

# QuadSearch: A novel metasearch engine

Leonidas Akritidis<sup>1</sup>, George Voutsakelis<sup>2</sup>, Dimitrios Katsaros<sup>1,2</sup>, and  
Panayiotis Bozanis<sup>2</sup>

<sup>1</sup> Data Engineering lab, Dept. of Informatics, Aristotle Univ., Thessaloniki, Hellas

<sup>2</sup> Computer & Communication Engineering Dept., Univ. of Thessaly, Volos, Hellas

lakritid@mywork.gr, gevoutsa@inf.uth.gr  
dkatsaro@csd.auth.gr, pbozanis@inf.uth.gr,  
<http://delab.csd.auth.gr/~dimitris>

**Abstract.** Metasearch engines are increasingly becoming a very useful tool for Web information retrieval. Their success depends mainly on their rank aggregation (fusion) method, their interface and their total “sustainability”, meaning that the engine is updated regularly and will still be around after some time.

In this paper we describe *QuadSearch*, an experimental metasearch engine that provides simultaneous access in four major conventional, crawler-based search engines. The heart of the new metasearch engine is based on two novel rank-based aggregation algorithms. Users can choose which rank aggregation algorithm to use and thus to adjust the results at their own needs. They can also alter the form of the result list in a more convenient and statistically enhanced way. The *QuadSearch* engine aims to combine speed, reliable rank aggregation method, “spam” free results, and detailed and enriched information. A publicly accessible interface for the new metasearch engine can be found at <http://delab.csd.auth.gr/~lakritid/metasearch/>.

## 1 Introduction

The Web is nowadays the main source of information; it is vast and doesn't have any specific structure, therefore it is extremely difficult for the user to find the information s/he desires without any external help. For this purpose, various systems have been developed to dig for information, but the most popular ones are the search engines, either general purpose search engines, like Google [18], or special purpose engines, like Medical World Search (<http://www.mwsearch.com/>). A search engine retrieves Web pages, relevant to a query which has been specified by the user. Although, search engines are extremely popular among Web users, they can not achieve large coverage and high scalability. It is a common belief among many people [21] that a single general purpose search engine for all Web data is unrealistic, because its processing power, no matter how large it is, can not scale to the rapidly increasing and unlimited amount of Web data.

The tool which rapidly gains acceptance by the users are the *metasearch engines* [17]. These systems work like a filter of the various crawler-based search engines that they combine. Metasearch engines run simultaneously a user query

across multiple *component search engines*, take the returned results and then aggregate them. The advantages of the metasearch engines are the following [17]: a) they increase the search coverage of the Web, b) they solve the scalability problem of searching the Web, c) they facilitate the exploitation of multiple search engines, and finally d) they improve the retrieval effectiveness.

The heart of any metasearch engine is the rank aggregation algorithm, which defines the final ranked result list from the individual results; this final list may be affected by the user's demands. For example, this can be done if the developers of the rank aggregation method provide parameters and choices so that the users can adjust some weights to any single conventional search engine depending on their confidence on these search engines.

Further process can be done in order to filter the results and allow the final result list of the metasearch engine to be relieved from unwanted, devious and undeservedly highly ranked Web pages. In a world which is frequently motivated by commercial interests, the user does not have a clear form of protection against the interests of individual search engines. Therefore, the metasearch engine should be capable to provide results to the user that are as free as they can be from paid listings and paid links.

This paper describes *QuadSearch*, an experimental metasearch engine that relies on two new rank aggregation algorithms. Although the default algorithm has some anti-spam properties, the user is granted the choice to select further antispam filtering. Moreover, users can define which search engines they want to be exploited by *QuadSearch*, the number of results which will be retrieved by each search engine, the number of results retrieved by *QuadSearch*, the form and appearance of result retrieval, and some other information enhancements.

## 1.1 Motivation and contributions

During the last years, the problem of paid listing within the retrieved results of search or metasearch engines has received a lot of attention. Generally, search engines present less paid listings than metasearch engines. In addition, paid links at the major search engines are, in some way, separated from the main result list. In contrast, the metasearch engines do not have such delineation, making unclear which links are paid. Related to this is the problem of "spam" by authors of Web pages who attempt to achieve undeservedly high rank for their Web pages by exploiting defects of the ranking functions of search engines.

Secondly, to the best of our knowledge, there is no rank aggregation algorithm that bears a wide variety of parameters like the number of the search engines where a particular item appeared, the total number of exploited search engines or the size of the top-k list returned from each search engine. We expect from a metasearch engine to:

- consider as many parameters of these as possible,
- have anti-spam and anti-paid list properties,
- provide extra information for demanding and experienced users, i.e., to support personalization properties,

- refrain from using any training data in order to perform the rank aggregation, because, there is usually no evidence about the underlying data properties and their distributions, and
- do not count upon the *scores* of the individual search engine rankings in order to perform the rank aggregation, because, most of the search engines do not provide such scores.

