# QueReSeek: Community-Based Web Navigation by Reverse Lookup of Search History

Hideyuki Tan
*Alpha Systems Inc.*
*tanh@alpha.co.jp*

Ikki Ohmukai
*National Institute of Informatics*
*i2k@nii.ac.jp*

Hideaki Takeda
*National Institute of Informatics*
*takeda@nii.ac.jp*

## Abstract

*In this paper, we propose a system called QueReSeek that realizes Web navigation by using search queries in a community. Web navigation is realized as follows: when a user browsing some Web content, if the Web content is included in the list of results of past search by people in the community, query strings used in the search are shown to the user. To realize this navigation, the system collects queries to search engines and their results, and builds the search query-URL index. It shows relevant queries from the URL of Web content which is browsed by users based on this index. By looking up this database reversely, it can show related query strings to Web contents. Since the search queries in the community are keywords related to information and knowledge of interest within the community, this navigation reflects implicit knowledge in the community. It is useful especially for community members who are not proficient in search. Such users can learn search expertise by following search strings provided by the system. We implemented this proposed method in two ways. We could display relevant queries for approximately 20% of the browsed Web content in this experiment.*

## 1. Introduction

Web browsing becomes an indispensable activity in our life and it is deeply supported by Search for the Internet. The rise of reliable search services substantially contributed to the spread of Web use, and users use search engines in their daily lives as tools in finding information and knowledge in the Internet. However, finding information with search engines requires expertise. A proficient user in search can reach the desired information with one or a few search actions. On the other hand, non-proficient user struggles to reach the desired information; she repeats search actions with different queries frequently. Even she sometimes gives it up. Their needs a method to share expertise of search.

In this paper, we propose QueReSeek, the method for sharing search behavior to Web content that is highly influential to the community. The system accumulates the search query-Web content relationship within the community, and shows query strings to individual Web content as keywords of interest in the community.

## 2. The Concept of proposed method

The key idea of QueReSeek is exploitation of restoring search queries collected within the community. Before going to the detail of the method of QueReSeek, we explain our background assumptions about search and browsing in communities, i.e., information gathering behavior, the relationship between search queries.

### 2.1. Information gathering behavior

A person's behavior in gathering information from many information sources consists of eight types of behavioral patterns, such as the start of the search, chained search, differentiating of the information source, browsing, monitoring of information, information extraction, verification of validity, and the close of the search [1]. This information search process model can be applied to search behavior using the Internet such as indefinite browsing for finding information without specific directions and chained searching for finding other relevant information based on the information.

A searcher enters a search query in the input box of the search engine in order to reach the goal of acquiring desired information. By reading the title of the content, the URL domain name, and snippets in the results view, she decides the next action that should be taken as follows. She may repeat the roaming behavior through Web content links if she cannot identify what she wants precisely. This roaming behavior as seeking

information and knowledge will not last forever because her knowledge changes along with browsing content while roaming. Knowledge after obtaining information is not a mere addition of knowledge before obtaining the information, but rather a dynamical change of the original knowledge through obtaining the information. By increasing the knowledge relevant to unknown matters, the search queries that she generates become closer to the desirable information. In short, the searcher gradually gets closer to the desirable Web content by roaming. When the searcher finally reaches the content with a description of obtained content and matters that the searcher wanted to know, the description of content obtained, matters that the searcher wanted to know, and the generated search queries will be firmly connected. In the course of approaching the desirable knowledge gradually, the query that the searcher generates functions as refining filter for the searcher's desired knowledge. In short, it represents the searcher's formalized search know-how. The proposed method utilizes the searcher's formalized search know-how.

## 2.2. The Relationship between search queries and search results content

The search queries by the user to the search engine are generated based on the user's background knowledge, which is a string related to the desired Web content. Then, the search engine executes a full-text search on the collected Web content set, determines the order by calculating scores with some ranking algorithm, and, finally, returns a list of URLs of the Web contents to the user. The search query is a fragment of the Web content that obtained from the full-text search. And the fragment is connected to Web content through a ranking algorithm in the search engine. When looking at this relationship in reverse, the search query is a summary of Web content. The proposed method adopts this idea.

