

Boris Lauser

**Semi-automatic ontology engineering and
ontology supported document indexing in
a multilingual environment**

Diploma Thesis

Bibliographic information published by the German National Library:

The German National Library lists this publication in the National Bibliography; detailed bibliographic data are available on the Internet at <http://dnb.dnb.de> .

This book is copyright material and must not be copied, reproduced, transferred, distributed, leased, licensed or publicly performed or used in any way except as specifically permitted in writing by the publishers, as allowed under the terms and conditions under which it was purchased or as strictly permitted by applicable copyright law. Any unauthorized distribution or use of this text may be a direct infringement of the author s and publisher s rights and those responsible may be liable in law accordingly.

Copyright © 2003 Diplomica Verlag GmbH
ISBN: 9783832469054

Boris Lauser

Semi-automatic ontology engineering and ontology supported document indexing in a multilingual environment

Boris Lauser

Semi-automatic ontology engineering and ontology supported document indexing in a multilingual environment

**Diplomarbeit
an der Universität Fridericiana Karlsruhe (TH)
Fachbereich Wirtschaftsingenieurwesen
Institut für Angewandte Informatik
Januar 2003 Abgabe**



Diplomica GmbH _____
Hermannstal 119k _____
22119 Hamburg _____

Fon: 040 / 655 99 20 _____
Fax: 040 / 655 99 222 _____

agentur@diplom.de _____
www.diplom.de _____

ID 6905

Lauser, Boris: Semi-automatic ontology engineering and ontology supported document indexing in a multilingual environment

Hamburg: Diplomica GmbH, 2003

Zugl.: Fachhochschule Südwestfalen, Technische Universität, Diplomarbeit, 2003

Dieses Werk ist urheberrechtlich geschützt. Die dadurch begründeten Rechte, insbesondere die der Übersetzung, des Nachdrucks, des Vortrags, der Entnahme von Abbildungen und Tabellen, der Funksendung, der Mikroverfilmung oder der Vervielfältigung auf anderen Wegen und der Speicherung in Datenverarbeitungsanlagen, bleiben, auch bei nur auszugsweiser Verwertung, vorbehalten. Eine Vervielfältigung dieses Werkes oder von Teilen dieses Werkes ist auch im Einzelfall nur in den Grenzen der gesetzlichen Bestimmungen des Urheberrechtsgesetzes der Bundesrepublik Deutschland in der jeweils geltenden Fassung zulässig. Sie ist grundsätzlich vergütungspflichtig. Zuwiderhandlungen unterliegen den Strafbestimmungen des Urheberrechtes.

Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Die Informationen in diesem Werk wurden mit Sorgfalt erarbeitet. Dennoch können Fehler nicht vollständig ausgeschlossen werden, und die Diplomarbeiten Agentur, die Autoren oder Übersetzer übernehmen keine juristische Verantwortung oder irgendeine Haftung für evtl. verbliebene fehlerhafte Angaben und deren Folgen.

