

SEMANTIC SIMILARITY MEASURES:  
A COMPARISON STUDY<sup>1</sup>

Paraskevi Raftopoulou and Euripides Petrakis

Technical Report  
TR-TUC-ISL-04-2005

TECHNICAL UNIVERSITY OF CRETE  
DEPARTMENT OF ELECTRONIC AND COMPUTER ENGINEERING

INTELLIGENT SYSTEMS LABORATORY

January 2005

---

<sup>1</sup>This work was carried out as part of the Heraclitus scholarship.

# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Writing and Using Ontologies</b>	<b>7</b>
2.1	Languages for Writing Ontologies . . . . .	8
2.2	Tools for Manipulating Ontologies . . . . .	9
<b>3</b>	<b>Comparing Concepts</b>	<b>10</b>
3.1	Ontology Approaches . . . . .	10
3.1.1	Single ontology approach . . . . .	10
3.1.2	Hybrid approach . . . . .	10
3.1.3	Multiple ontology approach . . . . .	11
3.2	Semantic Similarity Measures . . . . .	12
3.2.1	Edge-Counting Measures . . . . .	12
3.2.2	Information Content Measures . . . . .	14
3.2.3	Feature-Based Measure . . . . .	16
3.2.4	Combinational Measures . . . . .	16
<b>4</b>	<b>Conclusions</b>	<b>19</b>

# List of Figures

3.1	A fragment of the WordNet taxonomy . . . . .	14
-----	--	----

# List of Tables

3.1 Comparison between similarity measures . . . . .	17
--	----

# Chapter 1

## Introduction

The World Wide Web (WWW) can be thought of as a collection of distributed, autonomous and heterogeneous data sources. Large amounts of data are daily produced at an ever increasing rate. In many cases, data once stored is often never accessed again. Although we can imagine larger and more powerful databases and data warehouses in which to store data, humans and programs can access only a small portion of it. Wherefore, the existing technology offers the means to store and access data over the Web based on its simple form, but it seems to lack the ability to discover and extract knowledge from it, and to transform raw data into useful information.

A popular way for discovering data in the Web is the method of *information retrieval*. The user sends a query expressing his information needs and the query is suitably propagated in the net to meet the appropriate data source(s). Information is then retrieved and send back to the user. Depending on the way user expresses his information needs and the way resources in various data sources are described, special mechanisms need to be used in order the query to meet the appropriate data sources, to extract knowledge from the data and to match the query against the data. Our intuition says that if data is described in an “elegant” way, the knowledge extraction and therefore the retrieval of reliable results becomes more efficient.

The above presented scenario is an indicative example revealing that the presentation of data in the Web and the cooperation between Web sources requires sharing of knowledge. *Semantic Web* is defined as “an extension of the current Web in which information is given well-defined meaning, better enabling computers and people to work in cooperation” [5]. Thus, Semantic Web is regarded as the tool that will allow dynamic discovery of semantic meaning from raw data stored in Web resources. Achieving this goal requires tools for extraction, representation and manipulation of knowledge. In this context, *ontologies* are regarded as appropriate modelling structures for representing knowledge existing in or extracted from Web sources [7]. An ontology is a description of the logical structure of a domain, consisting of *terms*, representing domain concepts<sup>1</sup>, *relationships* between these terms, as well as *properties* of each concept, which describe various features and attributes of the concept. Concepts are often arranged in a *taxonomic hierarchy*, at the top of which the most general concept in the domain is defined, while moving downwards in the hierarchy more specific concepts are defined as subclasses of more general concepts or classes [38, 13]. In the general case however, no reference ontology is provided by the data source. In this case, we could automatically create an ontology or a semantic taxonomy using the available resources of the Web source and appropriate mechanisms to extract the knowledge (i.e., the concepts and the relationships) latent to the resources, like for example in [42, 53].

Then, assuming that knowledge is represented in ontologies, data sharing requires com-

---

<sup>1</sup>Words “terms” and “concepts” will be used interchangeably in the rest of this paper.

paring concepts in the same or different ontologies (e.g., comparing a concept across different ontologies). This is exactly the focus of this report. The same concept may be represented in different ways (i.e., may have different definitions in terms of names and properties) in different ontologies. As a result, direct comparison of concepts as keywords (i.e., syntactic comparison) is not always an efficient way for computing similarity between entities with semantic meaning. Obviously such methods cannot take concept properties or relationships between concepts into account. Alternatively, to relate concepts in different ontologies, it is more effective to focus on whether the concepts are semantically (rather than syntactically) similar by finding places in the ontologies where they overlap. This might also be an efficient way for performing tasks such as retrieving results to user queries, for checking ontologies for consistency or coherency (i.e., the same concept has the same meaning in different ontologies), for representation and for redundancy.

The purpose of this report is to present an overview of existing *similarity measures* to compare concepts by their semantic meaning, called henceforth *semantic similarity measures*. Such similarity measures are defined in an intuitive and algorithmic way and work by discovering linguistic relationships or affinities<sup>2</sup> between ontological terms across different ontologies.

In what follows, Section 2 presents the characteristics of and the languages used for an ontology followed by some examples. In Section 3 we present different ways to measure similarity between concepts followed by conclusions in Section 4.

---

<sup>2</sup>Concept used to identify terms with semantic relationships.

## Chapter 2

# Writing and Using Ontologies

Ontologies can be regarded as general tools of information representation on a subject. They can have different roles depending on the application domain and the level of specificity at which they are being used. In general, ontologies can be distinguished into *domain ontologies*, representing knowledge of a particular domain, and *generic ontologies* representing common sense knowledge about the world [55].

There are several examples of general purpose ontologies available including: (a) WordNet<sup>1</sup> [4, 33] attempts to model the lexical knowledge of a native speaker of English. It can be used as both a thesaurus and a dictionary. English nouns, verbs, adjectives, and adverbs are organized into synonym sets, called *synsets*, each representing a concept. (b) SENSUS<sup>2</sup> [24] is a 90,000-node concept thesaurus (ontology) derived as an extension and reorganization of WordNet. Each node in SENSUS represents one concept, i.e., one specific sense of a word, and the concepts are linked in a IS-A hierarchy, becoming more general towards the root of the ontology. (c) The Cyc<sup>3</sup> Knowledge Base (KB) [39, 46] consists of terms and assertions relating those terms, contains a vast quantity of fundamental human knowledge: facts, rules of thumb, and heuristics for reasoning about the objects and events of everyday life. At the present time, the Cyc KB contains nearly two hundred thousand terms and several dozen hand-entered assertions about/involving each term.

