

Ain Shams University

Faculty of Computer & Information Sciences

Prof. Mostafa Mahmoud Aref

Computer Science Department

# Developing a Framework For Building an Arabic Ontology

### **A THESIS**

Presented to Computer Science Department, Faculty of Computer and Information Sciences, Ain shams University.

Submitted in partial fulfillment of the requirements for the degree of doctor of philosophy in computer science.

Faculty of Computer and Information Sciences

Computer Science Department

Ain Shams University

## BY

## **Dalia Sayed Mohamed Hassan Fadl**

Assistant Lecturer Department of Computer Science, Akhbar El-Yom Academy.

M.Sc. of Computer and Information, Computer Science Department, Faculty of Computer and Information Sciences, Ain Shams University

## **Supervisors**

Dr. Safia Abbas

Professor of computer science, Associate Professor in Computer Science

Faculty of Computer and Information Sciences

Ain Shams University Ain Shams University

To My father, My Mother, and My family.

### **Abstract**

In Computer Science, ontology is a shared and common understanding of some domain that can be communicated across people and application systems or enabling knowledge sharing. It is a specification of a conceptualization. The rise of linguistic ontologies is a consequence of two simultaneous circumstances. Data organizing and description to encourage its usage by users later. Language is the best approach to vehicle data and information. So the requirement for semantic information is vital in all research fields. Web information always depends on the language which is written in; the accessibility of data identified with the language that would be much ideal as per the client would be an expanding need for today. Ontology is similar to a structure, dictionary or glossary, however, with more exceptional detail and structure. There is a robust requirement for Arabic language support since the ontology in English cannot be converted into Arabic. Distinctive languages have contained the particular semantic environment and the cultural context, which has brought on the need to build up the other ontology for various languages. Arabic ontology is an important natural language processing field it helps to enrich the Arabic language resources.

The aim of the framework is to automate the process of ontology generation, generating ontologies from pre-existing XML documents. It is an innovative framework, annotated as NAAO (Novel Automatic Arabic Ontology), which automates the ontology generation process from XML documents. The novelty of NAAO resides in generating the Arabic ontology, in the form of XML graph schema (XSG), from semi-structured data (XML documents associated with graph schema). The definition of this automation process was through four main steps necessary to achieve our goal. These steps conclude the main tasks of the automation process for building ontologies from XML documents. This thesis represents a framework that generates an Arabic Ontology from a semi-structured data (XML documents associated with graph schema), in which, XML schema is created and used in the graph schema development (XSG). The thesis provides two case studies, insectivore's case and mammal's case study where the developed Arabic ontology is applied. The results consist of 143 words, 10 concepts, 10 elements and 20 relationships. The generated ontology is evaluated using data-driven evaluation methods. 65% of the source XML documents have been included in the insectivore's case study.

Finally, the thesis provides the implementation of the framework for generating Arabic ontology containing the animal kingdom automatically. The ontology is divided into two parts to be more representative. The first part is the vertebrate's ontologies which provide 1576 concepts, 3836 element, and 2689 relations. Moreover, the second is the invertebrate's ontology which contains 320 concepts, 603 elements, and 783 relations. This result can be refined more than one time to reach satisfying results. The generated Arabic ontology is going to be evaluated using data-driven ontology measures cosine similarity measures and tree path mining. Finally, a comparison of the generated framework and three other system is provided.

## Acknowledgments

Thanks to Allah before and after. And heartily thankful to my supervisors; **Professor Mostafa Aref** for his encouragement guidance and support from the initial to the final level, enabled me to develop and understand Arabic Ontology Generation. He was not just a supervisor but a father in all the situations. Moreover, I would like to thank my supervisor; **Doctor Safia Abbas**. I would like to thank my family for all their love and encouragement; My parents who raised me with a love of science and supported me in all my pursuits; My brother Mohamed and my sister Dina. And, this thesis would not have been possible without my loving husband Tarek Khattab, who helped, supported, encouraged and believed in me all the time at all the situations. Thank you, is a very small word to express my feeling. Thanks to my son and my daughter .I am thankful to my friends and my colleges how supported me.

# **Publications**

Parts of this thesis have been published as original papers in the following references:

- Dalia Fadl, Safia Abbas, and Mostafa Aref, Automatic Arabic Ontology Construction framework: Insectivore's case study, AISI 2017: Proceedings of the International Conference on Advanced Intelligent Systems and Informatics 2017 pp 458-466, Cairo, Egypt.
- Dalia Fadl, Safia Abbas, and Mostafa Aref, Approach for Automatic Arabic Ontology Generation, International Journal of An Intelligent Computing and Information Sciences (IJCICIS) Faculty of Computer and Information Sciences, Ain Shams University, Cairo, Egypt, 2017.
- Dalia Fadl, Safia Abbas, and Mostafa Aref, Automatic Arabic Ontology Generation for the Animal Kingdom, accepted, Journal Of Theoretical And Applied Information Technology(JATIT) (E-ISSN 1817-3195 / ISSN 1992-8645),2017
- Dalia Fadl, Safia Abas, Mostafa Aref, "Arabic Ontology Using Different Ontology Learning Techniques," Proceedings of the Conference of Egypt Society of Language Engineering ESOLEC' 2014 Pages 68-67, Faculty of Engineering Ain Shams University, Cairo, Egypt.

