LexOnto is not exactly a lexicon, it maps domain knowledge represented in ontologies into linguistic structures which can be processed by NLP applications. In addition, LexOnto focuses on sub-categorization frames which are considered as a syntactic element of a lexical entry. However, our model focuses more on the semantics.

Another related system is the SIMPLE lexicon [29], which describes lexical entries in terms of their semantic properties. The model is designed to facilitate cross-language linking. The basic unit of the SIMPLE lexicon model is called Semantic Units (SemUs). Each SemanticUnit has an associated semantic type from the SIMPLE ontology. The main entities in the SIMPLE ontology are: SemanticUnits, Semantic Type, and Templates. SemanticUnits represent the primary means for encoding word senses. Each SemanticUnit is assigned a semantic type from the ontology, plus other information specified in the associated template, which contribute to the characterization of the word-sense (similar to features). SemanticUnits are assigned semantic types. Each type involves structured information represented as template. The semantic types themselves are organized into an ontology [30].

Our model shares with the SIMPLE lexicon the idea of defining semantic features for a word sense. However, one of the crucial differences between SIMPLE and our ontology is that our lexicon model is designed to aid lexicographers in the task of lexical change. SIMPLE, however, is intended to provide a standardized lexicon for cross language linking.

The Omega Ontology [31], a successor of SENSUS [32], is a large terminological ontology obtained by reemerging WordNet, and Mikrokosmos [33], a conceptual resource originally generated to support machine translation. It combines information from various other sources. It is a feature-oriented ontology that constitutes an upper model which merges the two models. It is considered a shallow lexical term taxonomy. Omega contains no formal concept definitions and only relatively few semantic relations between concepts. Omega does not make commitments to any specific theories of semantics or particular representations; it does not cover any semantic aspects of the language.

From our previous review of related work, it is clear that the main objective of computational lexicon development efforts were either used to aid in natural language processing tasks, or for cross language linking. None of the computational models reviewed addressed certain language characteristics which we focus on in our project. Our approach focuses on designing a model which supports the lexicographer's task of lexical change. Our knowledge representation model provides a structured approach to lexical and semantic change.

#### 4. THE ONTOLOGICAL MODEL

The approach we follow in ontology development is the UPON (Unified Process for ONtology) ontological engineering approach [34]. We had an expert in Arabic linguistics involved throughout the development process.

##### 4.1. Goal

The goal of the ontology is to provide a computational model capable of representing word meaning in a way which is suitable for the nature of language. It should facilitate lexicon growth, and provide a transparent and formal meaning representation. Ontology users are mainly linguists; however researchers in NLP and Semantic Web applications may find it useful since it provides a formal representation of word meaning using W3C standards. Within the scope of our research, three basic uses of the ontology are considered:

- Facilitate lexical and semantic change (generating new words/using old ones for new phenomena or concepts).
- Facilitate Semantic analysis (Componential analysis of word meanings).
- Facilitate collaboration among linguists (web based lexicon using open standards rather than proprietary systems).

The UPON ontology design methodology suggests setting a number of competency questions at a conceptual level which an ontology must be able to answer. The questions are used to evaluate the resulting ontology. The following is a set of competency questions that drives the design of our ontology:

- What is the meaning of a given word?
- A new phenomena or concept exists, and the lexicographer would like to name it, how can the ontological model support the lexicographer's task?

##### 4.2. Scope

The ontology we designed focuses on Arabic language vocabulary (nouns only) associated with the semantic field of "Time". The reason is that we already have readily available a detailed componential analysis of the associated vocabulary of this semantic field [35]. We also limited the vocabulary to those words which exists in the Holy Quran, the reason is that the Quran represents the purist and most authentic form of the classical Arabic language. This set of vocabulary constitutes a representative sample. The specific aspects of meaning we focus on in this model are componential analysis of word features (atomic components of word meaning).