Examples of domain specific ontologies include among others ontologies designed around (a) medical concepts such as UMLS<sup>4</sup> [36], SNOMED<sup>5</sup>, MESH<sup>6</sup> [35], (b) genomic data such as GO<sup>7</sup> [6, 15] and (c) spatial data such as SDTS<sup>8</sup>. The Unified Medical Language System (UMLS) contains a very large, multi-purpose and multi-lingual thesaurus concerning biomedical and health related concepts. In particular, it contains information about over 1 million biomedical concepts and 2.8 million concept names from more than 100 controlled vocabularies and classifications (some in multiple languages) used in patient records, administrative health data, bibliographic and full-text databases and expert systems. Furthermore, all the names and meanings are enhanced with attributes and inter-term relationships. UMLS includes other metathesaurus source vocabularies, such as Medical Subject Headings (MeSH) that is the National Library of Medicine's vocabulary thesaurus. MeSH consists of sets of terms naming *descriptors* in a hierarchical structure. Gene Ontology (GO) is a structured network of defined terms that describe gene proteins and concerns all organisms. The Spatial Data Transfer Standard (SDTS) contains an ontology used to describe the underlying conceptual model and the detailed specifications for the content,

---

<sup>1</sup><http://www.cogsci.princeton.edu/~wn/>

<sup>2</sup><http://mozart.isi.edu:8003/sensus2/>

<sup>3</sup><http://www.cyc.com/>, <http://www.opencyc.org/>

<sup>4</sup><http://www.nlm.nih.gov/research/umls>

<sup>5</sup><http://www.snomed.org>

<sup>6</sup><http://www.nlm.nih.gov/mesh>

<sup>7</sup><http://www.geneontology.org>

<sup>8</sup><http://mcmcweb.er.usgs.gov/sdts/>

structure, and format of spatial data, their features and associated attributes. Concepts in SDTS are commonly used on topographic quadrangle maps and hydrographic charts.

Intensive research efforts during the last few years have focused on providing tools for coherent, unambiguous and easy manipulation of information represented as ontologies. Such tools include languages providing the necessary syntax for the efficient representation of concepts and of their semantics as well as tools in the form of algorithms and graphic interfaces for viewing and manipulating the content of ontologies.

## 2.1 Languages for Writing Ontologies

The Resource Description Framework (RDF<sup>9</sup>) is a language for representing information about resources in the Web [2, 21]. It is particularly intended for representing metadata about Web resources, such as the title, author, and modification date of a document. RDF can also be used to represent information about things that can be identified on the Web, even when they cannot be directly retrieved, as for example information about items available from on-line shopping facilities (e.g., information about specifications, prices, and availability). RDF is intended for situations in which this information needs to be processed by applications, as it provides a common framework for expressing this information so it can be exchanged between applications without loss of meaning. RDF is based on the idea of identifying things using Web identifiers (called Uniform Resource Identifiers, or URIs), and describing resources in terms of simple properties and property values, which enables RDF to represent simple statements about resources as a graph of nodes and arcs representing the resources, and their properties and values. RDF also provides an XML-based syntax (called RDF/XML) for recording and exchanging these graphs. Although, RDF provides a way to express simple statements about resources, using named properties and values, it does not define the terms used in those statements. That is the role of RDF Schema (RDF-S<sup>10</sup>) that provides the facilities needed to describe such classes and properties, and to indicate which classes and properties are expected to be used together [28]. The RDF-S facilities are themselves provided in the form of an RDF vocabulary; that is, as a specialized set of predefined RDF resources with their own special meanings.

DAML+OIL<sup>11</sup>, which was the result of an initial joint effort by US and European researchers, is a semantic markup language for Web resources [18, 17]. It builds on RDF and RDF-S, and extends these languages with richer modelling primitives. In particular, DAML+OIL assigns a specific meaning to certain RDF triples. The model-theoretic semantics<sup>12</sup> specify exactly which triples are assigned a specific meaning, and what this meaning is.

The WWW Consortium (W3C) created the Web Ontology Working Group to develop a semantic markup language for publishing and sharing ontologies and the resulting language is Web Ontology Language (OWL<sup>13</sup>). OWL can be used to explicitly represent the meaning of terms in vocabularies and the relationships between those terms. OWL has more facilities for expressing meaning and semantics than XML, RDF, and RDF-S, and thus OWL goes beyond these languages in its ability to represent content on the Web. OWL is a revision of the DAML+OIL Web ontology language, adding more relations between classes (e.g., disjointness), cardinality (e.g., “exactly one”), equality, more properties, more characteristics of properties (e.g., symmetry), and enumerated classes.

To conclude, if machines are expected to perform useful reasoning tasks on Web resources, some language must be used in order to go beyond raw data, to express the semantics of the data and to extract knowledge from it. A summary of the existent recommenda-

---

<sup>9</sup><http://www.w3.org/RDF>

<sup>10</sup><http://www.w3.org/TR/rdf-schema>

<sup>11</sup><http://www.daml.org/language/>

<sup>12</sup><http://www.daml.org/2000/12/daml+oil.daml>

<sup>13</sup><http://www.w3.org/TR/owl-features>



tions related to the Semantic Web follows.

- XML provides a syntax for structured documents, but imposes no semantic constraints on the meaning of these documents.
- RDF is a datamodel describing resources and relations between them and provides a simple semantics for this datamodel. The datamodels can be represented in an XML syntax.
- RDF-S is a vocabulary for describing properties and classes of RDF resources.
- DAML+OIL assigns specific meaning to certain RDF triples.
- OWL adds more vocabulary for describing properties and classes.

There are also efforts for describing the semantics of Web services, resulting in the DAML-S<sup>14</sup> [51] and OWL-S<sup>15</sup> [20] languages.

## 2.2 Tools for Manipulating Ontologies

Examples of tools for manipulating ontologies include Protege-2000<sup>16</sup> [37] and Chimaera<sup>17</sup> [29, 30, 31]. Protege-2000 that allows users to construct domain ontologies, contains a platform that can be extended with graphical widgets for tables, diagrams, animation components to access other knowledge-based systems embedded applications, and has a library that other applications can use to access and display knowledge bases. Chimaera is a software system that supports users in creating and maintaining distributed ontologies on the Web. It supports two major functions that is merging multiple ontologies together and diagnosing<sup>18</sup> individual or multiple ontologies. It also provides users with tasks such as loading knowledge bases in different formats, reorganizing taxonomies, resolving name conflicts, browsing ontologies and editing terms.

---

<sup>14</sup><http://www.daml.org/services>

<sup>15</sup><http://www.mindswap.org/2004/owl-s/>

<sup>16</sup><http://protege.stanford.edu>

<sup>17</sup><http://www.ksl.stanford.edu/software/chimaera/>

<sup>18</sup>Tool used as an ontological sketchpad, and creating classes for example.

