

SEAN: A System for Semantic Annotation of Web Documents

Amarjeet Singh* Saikat Mukherjee* I. V. Ramakrishnan* Guizhen Yang† Zarana Shah*

*Department of Computer Science
Stony Brook University, Stony Brook, NY 11794

†Department of Computer Science and Engineering
University at Buffalo, Buffalo, NY 14260

Semantic Web documents use metadata to express the meaning of the content encapsulated within them. Although RDF/XML has been widely recognized as the standard vehicle for describing metadata, an enormous amount of semantic data is still being encoded in HTML documents that are designed primarily for human consumption. Tools such as those pioneered by SHOE [Heflin *et al.*, 2003] and Ontobroker [Fensel *et al.*, 1998] facilitate manual annotation of HTML documents with semantic markups.

the Web page in Figure 1 has a news taxonomy (on the left), which does not change, and a template for major headline news items. Each of these items begins with a hyperlink labeled with the news headline (*e.g.*, “White House ...”), followed by the news source (*e.g.*, “By REUTERS ...”), followed by a timestamp and a text summary of the article (*e.g.*, “The White House today ...”) and a (variable) number of pointers to related news. These concepts and concept instances can be organized into a semantic partition tree (such as the one shown in Figure 2, which represents the “semantics” of the HTML document.



Figure 1: New York Times front page

In this demo we will present SEAN, a system for *automatically annotating* HTML documents. It is based on the idea that well-organized HTML documents, especially those that are machine generated from templates, contain rich data denoting semantic concepts (*e.g.*, “News Taxonomy” and “Major Headline News”) and concept instances. These kinds of documents are increasingly common nowadays since most Web sites (*e.g.*, news portals, product portals, etc.) are typically maintained using content management software that creates HTML documents by populating templates from backend databases. For instance, observe that

In a semantic partition tree each partition (subtree) consists of items related to a semantic concept. For example, in Figure 2 all the major headline news items are grouped under the subtree labeled “Major Headline News”. There are two main tasks underlying the creation of a semantic partition tree from an HTML document: (i) identify segments of the document that correspond to semantic concepts; and (ii) assign labels to these segments. Informally, we say that several items are semantically related if they all belong to the same concept.

SEAN automatically transforms well-structured HTML documents into their semantic partition trees by exploiting two key observations. First, *semantically related items exhibit consistency in presentation style*. For instance, observe in Figure 1 the presentation styles of those items in the news taxonomy on the left. The main taxonomic items, “NEWS”, “OPINION”, “FEATURES”, etc., are all presented in bold font. All the subtaxonomic items (*e.g.*, “International”, “National”, “Washington”, etc.) under the main taxonomic item (*e.g.*, “NEWS”) are hyperlinks. A similar observation can also be made on all the major headline news items. Second, *semantically related items exhibit spatial locality*. For example, when rendered in a browser, all the taxonomic items are placed in close vicinity occupying the left portion of the page. Specifically, in the DOM tree corresponding to the HTML document in Figure 1 all the items in the news taxonomy will be grouped together under one single subtree.

The first observation leads to the idea of associating a type with every leaf node in the DOM tree. The type of a leaf node consists of the root-to-leaf path of this node in the DOM tree and captures the notion of consistency in presentation style. The second observation motivates the idea of propagating types bottom-up in the DOM tree and discovering structural recurrence patterns for semantically related items at the root

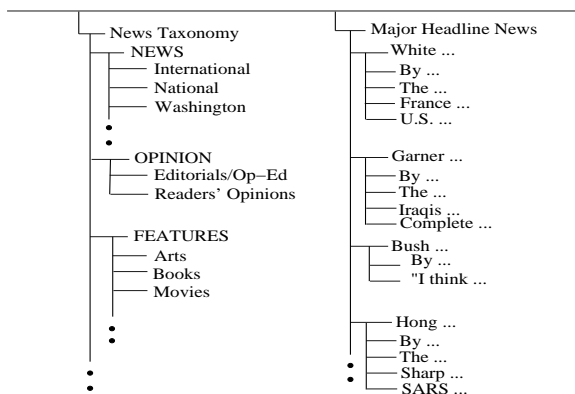


Figure 2: Partition tree of the New York Times front page

of each subtree. Based on the idea of types and type propagation, SEAN does structural analysis of an HTML document for automatically partitioning it into semantic structures. In this process it also discovers semantic labels and associates them with partitions when they are present in the document (e.g., “NATIONAL”, “INTERNATIONAL”, etc. appearing in the third column in Figure 1).