Motivated from these requirements, we developed the *QuadSearch* metasearch engine (named after the fact that it currently capitalizes on four most popular search engines), which satisfies the above criteria. Firstly, it has two fast rank aggregation algorithms; a default algorithm and an improved version of it, enhanced with more antispam properties. Secondly, it allows the user to exploit whichever version of the rank algorithm he desires, and adjust a lot of the interface and the appearance parameters. We have implemented an experimental version of this metasearch engine, which although not fully-fledged yet, it can be accessed at <http://delab.csd.auth.gr/~lakritid/metasearch/>.

The rest of this article is organized as follows: in Section 2 we briefly review the relevant work on metasearch engines; in Section 3, which presents the main article ideas, we describe the new rank aggregation methods and in Section 4 we present the implementation issues behind the developed metasearch engine. Finally, in Section 5 we highlight the main fetures of the new metasearch engine and in Section 6 we conclude the paper.

## 2 Existing metasearch engines

The first metasearch engines were established back in 1996. The fact that they allowed searches to be sent, at the same time, to various search engines gave the users the impression that they were getting more comprehensive results and so they gained popularity. The interested reader can find out the most popular metasearch engines of that period at [8, p. 388].

However, many problems occurred such as a lot of paid links inside organic results, the refusal of Google to cooperate with them and some fraud problems of pay-per-click search engines that led the metasearch engines to decline [5].

Now metasearch engines are coming back [4, 15] and a significant part of work is conducted for them [11, 12]. Researchers and developers work hard to prove that the results returned are defined by search algorithms and not by advertisers [5]. They use classification and personalization techniques that conventional search engines do not have. This does not mean that metasearch engines will overrun the major search engines, like Google, anytime soon, but it means that they are gradually gaining the position they deserve in the search market. In the sequel, we will simply list, among the many metasearch engines [1, 2, 14, 22], a few remarkable ones. An almost complete list of metasearch engines along with their main features and shortcomings can be found at [3, 6].

Vivisimo [7] and Jux2 are clustering engines, which automatically organize the retrieved pages on-the-fly into categories (groups). lxQuick ranks the results based on the top 10 rankings a site receives from various search engines,

iBoogie creates a list of categories related to search terms and InfoGrid provides direct links to major search engines and topical Web sites in different categories. SearchOnline offers a highly customizable interface, while Kartoo presents the results within a map that shows the most important sites and the linkage relationship between the results.

### 3 The heart of the proposed novel result merging

#### 3.1 Preliminaries

The ideal scenario for result merging is when each search engine gives a complete result list of all the alternative items, related to the keyword terms of a given query, in the universe of alternatives. This can not be done and it is far too unrealistic for two main reasons: (i) search engines' coverage is different, and (ii) search engines limit access only to a portion of the complete result list. The worst scenario is when the result lists of component search engines don't have overlapping elements between them. In this case there is nothing that a rank aggregation algorithm can do.

Several rank aggregation methods have been used by metasearch engines [16, 17, 19]. In the 1990s most of the metasearch engines used *score-based* ranking methods to produce their results, i.e., they utilized the scores (weights) returned by the component search engines in order to fuse the component rankings. Moreover, many metasearch techniques applied normalization on these ranking scores in order to make them comparable. Nowadays, no search engine provides the ranking scores, however it is possible to convert local ranks into ranking scores.

Although score-based methods appear to be more effective for rank fusion, the absence of scores (or denial to reveal) from many search engines' rankings turned these methods problematic [19]; thus the *rank-based* fusion became the mainstream in present metasearch engines [19]. For instance, the Borda Count [9, 19], which is a *voting-based* fusion method, is very popular among metasearch engines. Each result is a candidate and each search engine is the voter. Each candidate receives points from each voter according to its rank in the voter's list. For example, the top ranked candidate will receive  $n$  points, where  $n$  is the number of candidates. If a candidate is not in the top- $k$  list of some voter then it will receive a portion of the remaining points of the voter (each voter has a fixed number of points available for distribution) or a constant number (0 or 1), depending on the variation of the method. The Borda Count method can be found in different versions, like the weighted Borda Count method [20], where each voter also takes a score and therefore his opinion for a candidate is not treated equally against other voters. Improved methods for ranking comparison and merging in the case of ties can be found at [10, 11].

#### 3.2 The *ke* method

In the sequel, we will present the rank fusion method of *QuadSearch* using only four component search engines (Google, Yahoo!, Live Search, Ask Jeeves/Teoma)

since, for the present, our *QuadSearch* engine incoorporates only these engines. The consideration of more engines is straightforward though.