## Chapter 3

# Comparing Concepts

This section presents methods of computing the similarity between entities with some semantic meaning, such as concepts (i.e., classes) represented in ontologies, or elements (i.e., resources) represented in schemas. These methods, referred to as *semantic similarity methods*, exploit the fact that the entities which are compared may have (in addition to their name) properties (e.g., in the form of attributes) associated with them, taking also into account the level of generality (or specificity) of each entity within the ontology as well as their relationships with other concepts. Notice that, keyword-based similarity measures cannot use this information. Semantic similarity measures methods might be used for performing tasks such as retrieving results to user queries, for representation and for redundancy of retrieved resources, and for checking ontologies for consistency or coherency.

### 3.1 Ontology Approaches

Ontologies as tools for representing domain knowledge can be used in many different ways. Accordingly, different approaches for comparing concepts within or across ontologies can be defined [58].

#### 3.1.1 Single ontology approach

All information sources are related to one global (unique) ontology providing a common vocabulary for the specification of the entity semantics. A prominent approach is SIMS [3], which includes a hierarchical terminological knowledge base with nodes representing objects, actions and states. Each independent information source is described by relating its objects to the global domain model. Single ontology approach can be applied to cases where all information sources share nearly the same view for a domain (e.g., applications using common sense knowledge may use the Wordnet ontology). Comparing a concept with the ontology is translated into searching for the same or similar concepts within the ontology. How this task can be performed efficiently? How is the notion of “similarity” defined? How close two concepts must be so as to be characterized as “similar”?

#### 3.1.2 Hybrid approach

The semantics of each source is described by its own ontology, but all ontologies are built on one common vocabulary. The shared vocabulary contains basic terms (the primitives) of a domain, upon which the source ontologies are based to built complex<sup>1</sup> terms. As each concept of a source ontology is described by the use of some primitives, the problem of comparing a concept with an ontology can be solved based on methods proposed for the

---

<sup>1</sup>When the primitives are combined by some operators, as for example the union or intersection operator.

case of the single ontology approach. The drawback however is that existing ontologies cannot be used easily, but have to be redeveloped from scratch as all source ontologies must refer to the shared vocabulary.

In hybrid approaches, the interesting point is how the local ontologies are described, i.e. how the terms of the source ontologies are described by the primitives of the shared ontology. For example, in COIN [12] the local description of the information (called *context*) is an attribute value vector. The terms for the context comes out from the common shared vocabulary. In MECOTA [57] each source information is annotated by a label that indicates the semantics of the information. The label combines the primitives from the shared vocabulary.

### 3.1.3 Multiple ontology approach

Different information sources (e.g., knowledge about the application) are described by different ontologies. Knowledge within each ontology may be represented without reference to the other information sources or their ontologies. This approach has no common ontology commitment, thus simplifying the adding or modification of information sources. However, the lack of a common vocabulary makes the comparison of different source ontologies a very complicated task.

The multiple ontology approach, although the most general, is the most difficult to handle involving high complexity algorithmic approaches. A straightforward approach is the hard-coded conversion of all data sources into a common ontology which can be stored in a central warehouse. This approach is costly, while it requires substantial efforts from human experts and is not easily extensible to changes of information sources. There are also approaches build around the idea of using wrappers for the automated or semi-automated generation of mappings from the data sources into the global ontology [54]. An attempt to provide intuitive semantics for mappings between concepts in different ontologies is made in [32], where relationships borrowed from linguistics are used to relate terms in various ontologies. In general, the ontology mapping identifies semantically corresponding terms of different source ontologies (e.g., which terms are semantically equal or similar) and has also to consider different views on a domain (e.g., different aggregation of the ontology concepts), becoming thus a non-trivial task.

The problem of comparing concepts between different ontologies could be affronted by borrowing approaches already used in the database community, i.e. schema mapping by discovering semantic correspondences of attributes or instances across heterogeneous sources. The fundamental approach used in this case is “matching”, which takes two schemas as input and produces a mapping between elements that semantically correspond to each other [41], or maps concepts to schema elements [49]. Approaches of schema matching can be categorized into *label-based* and *instance-based*, according to the different information on which they rely [45]. Label-based approaches consider only the similarity between schema definitions or attribute labels of two information sources. Instance-based methods rely on the content overlap or statistical properties to determine the similarity of two attributes. Studies concerning the ontology mappings that are based on and extend the schema-matching techniques have been proposed, as for example the development of generic match algorithms [34] and the use of mining techniques [14].

An affinity-based unification method for global mapping construction across ontologies is proposed in [10]. The concept of *affinity* is used to identify terms with semantic relationships in different ontologies. The different ontology terms are firstly analyzed to identify those terms with affinity in different ontologies, which are then identified using a hierarchical clustering procedure [10]. The integration across ontologies is finally achieved by using these clusters.

The advent of peer-to-peer (P2P) systems introduce a different view to the problem by taking a social perspective which heavily relies on self-organization. Mappings between different ontologies are done by special mediator agents which are specialized to trans-

late between different ontologies and different languages in [43]. In this approach agents start from simple one-to-one mappings between classes and continue with mappings between complex expressions. Similarly, data sources introduce their own ontologies and then agents can incrementally come up with a global ontology by exchanging translations between the local ones in [1]. Finally, global semantics are seen in [40] as a matter of continuing negotiation, allowing the creation of global mapping that emerges from peer interactions.

## 3.2 Semantic Similarity Measures

We suppose that data in information sources is described by properties and is organized in a taxonomic (subclass-superclass) hierarchy based upon the ontology of the source. When a user sends a query to the Web what is the mechanism for discovering and retrieving answers to the user's inquiry? How are the concepts in the user's query compared with concepts presented in the ontology hierarchies owned by the different information sources? We will give some alternatives to cope with the polysemous nature of natural words, the multiple ways in which the same concept can be described, and the complex terms of source ontologies.

Many measures of semantic similarity with a variety of interesting properties have been proposed. In what follows, we present measures of similarity followed by a short discussion of their properties. Semantic similarity measures can be generally partitioned in four categories: those based on how close the two concepts in the taxonomy are, those based on how much information the two concepts share, those based on the properties of the concepts, and those based on combinations of the previous options.

Let  $\mathcal{C}$  be the set of concepts in an IS-A taxonomy. We want to measure the similarity of two concepts  $c_1, c_2 \in \mathcal{C}$ .

### 3.2.1 Edge-Counting Measures

In the first category we place measures that consider *where* two concepts  $c_1$  and  $c_2$  are in the taxonomy. The following measures are based on a simplified version of *spreading activation theory* [11, 52]. One of the assumptions of the theory of spreading activation is that the hierarchy of concepts is organized along the lines of semantic similarity. Thus, the more similar two concepts are, the more links there are between the concepts and the more closely related they are [44].