### 4.3. Design

Since our theoretical foundation is based on the field theory of semantics, the ontology is structured according to the hyponymy/hypernymy relationship between vocabulary concepts. The conceptual classification forms the main classes of the ontology, and the vocabulary are represented as individuals associated with each class. The conceptual classification is derived by word features, which are represented as individuals linked to words via ontological relations. The structure of the ontological lexicon implies that the more you go deep into the hierarchy, the more arguments the componential formula will have, and therefore the meaning narrows. In contrast, words at higher levels have less arguments in their componential formula, and therefore the meaning broadens. Semantic units (words/vocabulary) are the primary components of a language, and are associated with a specific semantic domain. Each semantic unit has at most two distinctive features, and a number of general features. Using this approach to language semantics, the meaning of a word consists of the combination of all the features from the top of the semantic field down until the word itself.

#### 4.3.1. Classes

There are two main categories of classes in the ontology, top level classes, and lexical classes. The top level classes of the ontology are those which represent major concepts from the field theory of semantics, they are reusable across different semantic fields and among different languages. Lexical classes are those classes which represent actual words (vocabulary) within a specific semantic domain. The reason for regarding vocabulary as classes is that there exists sub-class relationship between words in a vocabulary. In addition, representing them as concepts is much more applicable to the nature of language, as new words may emerge and need to be included as a subclass of an already existing word (individual).

Top level classes are shown in Figure 1 using UML notations. Top level classes are written in English, they can be applied to any semantic field for any language. This design decision was based on the need for flexibility in porting the ontology for other semantic fields or other languages. Top level classes are all grouped under the concept (Linguistic\_Concept). These classes are based on the major concepts in the field theory of semantics. The following is a description of the classes:

- **Linguistic\_Concept**: a class which represents all terminology used in our ontology.
- **Semantic\_Unit**: a collection of features which differentiates the meaning of the word from others within the semantic domain or language
- **Semantic\_Field**: a collection of semantic units semantically related to each other
- **Semantic\_Domain**: a division within a semantic field which divides the semantic field into conceptual spheres.
- **General\_Feature**: a feature of a certain semantic unit, which can be shared with other semantic units. A single feature may be shared with more than one semantic unit.

- **Distinctive\_Feature**: a feature of a certain semantic unit, which is unique, and can only be shared where word synonymy exists. A single semantic unit can have at most two distinctive features.

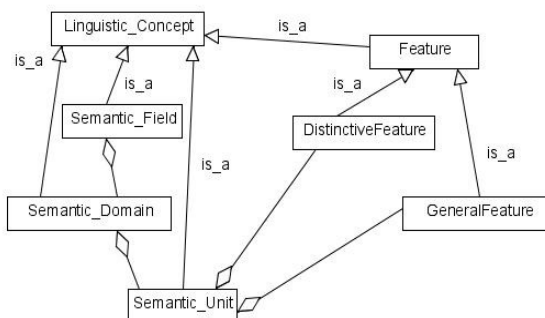


Figure 1. Ontology Top Level Classes

Lexical classes are related to the top classes by means of subclass relationship. In the ontology, all semantic units (language vocabulary) are represented as both classes and individuals. This is possible using OWL 2 because of the new feature called punning. This design decision was required, since a word in the vocabulary may act as a class (umbrella) which contains other vocabularies, and may act as a word in its own accord, which means it should be an individual which has features. Figure 2 shows lexical classes. It shows only part of the ontology which is concerned with vocabulary for “Day\_Domain” as opposed to “Night-Domain” semantic domain. The classes are translated into English if a matching word exists, otherwise a transliteration is provided.

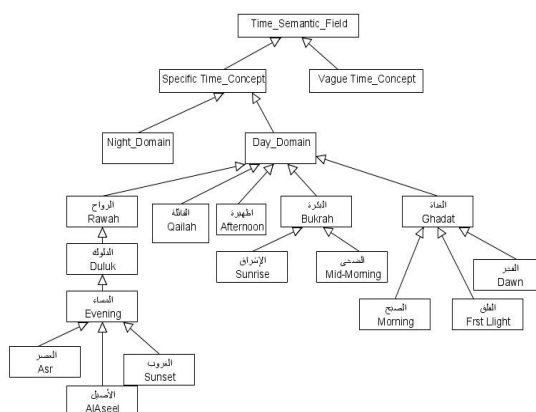


Figure 2. Ontology Lexical Classes for Part of the “Time” Ontology

#### 4.3.2. Individuals

Individuals in the ontology are of two types; words, and features. Words are represented as individuals since they must have features. Features are represented as individuals, since they describe the lexical unit and are not subject to further classification. Figure 3 shows a sample of the individuals in our ontology, rectangles represent classes from the ontology hierarchy, yellow ellipses represent individuals (words), and green ellipses represent individuals (features).



