

A Personalizable Agent for Semantic Taxonomy-Based Web Search

Larry Kerschberg, Wooju Kim, and Anthony Scime

E-Center for E-Business, George Mason University
4400 University Drive, Fairfax, VA 22030, USA

kersch@gmu.edu

<http://eceb.gmu.edu/>

Department of Industrial Engineering, Chonbuk National University, Korea

wjkim@chonbuk.ac.kr

Department of Computer Science, SUNY-Brockport

ascime@brockport.edu

Abstract. This paper addresses the problem of specifying Web searches and retrieving, filtering, and rating Web pages so as to improve the relevance and quality of hits, based on the user's search intent and preferences. We present a methodology and architecture for an agent-based system, called WebSifter II, that captures the semantics of a user's decision-oriented search intent, transforms the semantic query into target queries for existing search engines, and then ranks the resulting page hits according to a user-specified weighted-rating scheme. Users create personalized search taxonomies via our Weighted Semantic-Taxonomy Tree. Consulting a Web taxonomy agent such as WordNet helps refine the terms in the tree. The concepts represented in the tree are then transformed into a collection of queries processed by existing search engines. Each returned page is rated according to user-specified preferences such as semantic relevance, syntactic relevance, categorical match, page popularity and authority/hub rating.

1 Introduction

With the advent of Internet and WWW, the amount of information available on the Web grows daily. However, having too much information at one's fingertips does not always mean good quality information, in fact, it may often prevent a decision maker from making sound decisions, by degrading the quality of the decision. Helping decision makers to locate relevant information in an efficient manner is very important both to the person and to an organization in terms of time, cost, data quality and risk management.

Although search engines assist users in finding information, many of the results are irrelevant to the decision problem. This is due in part, to the keyword search approach, which does not capture the user's intent, what we call meta-knowledge. Another reason for irrelevant results from search engines is a "semantic gap" between the meanings of terms used by the user and those recognized by the search engines. In

addition, each search engine has its own uncustomizable ranking system, where users cannot “tell” the search engine what preferences to use for search criteria. For example, a shopping agent may go for the lowest price, while the user might want the “most flexible return policy.”

To overcome these three problems, we propose a semantic taxonomy-based personalizable meta-search agent approach. We build upon the ideas presented by Scime and Kerschberg [1, 2]. We develop a tree-structured representation scheme with which users specify their search intent. We call this representation scheme the “Weighted Semantic Taxonomy Tree (WSTT)”, in which each node denotes a concept that pertains to the user’s problem-domain. To address the second weakness, we present an elaborate user preference representation scheme based on various components, each of which represents a specific decision-criterion. Users can easily and precisely express their preference for a search using this representation scheme.

In order to rate the relevance of a page hit, we use a rating mechanism combining the WSTT and the component-based preference representation. Since Web page rating can itself be viewed as a decision-making problem, where a decision maker (a user) must evaluate various alternatives (Web pages) for his/her problem (user’s Web search intention), we use decision-analytic methods in the design of our rating mechanism.

Finally, we have designed and implemented a meta-search agent called WebSifter II that cooperates with WordNet for concept retrieval, and consults well-known search engines. For the empirical validation of our approach, we also present some real world examples of our system.

The remainder of the paper is organized as follows. Section 2 presents related research. Section 3 presents the major aspects of our semantic-based personalizable approach to represent user intention, and the multi-component rating of search hits. In Section 4, we discuss the system architecture of WebSifter II, the search agent that implements our methodology. We also deal with some collaboration issues too in this section. The results of empirical studies are presented in Section 5.

2 Related Work

Most of current Internet search engines such as Yahoo, Excite, Altavista, WebCrawler, Lycos, Google, etc. suffer from *Recall* and *Precision* problems [3]. The relatively low coverage of individual search engines leads to using meta-search engines to improve the recall of a query. Examples are MetaCrawler [4], SavvySearch [5], NECI Metasearch Engine [6], and Copernic (<http://www.copernic.com>). This meta-search engine approach partly addresses the recall problem but still suffers from the precision problem.

We can categorize research regarding the precision problem into three major themes: content-based, collaborative, and domain-knowledge approaches.

The content-based approach first represents a user’s explicit preferences and then evaluates Web page relevance in terms of its content and user preferences. Syskill & Webert [7], WebWatcher [8], WAWA [9], and WebSail [10] fall into this category.

Further, some research takes into account not only Web page content but also its structure (e.g. hyperlinks) to evaluate relevance [11, 12].

The collaborative approach determines information relevancy based on similarity among users rather than similarity of the information itself. Example systems are Firefly and Ringo [13], Phoaks [14], and Siteseer [15]. In addition, some hybrid approaches incorporate both approaches for example Fab [16], Lifestyle Finder [17], WebCobra [18].

The third category is the domain knowledge approach that uses user and organizational domain knowledge to improve the relevancy of search results. Yahoo! uses domain knowledge and provides a pre-defined taxonomy path. So, classifying Web pages automatically into a pre-defined, or a dynamically created taxonomy [19] is a related issue to this approach. NorthernLight (www.northernlight.com) is a search engine that supports this kind of dynamic taxonomy service. Using NorthernLight's *Custom Search Folder* service, users can refine their search query to a specific domain, when the search engine presents too much information.

Some research incorporates user domain knowledge in a more explicit way. For example, Aridor et al. [20] represent user domain knowledge as a small set of example Web pages provided by users. Chakrabarti et al. adopted both a pre-defined (but modifiable) taxonomy and a set of example user-provided Web pages as domain knowledge [21].

From this survey of related research, we have identified several aspects that merit further consideration. First, most approaches force users to use a search engine in a passive rather than active manner. Often, the user cannot understand why extraneous and irrelevant results are retrieved. There is a pressing need for users to be able to express their query intent in a more natural and structured manner. Second, current approaches lack sufficient expressive power to capture a users' search intent and preferences, because most of the representation schemes are based on a vector space model [22] or its variants. Third, most approaches do not take full advantage of domain-specific knowledge with which to scope the search, filter the hits, and classify the query result.

Regarding the first limitation, there is another related research category, the ontology-based approach by which users can express their search intent in a more semantic fashion. Domain-specific ontologies are being developed for commercial and public purposes [23] and OntoSeek [24], On2Broker [25], GETESS [26], and WebKB [27] are example systems.

Although the ontology-based approach is a promising way to solve some aspects of the precision problem, it still requires two major pre-requisites. First, the entire collection of Web pages must be transformed into ontological form. Second, there is as yet no common agreement on the representation of the ontology, nor the query or reasoning mechanisms. Even if these two prerequisites were satisfied, the precision problem in Web search would remain due to the huge amount of the information on the web. A user-centric information relevancy evaluation scheme will complement the above approaches.

3 Semantic Taxonomy-Tree-Based Approach for Personalized Information Retrieval

3.1 Weighted Semantic Taxonomy Tree

Usually a keyword-based search representation is insufficient to express a user's search intent. By postulating a user's decision-making process as depicted in Figure 1, we can support readily query formulation and search.

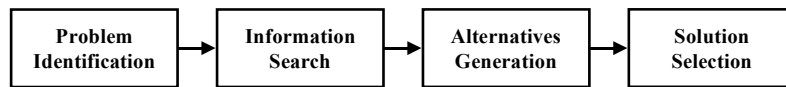


Fig. 1. Four Phases of Decision Making Process

This process starts with a problem identification phase and then a user seeks relevant information to solve the identified problem. Based on the collected information, listing alternatives, evaluating them, and selecting a solution are the subsequent steps. One implication of the decision-making process is that the more we understand a user's problems, the better we can support a user's information search. In our approach, we represent a user's search intent by a hierarchical concept tree with weights associated with each concept, thereby reflecting user-perceived relevance of concepts to the search.

Let's assume that a person has started a new business and is looking for office equipment. He wants to search for information about office equipment on the Web. Suppose he wants information about chairs, so he might build a query using a single term, "chair". If he is a more skilled user of Internet search engines, he might build a query using two terms, "office" and "chair" to obtain more precise results. He may also use the 'AND' or 'OR' operator between them. In this case, the term "office" provides added context for the search. However, this formulation is still very implicit and passive. As we mentioned earlier, one way to express this kind of context information is by using a taxonomy tree as shown in Figure 2. Figure 2(a) shows a simple taxonomy tree that represents a search intention to find a chair in the context of office, while a search for finding an office in the context of chair is expressed by Figure 2(b). The taxonomy tree provides more expressive semantics than simple keyword-based representations used by most current search engines.

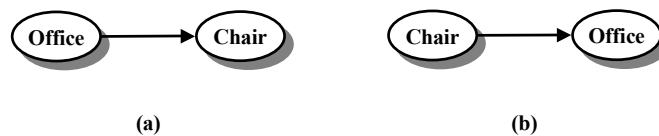


Fig. 2. A simple example of a taxonomy tree

The taxonomy tree approach is already used in many search engines such as Yahoo! We have devised a tree-based search representation model that allows users to present

their search intention by defining their own taxonomy topology. We call this the *Weighted Semantic Taxonomy Tree* (WSTT) model. Now, let us formally define this model. The WSTT consists of a set of nodes that is denoted as N in the sequel. Because it is a tree, all nodes, except the root node, must have one parent node. Every node should have one representative term and a weight that represents the importance of this node for a search. For a node $n \in N$, we denote a representative term, or label, and its weight as $rt(n)$ and $w(n)$, respectively. We restrict the feasible range of the value of $w(n)$ from 0 to 10. Figure 3 shows a realistic example of the businessman’s search intention using our WSTT scheme. Users can build their own hierarchical taxonomy tree, and assign importance levels to each term within the context of their antecedent terms. For example, we can translate the upper sub-tree as that a businessman wants to find information about chairs, desks, and phones within the context of office furniture and office equipment where the numbers that appear to the left to each term, 10, 9, and 6 denote the respective importance levels of chairs, desks, and phones.

