



# **AUTOMATIC ONTOLOGY CONSTRUCTION APPROACH FROM WEB PAGES**

By

**Naglaa Elsayed Ahmed Mohamed Elmesalmy**

A Thesis Submitted to the  
Faculty of Engineering at Cairo University  
in Partial Fulfillment of the  
Requirements for the Degree of  
**MASTER OF SCIENCE**  
in  
**COMPUTER ENGINEERING**

FACULTY OF ENGINEERING, CAIRO UNIVERSITY  
GIZA, EGYPT  
2019

# **AUTOMATIC ONTOLOGY CONSTRUCTION APPROACH FROM WEB PAGES**

By

**Naglaa Elsayed Ahmed Mohamed Elmesalmy**

A Thesis Submitted to the  
Faculty of Engineering at Cairo University  
in Partial Fulfillment of the  
Requirements for the Degree of  
**MASTER OF SCIENCE**  
in  
**COMPUTER ENGINEERING**

Under the Supervision of

**Prof. Magda B. Fayek**

Computer Engineering Department  
Faculty of Engineering, Cairo University

**Dr. Mayada M.Hadhoud**

Computer Engineering Department  
Faculty of Engineering, Cairo University

FACULTY OF ENGINEERING, CAIRO UNIVERSITY  
GIZA, EGYPT  
2019

# **AUTOMATIC ONTOLOGY CONSTRUCTION APPROACH FROM WEB PAGES**

By  
**Naglaa Elsayed Ahmed Mohamed Elmesalmy**

A Thesis Submitted to the  
Faculty of Engineering at Cairo University  
in Partial Fulfillment of the  
Requirements for the Degree of  
**MASTER OF SCIENCE**  
in  
**COMPUTER ENGINEERING**

Approved by the  
Examining Committee

---

Prof. Dr. **Magda B. Fayek**, Thesis Main Advisor

---

Prof. Dr. **Nevin M. Darwish**, Internal Examiner

- Professor in Computers Engineering Department, Faculty of Engineering, Cairo university

---

Prof. Dr. **Samia A. Mashali**, External Examiner

- Professor in Computers & Systems Department, Electronics Research Institute (ERI)

FACULTY OF ENGINEERING, CAIRO UNIVERSITY  
GIZA, EGYPT  
2019

**Engineer's Name:** Naglaa Elsayed Ahmed Mohamed  
**Date of Birth:** 25/11/1986  
**Nationality:** Egyptain  
**E-mail:** Naglaa.elmesalmy@gmail.com  
**Phone:** 01008879638  
**Address:** 21 Sheikh Al-Kafer st.- El Gamma Dept  
Zagazig- El shariqua  
**Registration Date:** 1/10/2011  
**Awarding Date:** ....../....../2019  
**Degree:** Master of Science  
**Department:** Computer Engineering



**Supervisors:**

Prof. Dr. Magda B. Fayek  
Dr. Mayada M. Hadhoud

**Examiners:**

Porf. Dr. Magda B. Fayek (Thesis main advisor)  
Prof. Dr. Nevin M. Darwish (Internal examiner)  
Prof. Dr. Samia A. Mashali (External examiner)  
- Professor in Computers & Systems Department,  
Electronics Research Institute (ERI)

**Title of Thesis:**

AUTOMATIC ONTOLOGY CONSTRUCTION APPROACH FROM WEB PAGES

**Key Words:**

Ontologies ; semantic web ; automatic ontology Construction ; web page ;Wordnet

**Summary:**

In the thesis, a new approach for automatic ontology construction from web pages is proposed. The proposed approach works as follows: first the triples of the document are extracted, next it utilizes natural language processing techniques to process the extracted triples, and finally it applies ontology design patterns before inserting triples into ontology. Our proposed approach is compared with the previous related work and the results show that we are closed to cover all the information in the documents correctly. Our system is better at reducing repeat rate where it do not insert the same object twice and do not repeat the relation between two concepts. Also our ontology is simpler representation than other work.

## **Disclaimer**

I hereby declare that this thesis is my own original work and that no part of it has been submitted for a degree qualification at any other university or institute.

I further declare that I have appropriately acknowledged all sources used and have cited them in the references section.