In this paper, we define "community" as a group of search engine users with a common purpose, such as people in a similar environment; for example, it is a group of a unit in a company, a school, or a research department. Accordingly, the search query issued by users who consist of a community is assumed to be potentially related to the common purpose of the community. In short, the search queries generated in the community potentially share the purpose of the community and related to Web content which people in the community may be interested in. The results returned from the search engine are usually links to multiple Web contents. It means that a keyword represented as a query is related to multiple Web

contents. Therefore, by threading out the relation of the search query-Web content reversely, we can collect pairs of Web content and keyword that would represent community interest. Web content that a user in the community is browsing can be abstracted as keywords of interest that are queries within the community. In other words, we provide links from Web contents to keywords that people in the community may be interested in.

## 2.3. Support for individual Web searches and browsing assisted by the community

QueReSeek shows search queries relevant to Web content during browsing by reverse lookup from the search history of the community.
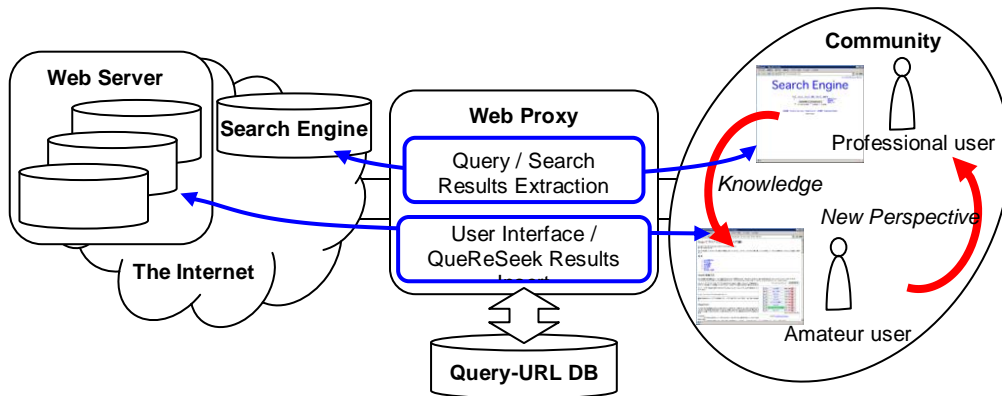
Suppose that a user browses Web content. If the Web content is included in the list of results of past search by people in the community, query strings used in the search are represented to the user. If she recognizes that some of these strings are useful for her search purpose, she may use these strings to obtain better ranked Web contents. Thus she can change browsing content to associative Web content substantially closer by finding a new search query.

While the user browsing Web content, QueReSeek presents her other's search queries, and promotes Web browsing navigation with the community as a base.

By sharing search query expertise in a community, any user who is not good at performing searches can obtain information as clues in the issuance of a search query. For a proficient searcher, it provides different point-of-views with which she can understand what other people in the community are thinking. The proposed method provides a new way of browsing, i.e., browsing with the context of the community. When browsing a Web content, a user can be aware how other people are interested in it.

Any profit-making organization, such as a company, is a group in which a member's role and function are specialized and integrated to achieve a certain common purpose. Each member seeks Web content on the Internet with purpose which reflects her role in the organization. Therefore aggregating search queries is to constituting knowledge that the organization needs to acquire. From this point of view, we expect that the more the community specializes in some specialized area, the more this proposed method will have an effect on knowledge sharing.

To summarize, QueReSeek facilitates the user in finding new search queries through browsing behavior and provides a new starting point for searches. Because the information search skills of others can be used and the proposed search queries reflect the interest and

**Figure 1.** Web Proxy type QueReSeek Implementation Structure and Utilization Image

concern of the community, it is expected to reduce contextual differences. The important point is that QueReSeek does not require any extra information to users. It just intercepts search queries and re-uses them for other people.

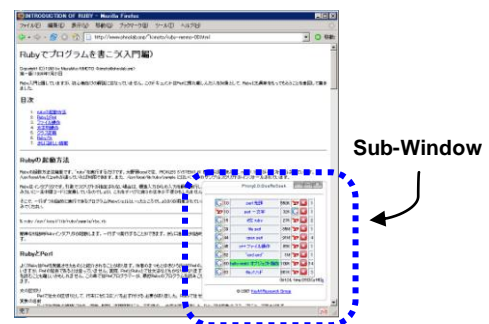## 3. The Implementation of QueReSeek

To realize a system that collects content search history of users belonging to the community and returns their processing results to the community, three function modules are needed. This chapter describes about the processing of three function modules which are needed to realize QueReSeek.