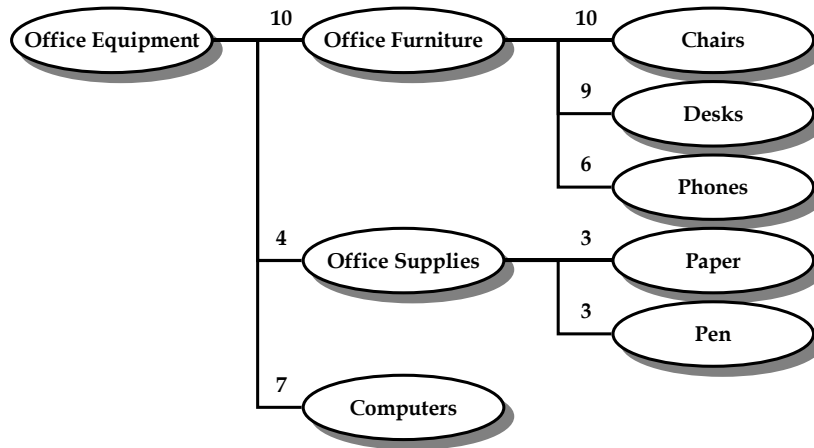


Fig. 3. An example of a WSTT representing a businessman’s search intention

One drawback is that the terms may have multiple meanings, and this is one of the major reasons that search engines return irrelevant search results. To address this limitation, we introduce the notion of “word senses” from WordNet [28] into our WSTT scheme to allow users to refine their search intention. WordNet is a linguistic database that uses sets of terms that have similar semantics (*synsets*) to represent word senses. Each synset corresponds to terms with synonymous meaning in English and so each word may be associated with multiple synsets. In this paper, we rename this synset as *Concept* for our own use and the user can choose one of the concepts available from WordNet for the term of a specific node in WSTT. We denote an available concept, that is, a set of terms for a node n as $c(n)$. For example, the “chair” term has the following four possible concepts from WordNet.

1. {chair, seat} // a seat for one person, with a support for the back,
2. {professorship, chair} // the position of professor, or a chaired professorship,

3. {president, chairman, chairwoman, chair, chairperson} // the officer who presides at the meetings of an organization, and
4. {electric chair, chair, death chair, hot seat} // an instrument of death by electrocution that resembles a chair.

If the user wants to search for a chair to sit on, he would choose the first concept. If the user selects the first concept, then without loss of generality, we can assume that the remaining concepts are not of interest, thereby obtaining both positive and negative indicators of his intent. Now, let's distinguish the set of terms of selected concept from the set of terms of the unselected concepts as *Positive Concept Terms* and *Negative Concept Terms*, and denote them as $pct(n)$ and $nct(n)$ for a node n , respectively. If we denote a term as t and assume that a user selects the k -th concept, then we can formalize the definitions of them for a given node n as follows:

$$pct(n) = \{t \mid t \in c_k(n)\} \quad (1)$$

$$nct(n) = \left\{ t \mid t \in \bigcup_{i \neq k} c_i(n) \right\} - \{rt(n)\} \quad (2)$$

where $c_i(n)$ denotes the i -th concept available from WordNet for a node n and $rt(n)$ denotes the representative term of n .

If a user selects the second concept from our example, according to the definitions from (1) and (2), $pct(n)$ and $nct(n)$ are as follows: $pct(n) = \{\text{professorship, chair}\}$ and $nct(n) = \{\text{seat, president, chairman, chairwoman, chairperson, electric chair, death chair, hot seat}\}$.

Figure 4 shows an internal representation of the user's intention via the WSTT schema, after the concept selection process has finished; the user however sees the tree of Figure 3. Another advantage using the tree structure is that it is possible to represent many concepts at the same time. This allows the user to specify a broad range of interests simultaneously.

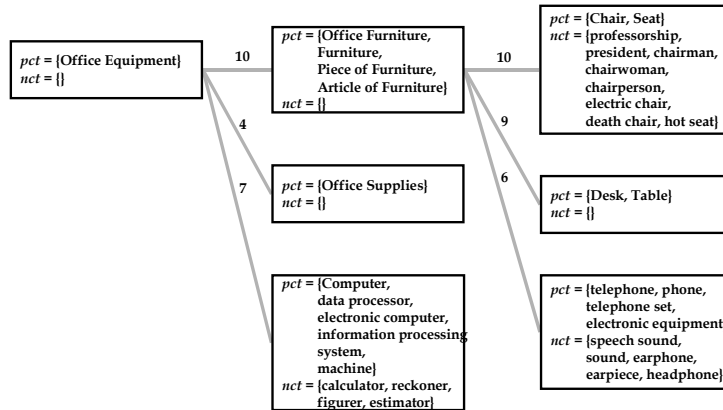


Fig. 4. An example of the internal representation of the user's search intention

3.2 Multi-Attribute-Based Search Preference Representation

The ranking of Web search hits by users involves the evaluation of multiple attributes, which reflect user preferences and their conception of the decision problem. In our approach, we pose the ranking problem as a multi-attribute decision problem. Thus, we examine the search results provided by multiple search engines, and rank the pages, according to multiple decision criteria. Both Multi-Attribute Utility Technology (MAUT) [29] and Repertory Grid [30] are two major approaches that address our information evaluation problem. Our ranking approach combines MAUT and the Repertory Grid. We define six search evaluation components as follows:

1. *Semantic* component: represents a Web page's relevance with respect to its content.
2. *Syntactic* component: represents the syntactic relevance with respect to its URL. This considers URL structure, the location of the document, the type of information provider, and the page type (e.g., home, directory, and content).
3. *Categorical Match* component: represents the similarity measure between the structure of the user-created taxonomy and the category information provided by search engines for the retrieved Web pages.
4. *Search Engine* component: represents the user's biases toward and confidence in search engine's results.
5. *Authority/Hub* component: represents the level of user preference for *Authority* or *Hub* sites and pages. Authority sites usually have larger in-degree from Hub sites and Hub sites usually have larger out-degree to Authority sites [31].
6. *Popularity* component: represents the user's preference for popular sites. The number of visitors or the number of requests for the specific page or site can measure popularity.

Further, in this multi-component-based preference representation scheme, the user can assign a preference level to each of these components, and also to each available search engine within the search engine component. Then, these components and the assigned preference level are eventually synthesized into a single unified value resulting in the relevance measure for a specific Web page. Figure 5 conceptually depicts our scheme. In this figure, each number assigned to an edge denotes the user's preference level for that component. This multi-component preference scheme allows users more control over their searches and the determination of a page's relevance.

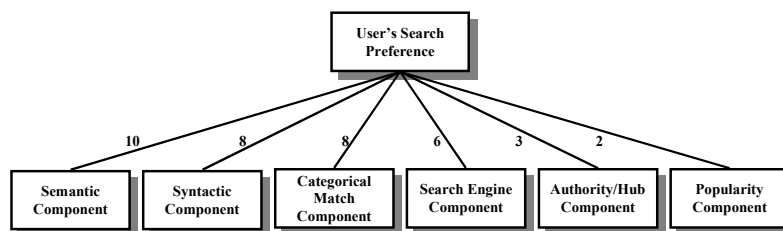


Fig. 5. A conceptual model of user's preference representation scheme

Thus far, we have discussed how to capture and represent semantically the user's search intention and search preferences. Now, we turn our attention to deriving a good

estimate of the relevancy of a Web page based on these semantics. In the following sections, we will discuss how to obtain Web information using existing search engines and then address the derivation of relevance estimates.

3.3 Gathering Web Information based on Search Intention

Since we adopt a meta-search approach to Web information gathering to preserve the benefits of meta-search engines discussed in [4, 5, 20], we neither create nor maintain our own index database of Web information. At present, there is no search engine that accepts a search request based on the WSTT. We have developed a translation mechanism from our WSTT-based query, to Boolean queries that most of current search engines can process.

As already mentioned, we represent a user's search intention as a tree, as shown in Figure 4. The leaf nodes denote the terms of interest to the user, and the antecedent nodes for each node form a search context. We transform the entire tree into a set of separate queries where each is acceptable to existing search engines. To do this, first we decompose the tree into a set of paths from the root to each leaf node. Then for each path, we generate all possible combinations of terms, by selecting one term from the positive concept terms of each node in the path from a root node to a leaf node. Finally, we pose each query to search engines to obtain query results.