In *QuadSearch*, we treat all four component search engines equally. The reason we do this is due to the following observations: (i) all of them are considered by experts as “major” search engines, (ii) during their lifetimes they have been proved reliable and (iii) most users and metasearch engines prefer them.

The default rank aggregation method of *QuadSearch* is *rank-based*. Each returned ranked item is assigned a score based on the following formula:

$$ke = \frac{S}{n^m * (\frac{k}{10} + 1)^n} \quad (1)$$

where  $S$  is the sum of all rankings that the item has taken,  $n$  is the number of search engine top-k lists the item is listed in,  $m$  is the total number of search engines exploited,  $k$  is the total number of ranked items that *QuadSearch* uses from each search engine. We named this weight as  $ke$ . The less the  $ke$  value for an item, the larger the final rank this item will take is. For example, consider the following listings of two search engines (see Table 1) for a particular query:

Rank	SE1	SE2
1	$U_1$	$U_{11}$
2	$U_2$	$U_{12}$
3	$U_3$	$U_{13}$
4	$U_4$	$U_{14}$
5	$U_5$	$U_4$
6	$U_6$	$U_{15}$
7	$U_7$	$U_{16}$
8	$U_8$	$U_{17}$
9	$U_9$	$U_{18}$
10	$U_{10}$	$U_{10}$

**Table 1.** Results of two search engines for a particular query.

Let us elaborate a bit more on this table. Firstly, we presume that we deal with the top-10 lists ( $k = 10$ ) from each conventional search engine (SE: search engine). Also we name the URLs of each result as  $U_i$ , in order to demonstrate the overlapping URLs more easily. As we can see, there are two overlapping URLs in the above listings, the  $U_4$  which was ranked 4-th by *SE1* and 5-th by *SE2* and the  $U_{10}$  which was ranked 10-th by both search engines. All the others are found only in one of the two search engine top-10 lists.

In Table 2 we can see the ranking scores of  $ke$  and Borda Count methods for each URL. In this point we should mention a compact. We assume that when two URLs have the same score, then the URL that is in both top-10 lists will be ranked first, otherwise the URL of the first search engine will be ranked first.

Finally, in Table 3 we can see the final top-10 lists of the two methods.

In this point, we must stress some differences between Borda Count and the  $ke$  method.

URL	$ke$	$ke$ rank	$BC_{18}$	BC rank
$U_1$	0.5	1	18	3
$U_2$	1	4	17	5
$U_3$	1.5	7	16	7
$U_4$	0.5625	3	29	1
$U_5$	2.5	10	14	10
$U_6$	3	11	13	11
$U_7$	3.5	13	12	13
$U_8$	4	15	11	15
$U_9$	4.5	17	10	17
$U_{10}$	1.25	6	18	2
$U_{11}$	0.5	2	18	4
$U_{12}$	1	5	17	6
$U_{13}$	1.5	8	16	8
$U_{14}$	2	9	15	9
$U_{15}$	3	12	13	12
$U_{16}$	3.5	14	12	14
$U_{17}$	4	16	11	16
$U_{18}$	4.5	18	10	18

**Table 2.** Ranking scores of  $ke$  and Borda Count methods.

Rank	SE1	SE2	$ke$ result list (top-10)	BC result list (top-10)
1	$U_1$	$U_{11}$	$U_1$	$U_4$
2	$U_2$	$U_{12}$	$U_{11}$	$U_{10}$
3	$U_3$	$U_{13}$	$U_4$	$U_1$
4	$U_4$	$U_{14}$	$U_2$	$U_{11}$
5	$U_5$	$U_4$	$U_{12}$	$U_2$
6	$U_6$	$U_{15}$	$U_{10}$	$U_{12}$
7	$U_7$	$U_{16}$	$U_3$	$U_3$
8	$U_8$	$U_{17}$	$U_{13}$	$U_{13}$
9	$U_9$	$U_{18}$	$U_{14}$	$U_{14}$
10	$U_{10}$	$U_{10}$	$U_5$	$U_5$

**Table 3.** Final top-10 lists of  $ke$  and BC methods.

- The Borda Count method takes into consideration the total number of candidates, while  $ke$  takes into consideration the number of voters.
- Some Borda Count variations assign scores to each and every candidate; a candidate which is not included in the top-k list of a particular search engine takes a part of the remaining points. This does not hold for the  $ke$  method. In the  $ke$  method, a candidate will be assigned a score only when it is contained in the top-k list of a particular search engine, otherwise its score is zero.
- The  $ke$  method takes into consideration the total number of exploited search engines, the number of search engines where a candidate has been appeared and the size of the top-k list.