### 3.1. Search query collection from Web browsing history

QueReSeek shares search query that each user turned on to the search engine in the community. A search result from the search engine is one type of Web content which can be displayed in a URL. In a word, retrieval query is included in the URL strings. Therefore, if the Web browsing history of the user who first belongs to the community can be collected by some methods, retrieval query can be obtained. Search engine sites for such objects include Google[1], Yahoo![2], Live Search[3], goo[4], Excite[5], and Infoseek[6].

### 3.2. Store collected search queries and results

The search queries obtained from the web browsing history shows the knowledge for which the community wanted. By using these search queries, this module

---

[1] http://www.google.{co.jp, com}
[2] http://www.yahoo.{co.jp, com}
[3] http://www.live.com
[4] http://www.goo.ne.jp
[5] http://www.excite.co.jp
[6] http://www.infoseekc.co.jp



**Figure 2.** Web browser screen presenting a search query of Web Proxy type QueReSeek

gets search results with Google and Yahoo! again, and obtains URLs, rankings, and estimated total hits to store them in the "Query-URL DB." The search query collection module stores the search results in the DB within 15 seconds in quasi-real-time after detecting search engine use.

### 3.3. Present search queries related to the browsing Web content

QueReSeek presents search queries related to the Web content that each user browsed in the community. For this, it only has to make URL that indicates the user browsing Web content a key and to refer to the "Query-URL DB." By referring to DB, this module can get search queries that are related to user browsing Web content, and that was used by community. These obtained search queries are sorted in order of rank in the search engine, the number of hits, and the number of times information that was input by the community.

## 4. The type of Implementation

There are three kinds of implementation, at content provider side, at a transfer point between the content provider and the user, and at a user side, are thought by the method of arranging those function modules. We
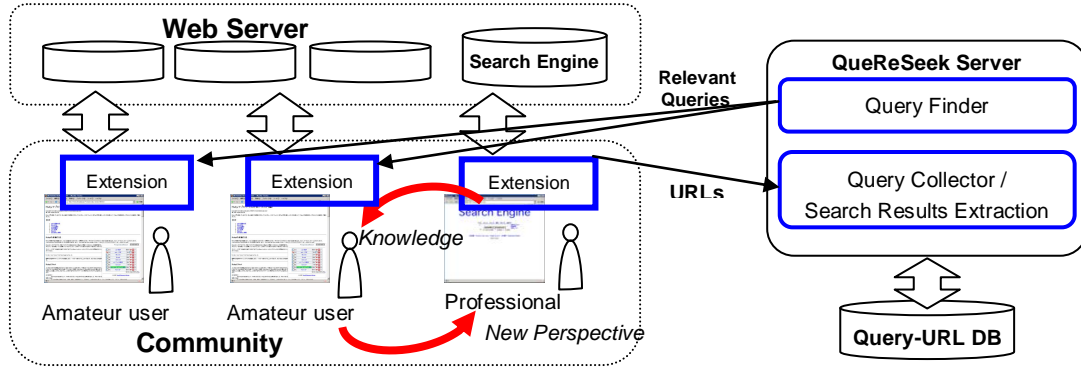
**Figure 3.** Client-Server type QueReSeek Implementation Structure and Utilization Image

mounted two types of implementations: Web Proxy type and Client-Server type.

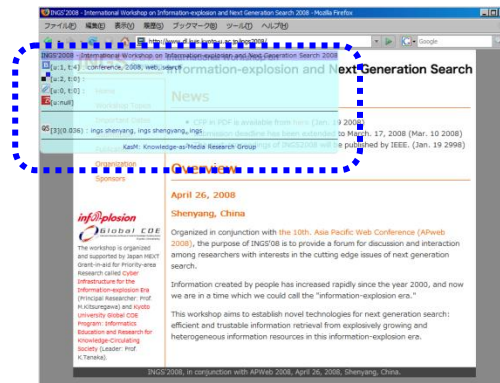## 4.1. Web Proxy type implementation

The Web proxy is located in network topology between a Web server and a Web browser, and it is introduced in organizations such as companies. Here, the Web browsing history of the community can be collected. In this case, the community becomes the group of the user who is using this proxy.