We now provide definitions to formalize the above discussion. Let's first define a *Leaf Path* as an ordered set of nodes, $\{n_0, n_1, n_2, \dots, n_{l-1}, n_l\}$, where n_0 is a root node, n_l is a leaf node, and n_1, n_2, \dots, n_{l-1} are consecutive intermediate nodes on the path from n_0 to n_l in the WSTT. We denote a leaf path as lp . We also define a set of all distinct leaf paths available from the WSTT as $lpset$. For example, we have six leaf paths from the example WSTT as in the Figure 3 and its $lpset$ becomes $\{\{\text{Office Equipment, Office Furniture, Chairs}\}, \{\text{Office Equipment, Office Furniture, Desks}\}, \{\text{Office Equipment, Office Furniture, Phones}\}, \{\text{Office Equipment, Office Supplies, Paper}\}, \{\text{Office Equipment, Office Supplies, Pen}\}, \{\text{Office Equipment, Computers}\}\}$. Now, let's define a *Term Combination Set* for a *Leaf Path* lp , as a set of all possible combinations of terms by selecting one term from each $pct(n)$, where $n \in lp$ and denote it as $tcslp(lp)$. We also denote a set of all term combinations available from a given WSTT and each of its elements as tcs and tc , respectively. Then, using the above definitions, a $tcslp(lp)$ and tcs can be formally represented respectively as follows:

$$tcslp(lp) = pct(n_0) \times pct(n_1) \times pct(n_2) \times \dots \times pct(n_l) \quad (3)$$

where symbol \times denotes the Cartesian product of sets.

$$tcs = \bigcup_{lp \in lpset} tcslp(lp) \quad (4)$$

If lp is the first element, that is, $\{\text{Office Equipment, Office Furniture, Chairs}\}$ of the $lpset$ in the case of Figure 3 and Figure 4, then according to equation (3), $tcslp(lp) = \{\{\text{Office Equipment, Office Furniture, Chair}\}, \{\text{Office Equipment, Office Furniture,$

Seat}, {Office Equipment, Furniture, Chair}, {Office Equipment, Furniture, Seat}, {Office Equipment, Piece of Furniture, Chair}, {Office Equipment, Piece of Furniture, Seat}, {Office Equipment, Article of Furniture, Chair}, {Office Equipment, Article of Furniture, Seat}}.

Once we get tcs , then we make each term combination, $tc \in tcs$ as a separate request and pose them to each search engine for Web information gathering. Now, the problem is how to generate actual query statements to each query engine based on each tc . We have trade-offs between *Precision* and *Coverage* depending on which logical operators we impose between terms. Actually, each tc is a set of terms and so, it can be represented as $\{t_1, t_2, \dots, t_n\}$ where $t_1, t_2, \dots, t_n \in tc$. To generate an actual query statement from a tc , we can have two different alternative choices, " $t_1 \wedge t_2 \wedge \dots \wedge t_n$ " and " $t_1 \vee t_2 \vee \dots \vee t_n$ " where \wedge denotes AND and \vee denotes OR. The first one provides more precise search results, while the second allows greater coverage.

Based on the fact that a general user tends to use the AND operator between terms when considering additional terms for the context of a search, we adopt the AND operator in generating actual query statements. We leave the more general scheme for future research. For the illustration of our query generation method, let's use the case depicted in Figure 4 again. According to the procedures mentioned thus far, the upper-most leaf path of the WSTT in Figure 4 is translated into eight separate query statements as follow. (1) "Office Equipment" AND "Office Furniture" AND "Chair", (2) "Office Equipment" AND "Office Furniture" AND "Seat", (3) "Office Equipment" AND "Furniture" AND "Chair", (4) "Office Equipment" AND "Furniture" AND "Seat", (5) "Office Equipment" AND "Piece of Furniture" AND "Chair", (6) "Office Equipment" AND "Piece of Furniture" AND "Seat", (7) "Office Equipment" AND "Article of Furniture" AND "Chair", and (8) "Office Equipment" AND "Article of Furniture" AND "Seat".

These queries can now be submitted to each target search engine, and the query results are stored for further processing, as discussed in the next section.

3.4 Unified Web Information Rating Mechanism

In this section, we discuss a rating mechanism to evaluate each resulting page hit from the target search engines for the generated query statements. Through this mechanism, each Web page will have its own value representing the relevance level from the user's viewpoint. To accomplish this goal, six relevance values of a Web page are computed, corresponding to each of the six components. Then a composite value of these six relevance values is computed based on a function of the multi-attribute-based search preference representation scheme. In the following sub-sections, we will first discuss how this composite relevance value is computed, and then a set of methods to compute each of component's relevance values.

3.4.1 Composite Relevance Value Computation

Let's first assume we have evaluated the six components' relevance values for a Web page retrieved from search engines. Then we need to synthesize these six values into one single composite relevance value to compare Web pages to each other and to list

them to the user in an order of relevance. This problem can be viewed as a multi-attribute decision-making problem.

One of the popularly accepted approaches in decision science community is AHP (Analytic Hierarchy Process) [32]. It converts user's subjective assessments of relative importance between preference components into a linear set of weights, which is further used to rank alternatives. Although we adopt AHP approach as a basis of our synthesizing mechanism, we have modified the original AHP to fit to our weight acquisition scheme, because it requires pair-wise comparisons between all components to obtain importance ratios between each pair of them. Actually in our approach, a user assigns an absolute importance weight on each component rather than relative ratios between components. However, since we still need those relative ratios, we first approximate them by dividing absolute importance weights of components by each other. Then, we follow the same remaining steps of AHP to compute the composite relevance value for each Web page.

We now provide notations to formalize the above discussion as follows.

compset: denotes a set of preference components to be considered in our scheme.

$cw^U(x)$: denotes a weight provided by the user to represent the importance of a component x .

$rv(x, pg)$: denotes a relevance value of a Web page pg with respect to a component x .

$lr(x, y)$: denotes a relative importance ratio of component x compared to component y .

$ns(z)$: denotes a function that returns the number of elements in a set z .

We first approximate $lr(x, y)$ by (5) based on the user-provided importance weights for each pair of components:

$$lr(x, y) = cw^U(x) / cw^U(y) \quad (5)$$

where $x \in compset$ and $y \in compset$.

Then, the AHP computes normalized importance weights for each component based on these relative ratios. We denote the normalized importance weight for a component com and the composite relevance value of a Web page pg as $cw^N(com)$ and $rv(pg)$, respectively. According to AHP, these two values can be calculated respectively as follows:

$$cw^N(com) = \left[\sum_x \left(\frac{lr(com, x)}{\sum_y lr(x, y)} \right) \right] / ns(compset) \quad (6)$$

where $x \in compset$ and $y \in compset$.

$$rv(pg) = \sum_x cw^N(com) \cdot rv(com, pg) \quad (7)$$

where $com \in compset$.

Finally, Web pages are presented to users in descending order of rv . This, together with the page relevancy value indicates the relative importance of that page to the user.

Thus far, we have discussed how to synthesize the relevance values of a user's preference components into a single composite value, under the assumption that these relevance values of the components have already been computed. Now, we show how to compute relevance values of each of the six preference components based on the user's preference, as well as the user's search intent as represented by the WSTT.

3.4.2 Semantic Component Relevancy Computation

The semantic component represents relevancy of a Web page to a user's search intent represented by the WSTT with respect to its content. To compute this relevance, we conceptually follow the reverse steps that we performed in the section 3.3 to generate separate queries from the WSTT.

First, we evaluate the semantic relevancies of a retrieved Web page for each of the term combinations. We then combine the semantic measures for each leaf path and bind each of these semantic measures to the corresponding leaf node. Finally we compute a semantic component relevancy of the Web page using an AHP-based WSTT relevance value composition mechanism. This mechanism propagates the bound values on the leaf nodes toward the root node, thereby providing a single combined relevance value at the root node.

Now, let's explain the details of this procedure in a formal manner. We first define $rvtc^{SM}(tc, pg)$ as a semantic relevance value of a Web page pg to a term combination tc and it is computed by a simple counting method as follows:

$$rvtc^{SM}(tc, pg) = \frac{\sum_{t \in tc} appear(t, pg)}{ns(tc)} \quad (8)$$

where t is a term and the function $appear(t, pg)$ returns 1 if t appears in pg and 0, otherwise.

Based on these $rvtc^{SM}$ values, we define $rvlp^{SM}(lp, pg)$ as a semantic relevance value of a Web page pg to a leaf path lp . When we compute this $rvlp^{SM}$, we have to consider two aspects. First, we have to synthesize multiple $rvtc^{SM}$ values obtained from equation (8) for a leaf path into a single measure and we adopt a max function for this. Second, we have to consider negative concepts related to a leaf path. To incorporate these negative concepts into computing $rvlp^{SM}$, we first develop a measure to evaluate irrelevancy of a Web page pg in terms of negative concept terms related to a leaf path lp and we denote it as $irv(lp, pg)$. The following equation (9) shows its mathematical definition.

$$irv(lp, pg) = \sum_t appear(t, pg) \quad (9)$$

where t is a term and also $t \in \bigcup_{n \in lp} nct(n)$.

Now, we can compute $rvlp^{SM}$ using the following equation (10).

$$rvlp^{SM}(lp, pg) = \left(\frac{\sum rvtc^{SM}(tc, pg)}{ns(tcslp(lp))} \right) \cdot (1 - \theta)^{irv(lp, pg)} \quad (10)$$

where $tc \in tcslp(lp)$ and θ is a given $[0, 1]$ scale degradation rate.

In equation (10), θ denotes the level of degradation with respect to the irrelevance caused by negative concepts. So if θ is close to 1, then a little irrelevancy results in a big impact on $rvlp^{SM}$. On the other hand, if it is close to 0, the irrelevancy does not have any impact on the rvp value. The user can control this rate and we set it to a default of 0.1.

Now, we synthesize a single semantic relevancy value of a Web page according to the WSTT. Since AHP was originally developed to derive a unified measure to evaluate decision alternatives based on a tree like the WSTT, we apply this approach to our WSTT scheme by combining our $rvlp^{SM}$ values for each leaf path into a single semantic relevance value of a Web page. However, we need to normalize the user-provided weights for the nodes of WSTT, for reasons similar to those discussed in the previous section. For this normalization, we apply equation (5) to each hierarchical branch of the WSTT, and we obtain a set of normalized weights for each node within the scope of the branch to which the nodes belong. We denote this normalized weight for a node n , $w^N(n)$. With the normalized weights, let's formalize the AHP-based WSTT relevance value composition mechanism.