- The *ke* method has better “resolution”, in the sense that the possibility of two scores being the same is less than that of the Borda Count. For example, in Table 2 we can see that Borda Count assigned three URLs with the same score ( $U_1$ ,  $U_{10}$ ,  $U_{11}$  with score 18) while the *ke* method has given  $U_{10}$  a different score.
- The lower the *ke* weight an item has the higher will be ranked in the final result list. In Borda Count holds the opposite.

### 3.3 Antispam version of *ke* method

Informally, we say that a search engine has been spammed by a page in its result list when it ranks the page too highly with respect to the other pages, according to the view of a “typical” (average) user. This is unavoidable for search engines, because their ranking algorithms have “defects” that can be exploited by Web page developers in order to achieve an undeservedly high page rank. Thus, if a page spams all or even most of the search engines, then the metasearch engine could not defeat this problem as well, because the aggregation function would work with bad data.

In *QuadSearch* we gave to the users the option to use an antispam version of *ke* method. This method takes into consideration the *Condorcet Criteria* [13]. In the context of metasearching, these criteria tell us, in a few words, that an item which is enlisted in the top- $k$  lists of some search engines should be ranked above an item that is ranked in the top- $k$  lists of fewer search engines. The *QuadSearch* engine attempts to satisfy the intuition that if a page spams fewer than half of the search engines, then the majority of search engines will prefer a relatively good page to a spam page. The following pseudocode describes the antispam version:

1. Find which items appear in more than half pages (let the number of these items be  $c$ ).
2. Apply the *ke* method for these items.
3. Position them in *QuadSearch* result list, starting at rank 1.
4. Apply the *ke* method for the rest of the items.
5. Position them in *QuadSearch* result list, starting at rank  $1 + c$ .

The result list of *QuadSearch* is changed into that presented in Table 4. As we can see  $U_4$  and  $U_{10}$  went above  $U_1$  for the simple reason that they appeared in more than the half of the search engines, contrary to  $U_1$ , which appeared only in one. Also, in that point we should note the similarity of Antispam *ke* and Borda Count top-10 lists.

## 4 System Implementation

In the following subsections, we describe the technical issues regarding the implementation of the new metasearch engine.

Rank	SE1	SE2	$ke$ result list (top-10)	Antispam $ke$ result list (top-10)
1	$U_1$	$U_{11}$	$U_1$	$U_4$
2	$U_2$	$U_{12}$	$U_{11}$	$U_{10}$
3	$U_3$	$U_{13}$	$U_4$	$U_1$
4	$U_4$	$U_{14}$	$U_2$	$U_{11}$
5	$U_5$	$U_4$	$U_{12}$	$U_2$
6	$U_6$	$U_{15}$	$U_{10}$	$U_{12}$
7	$U_7$	$U_{16}$	$U_3$	$U_3$
8	$U_8$	$U_{17}$	$U_{13}$	$U_{13}$
9	$U_9$	$U_{18}$	$U_{14}$	$U_{14}$
10	$U_{10}$	$U_{10}$	$U_5$	$U_5$

Table 4. Antispam  $ke$  ranking.

#### 4.1 Architecture

The most significant modules of *QuadSearch* are the Quad Bot, the Object Builder, the Classification Module and the Presentation Module. These modules are described in the next subsections. A schematic diagram of the architecture is depicted in the left part of Figure 1.