Name:

Date:

Signature:

## **Acknowledgments**

All praises and thanks are due to Allah, the most gracious, the most merciful, for providing me with the strength, and patience to complete this work.

I am grateful to my supervisors, Prof. Dr. Magda Fayek, and Dr. Mayada for their guidance, advice, and encouragement toward successful completion of this work. They were very helpful indeed, reading and correcting me all the way.

I would like also to express my gratitude to my parents for their endless care and support from the very beginning and my husband for his wonderful support and encouragement.

Finally, I would like to dedicate this work to my husband, and my family.

# Table of Contents

DISCLAIMER .....	I
ACKNOWLEDGMENTS .....	II
TABLE OF CONTENTS .....	III
LIST OF TABLES .....	V
LIST OF FIGURES .....	VI
NOMENCLATURE .....	VII
ABSTRACT.....	VII
CHAPTER 1 : INTRODUCTION.....	1
1.1.    MOTIVATION.....	1
1.2.    THE STUDY PROBLEM.....	1
1.3.    THE THESIS ORGANIZATION .....	2
CHAPTER 2 : BACKGROUND .....	3
2.1.    INTRODUCTION .....	3
2.2.    WEB EVOLUTION AND SEMANTIC WEB .....	3
2.2.1.    Web 1.0.....	4
2.2.2.    Web 2.0.....	4
2.2.3.    Web 3.0.....	5
2.2.3.1.    Layered architecture of semantic web.....	5
2.2.3.2.    Difference between current web and semantic web .....	7
2.2.3.3.    The objectives of semantic web .....	7
2.2.3.4.    The challenges of semantic web .....	8
2.3.    SUMMARY .....	8
CHAPTER 3 : RELATED WORK .....	10
3.1.    INTRODUCTION .....	10
3.2.    ONTOLOGY EXTRACTION CATEGORIES .....	10
3.2.1.    Semi-automatic Approaches .....	10
3.2.2.    Full automatic Approaches .....	19
CHAPTER 4 :THE PROPOSED APPROACH .....	25
4.1.    INTRODUCTION .....	25
4.1.1.    Stage 1: HTML Parser .....	26
4.1.2.    Stage 2 : Co-reference resolution.....	26
4.1.3.    Stage 3 :Triple extraction.....	28
4.1.4.    Stage 4: Name Entity recognizer (NER).....	31
4.1.5.    Stage 5 : Formatting and Linguistic Unification .....	31
4.1.5.1.    Get Synonym list .....	32
4.1.5.2.    Refining and Formatting triples .....	32
4.1.5.3.    Applying Design Patterns .....	32
4.1.6.    Stage 6 – ontology Generation.....	34
4.2.    WEB2ONTO_V2.....	34
4.3.    SUMMARY .....	34
CHAPTER 5 : RESULTS AND DISCUSSION .....	36

5.1.	EXPERIMENTAL SETUP.....	36
5.2.	EXPERIMENT NO.1 .....	36
5.2.1	Used Dataset .....	36
5.2.2	Proposed Comparison Metrics.....	36
5.2.3.	Experimental RESULT .....	37
5.3.	EXPERIMENT NO.2.....	42
5.3.1	Used Dataset.....	42
5.4.	WEB 2 ONTO _v2 RESULT.....	47
5.5.	DISCUSSION.....	49
CHAPTER 6 CONCLUSIONS AND FUTURE WORK .....		51
6.1.	CONCLUSION.....	51
6.2.	FUTURE WORK .....	52
REFERENCES.....		53
APPENDIX A: CASE STUDY .....		54
A.1.	CASE STUDY.....	54



## List of Tables

Table 2-1: World internet usage and population statistics.....	4
Table 2-2: Comparison of Web 1.0, Web 2.0 and Web 3.0.....	8
Table 3-1: Samples of sentences that can match semantic pattern <Plant Part><Becomes. Verb><Color>.....	13
Table 3.2: the candidates for biosensor.....	21
Table 4.1: Co-reference result stage.....	28
Table 5.1: Used dataset description .....	36
Table 5.2: The result of comparison between FRED and WEB2ONTO.....	41
Table 5.3: The URLs of the dataset.....	42
Table 5.4: The number of extracted paragraphs.....	43
Table 5.5: The number of co-reference resolution.....	43
Table 5.6: The number of triples.....	43
Table 5.7: The number of NER .....	44
Table 5.8: The number of final triples.....	44
Table 5.9: The difference in the URL between WEB2ONTO and WEB2ONTO_V2..	47
Table A.1: The result of extracted sentences .....	57
Table A.2: The result of Co-reference resolution.....	59
Table A.3: The result of Name Entity recognizer.....	67
Table A.4: Triples Inserted Onto Ontology.....	69