Equation (11) shows a relevance value determination rule on each node of WSTT for a Web page pg and we denote a relevance value of a Web page pg on a node n as $rvn(n, pg)$.

$$rvn(n, pg) = \begin{cases} bndfn(lp, pg) & \text{if } n \text{ is a leaf node} \\ \sum_{x \in children(n)} w^N(x) \cdot rvn(x, pg) & \text{of a leaf path } lp. \\ & \text{otherwise.} \end{cases} \quad (11)$$

where $children(n)$ is a set of nodes that is a child of n and $bndfn(lp, pg)$ is an arbitrary value binding function to leaf nodes.

To perform this mechanism, we first bind relevance values from $bndfn()$ to all corresponding leaf nodes and then these values are propagated from leaf nodes to the root node, finally obtaining a single composite relevance value of a Web page for the WSTT. In this semantic component case, by setting $bndfn(lp, pg)$ as $rvlp^{SM}(lp, pg)$ in the equation (10), we can obtain a single composite semantic relevance value of a Web page pg as $rvn(n_0, pg)$, where n_0 is the root node of the WSTT. This obtained value is then assigned to $rvc(\text{Semantic Component}, pg)$ for further computing of composite relevance value with other preference components, discussed in the previous section.

Figure 6 shows conceptually the entire flow of computation from relevancy computing for a term combination to relevancy computing across the WSTT, which is required to compute a semantic relevance value of a Web page. In this figure, $PageSet(tc_i, s_j)$ denotes a set of resulting pages from a search engine s_j for a term

combination tc_i . Actually, we will use a similar method when computing categorical match and search engine components' relevancies in the following sections.

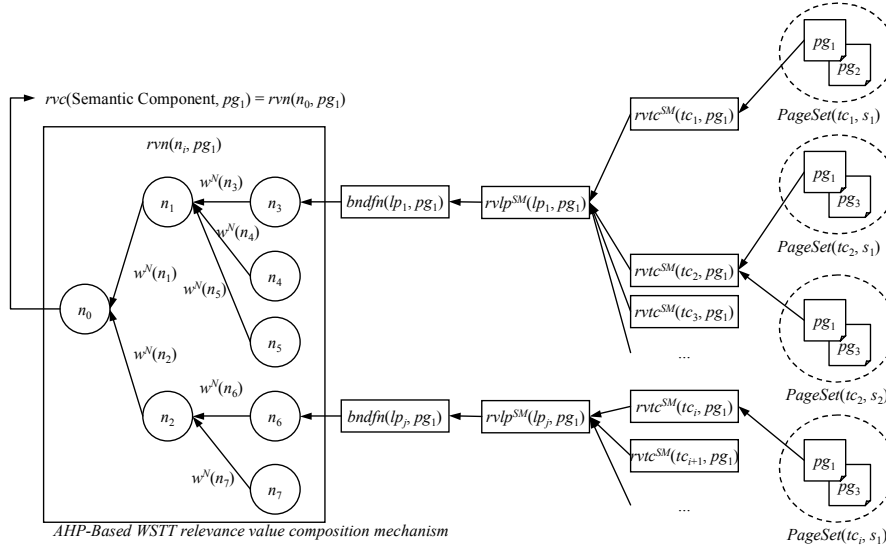


Fig. 6. Conceptual flow of computation of semantic component relevancy

3.4.3 Syntactic Component Relevancy Computation

The syntactic component of Web document measures the structural aspects of the page as a function of the role of that page within the structure of a Web site. Our approach takes into account the location of the document, its role (e.g., home, directory, and content), and the well formedness of its URL.

We define three types of Web pages:

- *Direct-Hit* – the page may be a home page or a page with significant content within its domain.
- *Directory-Hit* – this page has links to other pages in the domain of the Web site.
- *Page-Hit* – Web pages that are subordinate to direct-hit and directory-hit pages fall into this category. These pages contain partial information about the Web site domain.

Scime and Kerschberg [1, 2] define a set of heuristics to classify a Web page returned from a search engine as either a direct, directory, or page hit. Further, a page may have more than one classification. In order to manipulate syntactic relevancy, we assign a numeric value to each type as a real number in the interval [0, 1]. Default values for direct, directory, and page hits are 1.0, 0.6, and 0.4 respectively. The assumption is that users would prefer to view direct hits over the other two.

Since a Web page might be classified into more than one class, we need to synthesize those multiple matches into one measure. To do this, we introduce an averaging mechanism and define some necessary notations and a formula to compute

the syntactic relevance value of a Web page pg , $rvc(\text{Syntactic Component}, pg)$ as follows:

$rset(cl)$: denotes a set of rules to classify a Web page into the class cl .

$rsc(r)$: denotes a score of a rule r and it returns 1.0 if $r \in rset(\text{Direct Hit})$, 0.6 if $r \in rset(\text{Directory Hit})$, and 0.4 if $r \in rset(\text{Page Hit})$.

$mat(r, pg)$: denotes a function that returns 1 if a rule r is matched to a Web page pg and 0, otherwise.

$$rvc(\text{Syntactic Component}, pg) = \left(\sum_r rsc(r) \cdot mat(r, pg) \right) / \sum_r mat(r, pg) \quad (12)$$

3.4.4 Categorical Match Component Relevancy Computation

Categorical Match component represents the similarity measure between the structure of user-created taxonomy and the category information provided by search engines for the retrieved Web pages. Nowadays, many popular search engines respond to the users query not only with a list of URLs for Web pages but also with their own categorical information for each Web page. For example, the following is an extract of search results provided by Lycos for the query “chair”.

- (1) Donald B. Brown Research Chair on Obesity
Health > Mental Health > Disorders
> Eating Disorders > Obesity
- (2) Steel Chair Wrestling
Sports > Fantasy > Pro Wrestling
- ...
- (3) Chair Technologies
Business > Industries > Manufacturing
> Consumer Products > Furniture
> Seating > Office Chairs
- ...

In the search results, the numbers on the left hand side denote the ranks of the corresponding Web pages and the associated lines below each title show the related category information for those Web pages. Although different search engines associate different category information to

the same Web page, such categorical information helps users filter out some of the returned search results without actually visiting the URL. Actually, the categorical match component is designed to provide the benefits of manual filtering by automatic means; this is accomplished by comparing the WWST terms with the categorical information provided by search engines. This is one of the major contributions of this paper.

Now, let's discuss how to measure the relevancy between the WSTT and the categorical information in more detail. We first represent the category information for a Web page pg from a search engine s , as an ordered set of category terms in a form such as $\{cat_1, cat_2, \dots, cat_m\}$, where cat_i is the i -th category term and m is total number of category terms in the set and we denote it $catinfo(pg, s)$. For example, $catinfo(\text{Chair Technologies}, \text{Lycos})$ in the above case, can be represented as the ordered set of category terms, $\{\text{Business}, \text{Industries}, \text{Manufacturing}, \text{Consumer Products}, \text{Furniture}, \text{Seating}, \text{Office Chairs}\}$. However, since it is hard to directly compare such $catinfo$ to the entire WSTT, here we adopt an approach similar to that

applied to the Semantic Component case, where we first measure the relevance of a *catinfo* to a single term combination, and then, combine them up to a single composite measure with respect to the entire WSTT.

So now, the relevance between a *catinfo* and a term combination *tc* can be measured from two different aspects, co-occurrence of terms and order consistency of terms. To measure the co-occurrence, we use the following formula (13).

$$coccur(tc, catinfo) = \left(\frac{\sum_t member(t, catinfo)}{ns(tc)} \right) \cdot \left(\frac{\sum_{cat} member(cat, tc)}{ns(catinfo)} \right) \quad (13)$$

where $t \in tc$ is a term, $cat \in catinfo$ is a category term, and $member(x, y)$ is a function that returns 1 if x is a member of y and 0, otherwise.

To consider the order consistency, let's first denote the precedence relationship of two arbitrary terms, t_l and t_r as (t_l, t_r) , and that means t_l precedes t_r in an ordered terms set. We also define a set of all available precedence relationships from an ordered set of terms x , as $prelset(x)$. Then we measure the consistency of *catinfo* with respect to a precedence relationship, (t_l, t_r) as follows:

$$cons((t_l, t_r), catinfo) = \begin{cases} 1 & \text{if } t_l, t_r \in catinfo \text{ and } t_l \text{ precedes } t_r \text{ in } catinfo. \\ 0 & \text{otherwise} \end{cases} \quad (14)$$

Now, let's define a consistency of a category information *catinfo* to a term combination *tc* as $constc(tc, catinfo)$ and the equation (15) shows how to compute it. Because we want to focus only on order consistency between *catinfo* and *tc* not depending on co-occurrence between them, we additionally define an ordered intersection set of *tc* and *catinfo*, where order of its element terms follows *tc*, as $isset(tc, catinfo)$ and then we can remove co-occurrence effect by only considering the precedence relationships in that set.

$$constc(tc, catinfo) = \sum_{pr} cons(pr, catinfo) / \binom{ns(isset(tc, catinfo))}{2} \quad (15)$$

where $pr \in prelset(isset(tc, catinfo))$, is a precedence relationship.