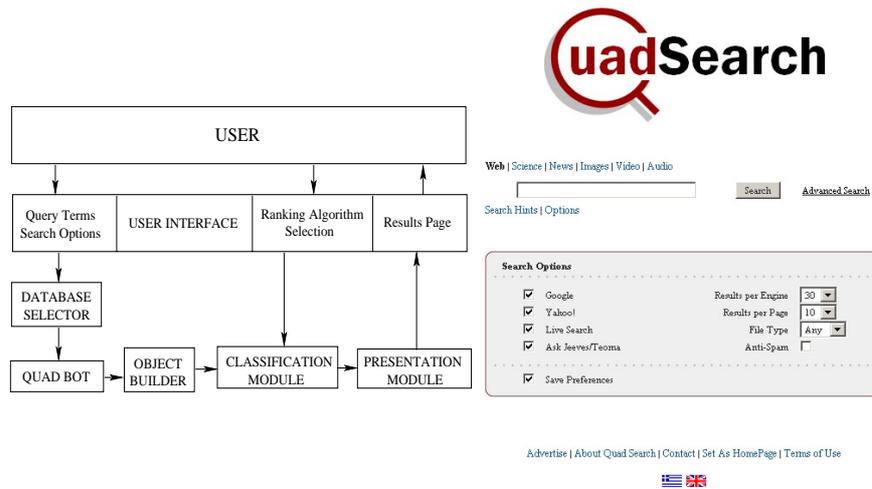


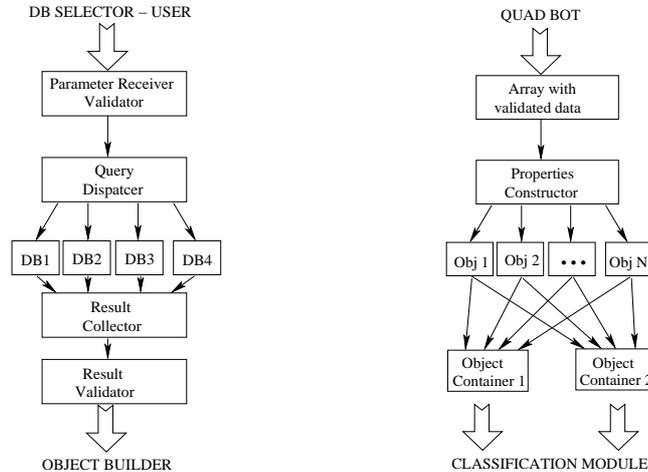
Fig. 1. (Left) Architecture of QuadSearch. (Right) Quad Search's homepage.

**User interface and database selector.** The static content of Quad Search's user interface is built with plain HTML, while the dynamic procedures are being processed by PHP. Additionally, we invoked Cascading Style Sheets (CSS) for page formatting. We decided that the Web pages' layout should be as simple as possible, in order to ensure: a) short download times, b) compatibility with all major browsers, and c) convenient usage.

For these reasons we avoided using large graphics files, or embedded objects like ActiveX Controls, or Flash presentations. We also rejected Javascript, because many experienced users tend to deactivate it, due to security reasons. Until now, Quad Search supports the classic Web search and also searching procedures for news, images, video and audio sources. The ability to search for scientific articles in the most popular scientific databases is also supported, but not very efficient. The user can switch among these features from either the home page or the results page. Regarding the database selector, the default search will be run using all four search engines. The user, however, has the option to choose which search engines will be exploited, as shown in the right part of Figure 1.

Apart from the classic text box, we included the most significant search preferences in the home page. The user can select the query resources (the search engines that will participate in the search process), the number of results to be retrieved per resource, the number of results that will be displayed per page etc. The interface provides the ability to store the values of these parameters, by setting cookies in the client's computer. Thus, the user is not obliged to define again these parameters for future queries. The interface includes an extra option to filter the results, to prevent spam records from entering the *ke* list. Finally, in the options page the user can select the ranking algorithm (*ke* or Borda Count).

**Quad bot.** The Quad Bot receives its inputs from both the database selector and the user interface. It is responsible for validating the input data and parameters, passing the query to the selected databases and collecting the results. Its internal structure is depicted in the left part of Figure 2.



**Fig. 2.** (Left) Quad Bot's structure. (Right) Object Builder's architecture.

*Parameter Receiver/Validator.* It accepts all the data coming from the database selector and the user. The validation process includes transformation of the

inputs in a way that can be sent to the search engines. For example, the procedure removes all leading and trailing spaces from the query string and replaces all spaces by the character “+”. The Result Validator also performs security checks to ensure that a “should-be” numeric parameter is really numeric, or that an attacking user does not send a malicious script instead of a query string. We developed this compartment by keeping in mind that all data coming from the Internet should be treated as suspicious. There are over a dozens of security checks that are performed for each parameter.

*Query Dispatcher.* The Query Dispatcher is the Quad Bot’s heart. It gets the validated data and creates http requests to the selected search engines. This is the slowest procedure of the whole system; its speed depends on the number of the invoked search engines, the requested results, the server’s Internet connection, etc. We have accelerated this procedure by submitting all the requests to the search engines simultaneously. To achieve that, we had to employ the *libcurl* library with cURL (client URL) extensions 7.16.0 for PHP 5.1, that support multiple connections at a time. By building the Query Dispatcher this way, we managed to shrink the idle time to no more than 1 or 2 seconds.

*Result Collector.* The Result Collector embraces the http responses transmitted by the search engines. Each involved search engine must respond to the Query Dispatcher’s request, by sending the source code of its result page. The source code is being filtered by using pattern matching techniques. The module retrieves the Rank, the URL, the Title and the Abstract for each candidate. When it receives all the information, it stores it in temporary arrays and sends them to the next module for validation.

*Result Validator.* The Result Validator is the most complex compartment of this module, as it performs multiple conversions to the collected data. The URL validation part is responsible for the appropriate formatting of the collected URLs, so that the overlapping candidates could be correctly detected later. At first, a UTF-8 decoding function converts all UTF-8 encoded characters to their ISO equivalent. For example, the %27 set of characters is being converted to an opening single apostrophe character ('). At next, a formatting function trims the trailing slash from a URL (if exists), a third procedure checks if an engine has returned two identical URLs etc. As already mentioned above, most of these conversions will be presented later on.