to our model, for example the antonymy relationship can be identified if there exists exactly one feature in each word formulae in which an opposition occurs for example male/female, adult/child. We plan to enhance the ontological structure with other lexical relations such as synonymy, antonymy, and polysemy. In addition, we plan to develop semantic web applications capable of exploiting the rich structure of the ontological model.

## ACKNOWLEDGEMENTS

This research is sponsored by KACST (King AbdulAziz City for Science and Technology)

## REFERENCES

- [1] De Saussure, F. *Cours de linguistique generale*. Payot, Paris (1949)
- [2] Web Ontology language (OWL): <http://www.w3.org/2004/OWL>
- [3] Hirst, G. Ontology and the Lexicon. In Staab, S. and Studer, R. (eds) *Handbook on Ontologies and Information Systems*. pp. 209-230. Springer, Heidelberg (2004)
- [4] Ullman, S. *Semantics: An Introduction to the Science of Meaning*. Blackwell, Oxford (1972)
- [5] Trier, J. *Der deutsche Wortschatz im Sinnbezirk des Verstandes*. Heidelberg (1931)
- [6] Lyons, J. *Semantics*. Cambridge University Press, Cambridge (1977)
- [7] Fodor, J. D. *Semantics: Theories of Meaning in Generative Grammar*. Harvester Press Limited, Sussex (1977)
- [8] Katz, J. J. *Semantic Theory*. Harper and Row, New York (1972)
- [9] Lakoff, G. On Generative Semantics. In Steinberg, D. D. and Jakobovits, L. A. (eds) *Semantics: An Interdisciplinary reader in philosophy, linguistics, and Psychology*. pp 232-296. Cambridge University Press, Cambridge (1971)
- [10] Leech, G. *Semantics*. Penguin Books, Harmondsworth (1974)
- [11] Hollmann, W. B. Semantic Change. In Culpeper, J. Katamba, F., Kerswill, P., McEnery, T. (eds) *The English Language*. Palgrave, Basingstoke, To Appear (n.d)
- [12] Grzega, J., Marion, S. *English and General Historical Lexicography*. <http://www1.ku-eichstaett.de/SLF/EnglVgISW/OnOnMon1.pdf> (2007)
- [13] Briscoe, T., de Paiva, V., Copestake, A. *Inheritance, Defaults, and the Lexicon*. Cambridge University Press, Cambridge (1993)
- [14] Macleod, C., Grishman, R., Meyers, A. *COMLEX Syntax: A Large Syntactic Dictionary for Natural Language Processing*. Computers and the Humanities. 31, 459-481. Springer, Netherlands (1997)
- [15] Fellbaum, C. *WordNet: An Electronic Lexical Database*. MIT Press, Cambridge, MA (1998)
- [16] Vossen, P. Eurowordnet: a multilingual database for information retrieval. In *Delos Workshop on Cross-language Information Retrieval* (1997)
- [17] Baker, C. F., Fillmore, C. J., Lowe, J. B. The Berkeley FrameNet Project. In the 17th International Conference on Computational linguistics. Montreal, Quebec, Canada (1998)
- [18] Dzikovska, M. O., Swift, M. D., Allen, J. F. Building a Computational Lexicon and Ontology with FrameNet. In *LREC workshop on Building Lexical Resources from Semantically Annotated Corpora*. Lisbon, Portugal (2004)
- [19] Kipper, K., Trang Dang, H., Palmer, M. Class-based construction of a verb lexicon. In the 7th Conference on Artificial Intelligence (AAAI-00) and of the 12th Conference on Innovative Applications of Artificial Intelligence (IAAI-00). AAAI Press, Menlo Park, CA (2000)