**Shortest path [44, 8]:** The first measure has to do with how close in the taxonomy the two concepts are.

$$sim_{sp} = 2MAX - L \quad (3.1)$$

where  $MAX$  is the maximum path length between two concepts in the taxonomy and  $L$  is the minimum number of links between concepts  $c_1$  and  $c_2$ . This measure is a variant on the *distance* method [44] and is principally designed to work with hierarchies. It is motivated by two observations: the behavior of conceptual distance resembles that of a metric, and the conceptual distance between two nodes is often proportional to the number of edges separating the two nodes in the hierarchy. A measure like this might be implemented in an information retrieval system that is based on indexing documents and queries into terms from a semantic hierarchy, or might be applied to help rank the documents to the query. There are many specific questions about the cognitive realism of shortest path measure, however it is a simple and powerful measure in hierarchical semantic nets.

**Weighted links [48]:** Extending the above measure, the use of weighted links is proposed to compute the similarity between two concepts. The weight of a link may be affected

by: (a) the density of the taxonomy at that point, (b) the depth in the hierarchy, and (c) the strength of connotation<sup>2</sup> between parent and child nodes. Then, computing the distance between two concepts is translated into summing up the weights of the traversed links instead of counting them.

**Hirst and St-Onge [16]:** The idea behind this measure is that two concepts  $c_1$  and  $c_2$  are semantically close if they are connected by a path that is not too long and that does not change direction too often.

$$sim_{H\&S}(c_1, c_2) = C - L - kd \quad (3.2)$$

where  $d$  is the number of changes of direction in the path, and  $C$ ,  $k$  are constants. Although this measure gives a different perspective of similarity between two concepts, it seems to poorly perform (see [9]) mainly because it lies in its tendency to wander than in the use of concept relationships.

**Wu and Palmer [59]:** This similarity measure considers the position of concepts  $c_1$  and  $c_2$  in the taxonomy relatively to the position of the most specific common concept  $c$ . As there may be multiple parents for each concept, two concepts can share parents by multiple paths. The most specific common concept  $c$  is the common parent related with the minimum number of IS-A links with concepts  $c_1$  and  $c_2$ .

$$sim_{W\&P}(c_1, c_2) = \frac{2H}{N_1 + N_2 + 2H} \quad (3.3)$$

where  $N_1$  and  $N_2$  is the number of IS-A links from  $c_1$  and  $c_2$  respectively to the most specific common concept  $c$ , and  $H$  is the number of IS-A links from  $c$  to the root of the taxonomy. It scores between 1 (for similar concepts) and 0.

**Li et al. [25]:** The following similarity measure, which was intuitively and empirically derived, combines the shortest path length between two concepts  $c_1$  and  $c_2$ ,  $L$ , and the depth in the taxonomy of the most specific common concept  $c$ ,  $H$ , in a non-linear function.

$$sim_{Li}(c_1, c_2) = e^{-\alpha L} \cdot \frac{e^{\beta H} - e^{-\beta H}}{e^{\beta H} + e^{-\beta H}} \quad (3.4)$$

where  $\alpha \geq 0$  and  $\beta > 0$  are parameters scaling the contribution of shortest path length and depth respectively. Based on [25] the optimal parameters are  $\alpha = 0.2$  and  $\beta = 0.6$ . This measure is motivated by the fact that information sources are infinite to some extent while humans compare word similarity with a finite interval between completely similar and nothing similar. Intuitively the transformation between an infinite interval to a finite one is non-linear. It is thus obvious that this measure scores between 1 (for similar concepts) and 0.

The above mentioned measures are based only on taxonomic (IS-A) links between concepts, assuming that links in the taxonomy represent distances. However, the density of terms throughout the taxonomy is generally not constant. Typically, more general terms exist higher in the hierarchy and represent a smaller set of nodes than the larger number of more specific terms that populate a much denser space lower in the hierarchy. For example, specify that distance between *plant* and *animal* is 2 in WordNet (their common parent is *living thing*), and the distance between *zebra* and *horse* is also 2 (their common parent is *equine*). Intuitively *horse* and *zebra* seem more closely related than *plant* and *animal*. Using in our example either the *Wu & Palmer* measure, the measure based on the *Weighted Links* (if the link weights are fixed accordingly) or the *Li et al.* measure, we take

<sup>2</sup>The *connotation* of a term is the list of membership conditions for the denotation. The *denotation* of a term is the class of things to which the term correctly applies. For example, the connotation of the general term “square” is “rectangular and equilateral”, while its denotation is all squares.

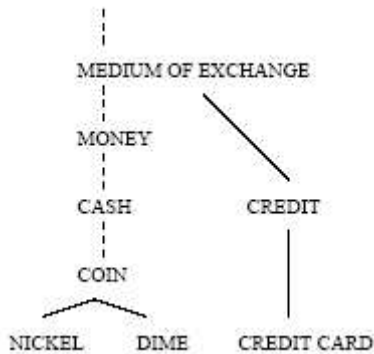


Figure 3.1: A fragment of the WordNet taxonomy

into account the fact that the first two terms occupy a much higher place in the hierarchy than the latter two terms and the results will be more realistic. Furthermore, in taxonomies there is wide variability in what is covered by a single taxonomic link. For example, *safety valve* IS-A *valve* seems much narrower than *knitting machine* IS-A *machine*. The *Weighted Links* measure may take into account the strength of links if link weights are computed accordingly. Finally, experimental results presented in [25] have demonstrated that the *Li et al.* measure significantly outperforms previous measures.

In what follows, we present measures involving information content, which seem to perform better than edge-counting measures.

### 3.2.2 Information Content Measures

In this category, similarity measures are based on the *information content* of each concept. The notion of information content of the concept practically has to do with the frequency of the term in a given document collection. The frequencies of terms in the taxonomy are estimated using noun frequencies in some large (1,000,000 word) collection of texts [47]. Furthermore, the key to the similarity of two concepts is the extent to which they share information in common, indicated by a highly specific concept that subsumes them both.

Associating probabilities with concepts in the taxonomy, let the taxonomy be augmented by the function  $p : \mathcal{C} \rightarrow [0, 1]$ , such that for any concept  $c \in \mathcal{C}$ ,  $p(c)$  is the probability of encountering an instance of concept  $c$ . The concept probability is defined as  $p(c) = \text{freq}(c)/N$ , where  $N$  is the total number of terms in the taxonomy,  $\text{freq}(c) = \sum_{n \in \text{words}(c)} n$  and  $\text{words}(c)$  is the set of terms subsumed by  $c$ . This function implies that if  $c_1$  IS-A  $c_2$ , then  $p(c_1) \leq p(c_2)$ , which intuitively means that the more general the concept is, the higher its associated probability. Then, the information content of a concept  $c$  can be quantified as the log likelihood,  $-\ln p(c)$ , which means that as probability increases, informativeness decreases, so the more abstract a concept, the lower its information content.