# TABLE OF CONTENTS

		Page
ABST	TRACT	III
ACK	NOWLEDGEMENTS	IV
TABL	LE OF CONTENTS	V
LIST	OF FIGURES	IX
LIST	OF TABLES	X
LIST	OF ABBREVIATIONS	XII
СНА	PTER 1 INTRODUCTION	1
1.1	Overview	1
1.2	Aims and Objectives.	4
1.3	Contributions	4
1.4	The Structure of The Thesis.	5
СНА	PTER 2 BACKGROUND AND RELATED WORK	6
2.1	Types of Ontologies.	6
2.2	Upper Ontology	7
	2.2.1 Basic Formal Ontology (BFO)	7
	2.2.2 Cyc	7
	2.2.3 DOLCE	7
	2.2.4 GFO (General Formal Ontology).	7
	2.2.5 PROTON (PROTo ONtology)	8

	2.2.6 SUO (IEEE Standard upper ontology)	8
	2.2.7 Sowa's Ontology	8
	2.2.8 Merging Upper Level Ontologies.	8
2.3	Ontology Development Methodologies	9
	2.3.1 Manually Development Methodologies	10
	2.3.2 Semi-automatic Development Methodologies	12
	2.3.3 Automatic Ontology Development	16
2.4	Arabic Ontology	17
	2.4.1 Arabic Language and The Semantic Web Research	18
	2.4.2 Arabic Ontology Generation Issues.	18
2.5	Ontology Construction Requirements.	18
2.6	XML Schema and Ontology Comparison.	19
2.7	XML Document Mining.	20
2.8	Related Work	21
	2.8.1 A Linguistic Ontology for the Semantic Web	21
	2.8.2 Developing Ontology for Arabic Blogs Retrieval	22
	2.8.3 Building a Framework for Arabic Ontology Learning	22
	2.8.4 Automatic Construction of Ontology from Arabic Texts	22
	2.8.5 A Compact Arabic Lexical Semantics Language Resource	23
	2.8.6 Arabic Word Net.	23
	2.8.7 AN Application on Time Nouns in the HOLY OURAN	25

	2.8.8 Al –Khalil: The Arabic Linguistic Ontology Project	26
	2.8.9 Automatic Ontology Construction.	27
	2.8.10 Ontology Learning from Textual Web Documents	28
	2.8.11 Mining Multiword Terms from Wikipedia	28
	2.8.12 Building Ontologies from XML Data Sources	29
СН	APTER 3 DESIGN OF ARABIC ONTOLOGY FRAMEWORK	30
3.1	Extraction.	32
3.2	XML Schema Parsing.	33
3.3	Ontology Generation.	34
	3.3.1 Ontology Generation Rules.	35
	3.3.2 Generation of Properties	36
	3.3.3 Generation of Mapping Connections	38
3.4	Refinements and Evaluation.	39
3.5	Case Studies.	42
	3.5.1 Insectivore's Animal Case Study	42
	3.5.2 Mammals Case Study	48
СН	APTER 4 ARABIC ONTOLOGY MODEL IMPLEMENTATION	54
4.1	Input Documents	54
4.2	Ontology Generation.	55
4.3	NAAO Output and Analysis	56
4 4	Evaluation of The Results	59

4.5	Comparison with Other Systems.	65
СНАР	PTER 5 CONCLUSION AND FUTURE WORK	68
5.1	Conclusion.	68
5.2	Contributions.	69
5 3	Future Work	70

# LIST OF TABLES

Table		Page
4.1	Animal XML Documents Analysis before Extraction	55
4.2	Ontology Input Analysis.	57
4.3	Ontology Output Analysis.	60
4.4	Arabic Ontologies Output Analysis	60

# LIST OF FIGURES

Figure		Page
2.1	The Ontology Learning Process.	15
3.1	Arabic Ontology FrameWork	30
3.2	Detailed Arabic Ontology FrameWork.	32
3.3	Steps of Extraction	33
3.4	XML Schema Parsing	34
3.5	XML Schema Parsing Steps.	34
3.6	OWL Classes	36
3.7	OWL Classes Generation Steps	36
3.8	Rules for Object Properties	37
3.9	Rules of Data Type Properties	37
3.10	Class Connections Rules	38
3.11	Element Connections Rules.	38
3.12	Object Properties Connections Rules.	39
3.13	Tree Path Evaluation Algorithm.	41