For example, let a term combination *tc* be $\{a, b, c, d, e\}$ and a category information *catinfo* be $\{a, e, c, f\}$. Then $isset(tc, catinfo)$ becomes $\{a, c, e\}$ and also $prelset(isset(tc, catinfo))$ becomes $\{(a, c), (a, e), (c, e)\}$. According to the formula (14), $cons((a, c), catinfo)$, $cons((a, e), catinfo)$, and $cons((c, e), catinfo)$ have their value as 1, 1, and 0, respectively. Since $ns(isset(tc, catinfo))$ is 3 in this case, the denominator of the equation (15) becomes 3, and finally $constc(tc, catinfo)$ becomes $(1+1+0)/3 = 2/3$. Also in this case, $coccur(tc, catinfo)$ becomes $3/5 \times 3/4 = 9/20$, because 3 of 5 terms of *tc* appear in *catinfo* and 3 of 4 terms of *catinfo* appear in *tc*.

To synthesize both the above aspects of categorical match between a term combination and a category information, we define the following measure, $rvtcc(tc, catinfo)$.

$$rvtcc(tc, catinfo) = \alpha \cdot coccur(tc, catinfo) + (1 - \alpha) \cdot constc(tc, catinfo) \quad (16)$$

where α is a $[0, 1]$ scale factor to represent the relative importance of co-occurrence to order consistency and it is set to 0.5 by default.

Actually since a Web page can have several category labels from different search engines for a given term combination, we need to further synthesize to obtain a single categorical match relevance value of a Web page pg for a term combination tc , $rvtc^{CM}(tc, pg)$ and it is formalized in (17).

$$rvtc^{CM}(tc, pg) = \frac{\sum_{s \in SC} rvtcc(tc, catinfo(pg, s))}{ns(SC)} \quad (17)$$

where s is a search engine and SC is a set of search engines that have categorical information for the page pg .

As in the case of Semantic Component, we adopt the max function to synthesize $rvtc^{CM}$ s to obtain a categorical match relevance value of a Web page pg for a leaf path lp , $rvlp^{CM}(lp, pg)$ as follows:

$$rvlp^{CM}(lp, pg) = \max_{tc \in tcs(lp)} rvtc^{CM}(tc, pg) \quad (18)$$

We also can obtain a single composite categorical match relevance value of a Web page pg , $rvc(\text{Categorical Match Component}, pg)$ using the AHP-based WSTT relevance value composition mechanism that is formalized in (11). To do this, we first set $bndfn(lp, pg)$ in the equation (11) as $rvlp^{CM}(lp, pg)$, then we propagate values from leaf nodes to the root node. At the root node n_0 , we obtain a single composite categorical match relevance value of a Web page pg as $rvn(n_0, pg)$ and we finally assign this value to $rvc(\text{Categorical Match Component}, pg)$, which will be used to obtain a composite relevance value with other preference components.

3.4.5 Search Engine Component Relevancy Computation

The Search Engine component represents the user's biases toward and confidence in a search engine's results. To measure this search engine component, let's first define a basic unit information, that is, rank of a Web page pg by search engine s for the request from term combination tc as $rank(tc, pg, s)$ and also define the number of resulting Web pages from search engine s for term combination tc as $npg(tc, s)$. In order to synthesize the search engine component with other components, we transform the rank information to a $[0, 1]$ scale normalized rank, $rank^N(tc, pg, s)$ according to the following equation.

$$rank^N(tc, pg, s) = 1 - \frac{(rank(tc, pg, s) - 1)}{npg(tc, s)} \quad (19)$$

The above normalization implies our intention to further discriminate the similarly ranked pages depending on the size of populations of those pages. For example, it transforms the second ranked page of ten result pages to a larger value than the same second of five results. Now, to obtain a composite search engine relevance value of a Web page pg for a term combination tc , $rvtc^{SE}(tc, pg)$, we adopt a weighted average method based on user's search engine preference as follows:

$$rvtc^{SE}(tc, pg) = \sum_s sw(s) \cdot rank^N(tc, pg, s) \quad (20)$$

To synthesize this in terms of a leaf path, we also define a search engine relevance measure of a Web page pg for a leaf path lp as $rvlp^{SE}(lp, pg)$ and formalize it as the equation (21).

$$rvlp^{SE}(lp, pg) = \frac{\sum_{tc \in tcslp(lp)} rvtc^{SE}(tc, pg)}{ns(tcslp(lp))} \quad (21)$$

Finally to obtain a search engine relevance value of a Web page with respect to WSTT, we also adopt AHP-based WSTT relevance value composition mechanism and so, we set $bndfn(lp, pg)$ in the equation (11) as $rvlp^{SE}(lp, pg)$. After value propagation process, we obtain a single synthesized search engine relevance value at the root node n_0 and assign its value, $rvn(n_0, pg)$ to $rvc(\text{Search Engine Component}, pg)$.

3.4.6 Authority/Hub Component Relevancy Computation

Authority/Hub component: represents the level of user preference for *Authority* or *Hub* sites and pages [31]. At present, no such authority or hub ranking service exists on the Web. Therefore, we have not incorporated this component into our proof-of-concept prototype.

3.4.7 Popularity Component Relevancy Computation

Our final component to be considered is the Popularity component and it represents the user's preference for popular sites. Popularity can be measured by the number of visitors or the number of requests for the specific page or site and there exist some publicly available services for this popularity information like www.yep.com. To compute the relevance value of a Web page pg in terms of the popularity component, let's introduce some definitions as follows.

$pop(pg)$: denotes the average number of daily visitors to the Web page pg .

$pgset$: denotes the set of whole Web pages retrieved.

Based on the definitions, we formalize the popularity relevance measure of a Web page pg as follows:

$$rvc(\text{Popularity Component}, pg) = \frac{pop(pg)}{\max_{x \in pgset} pop(x)} \quad (22)$$

So far, we have presented our approach for users to express their search intent, their search preference in terms of six preference components, have proposed a series of rating methods to compute for each a relevance value of the component, and provided a mechanism to combine them into a single measure of relevance. Finally we use this single measure to provide the users more relevant information with a list of resulting Web pages in a descending order of relevance value.

4 WebSifter II System Architecture

In this section we present the architecture of WebSifter II. Figure 7 shows the overall architecture of WebSifter II and its components. Major information flows are also depicted. WebSifter II consists of eight subsystems and four major information stores.

Now let's briefly introduce each of the components, their roles, and related architectural issues.

1) *WSTT Elicitor*

The WSTT elicitor supports the entire process (see section 3.1) of specifying a WSTT in a GUI environment. A user can express his search intent as a WSTT through interactions with the WSTT elicitor. This includes building a taxonomy tree, assigning weights to each node, and choosing a concept from an available list of WordNet concepts. To achieve this goal, the WSTT elicitor also cooperates with an Ontology agent, a Stemming agent, and a Spell Check agent. Once a user finishes building a WSTT, then the WSTT elicitor stores the WSTT information into the WSTT base in XML format.

2) *Ontology Agent*

The ontology agent is responsible for requesting available concepts of a given term via a Web version of WordNet (<http://www.cogsci.princeton.edu/cgi-bin/webwn/>) and also for interpreting the corresponding HTTP-based results. The agent receives requests for the concepts from WSTT elicitor and returns available concepts in an understandable form. Although WebSifter presently supports cooperation only with WordNet, its design can be easily extended to cooperate with other ontology servers such as CYC [33] and EDR [34].

3) *Stemming Agent*

Our stemming agent is based on Porter's algorithm [35]. It has two major roles: 1) to cooperate with the WSTT elicitor in transforming the terms in a concept to stemmed terms, and 2) to transform the content of Web pages into the stemmed terms internally through cooperation with a page request broker. As a result, the terms in concepts and the terms in Web pages can be compared to each other via their stemmed versions.

4) *Spell Check Agent*

The spell check agent monitors user's text input to the WSTT elicitor and checks and suggests correct words to the user in real time.

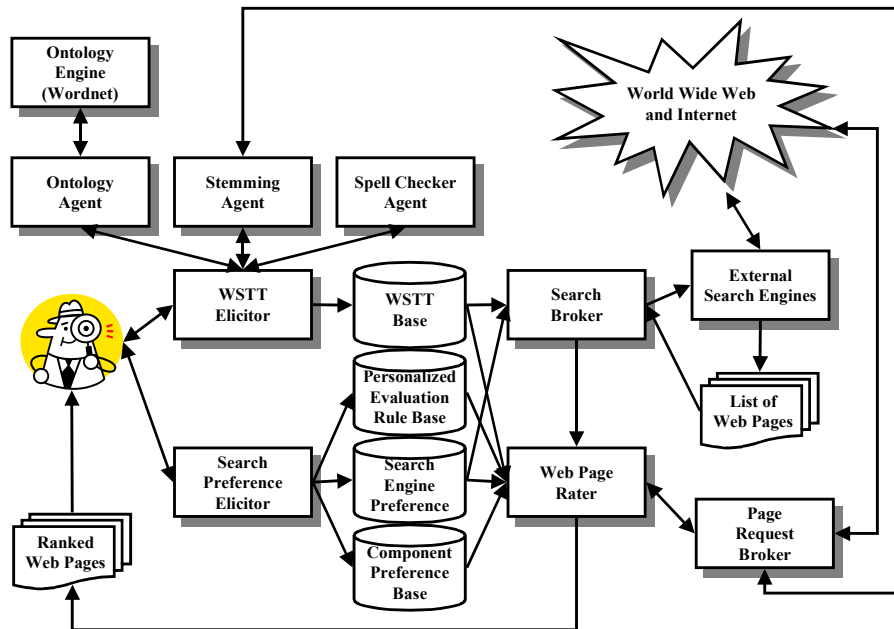


Fig. 7. System architecture of WebSifter II

5) *Search Preference Elicitor*

The search preference elicitor, via a GUI, supports the process (cf. section 3.2) to capture the user's search preferences. A user can express his search preference by assigning their preference weights to each of the preference components and also to their favorite search engines. Moreover, it allows the user to modify the default values assigned to each syntactic URL class such as Direct Hit, Directory Hit and Page Hit. Whenever the user modifies them, it updates the related information stored in the Personalized Evaluation Rule Base, the Search Engine Preference Base, and the Component Preference Base.