SEAN augments structural analysis with semantic analysis to factor in structural variations in concept instances (e.g., absence of pointers to related news in the third major headline news item in Figure 1 in contrast to others). Semantic analysis makes lexical associations via WordNet to more accurately put the pieces of a concept instance together. To assign informative labels that are not present in an HTML document (e.g., “Major Headline News” in Figure 1) to partitions, semantic analysis makes concept associations by classifying the content of a partition using an ontology encoding domain knowledge.

By combining structural and semantic analysis SEAN automatically discovers and labels concept instances in template-based, content-rich HTML documents w.r.t. a domain ontology. Details appear in [Mukherjee *et al.*, 2003]. Our demo will illustrate how SEAN is used to assign semantic labels to HTML documents. For semantic analysis SEAN currently provides a simple editor for creating/editing ontologies for domains of interest. The generated semantic partitions are assigned concept labels by either matching keywords in the partition’s content to those associated with concepts in the ontology or by applying concept classification rules to features extracted from the content. The keywords as well as the rules used for classification can both be edited. We point out that there has been extensive work on ontology tools and classifiers. AS future work we plan on designing a plug-in architecture for SEAN that will support the use of any sophisticated ontology editing tools such as Protege [Protege, 2000], SHOE [Heflin *et al.*, 2003], Ontobroker [Fensel *et al.*, 1998], etc. and powerful statistical and rule-based classifiers such as Naive Bayes and decision trees [Mitchell, 1997] for semantic analysis.

Recently, a number of works proposed to partition and annotate HTML documents based on structural analysis and tools based on the idea of combining structural analysis

with domain ontologies for semantic annotation were described in [Dill *et al.*, 2003; Handschuh and Staab, 2002; Handschuh *et al.*, 2003; Heflin *et al.*, 2003]. In [Handschuh and Staab, 2002; Handschuh *et al.*, 2003; Heflin *et al.*, 2003] powerful ontology management systems form the backbone that supports interactive annotation of HTML documents. The observation that semantically related items exhibit spatial locality in the DOM tree of an HTML document is not exploited in [Dill *et al.*, 2003]. As a result, their partitioning algorithm may fail to identify proper concept instances in HTML documents that are generated from templates.

References

- [Dill *et al.*, 2003] Stephen Dill, Nadav Eiron, Daniel Gibson, Daniel Gruhl, R. Guha, Anant Jhingran, Tapas Kanungo, Sridhar Rajagopalan, Andrew Tomkins, John Tomlin, and Jason Yien. SemTag and Seeker: Bootstrapping the Semantic Web via automated semantic annotation. In *International World Wide Web Conference*, 2003.
- [Fensel *et al.*, 1998] Dieter Fensel, Stefan Decker, Michael Erdmann, and Rudi Studer. Ontobroker: Or how to enable intelligent access to the WWW. In *11th Banff Knowledge Acquisition for Knowledge-Based Systems Workshop*, Banff, Canada, 1998.
- [Handschuh and Staab, 2002] Siegfried Handschuh and Steffen Staab. Authoring and annotation of web pages in CREAM. In *International World Wide Web Conference*, 2002.
- [Handschuh *et al.*, 2003] Siegfried Handschuh, Steffen Staab, and Raphael Volz. On deep annotation. In *International World Wide Web Conference*, 2003.
- [Heflin *et al.*, 2003] Jeff Heflin, James A. Hendler, and Sean Luke. SHOE: A blueprint for the Semantic Web. In Dieter Fensel, James A. Hendler, Henry Lieberman, and Wolfgang Wahlster, editors, *Spinning the Semantic Web*, pages 29–63. MIT Press, 2003.
- [Mitchell, 1997] Tom M. Mitchell. *Machine Learning*. McGraw Hill, 1997.
- [Mukherjee *et al.*, 2003] Saikat Mukherjee, Guizhen Yang, and I. V. Ramakrishnan. Automatic annotation of content-rich HTML documents: Structural and semantic analysis. In *International Semantic Web Conference (ISWC)*, 2003.
- [Protege, 2000] Protege, 2000. <http://protege.stanford.edu>.