Fig. 1 shows the structure and image of the utilization of a system with Web Proxy type. The three functions are mounted in Web Proxy. The search query display module inserts JavaScript in Web content to generate an interface to display the results of QueReSeek. At this time, the system refers to "Query-URL DB" with URL of browsing Web contents, and takes out search query that the community used, and presents it to the user. While browsing Web content, the Web browser where a relevant search query exists, is shown as a screen in Fig. 2, and an interface is shown as a sub-window in the bottom right corner of the content display area. In the sub-window, the search query is displayed in the table with the icon representing a search engine, the character strings of the search query, the estimated total hits, and the number of uses as a search query in the community. Top 10 is displayed among the search query obtained by processing described in 3.3. These search query strings are added links to search results with a large number of hits on Google or Yahoo!.

The advantages of Web Proxy type include realization of the proposed method without obtaining a prevailed Web server and Web browser, and without distributing a new application to users.

## 4.2. Client-Server type implementation

If a client can collect URLs of Web content that a user is browsing, the client may collect the search



**Figure 4.** Pop-up window of Client-Server type

query input in the search engine by the user, and the search query related to Web content may be returned to the client. Fig. 3 shows the system configuration implemented in the Client-Server type. The client consists of the portion that sends URLs of Web content while browsing to a dedicated server and the portion that receives the search query relevant to Web content and shows it to the user. The server is implemented with a function to analyze URL from the client of browsing Web content and collect the search query in the same way as the Web Proxy type. Moreover, the function to send the search queries input by the community through searching the "Query-URL DB" from URLs of Web content while browsing to the client is mounted on the server. The client utilized Greasemonkey[7], an add-on in the Firefox Web browser, and carried out the process with user script. By introducing this add-on and user script, the icon is always displayed on the left top of the content display area while browsing Web content. And if the search query related to Web content in browsing exists, the color of the icon changes and notifies the user of the presence of relevant information. The pop-up window, Fig. 4, posting information relevant to the browsing content is displayed by putting a mouse pointer over

---

[7] https://addons.mozilla.org/en-US/firefox/addon/748

this icon. In the pop-up window, up to 100 search queries got from server are lined up and displayed. Just as the Web Proxy type, links to search results are added to search query strings.

For the use of Client-Server type, the user must install add-ons and user script. However, in this type, the server can be set up in outside network. Therefore, the user can use this service by exceeding the wall of differences between the Internet and Intranet.

## 5. Evaluation

Here, we would like to speak about a brief evaluation performed on two systems implemented for verification of the effectiveness of QueReSeek.

### 5.1. Evaluation of Web Proxy type

At the initial stage of implemented system operation, the "Query-URL DB" is empty. If the community has already used the Web Proxy, the "Query-URL DB" can be prepared even in the initial stages of system operation by using its log. In this research, we measured how much we could exhibit the search query previously extracted from the log for Web content in browsing, and observed the guide for Web browsing going through the search query exhibited.

**5.1.1. Experiment Preparations.** We extracted the search query input into the search engine of the log of approximately 1,500 users among the organization that we belong to who browsed the Web for a month. And we configured "Query-URL DB" where about 84,000 search queries and about 8 million URLs are registered. For details of this process, we summarized the number of search queries, total URLs and unique URLs in Table 1.

**5.1.2. Experimental Conditions.** We invited people within the department that we belong to cooperate, and then asked them to perform Web browsing during business hours, including break times as usual after giving a brief explanation of the functions provided by the system. The measurement period lasted for five days and the number of users was twenty in total.

**5.1.3. Results of the experiment.** As results obtained, we summarized the amount of Web content in browsing, Web content that displayed the search query, the search query added from the start of measurement, and numbers of tracks of links added to the search query exhibited in Table 2. Web content displayed could display the relevant search query at the rate of approximately 33% for the content browsed. While

**Table 1.** Details of the prepared "Query-URL DB" for evaluation of Web Proxy type

| Search Queries | 84,135 |
|---|---|
| Total URLs | 11,366,183 |
| Unique URLs | 8,106,381 |