6) *Search Broker*

The search broker performs the processes specified in section 3.3. It first interprets the XML-based WSTT and then generates all corresponding query statements. Using this set of queries, it requests information from a set of popular search engines simultaneously. Finally, it interprets the results returned from the search engines and then stores parsed information in a temporary data store. When it finishes its work, it activates the Web page rater to begin the rating process.

7) *Page Request Broker*

Page request broker is responsible for requesting the content of a specific URL and it cooperates with both the stemming agent and the Web page rater.

8) *Web Page Rater*

Web page rater supports the entire Web page evaluation process specified in section 3.4 and also is responsible for displaying the results to users. This subsystem is the most complex and computationally intensive module of WebSifter II, and it uses all four major information stores and also communicates with search broker and page request broker.

5 Empirical Results

5.1 Implementation

We have implemented the semantic taxonomy-based personalizable meta-search agent in a working prototype using Java, with the exception of the spell check agent. Now, we plan to incorporate a commercial spell check agent into this system.

Figure 8 shows an illustrative screen where the user builds a WSTT using the WSTT elicitor. Figure 9 shows another screen of the WSTT elicitor supporting the selection of an intended concept from available concepts for a given term, obtained through cooperation between the ontology agent and WordNet.

Figure 10 shows a panel by which a user specifies his search preference using the search preference elicitor. The four tab windows in Figure 10 accept user preference for the relevance components, search engines, advanced parameters, and classification rules for Web pages, respectively. However, only the tab window for preference components is shown in Figure 10.

Finally, Figure 11 is a main screen of our WebSifter II system and it shows illustrative results for a given search query generated through the above steps.

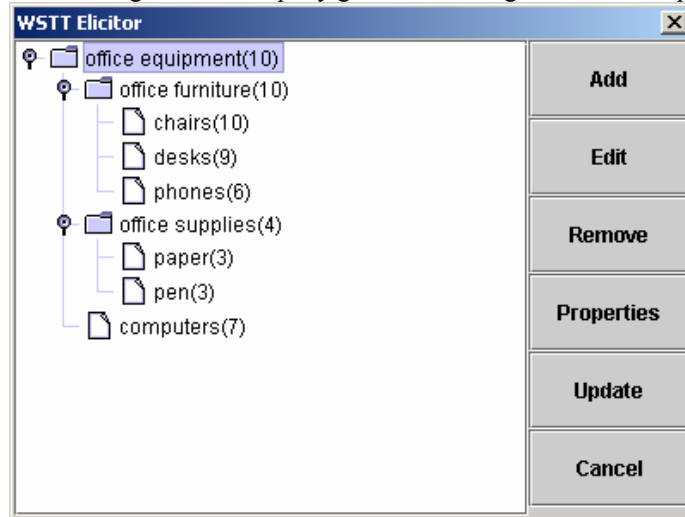


Fig. 8. An illustrative screen of the WSTT Elicitor

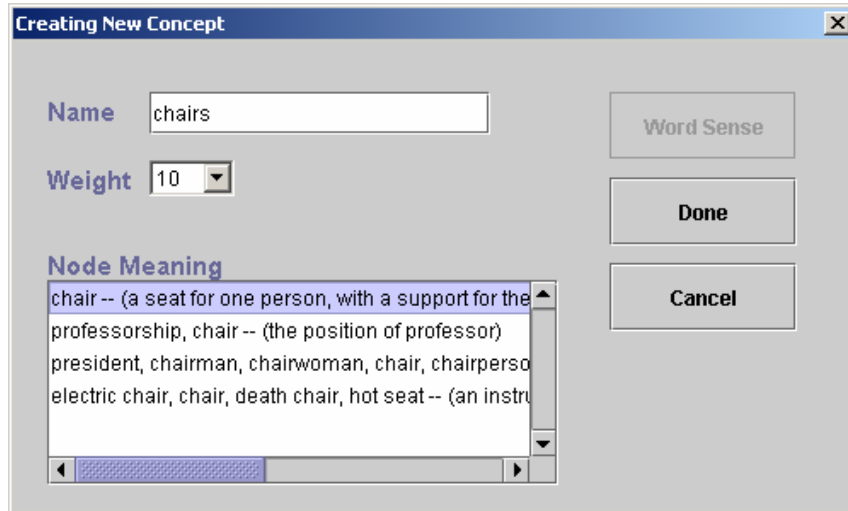


Fig. 9. An illustrative screen for concept selection

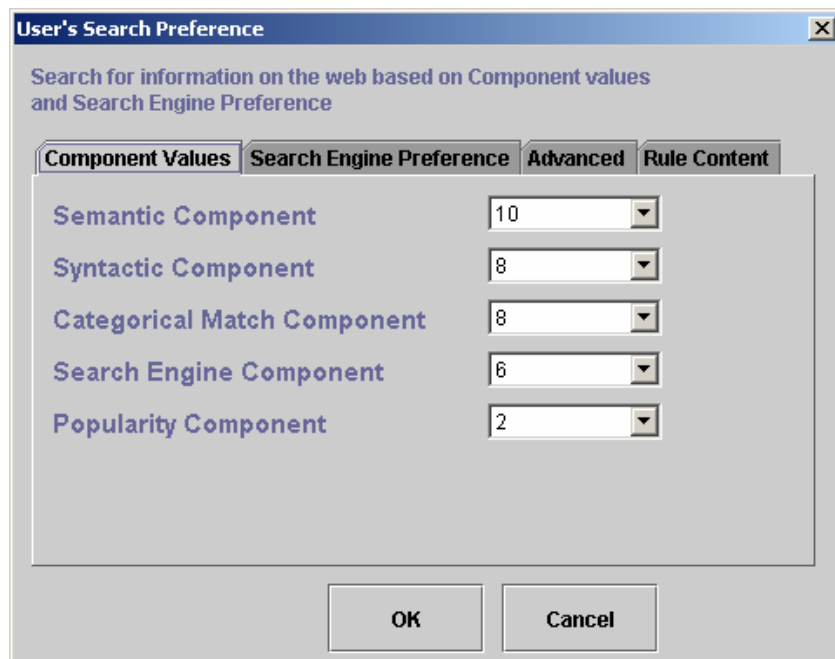


Fig. 10. A tab window of the Search Preference Elicitor

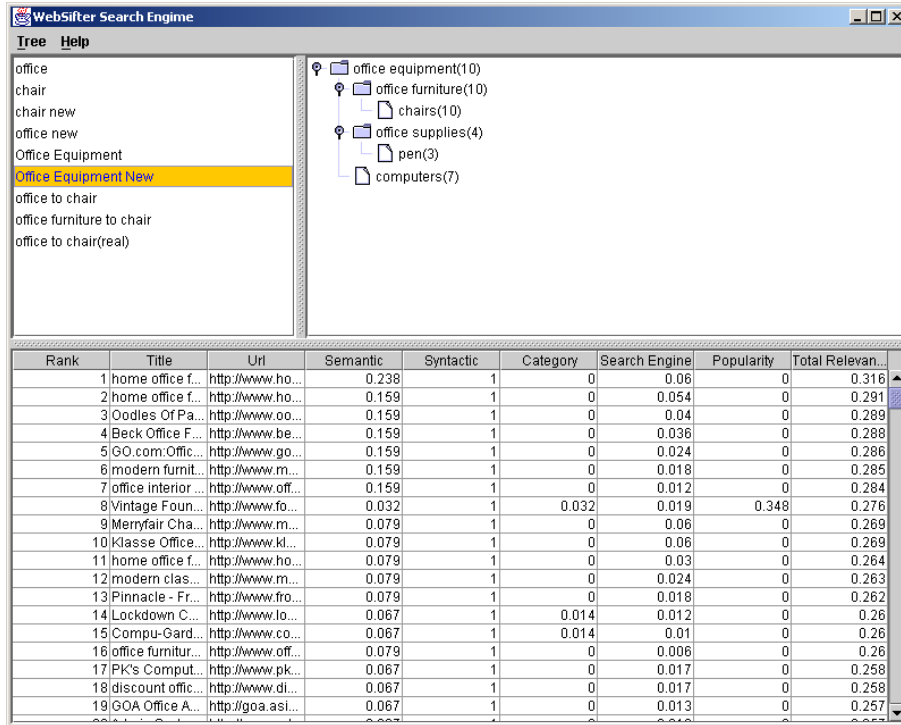


Fig. 11. An illustrative screen of the results from the Web Page Rater

Note that in Figure 11 the top-left pane contains specified WSTT queries with “Office Equipment New” highlighted, the top-right pane shows the WSTT of “Office Equipment New” and the bottom pane shows the search results ranked by the Total Relevance component.

5.2 Experimental Results

We are currently doing empirical experiments on our approach and some of them are presented. We first show results from the experiments related to the relevancy enhancement in the page hits by considering both positive and negative concepts obtained through user’s concept selection as shown in Figure 9.