Given these probabilities, several measures of semantic similarity, presented later in the section, have been defined. All these measures use the information content of the shared parents of two terms  $c_1$  and  $c_2$  (see Equation 3.5), where  $S(c_1, c_2)$  is the set of concepts that subsume  $c_1$  and  $c_2$ . As there may be multiple parents for each concept, two concepts can share parents by multiple paths. We take the minimum  $p(c)$  when there is more than one shared parents, and then we call concept  $c$  the *most informative subsumer*.

$$p_{\text{mis}}(c_1, c_2) = \min_{c \in S(c_1, c_2)} \{p(c)\} \quad (3.5)$$

For example, in Figure 3.1 *coin*, *cash*, etc. are all members of  $S(\text{nickel}, \text{dime})$ , but the term that is structurally the minimal upper bound is *coin*, and will also be the most

informative subsumer. The information content of the most informative subsumer will be used to quantify the similarity of the two words.

**Lord et al. [27]:** The first way to compare two terms is by using a measure that simply uses the probability of the most specific shared parent.

$$sim_{Lord}(c_1, c_2) = 1 - p_{mis} \quad (3.6)$$

The probability-based similarity score takes values between 1 (for the very similar concepts) and 0. It is used in order to access the extent to which similarity judgements might be sensitive to frequency per se, rather than information content.

**Resnik [47]:** The next measure uses the information content of the shared parents.

$$sim_{Resnik}(c_1, c_2) = -\ln p_{mis} \quad (3.7)$$

This measure signifies that the more information two terms share in common, the more similar they are, and the information shared by two terms is indicated by the information content of the term that subsume them in the taxonomy. As  $p_{mis}$  can vary between 0 and 1, this measure varies between infinity (for very similar terms) to 0. In practice, if  $N$  is the number of terms in the taxonomy, the maximum value of  $p_{mis}$  is  $1/N$  (see Equation 3.5), and the maximum value of the measure is defined by  $-\ln(1/N) = \ln(N)$ . Thus, this measure provides us with information such as the size of the corpus; a large numerical value indicates a large corpus. Furthermore, the score from comparing a term with itself depends on where in the taxonomy the term is, with less frequently occurring terms having higher scores, and thus the measure reveals information about the usage within corpus of the part of the ontology queried.

**Lin [26]:** This measure uses both the amount of information needed to state the *commonality* of two terms and the information needed to fully *describe* these terms.

$$sim_{Lin}(c_1, c_2) = \frac{2 \ln p_{mis}(c_1, c_2)}{\ln p(c_1) + \ln p(c_2)} \quad (3.8)$$

As  $p_{mis} \geq p(c_1)$  and  $p_{mis} \geq p(c_2)$ , the values of this measure vary between 1 (for similar concepts) and 0. In this case, a term compared with itself will always score 1, hiding the information revealed by the *Resnik* measure. However, the *Resnik* measure depends solely on the information content of the shared parents, and there are only as many discrete scores as there are ontology terms. By using the information content of both the compared terms and the shared parent the number of discrete scores is quadratic in the number of terms appearing in the ontology [27], thus augmenting the probability to have different scores for different pairs of terms. Consequently, using this measure to compare the terms of an ontology can have a better ranking of similarity than the *Resnik* measure.

**Jiang et al. [19]:** Contrary to the above similarity measures, this measure is of *semantic distance*.

$$dist_{Jiang}(c_1, c_2) = -2 \ln p_{mis}(c_1, c_2) - (\ln p(c_1) + \ln p(c_2)) \quad (3.9)$$

Thus, the similarity between two concepts  $c_1$  and  $c_2$ ,  $sim_{Jiang}(c_1, c_2)$ , is computed as  $1 - dist_{Jiang}(c_1, c_2)$ . This measure can give arbitrarily large values, like the *Resnik* measure, although in practice has a maximum value of  $2 \ln(N)$ , where  $N$  is the size of the corpus. Furthermore, it combines information content from the shared parent and the compared concepts, as the *Lin* measure. Thus, this measure seems to combine the properties of the above presented measures, i.e. provides us with both information about the size of the ontology and ranking of different term pairs.

### 3.2.3 Feature-Based Measure

Up to now, the features of the terms in the ontology are not taken into account. However, the features of a term contain valuable information concerning knowledge about the term. The following measure considers also the features of terms in order to compute similarity between different concept, while it ignores the position of the terms in the taxonomy and the information content of the term.

**Tversky [56]:** This measure is based on the *description sets* of the terms. We suppose that each term is described by a set of words indicating its properties or features. Then, the *more common* characteristics two terms have and the *less non-common* characteristics they have, the more similar the terms are.

$$sim_{Tversky}(c_1, c_2) = \frac{|C_1 \cap C_2|}{|C_1 \cap C_2| + \kappa|C_1 \setminus C_2| + (\kappa - 1)|C_2 \setminus C_1|} \quad (3.10)$$

where  $C_1, C_2$  correspond to description sets of terms  $c_1$  and  $c_2$  respectively and  $\kappa \in [0, 1]$  defines the relative importance of the non-common characteristics. This measure scores between 1 (for similar concepts) and 0, it increases with commonality and decreases with the difference between the two concepts. In reverse to all the above presented measures, it has nothing to do with the taxonomy and the subsumers of the terms, and seems to better exploit the properties of the ontology used.

In the above presented measure, the determination of  $\kappa$  is based on the observation that similarity is not necessarily a symmetric relation: the common, as opposed to the different, features between a subclass and its superclass have a larger contribution to the similarity evaluation than the common features in the inverse direction. Given this assumption, it provides a systematic approach to determine the asymmetry of a similarity evaluation.

### 3.2.4 Combinational Measures

The next approaches used to compare two concepts  $c_1$  and  $c_2$  combine some of the above presented approaches, considering the path connecting the two terms in the taxonomy, the IS-A links of the terms with their parents in the graph and the features of the terms.

**Rodriguez et al. [50]:** This similarity measure suggests a different identification of distinguishing features than the typical single classification of features into attributes by classifying them into *functions*, *parts* and *attributes*. Functions represent what is done to or with instances of a class, parts are structural elements of a class, and attributes correspond to additional characteristics. For example, considering the term *college*, a function is *to educate*, its parts may be *roof* and *floor*, and other attributes can be *architectural properties*. Then, we consider all the distinguishing features of a term and the global similarity function is a weighted sum of the similarity values for parts, functions, and attributes, where  $\omega_p, \omega_f$  and  $\omega_a$  are the corresponding weights.

$$sim_{Rodriguez}(c_1, c_2) = \omega_p \cdot S_p(c_1, c_2) + \omega_f \cdot S_f(c_1, c_2) + \omega_a \cdot S_a(c_1, c_2) \quad (3.11)$$

where  $\omega_p, \omega_f$  and  $\omega_a \geq 0$  and  $\omega_p + \omega_f + \omega_a = 1$ . For each type of distinguishing features,  $S_p, S_f$  and  $S_a$  a similarity function  $sim_{Tversky}(c_1, c_2)$  is used based on the Tversky feature-matching model.