**Table 2.** Ratio of search queries which could be displayed by Web Proxy type

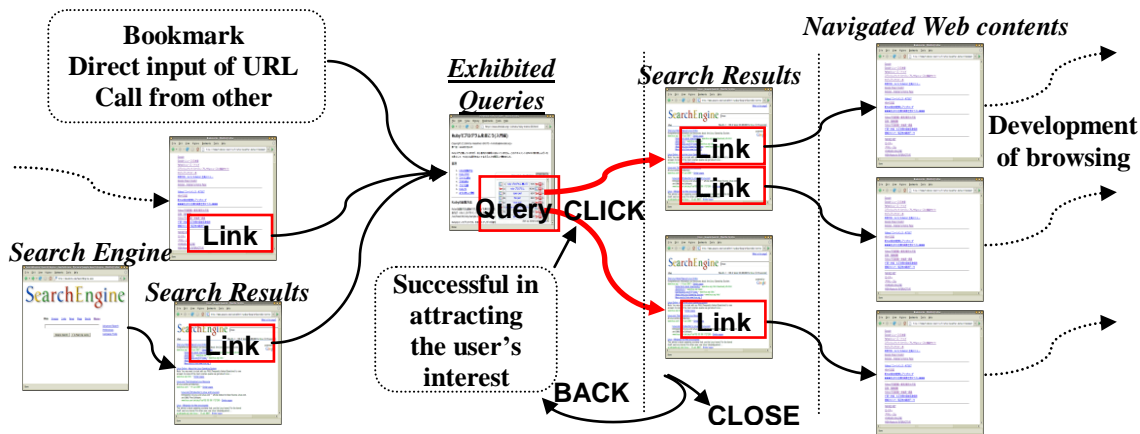| Browsed contents | 6,047 | |
|---|---|---|
| Displayed the search query contents | 2,019 | 33.4% |
| Added search queries | 593 | |
| Tracks of links | 49 | |

subjects were expecting to browse the Web content with a high upgrade frequency, including news and weblogs, more than 30% of Web content could display the search queries, which indicate about 8 million unique URLs in the "Query-URL DB," may be more than enough for around 6,000 Web pages that were browsed. However, although it is browsed after selected from many Web contents, at least one display of information related to the content was available during three sessions of Web content browsing. Therefore, this fact makes us think that background knowledge contained in the community and the domain of information and knowledge required by the community are overlapping.

The number of tracks of links added to the search query exhibited was 49. Fig. 5 shows the browsing transition of Web content by QueReSeek. If a user clicks the link of the search query displayed, the display of such search query was successful in attracting the user's interest. In this experiment, displays successful in attracting the user's interest were for approximately 2% of Web content browsed. Among them, about 15% of search queries clicked had site names, including Slashdot [8], which was a "Navigational Query" [2] used by a user with a predetermined destination. About 30% of the search query whose links were tracked was the ones added after starting the experiment. From this, we can assume that a new concept is introduced in the community, and the user is trying to find new content based on his/her behavior.

In the observed browsing transition, it is observed that the user searched "iptables [9]" from the search window of the Web browser, then browsed the search results of the search query, such as "Firewall," "Resetting" and "iptables packet" displayed in the Web

---

[8] http://slashdot.jp
[9] This is an administration tool for IPv4 packet filtering on Linux

**Figure 5.** Browsing transition with QueReSeek

content reached based on the search results of "iptables". This shows that it could display the search queries that are substantially closer to the search query that the user voluntarily generated. However, as the measurement was conducted in a short period of time, the number of tracks of links from the search query exhibited was few, and sufficient verification is required for the effectiveness of the system. Considering the opinion from the subject that browsing is interfered because the sub-window for showing relevant information hides Web content, a reexamination of the user interface is also necessary.

## 5.2. Evaluation of Client-Server type

Same as the evaluation for Web Proxy type, we measured how much the search query could be exhibited to users for Web content in browsing as an index to measure its usability. In this time, we only treated the search query added whenever detecting the subject's search engines used as they were.

**5.2.1. Experiment Conditions.** We invited people to cooperate in the experiment from the research room and department that we belong to, and asked them to perform Web browsing as normal after giving a brief explanation on functions provided by the system. We also told the people cooperating in the experiment about a method to stop the function of the user script as a measure for cases when they do not want to provide notification of Web content URLs while browsing. Data measurement took ten days and the number of users was twenty-four in total.

**5.2.2. Results of the experiment.** We summarized the obtained results such as the amount of Web content browsed, Web content that exhibited search queries, the search queries added from the start of measurement, added numbers of URLs and unique URLs in Table 3.