3.14	XSG for XML File	42
3.15	English XSG of The Running Insectivores' Animal Case Study	43
3.16	Arabic XSG of The Running Insectivores' Animal Case Study	44
3.17	Arabic Ontology Connections.	45
3.18	Insectivores' Ontology	46
3.19	XSG for XML File	48
3.20	English XSG of The Running Insectivores' Animal Case Study	49
3.21	Arabic XSG of The Running Insectivores' Animal Case Study	50
3.22	Arabic Ontology Connections.	51
3.23	Mammals Ontology	52
4.1	Arabic Ontology Output Analysis	56
4.2	Graphical Representation of Part of The Generated Arabic Ontology	58
4.3	The Mammal's Ontology Output Analysis	61
4.4	The Bird's Ontology Output Analysis	62
4.5	The Reptile's Ontology Output Analysis	62
4.6	The Amphibian's Ontology Output Analysis	63
4.7	The Invertebrate's Ontology Output Analysis	64

## LIST OF ABBREVIATIONS

CR Context Resemblance 1 2 IDF Inverse Document Frequency 3 Novel Automatic Arabic Ontology NAAO OWL Ontology Web Language 5 OOT OWL Ontology Term 6 TF Term Frequency 7 TO Term Occurrence 8 Tree Path Mining TPM Extensible Marked Up Language 9 XML

## **Chapter 1**

## Introduction

#### 1.1. Overview

The term "Ontology" has been introduced with the data sciences and research fields amid the 1990's by a few Artificial Intelligence (AI) research groups. AI analysts received the expression "Ontology" basically to depict what they thought would be (from the outlook of computational aspects) a legitimate representation of the world in a program code. Ontologies are of essential enthusiasm for a broad range of fields, to a great extent because of what they guarantee: a mutual and normal comprehension of some area that can be the reason for correspondence ground over the crevices amongst individuals and PCs. They (Ontology approaches) take into account sharing and reuse of knowledge bodies in computational shape. The same numbers of traditional activities are changing their way in the realm of today because of the accessibility of data brought by the World-Wide-Web (WWW), Ontologies are probably going to turn increasingly when the learning is organized in a machine-readable way, and the abstract concepts it contains are shared.

Of the many definitions which have aroused for Ontology the following "Ontology is a formal, explicit specification of a shared conceptualization." A "conceptualization" is a theoretical model of a phenomenon, made by ID of the related concepts of the phenomenon. The concepts, the relations amongst them and the imperatives on their usege are unequivocally characterized. "Formal" implies that Ontology is machine-readable and rejects the use of natural languages. For instance, in restorative spaces, the concepts are illnesses and manifestations, the relations between them are causal, and an imperative is that an infection cannot bring itself.

Ontology is a "shared conceptualization" states that Ontologies expect to speak to consensual information proposed for the use of the group. In a perfect world, the Ontology catches knowledge autonomously of its usege and in a way that can be shared all around, however for all intents and purposes unique errands and utilizations call for various portrayals of the knowledge in Ontology. Ontology is some of the time mistook for taxonomy, which is an arrangement of the information in a domain.

The distinction between them is in two vital contexts:

- 1. Ontology has a richer internal structure as it includes relations and constraints between the concepts.
- 2. Ontology claims to represent a certain consensus about the knowledge in the domain.

This consensus is among the intended users of the knowledge, e.g. doctors using a hospital Ontology regarding a certain disease, artists relating to historical art and so on. Ontologies are divided into types in accord with the degree of generality of the principles they contain. Using the distinctive languages in the investigation of Ontology can be a challenge to the many endeavors of the Web designs to cater the huge number of users on the World Wide Web. Web data is typically language dependent. The accessibility of data identified with the language that would be much ideal as indicated by the client would be an expanding need of today.

The developing enthusiasm for ontologies for some common language applications in the current years has prompted to the production of ontologies. These ontologies are for various purposes and with various elements frameworks. Additionally, the current work in Artificial Intelligence is investigating the usege of formal ontologies.

Its use is a way of specifying content-specific agreements for the sharing and reuse of knowledge among software entities. There are various studies conducted in the Arabic language in Semantic Web. The propose of this studies is to improve the Arabic information retrieval on the web [1]. The ontology development life cycle had many questions around it in the last few years. Are there common designing criteria or not? Although Arabic is the language of hundred millions of people over the world, little has been done regarding computerized linguistic resources, tools or applications.

There are six parts in the life cycle in the development of ontology: Creation, Population, Validation, Deployment, Maintenance, and Evolution [3]. Manual ontology building is a time-consuming activity that requires many efforts for knowledge domain acquisition and knowledge domain modeling. To overcome these problems; many methods have been developed, including systems and tools that automatically or semi-automatically, using text mining and machine learning techniques, allows generating ontologies. The research fields which study this issues is