Table 1 shows the page hits retrieved by WebSifter II for the search of a just a single term, “chair”, associated with the selected concept as a seat for one person, which appears as the first concept in Figure 9. In this table, the ranks of web pages provided by WebSifter appear in the first left column and the next column shows their corresponding URLs.

Table 1. Result comparison of WebSifter query “chair” with other search engines

Rank	URL	Relevancy	Copernic	Altavista	Google	Yahoo	Excite
1	http://www.countryseat.com	Y	-	-	-	-	-
2	http://www.infant-car-seat.com/	N	-	-	-	-	-
3	http://www.chairmaker.co.uk/	Y	-	-	-	19	-
4	http://www.convertible-car-seat.com/	N	-	-	-	-	-
5	http://www.booster-car-seats.com/	N	-	-	-	-	-
6	http://www.booster-seats-online.com/	N	-	-	-	-	-
7	http://www.booster-car-seat.com/	N	-	-	-	-	-
8	http://www.podiatrychair.com/	N	-	-	-	-	9
9	http://www.carolnchair.com/	Y	-	-	-	9	-
10	http://www.chairdancing.com/	N	-	-	-	-	-
11	http://www.message-chairs-online.com/	N	-	-	-	-	13
12	http://www.panasonic-massage-loungers.com/	N	-	-	-	-	14
13	http://www.fairfieldchair.com/	Y	-	-	15	-	-
14	http://www.gasserchair.com/	Y	-	16	-	-	-
15	http://www.chairtech.com/	Y	-	18	-	-	-
16	http://www.snugseat.com/	N	-	-	-	-	-
17	http://www.seat.com/	N	-	-	-	-	-
18	http://www.fifthchair.org/	N	3	2	5	8	3
19	house.com/mainframes/About_each_show/show_list/	N	19	-	9	-	-
20	http://www.jeanmonnetprogram.org/	N	5	1	1	-	5

The right five columns show a comparison of our result to the results by the search engines, Copernic, Altavista, Google, Yahoo, and Excite. The number in each cell of those columns indicates the rank assigned by the corresponding search engine (column) for the given page-hit (row). Especially, the ‘-’ sign in the table indicates the corresponding Web page was not retrieved, or was ranked lower than 20th by the corresponding search engine. Actually, most of the page-hits, which do not have any corresponding ranks from other search engines, are due to WebSifter’s use of *semantic concepts* rather than *terms*. A concept usually consists of multiple terms as mentioned in Section 3.1, and so, our approach generates multiple queries based on such multiple terms. For this “chair” case, not only “chair” but also “seat” is used in the search. But when we just use the search engines, “chair” is the only term posed to them because they do not support query extensions using concepts.

Relevancy column shows the relevancies of the corresponding page-hits, where ‘Y’ means the page is relevant to the query, while ‘N’ means it is irrelevant. The decisions are based on our subjective but clear evaluation criterion as to whether the corresponding Web page is relevant to a real chair or seat. As shown in Table 1, we found six relevant Web pages in our page-hit result based on this criterion and most of them are also relatively high-ranked. The results also show most of high-ranked pages in the search engines appeared to be irrelevant and were low-ranked in our page-hits list.

In Table 2, we also compare our results with Copernic, which is one of the leading commercial meta-search engines. Copernic retrieved only one relevant Web page and the reader can easily find its rank information is strongly related to the ranks from the search engines in this table. Evidence for this is the fact that the first-ranked Web page is “www.theelectricchair.com” and has strong support from Google and Yahoo, where its ranks are 2 and 1, respectively. However, it is about the real electric chair and totally irrelevant to the chair that we are looking for. Our approach ranked it as lower than 20th and so, it does not appear on our list. Even though our approach also considers the rank information from the search engines, WebSifter ranks that electric chair page lower using negative concept.

Table 2. Comparison of Copernic with WebSifter and other search engines

RANK	URL	Relevancy	WebSifter	Altavista	Google	Yahoo	Excite
1	http://www.theelectricchair.com/	N	-	-	2	1	-
2	http://www.chair-sales.com/	N	-	-	-	-	1
3	http://www.fifthchair.org/	N	18	2	5	8	3
4	http://www.urbanlegends.com/death/electric.chair/electric_chair	N	-	-	-	15	-
5	http://www.jeanmonnetprogram.org/	N	20	1	1	-	5
6	http://www.the-perfect-chair.com/	N	-	-	-	-	2
7	http://www.widc.org	N	-	3	4	2	4
8	http://www.law.harvard.edu/programs/JeanMonnet/	N	-	-	3	-	-
9	http://www.chairpage.com/	N	-	-	-	-	6
10	http://www.obesity.chair.ulaval.ca/	N	-	-	7	12	-
11	http://www.tderc.com	N	-	-	-	3	-
12	http://www.titanicdeckchair.com	N	-	-	-	6	-
13	http://www.electricchair.com/	N	-	-	8	-	11
14	http://www.gsb.stthomas.edu/ethics/	N	-	-	-	4,5	-
15	http://www.producerschair.com/	N	-	-	6	-	-
16	http://www.nantucketbeachchair.com	N	-	-	17	20	-
17	http://www.emf.net/~troop24/scouting/scouter.html	N	-	-	-	-	7
18	http://www.examchair.com/	N	-	-	-	-	8
19	http://www.painted-house.com/mainframes/About_each_show/sh	N	19	-	9	-	-
20	http://www.windsorchairresources.com/	Y	-	-	10	-	-

Overall search performance comparison between WebSifter and other competing search engines are summarized in Table 3. In this table, hit ratio means the percentage of the relevant page-hits to 20 high ranked pages and average rank of relevant pages partially measures the quality of ranking. That is, as far as the hit ratios are same or similar, the search method that produces the smallest average rank of relevant pages is the best ranked.

Table 3. Overall search performance for the query *chair*

Search Engines	Hit Ratio	Average Rank of Relevant Pages
WebSifter	30%	9.17
Copernic	5%	20.00
Altavista	20%	13.50
Google	20%	14.00
Yahoo	10%	14.00
Excite	20%	15.00

As shown in the table, the WebSifter approach outperforms other approaches in both measures. By the way, Copernic shows the poorest performance and this seems to be caused because most of relevant page hits were low-ranked by the search engines in this case, but meta-search engine like Copernic, tends to consider high-ranked page hits, first.

The above results show that consideration of the positive and negative concepts can contribute greatly to the precision of the ranking. To further validate the effect of the hierarchical concept tree in search, we extend the test query to a chair for office use. This query is represented as the case (a) in Figure 2 in our approach and we use “office” and “chair” terms for other search engines.

Table 4 shows the retrieved page-hits by WebSifter II for this query and their ranks. WebSifter shows 95%-hit ratio in this case and the only irrelevant page-hit is ranked as 20th. In this experiment, we also tested the WebSifter approach for the case where we suppress the effect of the categorical match component mentioned in 3.4.4, by setting its weight to zero. Through this additional test, we can evaluate the contribution of both the hierarchical concept tree and the categorical match

component to relevancy precision by comparing the normal case with the suppressed one. The corresponding ranks generated by the suppressed case for each page-hit are shown in the right most column in Table 4.

Table 4. WebSifter results for the query *office* and *chair*

Rank	URL	Relevancy	w/o Categ
1	http://www.seatingvfm.com/	Y	1
2	http://www.officechair.co.uk/	Y	2
3	http://www.AmericanErgonomics.com/	Y	9
4	http://www.ompchairs.com/	Y	22
5	http://www.klasse.com.au/	Y	4
6	http://www.cyberchair.com/	Y	46
7	http://www.leap-chair.com	Y	47
8	http://www.seizaseat.com/	Y	50
9	http://www.zackback.com	Y	49
10	http://www.fairfieldchair.com	Y	2
11	http://www.chair-ergonomics.com/	Y	5
12	http://www.buy-ergonomic-chairs.com/	Y	6
13	http://www.jfainc.com/	Y	7
14	http://www.chairtech.com/	Y	8
15	http://www.plasticfoldingchairs.com/	Y	13
16	http://www.kneelsit.com/	Y	10
17	http://www.home-office-furniture-store.com/	Y	11
18	http://www.home-office-furniture-site.com/	Y	12
19	http://www.amadio.it/uk/	Y	19
20	http://www.newtrim.co.uk/	N	15

Table 5 also shows the retrieved page-hits and their ranks from the case where we turn off the categorical match component. In this table, the column labelled by ‘with Categ’ also shows the corresponding ranks of the page-hits from the case when we include the categorical match component into the WebSifter II rating mechanism. By comparing both Tables 4 and 5, we can find that the consideration of concept hierarchy and the categorical match component affects greatly the resulting rankings. As shown in the results, three of four irrelevant page-hits were ranked lower than 20th, and were replaced with relevant page hits. One remaining irrelevant page hit is also downgraded from 15th rank to 20th rank. In summary, the hit ratio from Table 5 (80%) is enhanced to 95% in Table 6. This implies consideration of concept hierarchy and categorical match component contributes a 15% performance enhancement in this case and they are very important factors to be considered in retrieving the relevant page-hits.