**Table 3.** Ratio of search queries which could be displayed by the Client-Server type and details of "Query-URL DB"

| Browsed contents | 7,153 | |
| --- | --- | --- |
| Displayed the search query contents | 1,464 | 20.5% |
| Added search queries | 788 | |
| Tracks of links | 15 | |
| Added total URLs | 133,657 | |
| Added unique URLs | 109,098 | |

Web content that displayed search queries could exhibit the relevant search queries at the rate of 20% for the content browsed.

The results show that the number of search queries input into the search engine was about 800 and about 110,000 unique URLs were stored in the "Query-URL DB." This is 1.4% of unique URLs prepared in the evaluation of the Web Proxy type, a very small number. However, the fact that the search queries were exhibited in about 20% of Web content browsed shows, just as the Web Proxy type, that participants browsed Web content guided from the results of the search engine.

The ratio of Web contents that were able to present search query is compared with the community that applies by evaluating the Web Proxy type implementation by measuring Web browsing. As a result, including break times, and defining the users belonging to a department of a company as a community, there seemed to be a higher possibility that participants were browsing Web content mainly required for their jobs. On the other hand, as the Client-Server type provided a community mainly consisting of members of the research room of the university, connection via a common purpose was weak, areas of interest were broad. So, the Web content that is the domain of each member's interest,

browsed were hard to overlap, and the ratio of the Web content number displayed seemed to be reduced.

# 6. Related Work

The area that related to this work is Web content annotation, query recommendation, and query log analysis, knowledge sharing. In this paper we have identified three important topics of related work.

## 6.1. Recommendation of search query

In the context when search query is input, there are services to recommend search query [3][4]. These detect the event of the keyboard generated in the search query input box. And, the recommended candidate incremental retrieved whenever one character is input is presented to the user. This is the support function for inputting search queries by the extension of predictive transform input method that is realized on the Web browser. It generally improves usability in searching Web content related to a popular theme. As this proposal method exhibits search query in Web content browsing all the time, search query display timing for the user is not the same.

In relation to the retrieval support, there is an attempt to personalize and to do the search result by using the search history of the community [5]. This pays attention from the retrieval result to the part where previous contents have been selected referring to the title and snippet.

## 6.2. Exhibition of information related to Web content in browsing

There are some attempts to exhibit relevant information based on Web content URLs in browsing. For example, [6] is implemented as a functionality extension plug-in, and it displays the evaluation rate of the index within the search engine related to Web content browsing. [7] gives the user information in the surrounding generated from the link relation of the Web page for guide Web browsing.

There is also an application that uses the social bookmark (SBM) service as a source of exhibition information related to Web content. As the proposed method exhibits the search queries input from the community as relevant information to the user, the source of the handled information is different. SBM service is a service that releases and shares bookmarks with an unspecified majority of users, and users may add words for classification called free description tags for content when registering Web content [8]. [9] displays the number of bookmarked sites and added tags in relation to Web content while browsing. As these exhibited tags are linked to tag-added Web contents list page that is provided by the SBM service, it can be called guidance for Web browsing through SBM. Moreover, [10] is an attempt to improve the search area based on the tag given to contents by SBM.

## 6.3. The Relationship between search queries and Web page

There is an attempt to a typical query of Web pages from query logs [11]. Identical to our way of thinking about the proposed method, candidate search queries related to Web content are extracted from query logs by browsing transition from the search result using the relationship between search queries and Web pages. This method does not mention the community as an aggregation of users, and although it may allow us to extract typical and well-understood search queries, it does not intend to extract valuable queries for the community sharing a purpose. The proposed method is a method for sharing search queries as background knowledge, assuming the users gathered under the same purpose to be in a community.

[12] is a method for linking the world view of the search engine user community with that of the search engine itself. It is a method of accumulating the process of query reformulation in the search session, and giving weight to the relation between contents and search query. In our approach, the ranking of the search result and the use frequency in the community are used as it is without new relation weight.

# 7. Summary and future issues

In this research, we proposed QueReSeek as a reverse lookup engine that allows us to share search query input in the search engine within the community, display the search queries in Web content browsed by the users as relevant information, and support Web browsing and searches. We implemented the system to realize the above, and discussed one evaluation about its usability.