Like the *Tversky* measure, this measure also uses the number of common and different features between two terms to compute their similarity. However, this measure differs from the *Tversky* measure in the following ways: (a) it considers a non-typical classification of features in functions, parts and attributes and (b) it defines  $\kappa$  in terms



Property	Knappe	Rodriguez	Tversky	Jiang	Lin	Resnik	Lord	Li	Wu & Palmer	Hirst & St-Onge	shortest path
increase with commonality	yes	yes	yes	yes	yes	yes	yes	yes	yes	no	no
decrease with difference	no	yes	yes	yes	yes	no	no	yes	yes	yes	yes
information content	no	no	no	yes	yes	yes	yes	no	no	no	no
position in hierarchy	yes	no	no	yes	yes	yes	yes	yes	yes	yes	yes
path length	no	yes	no	no	no	no	no	yes	yes	yes	yes
max value = 1	yes	yes	yes	no	yes	no	yes	yes	yes	no	yes
symmetric	no	no	no	yes	yes	yes	yes	yes	yes	no	yes
different perspectives	yes	yes	yes	yes	yes	no	no	yes	yes	yes	no

Table 3.1: Comparison between similarity measures

of the distance among terms in the hierarchy; function  $\kappa$  has to do with the distance between terms  $c_1, c_2$  and the most informative subsumer  $c_{mis}$ :

$$\kappa(c_1, c_2) = \begin{cases} \frac{d(c_1, c_{mis})}{d(c_1, c_2)}, & d(c_1, c_{mis}) \leq d(c_2, c_{mis}); \\ 1 - \frac{d(c_1, c_{mis})}{d(c_1, c_2)}, & d(c_1, c_{mis}) > d(c_2, c_{mis}). \end{cases} \quad (3.12)$$

where  $d(c_1, c_2) = d(c_1, c_{mis}) + d(c_2, c_{mis})$ .

**Knappe [22]:** This measure is primarily based on the aspect that there may be *multiple paths* connecting two concepts. Taking all possible paths involves a substantial increase in complexity. Thus, the general idea puts emphasis on the “shared” concepts and a similarity measure representing the part of the ontology covering the compared concepts is defined. Furthermore, there is the notion of complex concepts that allows a concept to be constituted by more than one term.

Initially, the *term decomposition*  $\tau(c)$  of a concept  $c$  into a set  $C$  is defined, and then the *upwards expansion*  $\varpi(C)$  of  $C$  is performed. The term decomposition of  $c$  is defined as the set of all concepts included in  $c$  (if  $c$  is a complex concept, otherwise this set includes only  $c$ ) and all attributes of these concepts. If for example, the initial concept  $c$  is “dog” the term decomposition could be the set  $C = \{dog, colour\}$ . The upwards expansion,  $\varpi(C)$ , involves the IS-A links of all elements in  $C$ .

Let  $u(c)$  be the set of nodes upwards reachable from  $c$ , that is  $u(c) = \varpi(\tau(c))$ . The reachable nodes shared by both  $c_1$  and  $c_2$  are  $u(c_1) \cap u(c_2)$ . Then, we consider the upward and downward directions in the graph as *generalization* and *specialization* respectively. Three major desirable properties are considered in defining the similarity function: (a) the cost of generalization should be significantly higher than the cost of specialization, indicating that the similarity function cannot be symmetrical, (b) the cost for traversing edges should be lower when nodes are more specific and (c) further specialization implies reduced similarity.

$$sim_{Knappe}(c_1, c_2) = \rho \frac{|u(c_1) \cap u(c_2)|}{|u(c_1)|} + (1 - \rho) \frac{|u(c_1) \cap u(c_2)|}{|u(c_2)|} \quad (3.13)$$

where  $\rho \in [0, 1]$  determines the degree of influence of generalizations.

This measure scores between 1 (for matching concepts) and 0. The purpose of this similarity measure is to introduce soft rather than crisp evaluation, since we usually want to look for similar rather than exactly matching values. Furthermore, the idea of concept expansion leads the similarity matching towards a set comparison, incorporating in the similarity measure the knowledge represented by the ontology.

The key properties of the similarity measures presented in the previous sections are summarized in Table 3.1. We consider whether the similarity measures are affected by the common characteristics of the compared concepts and whether the differences between the concepts cause the measures to decrease. Furthermore, we think the relation of the similarity measures with the taxonomy and the taxonomic relations, i.e. whether the position of the concepts in the taxonomy and the number of IS-A links are considered. It is also

presented whether the similarity measures are taking into account the information content of the concepts, whether they are bounded or return infinite values, whether they are symmetric (i.e.,  $sim(c_1, c_2) = sim(c_2, c_1)$ ), and whether they give different perspectives.

## Chapter 4

# Conclusions

In this report, we discussed the shift towards the Semantic Web which justifies the need to share knowledge between different information sources and thus the need to organize data into ontologies. We then introduced the definition and the basic elements of an ontology, along with some examples of ontologies and tools for their manipulation, followed by a brief discussion of the way ontologies can be used. Furthermore, we presented aspects and measures for comparing concepts.

We initially presented measures that consider only the taxonomic positions of the concepts compared, and then we moved towards more sophisticated measures, which exploit the information content of the concepts and the properties given to the concepts by the use of the ontologies. These measures either use the subclass-superclass relation, the properties describing the concepts, or both of them.

We think that the measures proposed in [50] and in [23, 22] seem to be the most efficient, as they combine the attributes of the concepts with their taxonomic relations. Specially, the *Knappe* measure is based on knowledge introduced in data by the use of ontologies and better adopt the idea of looking similar concepts, as it expands crisp values into sets including more concepts. However, the efficiency of this measure remains to be explored experimentally.

