

Swoogle: Searching for Knowledge on the Semantic Web *

Tim Finin, Li Ding, Rong Pan, Anupam Joshi, Pranam Kolari, Akshay Java and Yun Peng
University of Maryland Baltimore County, Baltimore, MD

Introduction

Most knowledge on the Web is encoded as natural language text, which is convenient for human users but very difficult for software agents to understand. Even with increased use of XML-encoded information, software agents still need to process the tags and literal symbols using application dependent semantics. The Semantic Web offers an approach in which knowledge can be published by and shared among agents using symbols with a well defined, machine-interpretable semantics.

The Semantic Web is a “web of data” in that (i) both ontologies and instance data are published in a distributed fashion; (ii) symbols are either ‘literals’ or universally addressable ‘resources’ (URI references) each of which comes with unique semantics; and (iii) information is semi-structured. The Friend-of-a-Friend (FOAF) project (<http://www.foaf-project.org/>) is a good application of the Semantic Web in which users publish their personal profiles by instantiating the *foaf:Person* class and adding various properties drawn from any number of ontologies.

The Semantic Web’s distributed nature raises significant data access problems – how can an agent discover, index, search and navigate knowledge on the Semantic Web? Swoogle (Ding *et al.* 2004) was developed to facilitate web-scale semantic web data access by providing these services to both human and software agents. It focuses on two levels of knowledge granularity: URI based *semantic web vocabulary* and *semantic web documents* (SWDs), i.e., RDF and OWL documents encoded in XML, NTriples or N3.

Figure 1 shows Swoogle’s architecture. The **discovery** component automatically discovers and revisits SWDs using a set of integrated web crawlers. The **digest** component computes metadata for SWDs and *semantic web terms* (SWTs) as well as identifies relations among them, e.g., “an SWD instantiates an SWT class”, and “an SWT class is the domain of an SWT property”. The **analysis** component uses cached SWDs and their metadata to derive analytical reports, such as classifying ontologies among SWDs and ranking SWDs by their importance. The **service** component sup-

ports both human and software agents through conventional web interfaces and SOAP-based web service APIs. Two key services are (i) a *swoogle search* service that searches for SWDs by constraints on their URLs, the sites which host them, and the classes/properties used or defined by them and (ii) an *ontology dictionary* service that searches for SWTs and their relationships with other SWTs and SWDs.

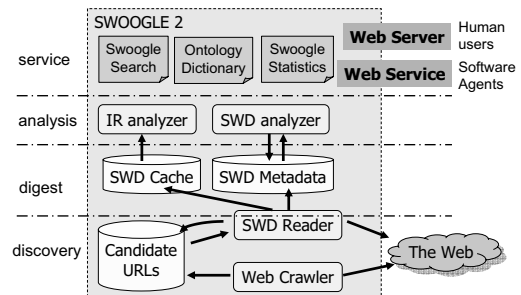


Figure 1: Swoogle has four components that discover, digest, analyze and serve semantic web data.

Discovering Semantic Web Documents

The size of the Semantic Web is measured by the number of discovered SWDs. (Eberhart 2002) reported finding 1,479 SWDs with about 255K triples out of nearly 3M web pages. As of May 2005, Swoogle has found over 368K SWDs with more than 70M triples. Although this number is far less than Google’s eight billion web pages, it represents a non-trivial collection of semantic web data (Guo, Pan, & Heflin 2004).

The Semantic Web’s content can be divided into two categories – program generated instance data and (mostly) hand crafted ontologies. The first category is the larger and includes FOAF personal profiles, RSS news feeds, RDF metadata embedded in PDF files, Dublin Core digital library metadata, Creative Commons’ copyright statements, and assertions extracted from structured data sources such as WordNet and the CIA fact book. While some ontologies have been derived from structured sources, most appear to be designed by semantic web researchers. Although these ontology documents are far outnumbered by instance data documents, they are critically important since they convey symbol semantics.

*Research support was provided by DARPA contract F30602-00-0591 and NSF awards NSF-ITR-IIS-0326460 and NSF-ITR-IDM-0219649.

Copyright © 2005, American Association for Artificial Intelligence (www.aaai.org). All rights reserved.