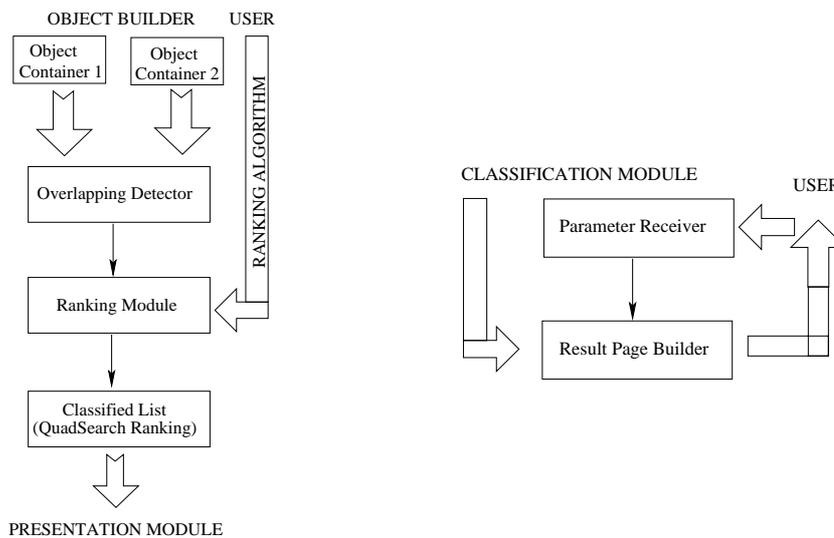
**Object Builder.** The Object Builder is a connecting bridge between the Quad Bot and the Classification module. We concluded that it is a good idea to treat our data as objects, so that we can use all the object-oriented programming features (such as inheritance, or support for multiple instances). This coding approach makes things very easy for the classification and presentation modules. In this section we describe how the collected results are being converted to objects. The Object Builder’s architecture is depicted in the right part of Figure 2.

*Array with validated data.* The Object Builder’s input is the array that the Quad Bot produces. It contains all the collected results that passed the Result Validator’s checks.

*Property Constructor.* This module implements a class that describes the properties of our objects. The properties that are being assigned to the objects are the URL, the Title, the Abstract and one to four Rankings (depending on the number of the selected search engines).

*Object Containers 1 and 2.* The first advantage of using objects in the system's implementation comes almost immediately, as we are provided with the ability to create multiple copies (not just references) of the objects. In this compartment, all the objects (the results) are being transferred to two new, identical object containers. The results enter the containers in groups. The first group consists of the results that the first search engine returns, the second group consists of the results that the second engine returns, etc. These containers will be the main tool in our effort to compare the search engine rankings and generate the final ranked list.

**Classification Module.** The Classification Module accepts the two result containers from the Object Builder and performs the result ranking according to the selected ranking algorithm. Its architecture is illustrated in the left part of Figure 3.



**Fig. 3.** (Left) Classification Module. (Right) Presentation Module.

*Overlapping Detector.* This section is responsible for detecting the overlapping candidates and for creating the final candidate list. It receives input from the two object containers and compares each object from the first container, to all objects from the second container. When the URL properties of two objects are identical, this object is marked as overlapping. Finally, the procedure constructs one container that holds all candidates, overlapping or not. The overlapping candidates appear only once in this container.

*Ranking Module.* The Ranking Module accepts the candidate container that the Overlapping Detector constructs, but it also receives the ranking algorithm that the user selected. The Ranking Module will apply the *ke* algorithm by default, unless the user selects another supported algorithm. Next, it computes the weight factors and/or the Borda Scores. Finally, it sorts the candidate list on ascending (for weight factors) or descending (for Borda Scores) order and passes the classified list to the Presentation Module.

**Presentation Module.** The task of this module is to construct the result page that will be presented to the user. In comparison to the other system compartments, this one has the simplest architecture. In the right part of Figure 3, we illustrate a schematic diagram of its internal structure.

*Parameter Receiver.* In this section, a set of parameters and user preferences are being transferred to the module. These preferences may either derive from user's direct selections, or from a previously stored cookie. The innovative element here is the view selector. QuadSearch is capable of displaying the results with the classic way that a search engine displays its results, but also can present the results by using the array view. The array view will present a matrix that shows only the titles of the candidates and the rankings they received from each search engine. This feature has been developed because it offers an easier way to compare the candidates.