Diplomica GmbH

<http://www.diplom.de>, Hamburg 2003

Printed in Germany

TABLE OF CONTENTS

1	INTRODUCTION.....	1
1.1	MOTIVATION.....	1
1.2	APPROACH.....	3
1.3	OUTLINE.....	4
2	THE PROJECT ENVIRONMENT.....	5
2.1	FAO AND THE AOS.....	5
2.2	INFORMATION MANAGEMENT AT THE FAO.....	7
2.2.1	<i>Resources and metadata</i>	7
2.2.2	<i>The information management system</i>	8
2.2.3	<i>AGROVOC Thesaurus and Document Indexing</i>	10
2.3	PROBLEMS WITH THE CURRENT SYSTEM AND PROPOSAL.....	13
3	SEMANTIC WEB.....	15
3.1	THE IDEA.....	15
3.2	ONTOLOGIES.....	17
3.2.1	<i>Introduction</i>	17
3.2.2	<i>Types of ontologies</i>	20
3.2.3	<i>Ontology representation languages</i>	22
3.2.4	<i>KAON</i>	25
3.2.5	<i>Ontology Engineering</i>	27
4	INTRODUCTION OF ONTOLOGY BASED INFORMATION MANAGEMENT SYSTEM AT THE FAO.....	29
4.1	THE PROTOTYPE PROJECT.....	29
4.2	REQUIREMENTS REGARDING THE AOS.....	30
4.3	ONTOLOGY ENGINEERING FRAMEWORK.....	32
4.3.1	<i>Overview</i>	32
4.3.2	<i>Initialisation of the cycle</i>	33
4.3.3	<i>The 5 phases of the framework</i>	35
4.4	THE ONTOLOGY BROWSER.....	40
4.5	REPRESENTATION OF AGROVOC IN KAON.....	42
4.6	RELATED WORK AND POSITIONING:.....	46
4.7	CURRENT STATUS AND FURTHER WORK:.....	48
5	THE ONTOLOGY PRUNER.....	50
5.1	INTRODUCTION TO THE PRUNING APPROACH.....	50
5.2	ADAPTATION OF THE ONTOLOGY PRUNER.....	53
5.3	EVALUATION.....	56
5.3.1	<i>Resources: Document corpus and source ontology</i>	56
5.3.2	<i>Hypotheses for evaluation</i>	58
5.3.3	<i>Evaluation plan</i> :.....	59
5.4	RESULTS AND DISCUSSION:.....	60
5.4.1	<i>Pruner Trie vs. Pruner</i> :.....	61
5.4.2	<i>Dependency of the statistics on different parameter settings</i> :.....	61
5.4.3	<i>Generic Document Set 1 (Gen) vs. Generic Document Set 2 (AG)</i> :.....	62
5.4.4	<i>Empirical evaluation</i> :.....	63
5.5	SUMMARY.....	67
6	AUTOMATIC CLASSIFICATION.....	69
6.1	INTRODUCTION.....	69
6.1.1	<i>What is text categorisation?</i>	69
6.1.2	<i>Motivation within the project context</i>	69
6.2	BASIC DEFINITIONS.....	70

6.2.1	<i>Using Support Vector Machines for Multi-label Document Indexing</i>	70
6.2.2	<i>Evaluation measures:</i>	74
6.3	ADAPTATION OF THE CLASSIFIER	78
6.3.1	<i>Multi-label vs. single-label Indexing</i>	78
6.3.2	<i>Multiple Languages</i>	80
6.3.3	<i>Integration of background knowledge</i>	80
6.3.4	<i>Multi-class problem and class hierarchy</i>	83
6.4	SET OF TRAINING AND TEST DOCUMENTS	85
6.5	EVALUATION.....	89
6.5.1	<i>Single-label vs. multi-label classification</i>	89
6.5.2	<i>Multilingual classification</i>	96
6.5.3	<i>Integration of domain specific background knowledge</i>	98
6.6	RELATED WORK	100
6.7	SUMMARY AND OUTLOOK	101
7	CONCLUSION	103
7.1	SUMMARY	103
7.2	OUTLOOK.....	105
	REFERENCES.....	106
A	KAON RDFS REPRESENTATION OF THE ONTOLOGY ON FOOD SAFETY, ANIMAL AND PLANT HEALTH (EXTRACT).....	113
B	COMPLETE LIST OF WEB SITES OUTPUT BY THE FOCUSED CRAWLER.....	114
C	AGROVOC CATEGORIES	119
D	RESULTS OF ONTOLOGY INTEGRATION INTO AUTOMATIC TEXT CLASSIFICATION.....	123