- [20] Schuler, K. K. VerbNet Overview. In *NAACL HLT: Tutorials*, pp13-14. Association for Computational Linguistics, Boulder, Colorado (2009)
- [21] Niles, I., Pease, A. Towards a Standard Upper Ontology.: In the 2nd International Conference on Formal Ontology in Information Systems (FOIS-2001). Ogunquit, Maine (2001)
- [22] Lenat, D. B. CYC: A large-scale investment in knowledge infrastructure. *Communications of the ACM*, 38, pp. 33-38 (1995)
- [23] Gangemi, A., Guarino, N., Masolo, C., Oltamari, A. Sweetening WordNet with DOLCE. *AI Magazine*, 24, pp. 13-24 (2003)
- [24] Gangemi, A., Navigli, R., Velardi, P. The OntoWordNet Project: extension and axiomatization of conceptual relations in WordNet. In *On the Move to Meaningful Internet Systems OTM2003*, pp 820-838. Springer-Verlag, Catania, Italy (2003)
- [25] Scheffczyk, J., Pease, A., Ellsworth, M. Linking FrameNet to the Suggested Upper Merged Ontology. In the conference on Formal Ontology in Information Systems (FOIS). Baltimore, USA (2006)
- [26] Scheffczyk, J., Baker, C. F., Narayanan, S. Ontology-Based Reasoning about Lexical Resources. In the Workshop on Interfacing Ontologies and Lexical Resources for Semantic Web Technologies (OntoLex 2006), pp. 1-8, Genoa, Italy (2006)
- [27] Calzolari, N., Zampolli, A., Lenci, A. Towards a Standard for a Multilingual Lexical Entry: The EAGLES/ISLE Initiative. In Hartmanis, G. G., van Leeuwen, J. (eds) *CICLing 2002 Computational Linguistics and Intelligent Text Processing*. LNCS, vol. 2276, pp. 65-80. Springer, Berlin (2002)
- [28] Cimiano, P., Haase, P., Herold, M., Mantel, M., Buite, P. LexOnto: A Model for Ontology Lexicons for Ontology-based NLP. In the workshop on Lexicon/Ontology Interface (OntoLex 2007), Busan, South Korea (2007)
- [29] Lenci, A., Bell, N., Busa, F., Calzolari, N., Gola, E., Monachini, M., Ogonowsky, A., Peters, I., Peters, W., Ruimy, N., Villegas, M., Zampolli, A. SIMPLE: A General Framework for the development of Multilingual Lexicons. *International Journal of Lexicography*, 13, 249-263 (2000)
- [30] Pustejovsky, J. *The Generative Lexicon*. MIT Press, Cambridge (1995)
- [31] Philpot, A., Hovy, E., Patrick, P. The Omega Ontology. In *OntoLex Workshop: Ontologies and Lexical Resources*, Jeju Island, South Korea (2005)
- [32] Knight, K., Luk, S.K. Building a Large- Scale Knowledge Base for Machine Translation. In the AAAI-94, Seattle, WA (1994)
- [33] O'Hara, T., Mahesh, K., Nirenburg, S. Lexical acquisition with WordNet and the Mikrokosmos Ontology. In the COLING/ACL Workshop on Usage of WordNet in Natural Language Processing Systems, Montreal, Canada (1998)
- [34] De Nicola, A., Missikoff, M., Navigli, R. A software engineering approach to ontology building. *Information Systems*. 34, 258-275 (2009)
- [35] Al-Helwah, Nawal. *Time Vocabulary in the Holy Quran: A Semantic Analysis Approach*. Princess Norah University, Riyadh, Saudi Arabia (2006)
- [36] Burton-Jones, A., Storey, V., Ahluwalia, P. A Semiotic metrics suit for assessing the quality of ontologies. *Data and Knowledge Engineering*. 55, 84-102 (2005)