# Bibliography

- [1] K. Aberer, P. Cudre-Mauroux, and M. Hauswirth. The Chatty Web: Emergent Semantics Through Gossiping. In *Proceedings of the 12th International World Wide Web Conference (WWW'03)*, Budapest, Hungary, 20-24 May 2003. ACM.
- [2] S. Alexaki, V. Christophides, G. Karvounarakis, D. Plexousakis, K. Tolle, B. Amann, I. Fundulaki, M. Scholl, and A.-M. Vercoustre. Managing RDF Metadata for Community Webs. In *Proceedings of the ER'00 2nd International Workshop on the World Wide Web and Conceptual Modeling (WCM'00)*, pages 140–151, Salt Lake City, Utah, 9-12 October 2000.
- [3] Yigal Arens, Chun-Nan Hsu, and Craig A. Knoblock. Query Processing in the SIMS Information Mediator. In *Advanced Planning Technology*, California, USA, 1996. AAAI Press.
- [4] R. Beckwith and G. A. Miller. Implementing a Lexical Network. Technical Report 43, Princeton University, 1992.
- [5] T. Berners-Lee, J. Hendler, and O. Lassila. The Semantic Web. *Scientific American*, 279(5):34–43, May 2001.
- [6] J.A. Blake and M. Harris. The Gene Ontology Project: Structured vocabularies for molecular biology and their application to genome and expression analysis. In A.D. Baxevanis, D.B. Davison, R. Page, G. Stormo, and L. Stein, editors, *Current Protocols in Bioinformatics*. Wiley and Sons, Inc., New York, 2003.
- [7] C. Brewster and K. O'Hara. Knowledge Representation with Ontologies: The Present and Future. *IEEE Intelligent Systems*, 19(1):72–81, January/February 2004.
- [8] H. Bulskov, R. Knappe, and T. Andreasen. On Measuring Similarity for Conceptual Querying. In T. Andreasen, A. Motro, H. Christiansen, and H.L. Larsen, editors, *Proceedings of the 5th International Conference on Flexible Query Answering Systems (FQAS'02)*, volume 2522 of *LNAI*, pages 100–111, Copenhagen, Denmark, 27-29 October 2002.
- [9] Alexander Butanisky and Graeme Hirst. Semantic distance in WordNet: An experimental, application-oriented evaluation of five measures. In ..., 1999.
- [10] S. Castano, V. Antonelis, and S. De Capitani di Vimercati. Global Viewing of Heterogeneous Data Sources. *IEEE Transactions on Knowledge and Data Engineering*, 13(2):277–297, March/April 2001.
- [11] P.R. Cohen and R. Kjeldsen. Information Retrieval by Constrained Spreading Activation in Semantic Networks. *Information Processing and Management*, 23(4):255–268, 1987.
- [12] Cheng Hian Goh. *Representing and Reasoning about Semantic Conflicts in Heterogeneous Information Sources*. Phd, MIT, 1997.

- [13] T.R. Gruber. A Translation Approach to Portable Ontologies. *Knowledge Acquisition*, 5(2):199–220, 1993.
- [14] Bin He, Kevin ChenChuan Chang, and Jiawei Han. Discovering Complex Matchings across Web Query Interfaces: A Correlation Mining Approach. In *Proceedings of the 2004 ACM SIGKDD Conference (KDD'04)*, Seattle, Washington, August 2004.
- [15] D.P. Hill, J.A. Blake, J.E. Richardson, and M. Ringwald. Extension and Integration of the Gene Ontology (GO): Combining GO vocabularies with external vocabularies. *Genome Res*, 12:1982–1991, 2002.
- [16] Graeme Hirst and David St-Onge. Lexical Chains as Representations of Context for the Detection and Correction of Malapropisms. In *Proceedings of Fellbaum*, pages 305–332, 1998.
- [17] I. Horrocks, P.F. Patel-Scheiner, and F. van Harmelen. Reviewing the Design of DAML+OIL: An Ontology Language for the Semantic Web. In *Proceedings of 18th National Conference on Artificial Intelligence (AAAI'02)*, 2002.
- [18] Ian Horrocks. DAML+OIL: a Reasonable Web Ontology Language. In *Proceedings of the International Conference on Extending DataBase Technology*, March 2002.
- [19] J.J. Jiang and D.W. Conrath. Semantic Similarity Based on Corpus Statistics and Lexical Taxonomy. In *Proceedings of the International Conference on Research in Computational Linguistic*, Taiwan, 1998.
- [20] Lalana Kagal, Grit Denker, Tim Finin, Massimo Paolucci, Naveen Srinivasan, and Katia Sycara. An Approach to Confidentiality and Integrity for OWL-S. In *Proceedings of the 1st International Semantic Web Services Symposium (ISWSS'04)*, AAAI'04 Spring Symposium Series, 22-24 March 2004.
- [21] G. Karvounarakis, V. Christophides, D. Plexousakis, and S. Alexaki. Querying RDF Descriptions for Community Web Portals. In *Proceedings of the 17ièmes Journées Bases de Données Avancées (BDA'01)*, pages 133–144, Agadir, Maroc, 29 October - 2 November 2001.
- [22] R. Knappe, H. Bulskov, and T. Andreasen. On Similarity Measures for Content-Based Querying. In O. Kaynak, editor, *Proceedings of the 10th International Fuzzy Systems Association World Congress (IFSA'03)*, pages 400–403, Instsnbul, Turkey, 29 June - 2 July 2003.
- [23] R. Knappe, H. Bulskov, and T. Andreasen. Similarity Graphs. In N. Zhong, Z.W. Ras, S. Tsumoto, and E. Suzuki, editors, *Proceedings of the 14th International Symposium on Methodologies for Intelligent Systems (ISMIS'03)*, number 2871 in LNAI, pages 668–672, Maebashi, Japan, 28-31 October 2003.
- [24] K. Knight and S. Luk. Building a Large-Scale Knowledge Base for Machine Translation. In *Proceedings of thr National Conference on Artificial Intelligence (AAAI'94)*, Seattle, WA, 1994.
- [25] Yuhua Li, Zuhair A. Bandar, and David McLean. An Approach for Measuring Semantic Similarity between Words Using Multiple Information Sources. *IEEE Transactions on Knowledge and Data Engineering*, 15(4):871–882, July/August 2003.
- [26] D. Lin. Principle-Based Parsing Without Overgeneration. In *Proceedings of the 31st Annual Meeting of the Association for Computational Linguistics (ACL'93)*, pages 112–120, Columbus, Ohio, 1993.