*Result Page Builder.* The Result Page Builder is a HTML code production factory. It accepts the ranked result list and the user preferences and constructs the source code of the result page from the scratch. Finally, the page is displayed to the user through the user interface.

## 5 Innovative Features.

In this section, a quick walkthrough of Quad Search's innovative features is presented.

1. *Classic/Array View Switch.* This feature has already been mentioned earlier. The user is able to view the results in the classic way, but can also select the array view that provides an easier way of comparing the collected results.

2. *Related Searches.* Apart from the desired results, the Quad Bot is capable of grabbing almost everything from the results' pages that the exploited search engines transmit. In order to provide more specific results, the search engines prompt their users to submit the queries that they propose. The Quad Bot can fetch these query strings and present it to the result page through the Presentation Module.

3. *File Type Filter.* Many users tend to search the Web for specific file types (e.g., Adobe Acrobat or Microsoft Word files) and QuadSearch includes a similar feature. The user can select one of the most popular file extensions and perform a Web search. At this time, the QuadSearch engine supports searches for the following file formats: PDF, DOC, XLS, PS, RTF and PPT.

4. *Search for Scientific Articles.* QuadSearch supports searches for scientific articles in the richest scientific databases. Google Scholar is also included in these

The screenshot shows the QuadSearch interface. At the top left is the 'quadSearch' logo. To its right are navigation links: 'Web' (4), 'Science', 'News', 'Images', 'Video', and 'Audio'. Below these is a search bar containing 'electronic engineering' and a 'Search' button. To the right of the search bar is a settings menu with checkboxes for 'Google', 'Yahoo!', 'Live Search', and 'Ask/Teoma'. It also includes dropdown menus for 'Results per Engine' (30), 'Results per Page' (20), and 'File Extension' (Any), along with an 'Anti-Span' checkbox. Below the search bar is a status bar indicating 'Displaying 1-20 of 90 results for query: electronic engineering (6.11 sec.)' and a 'Switch to array view' link (1). The main content area is titled '(Rankings based on KE Algorithm)' and lists three search results (5). Each result includes a title, a brief description, and a URL. To the right of the results is a 'Related Searches' section (2) with several links related to the query.

Fig. 4. The results' page with its innovative features.

databases. This type of search can be accessed from the “Science” link and will return papers, technical reports and books approved by the scientific community related to the query terms.

5. *Query String Explosion Feature.* This feature (see Figure 4) splits the query string to its search terms and gives the user the ability to perform ‘single term’ searches. For example, the query string ‘electronic engineering’ is being split to the terms ‘electronic’ and ‘engineering’. By clicking on any of these words QuadSearch will perform a Web search.

6. *Ranking Algorithm Selector.* This feature (see Figure 5) is only accessible from the options page and provides the user with the facility to determine how the collected results will be ranked, by employing one (or more) of the supported algorithms. At this time, QuadSearch supports our *ke* Algorithm and the Borda Count method. It also provides a third option that utilizes both algorithms and presents the results in array view (comparison mode). It is in our intentions to include more ranking algorithms in the system (e.g., Markov Chains).

7. *Engine Bombing Protection.* When various search resources are being exploited, a possibility that many similar results will enter the result’s list always exists. This phenomenon is called *engine bombing*. For example, it is not very informative and useful for a user to submit a query and receive five or more results from the same domain in the top, say, twenty listing. Thus, we developed a feature (which can be enabled or disabled) to prevent multiple results coming from the same domain to enter into the result list; alternatively the user can select the maximum number of such results.

**Results Filtering**

---

**Anti Spam Filter**

Enable Anti Spam Filtering

**Engine Bombing Protection**

Engine Bombing occurs when multiple results from the same domain enter the results list.

Allow  results from the same domain to enter the results list.

**Results Ranking and Presentation**

---

**Ranking Algorithm**

Determine how the collected results will be ranked, by employing one (or more) of the supported algorithms.

KE Algorithm (KEA)

Unweighted Borda Count (UBC).

KEA vs UBC (Comparison Mode - Experimental).

**Fig. 5.** Part of the options' page is where the Ranking Algorithm Selector and the Engine Bombing Protection lay.

## 6 Concluding remarks and future work

In this article, we considered the issue of developing a new metasearch engine to assist in the process of Web information retrieval. The main motivation to develop this novel metasearch engine was the common intuition that a rank aggregation algorithm should a) be related to the comparison of the top-k lists of each conventional search engine, and b) deal with the problem of the spam into metasearch result lists. Thus, we came up with a pair of new methods for rank aggregation, i.e., the *ke* method and its antispam version. We injected some new parameters, like the number of the top-k lists that a page appears, the total number of exploited search engines and the size of the top-k lists. The best way to experiment and test these two new methods was to develop a new metasearch engine, named *QuadSearch*, a name related to the current number of exploited engines. The new metasearch engine is publicly available at <http://delab.csd.auth.gr/~lakritid/metasearch/>.