TABLE OF FIGURES

FIGURE 1: ONTOLOGY EXAMPLE, EXCERPT.....	2
FIGURE 2: INFORMATION MANAGEMENT SYSTEM AT THE FAO	10
FIGURE 3: AGROVOC THESAURUS: A SAMPLE EXTRACT SHOWING A DESCRIPTOR AND A NON-DESCRIPTOR	12
FIGURE 4: XML SERIALISATION OF RDF, EXAMPLE.....	16
FIGURE 5: ONTOLOGY TYPES	21
FIGURE 6: ONTOLOGY REPRESENTATION LANGUAGES AND THEIR EXPRESSIVENESS TAKEN FROM [CG00].....	22
FIGURE 7: RDF SCHEMA EXAMPLE MODEL	23
FIGURE 8: LEXICAL OIMODEL.....	25
FIGURE 9: SPANNING OBJECT EXAMPLE.....	26
FIGURE 10: THE ONTOLOGY ENGINEERING FRAMEWORK	33
FIGURE 11: THE FOCUSED WEB CRAWLER.....	36
FIGURE 12: EVALUATION OF THE ONTOLOGY	39
FIGURE 13: COMMUNICATION BETWEEN THE CDS SYSTEM AND THE ONTOLOGY BROWSING INTERFACE	40
FIGURE 14: SCREENSHOT OF THE ADAPTED KAON PORTAL	41
FIGURE 15: MAPPING OF AGROVOC THESAURUS TO ONTOLOGY STRUCTURE.....	45
FIGURE 16: MODELLING OF AGROVOC CATEGORIES.....	46
FIGURE 17: THE ONTOLOGY PRUNING PROBLEM.....	51
FIGURE 18: PRUNING PROCESS – OLD VS. NEWLY ADAPTED VERSION	54
FIGURE 19: FREQUENCY PROPAGATION – FREQUENT CONCEPT WITH INFREQUENT SUPER CONCEPT	55
FIGURE 20: PRUNER VS. PRUNER TRIE, EVALUATION RESULTS.....	60
FIGURE 21: DEPENDENCY OF ALL STATISTICAL ONTOLOGY PARAMETERS ON VARIATION OF THE RATIO PARAMETER (EXEMPLARY FOR THE SETTING TFIDF ALL GEN WITH ONTOLOGY PRUNER TRIE)	62
FIGURE 22: DIFFERENCES IN SIZE BETWEEN LARGEST PRUNED ONTOLOGY AND ALL OTHERS (PRUNER TRIE)	65
FIGURE 23: NUMBER OF DOMAIN SPECIFIC CONCEPTS, WHICH HAVE NOT BEEN IDENTIFIED BY THE AUTOMATIC ONTOLOGY PRUNER	66
FIGURE 24: EXAMPLE MICRO-AVERAGING VS. MACRO-AVERAGING	77
FIGURE 25: DEVELOPMENT OF PRECISION, RECALL AND BREAKEVEN FOR TEST SET $X_{MULTI_EN_DESC}$	92
FIGURE 26: PRECISION VS. RECALL FOR TEST SET $X_{MULTI_EN_DESC}$	92
FIGURE 27: SINGLE-LABEL VS. MULTI-LABEL CLASSIFICATION: COMPARISON OF OVERALL PERFORMANCE	96
FIGURE 28: ONTOLOGY INTEGRATION VS. NO INTEGRATION OF BACKGROUND KNOWLEDGE, $X_{SINGLE_EN_DESC}$	99
FIGURE 29: INFLUENCE OF THE DIFFERENT MODES OF ONTOLOGY INTEGRATION ON THE OVERALL PERFORMANCE (EACH SERIES CORRESPONDS TO A SPECIFIC NUMBER OF TRAINING EXAMPLES PER CLASS, STARTING AT 5).....	100