# List of Figures

Figure 2.1: A global map of the web index for countries.....	3
Figure 2.2: Comparison between web1.0 and web2.0.....	5
Figure 2.3: Semantic web architecture .....	5
Figure 2.4: Difference Between Current Web (a) And Semantic Web (b) .....	8
Figure 3.1: Protégé interface[13].....	11
Figure 3.2: Text2Onto interface[14].....	12
Figure 3.3: General outline of Simple Method for Ontology Automatic Extraction from Documents[16] .....	14
Figure 3.4: Example of Knowledge extraction :(a) extraction result (b) knowledge triples[17] .....	15
Figure 3.5: Sample of Artequakt result: (a) xml file of extracted information (b) ontology .....	15
Figure 3.6: Structure of ontology in Research on Ontology Construction and Information Extraction Technology Based on WordNet[18].....	16
Figure 3.7: Structure of ontology in MOBM[19].....	17
Figure 3.8: The proposed system architecture in Domain Ontology Construction by Partial Import for Document Annotation[20].....	18
Figure 3.9: The architecture of Upper-Ontology-Based Approach for Automatic Construction of IOT Ontology[21].....	19
Figure 3.10: The concepts in triple extracted from unstructured documents[21] .....	19
Figure 3.11: The concepts in triple extracted from structured documents[21] .....	19
Figure 3.12: Proposed method of Creating ontologies from web documents[22] .....	20
Figure 3.13: The main components in FRED[23].....	22
Figure 3.14: Texting processing Module in FRED[23].....	23
Figure 3.15: Heuristic-based triplification module in FRED[23] .....	23
Figure 3.16: RDF graph Enrichment in FRED[23] .....	23
Figure 4.1: The proposed approach block diagram (WEB2ONTO ) .....	25
Figure 4.2: HTML Parser .....	26
Figure 4.3: Implement of negation .....	31
Figure 4.4: Sub-modules in stage5 .....	32
Figure 4.5: Applying N-ary relation design pattern .....	33
Figure 4.6: Implement the Adjectives .....	33
Figure 5.1: The result of FRED for input "John travels by car. John travels by car." ...	37
Figure 5.2: The result of WEB2ONTO for input "John travels by car. John travels by car." .....	37
Figure 5.3: WEB2ONTO representation for text No.2 (T2).....	38
Figure 5.4: FRED representation for Text No2 (T2).....	39
Figure 5.5: Part of the FRED Representation for text No.2 (T2).....	40
Figure 5.6: The Error in the representation of Text No.2 (T2) using FRED.....	40
Figure 5.7: Incorrect extrated triples in WEB2ONTO versus FRED.....	41
Figure 5.8: Duplicates in WEB2ONTO versus FRED.....	42
Figure 5.9: The graph for DOC2 .....	45
Figure 5.10: The graph for DOC2 and DOC1 .....	46
Figure 5.11: Result by WEB2ONTO .....	49
Figure 5.12: Result by WEB2ONTO_V2 .....	49
Figure A.1: Screenshot of test web page .....	54