The proposed method may provide the user with links to undiscovered and unknown Web content that is relevant to the Web content that are being browsed. Web networks are growing on a daily basis, and it is hard to cover everything with one search engine in present conditions. This method aggregates the results of multiple search engines. Therefore, this system makes it possible for the user using only a specific search engine to increase chances to discover new Web content. Also, this system only uses information obtainable from normal search behavior and does not

require the first-time user of this system to do different things other than Web browsing, including the past use of search engines. A system with the purpose of knowledge sharing is sometimes designed to make the user, who obtained knowledge, forcibly output the knowledge, and ends in losing users in many cases. On the other hand, this proposed method only requires the user to input the search query into the window, which is a normal search behavior, to make the user's search skills explicit knowledge for reuse.

In this proposed method, the more the user uses this, the more the entry of DB increases naturally. Considering that periodical crawling by robots results in keeping the freshness of Web content exhibited in the search engine site, the "Query-URL DB" is also required to update periodically. There are search queries with search results with zero hits and character string input with no meaning due to wrong input. However, since there is a connection to a certain extent in relations between Web content and the search queries and because knowledge domains can be narrowed down, there may be a small possibility that the increase in the amount of data is harmful to the obtained results. But, as large scale implementation causes slow DB response, we have to exercise ingenuity related to implementation. Web content, such as top pages of portal sites had a large amount of presence in search queries while browsing, so we have to review the method for ranking search queries to be exhibited. In this proposed method, we utilized the fact that the search queries are strongly connected with the matters that a searcher wants to know, but conversely this may become a privacy problem. For example, when referring to a patent database, the searcher's object of interest is aggregated to the search query, so we need appropriate access control.

The proposed method is considered to be one of the knowledge sharing tools to raise the level of information and knowledge held by entire aggregation of people with the same purpose. In that respect, I could not fully evaluate from the view of effective Web navigation. There are issues such as the selection method of search queries to be exhibited when plural search queries related to Web content are found and the user interface which displays the search queries. We would like to reexamine the method to display search queries hereafter and to proceed with effective evaluation and verification of the proposed method from now on.

## 8. Reference

[1] D. Ellis, "A behavioral approach to information retrieval design", *Journal of Documentation*, Vol. 45, No. 3, 1989, pp. 171-212.

[2] A. Broder, "A taxonomy of web search", *ACM SIGIR Forum*, Vol. 36, No. 2, ACM, 2002, pp. 3-10.

[3] Google, Inc. Google suggest Web Page, http://www.google.com/webhp?complete=1, 2004.

[4] NTT Resonant Inc. and JustSystems, Inc. goo suggest β with ATOK Web Page, http://suggest.search.goo.ne.jp/suggest/, 2005.

[5] O. Boydell, and B. Smyth, "Capturing community search expertise for personalized web search using snippet-indexes", *In Proceedings of the 15th ACM international Conference on information and Knowledge Management*, CIKM'06. ACM, 2006, New York, pp. 277-286.

[6] Google, Inc., Google toolbar Web Page, http://toolbar.google.com, 2000.

[7] S. Ikeda, K. Zettsu, S. Oyama, and K. Tanaka. "Supporting Web Navigation By Circumference Information", *The Database society of Japan Letters*, Vol. 2, No. 1, 2003, pp. 135-138.

[8] L. Daminos, J. Griffith, and D. Cuomo, "Onomi: Social Booking on a Corporate Intranet", Collaborative Web Tagging Workshop. WWW2006, Edinburgh, 2006.

[9] Glucose, Inc., Glucose nano Web Page, http://dev.glucose.jp/wiki/index.php/Glucose_nano, 2007.

[9] D. R. Millen, M. Yang, S. Whittaker, and J. Feinberg, "Social bookmarking and exploratory search", ECSCW 2007, Limerick, Ireland, 2007.

[10] Y. Kabutoya, T. Tumoto, S. Oyama, and K. Tanaka "Extracting Typical Queries for Web Pages Using Query Log and Its Applications", Symposium on Data Base and Web Information System (DBWeb) 2007, 2007.

[11] E. Amitay, A. Darlow, D. Konopnicki, and U. Weiss, "Queries as anchors: selection by association", *In Proceedings of the Sixteenth ACM Conference on Hypertext and Hypermedia*, HYPERTEXT '05. ACM, 2006, New York, pp. 193-201.