LIST OF TABLES

TABLE 1: AGROVOC ONTOLOGY STATISTICS	57
TABLE 2: ONTOLOGY PRUNER OUTPUT VS. SUBJECT ASSESSMENT OF THIS OUTPUT	64
TABLE 3: SPECIFICATION CORRECTNESS AND SPECIFICATION RECALL FOR AUTOMATICALLY PRUNED ONTOLOGIES	67
TABLE 4: CONTINGENCY TABLE FOR DOCUMENT x_i	75
TABLE 5: CONTINGENCY TABLE FOR CLASS c_i	76
TABLE 6: GLOBAL CONTINGENCY TABLE.....	76
TABLE 7: RAW TEST DOCUMENT SET FOR AUTOMATIC TEXT CLASSIFICATION, X_{RAW}	86
TABLE 8: COMPILED TEST DOCUMENT SET X_{MULTI} (MULTI-LABEL).....	87
TABLE 9: COMPILED TEST DOCUMENT SET X_{SINGLE} (SINGLE-LABEL)	88
TABLE 10: OVERVIEW ABOUT THE CLASSES OF THE TEST DOCUMENT SETS.....	89
TABLE 11: SINGLE-LABEL CLASSIFICATION ON ENGLISH DOCUMENTS SETS; WORD PRUNING THRESHOLD VS. VARIATION OF TRAINING EXAMPLES PER CLASS; AVERAGE PRECISION OVER 15 TEST RUNS FOR EACH CONFIGURATION	90
TABLE 12: PERFORMANCE OF MULTI-LABEL CLASSIFICATION WITH ENGLISH DOCUMENT SET $X_{MULTI_EN_DESC}$, AVERAGE PERFORMANCE MEASURES OVER 30 TEST RUNS	91
TABLE 13: PERFORMANCE OF MULTI-LABEL CLASSIFICATION WITH ENGLISH DOCUMENT SET $X_{MULTI_EN_CAT}$, AVERAGE PERFORMANCE MEASURES OVER 15 TEST RUNS	93
TABLE 14: PERFORMANCE OF MULTI-LABEL CLASSIFICATION WITH SPANISH DOCUMENT SET $X_{MULTI_FR_DESC}$, AVERAGE PERFORMANCE MEASURES OVER 30 TEST RUNS	94
TABLE 15: PERFORMANCE OF MULTI-LABEL CLASSIFICATION WITH SPANISH DOCUMENT SET $X_{MULTI_ES_DESC}$, AVERAGE PERFORMANCE MEASURES OVER 30 TEST RUNS	95
TABLE 17: AVERAGE PRECISION RESULTS OF SIMPLE LANGUAGE CLASSIFIER	97
TABLE 18: AVERAGE PRECISION OF SINGLE LABEL TEST RUNS IN ALL 3 LANGUAGES.....	97
TABLE 19: PERFORMANCE OF $X_{SINGLE_EN_DESC}$ WITH ONTOLOGY BACKGROUND KNOWLEDGE, AVERAGED PRECISION OVER 30 RUNS.....	98
TABLE 20: PERFORMANCE OF $X_{SINGLE_EN_CAT}$ WITH ONTOLOGY BACKGROUND KNOWLEDGE, AVERAGED PRECISION OVER 30 RUNS.....	123
TABLE 21: PERFORMANCE OF $X_{SINGLE_FR_CAT}$ WITH ONTOLOGY BACKGROUND KNOWLEDGE, AVERAGED PRECISION OVER 30 RUNS.....	123
TABLE 22: PERFORMANCE OF $X_{SINGLE_ES_CAT}$ WITH ONTOLOGY BACKGROUND KNOWLEDGE, AVERAGED PRECISION OVER 15 RUNS.....	123

1 Introduction

1.1 Motivation

The management of large amounts of information and knowledge is of ever increasing importance in today's large organisations. With the ongoing ease of supplying information online, especially in corporate intranets and knowledge bases, finding the right information becomes an increasingly difficult task. Today's search tools perform rather poorly in the sense that information access is mostly based on keyword searching or even mere browsing of topic areas. This unfocused approach often leads to undesired results. The following example illustrates the problem more clearly:

An agriculture scientist would like to find out which organisation established the Agreement on Agriculture. A simple search for "establish Agreement on Agriculture" might result in a huge list of documents containing these words, but actually none of them containing the desired result: WTO or World Trade Organisation. The problem becomes even worse if the result searched for only appears in a foreign language document.

Figure 1 shows an extract of an ontology, which could solve this problem by following links in a graph. The grey ellipses represent generic concepts, whereas the white ones represent specific instances of these concepts. The two concepts shown here are linked by a relationship. An ontology-enabled search application would first identify "Agreement on Agriculture" as a "standard" and would then detect the relationship "establish" to "international organisation" and its instances, and hence solve the problem by extending the search query. This example shows how ontologies can help to improve the management of information. Furthermore, it could provide added value by detecting other relationships that provide the user with more possibilities: for example, standards of other organisations could be presented.

Semantically annotated documents, i.e. documents that are indexed with ontological terms and concepts instead of simple keywords, provide several advantages. First, the ontological abstraction provides robustness against changes in the document. In the above example, the document representation might change using the term 'Agricultural Agreement' instead of 'Agreement on Agriculture'. However, since the document has been annotated with the ontological semantics, this will not affect the search results. Second, since the ontology used

for annotating the document in this example is domain-specific, the semantic meanings and interpretations of keywords are bound to that domain and therefore the retrieval is likely to be more efficient. A term can have several meanings in different domains. By first mapping the keyword to its semantic representation in a specific ontology and using the ontology's linked knowledge structure, a much more focused search approach can be taken. Third, document specific representations no longer affect the search. This is extremely important in the case of multilingual representations. Keywords of several languages are mapped to the same concept in an ontology and are therefore given the same meaning. Multilingual search portals can be established to produce the same results, no matter which language is used for retrieval.