# Nomenclature

AJAX	Asynchronous Javascript And Xml
API	Application Programming Interface
ASCII	American Standard Code For Information Interchange
ATCIS	Army Tactical Command Information System
CERN	The European Organization For Nuclear Research
CIDOC	International Committee For Documentation Of The International Council Of Museums
CRM	Conceptual Reference Model
D-FDED	Number Of Duplicates In FRED Ontology
D-ONTO	Number Of Duplicates In Our Ontology
DRS	Discourse Representation Structures
DTiMS	Defense Transaction Interface Module Systems
GWT	Google Web Toolkit
HTML	Hyper Text Markup Language
HTTP	Hypertext Transfer Protocol
IC-FRED	Number Of Incorrect Results In FRED Ontology
IC-ONTO	Number Of Incorrect Results In Our Ontology
IDF	Inverse Document Frequency
Info- onto	Number Of Information which Covered By Our Ontology
Info-FRED	Number Of Information which Covered By FRED Ontology
Info-No	Number Of Information
MOBM	Mixed Ontology Building Methodology
NER	Name Entity Recognizer
NLP	Natural Language Processing
NotCON-ONTO	Non-Connected Class In Our Ontology
NotCON-FRED	Non-Connected Class In FRED
OWL	Web Ontology Language
POM	Probabilistic Ontology Model
PVC	Professional Virtual Communities
RDF	Resource Description Framework
RSS	Really Simple Syndication
SQL	Structured Query Language
SVD	Singular Value Decomposition
TF	Term Frequency
TNO	Text Number
URI	Uniform Resource Identifier
URL	Uniform Resource Locator
W3C	World Wide Web Consortium
WWW	World Wide Web
XML	Extensible Markup Language

# Abstract

Over the past decades, the Semantic Web got a great attention as the next generation Web. The Semantic Web adds semantics to the Web pages to become meaningful, machine-processable, and understandable. Ontology is a main component of the Semantic Web. In various fields, Ontology has been studied by many researchers.

Ontology plays an important role in exchanging knowledge between different organizations and within various Systems. It provides the description of data and the relations between them.

This thesis presents an approach for automatically constructing ontology from web pages thus transforming data from web into ontology. The proposed approach applies to web pages in any domain. It is a step towards the utilization of the huge amount of data published on the web to build ontology.

The proposed approach first extracts a set of triples from Web page, next it utilizes natural language processing techniques to process the extracted triples, and finally it applies ontology design patterns before inserting these triples into ontology.

We defined a number of performance metrics for automatic ontology generation. These metrics were used to compare our proposed system with other systems.

Compared to other related work and results, our proposed approach, show that we are closer to cover all the information in the documents correctly. Our system is good at reducing repeat rate where it does not insert the same object twice and does not repeat the relation between two concepts. Also our ontology is a simpler representation than other automatic systems. Total precision of our system ontology is 0.84 and our system recall is 0.79.

# Chapter 1 : Introduction

The Semantic Web is the next generation of the current Web in which smart devices can understand the meaning of web contents because it contains Semantic annotations which provide additional information in the form of markup about various concepts (such as places, organizations, people, things etc.). It exchanges principles of the WWW (World Wide Web) from shared documents to shared data through a common framework which permits data to be shared and reused across applications, enterprise, and community boundaries. Ontology is the backbone of Semantic Web. The definition of Ontology is "a formal conceptualization of a particular domain that can be shared by a group of people (in and between organizations)"[1].

In semantic web, Ontology provides a sound semantic basis for the definition of meaning. It is typically used to present the natural language for communication between machines and humans [1]. It is a collection of concepts and the relationships between the concepts within a domain. It organizes the concepts into categories of concepts, each with its attributes, and describes relationships between concepts. When data is marked up using Ontology, softbots can well understand the semantics and therefore more intelligently locate and integrate data for a wide variety of tasks.

## 1.1. Motivation

Beside the significance of semantic networks, there are two more reasons that have encouraged me to propose this research:

- 1- Ontology availability is the most challenge of semantic web. Most proposed methods of ontology construction were for a specific domain. Also, ontology should be constructed by experts.
- 2- Amount of available knowledge on the current web is huge, important and in various fields. Therefore, there is a need to prepare this knowledge for usage in the semantic web.

These reasons led me think to develop a totally automatic approach for transforming the huge amount of knowledge in the web to ontology that can be used in the semantic web.

## 1.2. The Study Problem

The Wide World Web is a vast and rapidly growing source of information in all fields of science, biology, business, medicine, military, etc. The success of the Semantic Web depends strongly on transforming this information into a huge ontology, which provides a controlled vocabulary or conceptualization in these fields. Ontology is the core of the Semantic Web. In fact, Ontology is used to provide semantics and present a comprehensible, common foundation for resources on the Semantic Web. Also, ontology can present a common vocabulary and a rule for publishing data. Moreover, ontology can provide data with the semantic description. In this case, the WWW will be transformed from being machine-readable to machine-understandable.