For the near future, we are going to implement anonymous personalization techniques, further result filtering, and thorough search hints. Furthermore, we are orientated towards making *QuadSearch* a scientific tool by implementing most of the rank aggregation methods and giving the user the opportunity to choose between them as well as to provide more statistics and a user grading system of result lists.

## References

1. The best and most popular meta search engines, 2006. Retrieved on December 31st from [www.searchenginewatch.com/showPage?html.page=2160791](http://www.searchenginewatch.com/showPage?html.page=2160791).

2. The big four meta search engines, 2006. Retrieved on December 31st from [www.searchenginewatch.com/showPage?html.page=2160781#scene\\_1](http://www.searchenginewatch.com/showPage?html.page=2160781#scene_1).
3. A meta search engine roundup, 2006. Retrieved on December 31st from [www.searchenginewatch.com/showPage?html.page=2160801](http://www.searchenginewatch.com/showPage?html.page=2160801).
4. Meta search engines are back, 2006. Retrieved on December 31st from [www.searchenginewatch.com/showPage?html.page=3109441](http://www.searchenginewatch.com/showPage?html.page=3109441).
5. Meta search or meta ads?, 2006. Retrieved on December 31st from [www.searchenginewatch.com/showPage?html.page=2163821](http://www.searchenginewatch.com/showPage?html.page=2163821).
6. Metacrawlers and metasearch engines, 2006. Retrieved on December 31st from [www.searchenginewatch.com/showPage?html.page=2156241](http://www.searchenginewatch.com/showPage?html.page=2156241).
7. Vivisimo clustering engine, 2006. Retrieved on December 31st from <http://vivisimo.com/>.
8. R. Baeza-Yates and B. Ribeiro Neto. *Modern Information Retrieval*. Addison-Wesley and ACM Press, 1999.
9. C. Dwork, R. Kumar, M. Naor, and D. Sivakumar. Rank aggregation methods for the Web. In *Proceedings of the ACM International Conference on World Wide Web (WWW)*, pages 613–622, 2001.
10. R. Fagin, R. Kumar, M. Mahdian, D. Sivakumar, and E. Vee. Comparing partial rankings. In *Proceedings of the ACM International Symposium on Principles Of Database Systems (PODS)*, pages 47–58, 2004.
11. R. Fagin, R. Kumar, M. Mahdian, D. Sivakumar, and E. Vee. Comparing partial rankings. *SIAM Journal on Discrete Mathematics*, 20(3):628–648, 2006.
12. R. Fagin, R. Kumar, and D. Sivakumar. Comparing top  $k$  lists. *SIAM Journal on Discrete Mathematics*, 17(1):134–160, 2003.
13. N. Francis. Voting as a method for rank aggregation and spam reduction on the Web. Undergraduate senior thesis (CPSC 490). Department of Computer Science, Yale University, May 9th, 2005.
14. A. Gulli and A. Signorini. Building an open source meta-search engine. In *Proceedings of the ACM International Conference on World Wide Web (WWW)*, pages 1004–1005, 2005.
15. M. Levene. *An Introduction to Search Engines and Web Navigation*. Addison-Wesley, 2006.
16. Y. Lu, W. Meng, L. Shu, C. Yu, and K.-L. Liu. Evaluation of result merging strategies for metasearch engines. In *Proceedings of the IEEE International Conference on Web Information Systems Engineering (WISE)*, pages 53–66, 2005.
17. W. Meng, C. Yu, and K.-L. Liu. Building efficient and effective metasearch engines. *ACM Computing Surveys*, 34(1):48–89, 2002.
18. L. Page, S. Brin, R. Motwani, and T. Winograd. The PageRank citation ranking: Bringing order to the Web. Technical Report TR-1999-66, Stanford University, 1999.
19. M. E. Renda and U. Straccia. Web metasearch : Rank vs score based rank aggregation methods. In *Proceedings of the 18th ACM International Symposium on Applied Computing (SAC)*, pages 841–846, 2003.
20. S. Soudatos, T. Dalamagas, and T. Sellis. Sailing the Web with Captain Nemo: A personalized metasearch engine. In *Proceedings of the ICML workshop: Learning in Web Search (LWS)*, Bonn, Germany, 2005.
21. A. Sugiura and O. Etzioni. Query routing for Web search engines: Architecture and experiments. *Computer Networks*, 33(1–6):417–429, 2000.
22. Z. Wu, W. Meng, and Z. Yu, C. Li. Towards a highly-scalable and effective metasearch engine. In *Proceedings of the ACM International Conference on World Wide Web (WWW)*, pages 386–395, 2001.