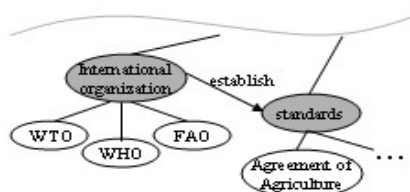


Figure 1: Ontology example, excerpt

An important task in knowledge management facilitating above described search scenario is the classification and indexing of documents. At present, subject specialists are responsible for this time consuming process. However, with today's vast amount of available information on the WWW, automatic support is needed to efficiently manage this task. Ontologies play a critical role in supporting the machine readable semantics needed to facilitate automation. They can be used for providing the categories and keywords needed to describe the content of documents. Automatic text classification tools still lack the necessary precision to replace human indexers and need to be extensively evaluated in different domains.

Before such powerful Semantic Web¹ applications can be built and used within certain domains of knowledge, the basic requirement - a machine readable vocabulary represented by a domain ontology - has to be established. The creation of ontologies is a time consuming task and often carried out in an ad-hoc manner. Only few methodologies exist and existing ones are often extremely complex and need extensive training and expertise. Even less automated tool support is available. Constituting the knowledge base for future Semantic Web applications, domain ontologies have to be created continuously in all possible areas and communities. The need for a reusable methodology is evident.

¹ Refer to [Pal01] for a short introduction to the Semantic Web.

1.2 Approach

The thesis introduces a comprehensive framework for building a domain-specific ontology. The approach combines classical methodologies for human-based ontology engineering with semiautomatic support of a heuristic toolkit. Two methods for ontology acquisition are applied in order to create the domain ontology. The first is to create a small, domain-specific core ontology from scratch. This step is supported by automatically extracting interesting concepts from a corpus of domain texts, which can be used to extend this base ontology. The second acquisition approach takes a well-established thesaurus as a basic vocabulary reference set, and converts it into an ontology representation. Then, a domain specific and a general corpus of texts are used to remove ontology concepts that are not descriptive for the domain from this converted representation. The rationale used here is that domain specific concepts are more frequent in the domain-specific text corpus. The results of these steps are assessed to assemble a first version of the domain specific ontology. This ontology is then accessible through a multilingual web portal to be incorporated into other applications, such as document indexing or keyword searching of indexed documents. It could eventually be used to automatically index documents available through this kind of search application.

Carried out in collaboration with the Food and Agriculture Organisation (FAO)² of the United Nations (UN), the main focus of this thesis is on the adoption of the proposed framework to the specific environment and needs of this large organisation. The framework has been applied to create a prototype biosecurity ontology for the domain of Food Safety, Animal and Plant Health to be incorporated into an Internet Portal to this domain. Within this context, the conversion of a thesaurus into an ontology and evaluations of two automatic tools especially, constitute the central parts of the academic research work. The first evaluation is on a tree-pruning algorithm used in the ontology creation process to retrieve domain specific concepts from the converted thesaurus. The second evaluation is on a text classification application based on support vector machines, enhanced by a domain specific ontology serving as background knowledge for the classification algorithm.

² [<http://www.fao.org>].

1.3 Outline

The next section gives an introduction and overview about the Food and Agriculture Organisation, and the Agricultural Ontology Service (AOS) Project, which provides the bigger context in which the research work of this thesis is embedded. The current information management structure will be introduced briefly, outlining the overall current status and problems within the organisation.

In section 3, I will give an introduction to the idea of the Semantic Web as well as to ontologies and their various representations and engineering approaches. The comprehensive framework for the creation of a multilingual domain ontology is covered in section 4. The application of the framework will be described in the context of the above-mentioned project to establish an International Portal on Food Safety, Animal and Plant Health. The conversion of an existing thesaurus into an ontology representation as well as the adaptation of a multilingual ontology web browser to be embedded into the system is discussed here in detail.

Sections 5 and 6 describe in detail the adaptation and evaluation of two automatic tools constituting parts of the framework. Section 5 describes the thesaurus pruning algorithm used within the ontology creation framework and discusses the results of an empirical evaluation carried out within the context of the project. Section 6 introduces the reader to the area of automatic text classification and describes the adaptation of an already existent automatic text classifier based on support vector machines to incorporate domain specific ontologies. Several evaluation results are discussed against the question of the applicability of the classifier in the context of the FAO and against results of earlier evaluations. Finally, section 7 summarises the findings and results and provides an overview on future work.