Due to ontology's significance, many efforts have been exerted to develop ontology. The manual and semi-automatic constructional approaches of large ontology from web pages will not be feasible because of the effort, time and costs required. This gives place to the initiation of proposing approaches for constructing automatic ontology generation.

In this thesis, a new approach is proposed for automatic ontology construction that extracts the concepts and their relationship from web page. Then, it stores the concepts with their relationship on the ontology after making sure that they do not exist in ontology with the same relationship. Also, it uses the ontology design patterns in constructing the ontology.

### **1.3. The Thesis Organization**

The remainder of the thesis is organized as follows. In chapter 2, an introduction about the basic concepts of semantic web and its layers is presented. In chapter 3, a literature survey about different ontology construction methods is presented. In chapter 4, the proposed approach for automatic ontology construction from web pages is explained in details. In chapter 5, the results of the comparison between the proposed approach and the previous work are pointed out and discussed. Finally, conclusion and future work are given in chapter 6.

## Chapter 2 : Background

### 2.1. Introduction

The World Wide Web (WWW) is a collection of web pages and other web resources which are linked together through the internet. It was invented in 1989 by the English scientist Tim Berners-Lee who is the director of the World Wide Web Consortium (W3C). He was an employee at the European Organization for Nuclear Research (CERN) in Switzerland and presented a proposal for CERN communication system. Thus, the idea was implemented throughout the world. At the beginning of 1990, Belgian with Tim proposed the hyperlink to navigate across web pages. Tim also released the first web browser computer program in 1990, and consequently WWW began to appear to the world [2, 3, 4 ].

The World Wide Web plays a real role in the evolution of the Information Age and is the main tool for billions of people to interact on the Internet. There are billions of web pages, published and indexed as shown in figure 2.1[5]. Web pages contain text that annotated and formatted with Hypertext Markup Language (HTML). In addition to formatted text, web pages can display images, audio, video, etc.

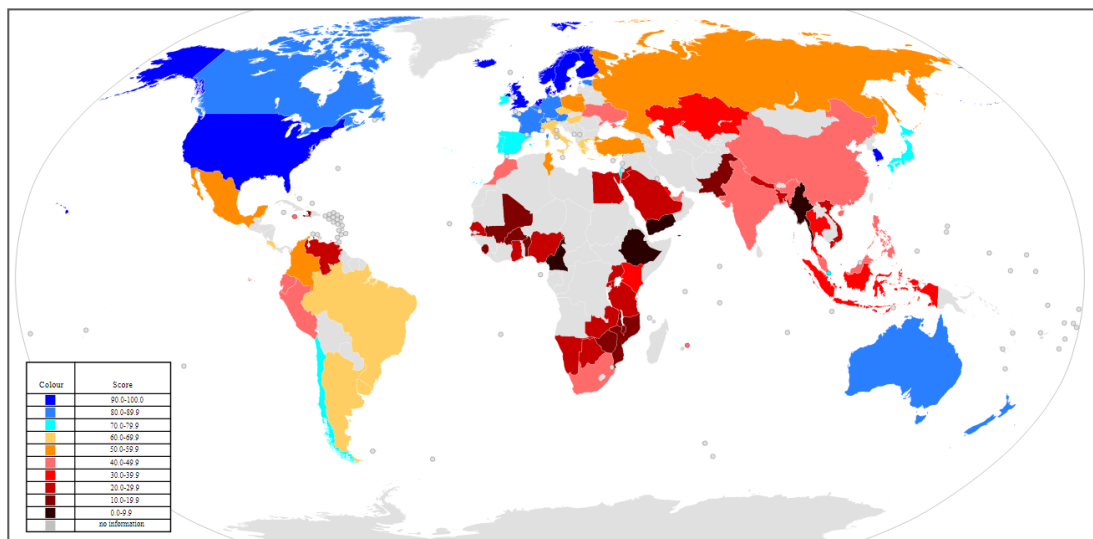


Figure 2.1: A global map of the web index for countries

### 2.2. Web Evolution And Semantic Web

Since the early 1990s, Web has been evolving in response to the never ending needs of users. Table 2-1 shows that the number of Internet users is 51.7% of the World population, and the growth rate is 976.4% in the period between 2000 to 2017[6]. This means that the Web should be more and more intelligent and sophisticated as the users' expectation is increasing.