- [27] P.W. Lord, R.D. Stevens, A. Brass, and C.A. Goble. Investigating Semantic Similarity Measures across the Gene Ontology: the Relationship between Sequence and Annotation. *Bioinformatics*, 19(10):1275–83, 2003.
- [28] A. Magkanaraki, S. Alexaki, V. Christophides, and D. Plexousakis. Benchmarking RDF Schemas for the Semantic Web. In *Proceedings of the 1st International Semantic Web Conference (ISWC'02)*, Sardinia, Italy, 9-12 June 2002.
- [29] D.L. McGuinness. Conceptual Modeling for Distributed Ontology Environments. In *Proceedings of the 8th International Conference on Conceptual Structures Logical, Linguistic, and Computational Issues (ICCS'00)*, Darmstadt, Germany, 14-18 August 2000.
- [30] D.L. McGuinness, R. Fikes, J. Rice, and S. Wilder. An Environment for Merging and Testing Large Ontologies. In *Proceedings of the 7th International Conference on Principles of Knowledge Representation and Reasoning (KR'00)*, Breckenridge, Colorado, USA, 12-15 April 2000.
- [31] D.L. McGuinness, R. Fikes, J. Rice, and S. Wilder. The Chimaera Ontology Environment. In *Proceedings of the 17th National Conference on Artificial Intelligence (AAAI'00)*, Austin, Texas, 30 July - 3 August 2000.
- [32] E. Mena, V. Kashyap, A. Sheth, and A. Illarramendi. OBSERVER: An Approach for Query Processing in Global Information Systems based on Interoperation across Pre-existing Ontologies. In *Proceedings of the 1st IFCIS International Conference on Cooperative Information Systems (CoopIS'96)*, Brussels, 1996.
- [33] G. A. Miller, R. Bechwith, C. Felbaum, D. Gross, and K. Miller. Introduction to WordNet: an on-line lexical database. *International Journal of Lexicography*, 3(4):235–244, 1990.
- [34] P. Mork and P.A. Bernstein. Adapting a Generic Match Algorithm to Align Ontologies of Human Anatomy. In *Proceedings of the 20th International Conference on Data Engineering*, pages 787–790, Boston, USA, 30 March - 2 April 2004.
- [35] S.J. Nelson, D. Johnston, and B.L. Humphreys. Relationships in Medical Subject Headings. In C.A. Bean and R. Green, editors, *Relationships in the Organization of Knowledge*, pages 171–184. Kluwer Academic Publishers, New York, 2001.
- [36] S.J. Nelson, T. Powell, and B.L. Humphreys. The Unified Medical Language System (UMLS) Project. In A. Kent and C.M. Hall, editors, *Encyclopedia of Library and Information Science*, pages 369–378. Marcel Dekker, Inc., New York, 2002.
- [37] N. F. Noy, M. Sintek, S. Decker, M. Crubezy, R. W. Fergerson, and M. A. Musen. Creating Semantic Web Contents with Protege-2000. *IEEE Intelligent Systems*, 16(2):60–71, 2001.
- [38] N.F. Noy and D.L. McGuinness. Ontology Development 101: A Guide to Creating Your First Ontology. Report SMI-20010880, Stanford University, Stanford, CA, 2001.
- [39] T. O'Hara, N. Salay, M. Witbrock, D. Schneider, B. Aldag, S. Bertolo, K. Panton, F. Lehmann, and et al. Inducing Criteria for Mass Noun Lexical Mappings using the Cyc KB, and its Extension to WordNet. In *Proceedings of the 5th International Workshop on Computational Semantics (IWCS-5)*, Tilburg, The Netherlands, 15-17 January 2003.
- [40] A. M. Ouksel. In-Context Peer-to-Peer Information Filtering on the Web. In *SIGMOD Record*, volume 32, pages 65–70, September 2003.

- [41] C. Parent and S. Spaccapietra. Issues and Approaches of Database Integration. *Communications of the ACM*, 41(5):166–178, 1998.
- [42] Andrew Pargellis, Eric Fosler-Lussier, Alexandros Potamianos, and Chin-Hui Lee. Auto-Induced Semantic Classes. *Speech Communication*, 43:183–203, August 2004.
- [43] A. D. Preece, K.-Y. Hui, W. Gray, P. Marti, Z. Cui, and D. M. Jones. Kraft: An Agent Architecture for Knowledge Fusion. *International Journal of Cooperative Information Systems (IJCIS)*, 10(1-2):171–195, 2001.
- [44] R. Rada, H. Mili, E. Bicknell, and M. Blettner. Development and Application of a Metric on Semantic Nets. *IEEE Transactions on Systems, Man, and Cybernetics*, 19(1):17–30, January/February 1989.
- [45] Erhard Rahm and Philip A. Bernstein. A Survey of Approaches to Automatic Schema Matching. *The VLDB Journal*, 10(4):334–350, 2001.
- [46] S. Reed and D. Lenat. Mapping Ontologies into Cyc. In *Proceedings of the AAAI’02 Conference Workshop on Ontologies For The Semantic Web*, Edmonton, Canada, July 2002.
- [47] O. Resnik. Semantic Similarity in a Taxonomy: An Information-Based Measure and its Application to Problems of Ambiguity and Natural Language. *Journal of Artificial Intelligence Research*, 11:95–130, 1999.
- [48] R. Richardson, A. Smeaton, and J. Murphy. Using WordNet as a Knowledge Base for Measuring Semantic Similarity Between Words. Technical Report Working paper CA-1294, School of Computer Applications, Dublin City University, Dublin, Ireland, 1994.
- [49] N. Rische, J. Yuan, R. Athauda, S.C. Chen, X. Lu, X. Ma, A. Vaschillo, A. Shaposhnikov, and D. Vasilevsky. Semantic Access: Semantic Interface for Querying Databases. In *Proceedings of the 26th International Conference On Very Large Data Bases*, pages 591–594, 2000.
- [50] M.A. Rodriguez and M.J. Egenhofer. Determining Semantic Similarity Among Entity Classes from Different Ontologies. *IEEE Transactions on Knowledge and Data Engineering*, 15(2):442–456, March/April 2003.
- [51] M. Sabou, D. Richards, and S. van Splunter. An Experience Report on using DAML-S. In *Proceedings of the 12th International World Wide Web Conference Workshop on E-Services and the Semantic Web (ESSW’03)*, Budapest, 2003.
- [52] G. Salton and C. Buckley. On the Use of Spreading Activation Methods in Automatic Information. In *Proceedings of the 11th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 147–160. ACM Press, 1988.
- [53] Peter Spyns, A. Johannes Pretorius, and Marie-Laure Reinberger. Evaluating DOGMA-lexons Generated Automatically from a Text Corpus. In *Proceedings of the 14th International Conference on Knowledge Engineering and Knowledge Management (EKAW’04)*, Whittlebury Hall, Northamptonshire, UK, 5-8 October 2004.
- [54] H. Stuckenschmidt and H. Wache. Context Modeling and Transformation for Semantic Interoperability. In *Knowledge Representation Meets Data Base (KRDB’00)*, 2000.

- [55] R. Studer. Knowledge Engineering and Agent Technology. In J. Cuenca and et al., editors, *Situation and Perspective of Knowledge Engineering*. IOS Press, Amsterdam, 2000.
- [56] A. Tversky. Features of Similarity. *Psychological Review*, 84(4):327–352, 1977.
- [57] H. Wache, T. Scholz, H. Stieghahn, and B. König-Ries. An Integration Method for the Specification of Rule-Oriented Mediators. In Yahiko Kambayashi and Hiroki Takakura, editors, *Proceedings of the international Symposium on Database Applications in Non-Traditional Environments (DANTE'99)*, pages 109–112, Kyoto, Japan, 28-30 November 1999.
- [58] H. Wache, T. Vögele, U. Visser, H. Stuckenschmidt, G. Schuster, H. Neumann, and H. Hubner. Ontology-Based Integration of Information - A Survey of Existing Approaches. In *Proceedings of the IJCAI'01 Workshop on Ontologies and Information Sharing*, Seattle, WA, 2001.
- [59] Z. Wu and M. Palmer. Verb Semantics and Lexical Selection. In *Proceedings of the 32nd Annual Meeting of the Association for Computational Linguistics (ACL'94)*, pages 133–138, Las Cruces, New Mexico, 1994.