Table 5. WebSifter result for query *office* and *chair* without category match component

Rank	URL	Relevancy	with Categ
1	http://www.seatingvfm.com/	Y	1
2	http://www.fairfieldchair.com	Y	10
3	http://www.officechair.co.uk/	Y	2
4	http://www.klasse.com.au/	Y	5
5	http://www.chair-ergonomics.com/	Y	11
6	http://www.buy-ergonomic-	Y	15
7	http://www.jfainc.com/	Y	13
8	http://www.chairtech.com/	Y	14
9	http://www.AmericanErgonomics.c	Y	3
10	http://www.kneelsit.com/	Y	16
11	http://www.home-office-furniture-	Y	17
12	http://www.home-office-furniture-	Y	18
13	http://www.plasticfoldingchairs.co	Y	15
14	http://www.office-interior-	N	21
15	http://www.newtrim.co.uk/	N	20
16	http://www.oa-chair.com/	N	23
17	http://www.buy-ergonomic-	Y	24
18	http://www.countryseat.com/	Y	26
19	http://www.amadio.it/uk/	Y	19
20	http://www.mobile-office-desk.com/	N	28

Overall search performance comparison results for the query, “office” and “chair” are shown in Table 6. WebSifter II approach also outperforms the other approaches and even our approach without consideration of the categorical match component is still better or quite competitive to other approaches.

Table 6. Overall search performance comparisons for the query *office* and *chair*

Search Engine	Hit Ratio
WebSifter	95%
WebSifter (w/o Categorical Match)	80%
Corpenic	75%
Altavista	65%
Google	60%
Yahoo	85%
Excite	65%

These results from the experiments so far show most promising evidence to validate our approach. We are still performing additional experiments to validate the performance of our approach in more broad cases.

6 Conclusions

We have proposed a semantic taxonomy-based personalizable meta-search agent approach to achieve two important and complementary goals: 1) allowing users more expressive power in formulating their Web searches, and 2) improving the relevancy of search results based on the user's real intent. In contrast to the previous research, we have focused not only on the search problem itself, but also on the decision-making problem that motivates users to search the Web.

Now, let's briefly summarize our contributions as follows. We have proposed a search-intention representation scheme, the Weighted Semantic-Taxonomy Tree, through which users express their real search intentions by specifying domain-specific concepts, assigning appropriate weights to each concept, and expressing their decision problem as a structured tree of concepts. We also allow users to express their search result evaluation preferences as a function of six preference components.

Second, to enhance the *precision* of the retrieved information, we present a hybrid rating mechanism which considers both the user's search intent represented by the WSTT and user's search preference represented by multi-preference components such as semantic relevance, syntactic relevance, categorical match, page popularity, and authority/hub rating.

Third, we have designed and have implemented a meta-search agent system called WebSifter II that cooperates with WordNet for concept retrieval, and most well known search engines for Web page retrieval. Our open and extensible architecture allows new services to be incorporated in WebSifter II, as they become available. For the empirical validation of our approach, we empirically validate our approach already for some limited cases and we are also doing some real world experiments of our system.

References

- 1 Scime, A. and L. Kerschberg, "WebSifter: An Ontology-Based Personalizable Search Agent for the Web," *International Conference on Digital Libraries: Research and Practice*, Kyoto Japan, 2000, pp. 493-446.
- 2 Scime, A. and L. Kerschberg, "Web Sifter: An Ontological Web-Mining Agent for E-Business," *Proceedings of the 9th IFIP 2.6 Working Conference on Database Semantics (DS-9): Semantic Issues in E-Commerce Systems*, Hong Kong, 2001.
- 3 Lawrence, S. and C. L. Giles, "Accessibility of Information on the Web," *Nature*, vol. 400, 1999, pp. 107-109.
- 4 Selberg, E. and O. Etzioni, "The MetaCrawler Architecture for Resource Aggregation on the Web," *IEEE Expert*, vol. 12, no. 1, 1997, pp. 11-14.
- 5 Howe, A. E. and D. Dreilinger, "Savvy Search: A Metasearch Engine that Learns which Search Engines to Query," *AI Magazine*, vol. 18, no. 2, 1997, pp. 19-25.
- 6 Lawrence, S. and C. L. Giles, "Context and Page Analysis for Improved Web Search," *IEEE Internet Computing*, vol. 2, no. 4, 1998, pp. 38-46.
- 7 Ackerman, M., et al., "Learning Probabilistic User Profiles - Applications for Finding Interesting Web Sites, Notifying Users of Relevant Changes to Web Pages, and Locating Grant Opportunities," *AI Magazine*, vol. 18, no. 2, 1997, pp. 47-56.

- 8 Armstrong, R., et al., "WebWatcher: A Learning Apprentice for the World Wide Web," *Proceedings of the 1995 AAAI Spring Symposium on Information Gathering from Heterogeneous, Distributed Environments*, 1995.
- 9 Shavlik, J. and T. Eliassi-Rad, "Building Intelligent Agents for Web-based Tasks: A Theory-Refinement Approach," *Proceedings of the Conference on Automated Learning and Discovery: Workshop on Learning from Text and the Web*, Pittsburgh, PA, 1998.
- 10 Chen, Z., et al., "WebSail: from On-line Learning to Web Search," *Proceedings of the First International Conference on Web Information Systems Engineering*, vol. 1, 2000, pp. 206-213.
- 11 Chakrabarti, S., et al., "Enhanced Hypertext Categorization using Hyperlinks," *Proceedings of ACM SIGMOD International Conference on Management of Data*, Seattle, Washington, 1998, pp. 307-318.
- 12 Li, Y., "Toward a Qualitative Search Engine," *IEEE Internet Computing*, vol. 2, no. 4, 1998, pp. 24-29.
- 13 Maes, P., "Agents that reduce work and information overload," *Communications of the ACM*, vol. 37, no. 7, 1994, pp. 30-40.
- 14 Terveen, L., et al., "PHOAKS: a System for Sharing Recommendations," *Communications of the ACM*, vol. 40, no. 3, 1997, pp. 59-62.
- 15 Bollacker, K. D., et al., "Discovering Relevant Scientific Literature on the Web," *IEEE Intelligent Systems*, vol. 15, no. 2, 2000, pp. 42-47.
- 16 Balabanovic, M. and Y. Shoham, "Content-Based, Collaborative Recommendation," *Communications of the ACM*, vol. 40, no. 3, 1997, pp. 66-72.
- 17 Krulwich, B., "Lifestyle Finder," *AI Magazine*, vol. 18, no. 2, 1997, pp. 37-46.
- 18 de Vel, O. and S. Nesbitt, "A Collaborative Filtering Agent System for Dynamic Virtual Communities on the Web," *Working notes of Learning from Text and the Web, Conference on Automated Learning and Discovery CONALD-98*, Carnegie Mellon University, Pittsburgh, 1998.
- 19 Chen, H. and S. Dumais, "Bringing Order to the Web: Automatically Categorizing Search Results," *Proceedings of the CHI 2000 conference on Human factors in computing systems*, The Hague Netherlands, 2000, pp. 145-152.
- 20 Aridor, Y., et al., "Knowledge Agent on the Web," *Proceedings of the 4th International Workshop on Cooperative Information Agents IV*, 2000, pp. 15-26.
- 21 Chakrabarti, S., et al., "Focused Crawling: A New Approach to Topic-Specific Web Resource Discovery," *Proceedings of the Eighth International WWW Conference*, 1999, pp. 545-562.
- 22 Salton, G., et al., "A Vector Space Model for Automatic Indexing," *Communications of the ACM*, vol. 18, no. 11, 1975, pp. 613-620.
- 23 Clark, D., "Mad Cows, Metathesauri, and Meaning," *IEEE Intelligent Systems*, vol. 14, no. 1, 1999, pp. 75-77.
- 24 Guarino, N., et al., "OntoSeek: Content-based Access to the Web," *IEEE Intelligent Systems*, vol. 14, no. 3, 1999, pp. 70-80.
- 25 Fensel, D., et al., "On2broker: Semantic-Based Access to Information Sources at the WWW," *Proceedings of the World Conference on the WWW and Internet (WebNet 99)*, Honolulu, Hawaii, USA, 1999, pp. 25-30.
- 26 Staab, S., et al., "A System for Facilitating and Enhancing Web Search," *Proceedings of IWANN '99 - International Working Conference on Artificial and Natural Neural Networks*, Berlin, Heidelberg, 1999.
- 27 Martin, P. and P. W. Eklund, "Knowledge Retrieval and the World Wide Web," *IEEE Intelligent Systems*, vol. 15, no. 3, 2000, pp. 18-25.
- 28 Miller, G. A., "WordNet a Lexical Database for English," *Communications of the ACM*, vol. 38, no. 11, 1995, pp. 39-41.

- 29 Klein, D. A., *Decision-Analytic Intelligent Systems: Automated Explanation and Knowledge Acquisition*, Lawrence Erlbaum Associates, 1994.
- 30 Boose, J. H. and J. M. Bradshaw, "Expertise Transfer and Complex Problems: Using AQUINAS as a Knowledge-acquisition Workbench for Knowledge-Based Systems," *Int. J. Man-Machine Studies*, vol. 26, 1987, pp. 3-28.
- 31 Kleinberg, J. M., "Authoritative Sources in a Hyperlinked Environment," *Journal of the ACM*, vol. 46, no. 5, 1999, pp. 604-632.
- 32 Saaty, T. L., *The Analytic Hierarchy Process*, New York, McGraw-Hill, 1980.
- 33 Lenat, D. B., "Cyc: A Large-Scale Investment in Knowledge Infrastructure," *Communications of the ACM*, vol. 38, no. 11, 1995, pp. 33-38.
- 34 Yokoi, T., "The EDR Electronic Dictionary," *Communications of the ACM*, vol. 38, no. 11, 1995, pp. 45-48.
- 35 Porter, M., "An Algorithm for Suffix Stripping," available at <http://www.muscat.co.uk/~martin/def.txt>.