



# **Secured Personalized Web Search using Ontology and Hashing Technique**

Santosh Sharma Kongbrailatpam, Dr. R.J. Anandhi

M.Tech Student, Dept of CSE, The Oxford College of Engineering, Bangalore, India

Vice-Principal, Professor & HOD (CSE), The Oxford College of Engineering, Bangalore, India

**ABSTRACT:** Personalized Web Search is a best way to improve search quality by customizing search results for the users with individual information goals. This is achieved through implicitly collection of user profiles, browsing history, clicked through data, bookmarks, location of the user. However, users are not willing to expose their private preference information to search engines. Many research are undergoing in improving Personalized Web Search (PWS), this system added another dimension that is Privacy Protection in PWS. This system has a framework by name personalized privacy preserving search framework, which is very helpful framework where eavesdroppers are not able to get any details of the users since privacy measures are taking place at client side. To achieve privacy protection, this system provides metadata and user's query encryption which is achieved using MD-5 hashing technique. In the beginning it will look like a keyword search but there is a behavioral observing system called the Spy-NB, which monitors the users' behavior and which will provide the results sets according to the users' interest in efficient way. To improve the personalized search, this system uses Taxonomy (ontology) concepts, RSMV, GreedyIL & GreedyDP algorithm. This system may be applied in real world applications like Google and Yahoo search engine, which allows users to describe their interests explicitly by selecting from pre-defined preference option, so that the results that match are re-rank according to the user's interests and also gives users the option to save web sites they like and block those they dislike.

**KEYWORDS:** Privacy protection, Hashing, Ontology (content & location), Risk, Profile

## **I. INTRODUCTION**

The measure of information on the web persistently develops; it has become an increasingly arduous for web search engines to return the information to the users that satisfy users' individual needs [2]. Web Search engines are the single largest source of web traffic, which allows people looking for useful information. When people are looking for something online they go to a search engine first and start searching directly to it by typing the domain name into the browser address bar.

However, users might face problem when web search engine returns irrelevant result that do not meet our real intentions [1]. This is because different users have different backgrounds and interests and also may have different information needs when providing exactly the same query. To overcome this problem, personalized web search came into picture since it can provide different results based on the users' preferences and information requirements of users. Consider the query "keyboard": Personalized web search can determine the query by collecting the following user information:

1. The user is a computer user, not a musician.
2. The user has just input a query "keyboard," but not "music".

The user would have viewed a web page with numerous words identified with computer keyboard before entering the inquiry, for example, computer, input-output device. Individual information, for example, bookmark, browsing history may be useful to recognize a client's certain purposes. On the other hand, clients have concerns about how their own data is used. As opposed to confidentiality and security, privacy depends on how the user may benefit from data and information sharing. The personalized web search proposed is an innovative method for personalizing the search result as it uses content mining and location concept for profiling the user by utilizing the location and content preferences.



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 5, May 2015

## II. RELATED WORK

In [1] authors have demonstrated its effectiveness in improving the quality of various search services on the Internet. However, evidences show that users' reluctance to disclose their private information during search has become a major barrier for the wide proliferation of PWS. Protection assurance in PWS applications that model client preference as various levelled client profiles has been considered. A PWS structure called UPS is proposed, that can adaptively sum up profiles by questions while regarding client indicated security prerequisites. The runtime speculation goes for striking a harmony between two predictive matrix that assess the utility of personalization and the protection danger of uncovering the summed up profile. Two voracious calculations, specifically GreedyDP and GreedyIL are exhibited, for runtime speculation. An online prediction method is accommodated choosing whether customizing a query is helpful. Its main limitation is query submitted by the user are not encrypted and also search results are not provided based on the location of the user.

When a same inquiry is presented by diverse clients [2], average web crawlers give back the same result paying little heed to who presented the question. Customized web hunt is a promising approach to enhance look quality by redoing query items for individuals with individual data objectives. However, customers are uncomfortable with uncovering private slant information to web crawlers. Then again, protection is not outright, and frequently can be traded off if there is an addition in administration or productivity to the client. A versatile way for customers to subsequently manufacture rich customer profiles is shown. An adaptable path for clients to consequently fabricate rich client profiles is displayed. These profiles abridge a client's advantage into a various levelled association as indicated by particular hobbies. Two parameters for indicating protection necessities are proposed to help the client to pick the content and level of subtle element of the profile data that is presented to the web search tool. Its main limitation is clicked through data are not collected for personalization and also not satisfactory in handling large number of information.

In [3] authors said although personalized search has been proposed for many years and many personalization strategies have been investigated, it is still unclear whether personalization is consistently effective on different queries for different users, and under different search contexts. This problem and provide some preliminary conclusions. A large-scale evaluation framework for personalized search based on query logs, and then evaluate five personalized search strategies (including two click-based and three profile-based ones) using 12-day MSN query logs [3]. By analyzing the results, it is reveal that personalized search has significant improvement over common web search on some queries but it has little effect on other queries (e.g., queries with small click entropy) and reveal that both long-term and short-term contexts are very important in improving search performance for profile-based personalized search strategies. Its main limitation is that the search result is not accurate in large database.

## III. PROBLEM STATEMENT

A UPS framework which is a client-side privacy protection framework has been presented in the existing system. UPS could conceivably be received by any PWS that catches client profiles in a progressive scientific classification. This framework allowed user to specify customized privacy requirements through the hierarchical profiles. Furthermore, UPS likewise executed online speculation on client profiles so as to secure the security of the client without trading off the hunt quality. Two algorithms namely GreedyDP and GreedyIL, have been proposed for the online-generalization which achieves better search results while user's customized privacy requirements are preserved. It also improves efficiency and effectiveness.

At the same time in the current framework, it utilizes just the speculation idea. It corrupts the execution of existing framework. For this we are going execute and augment the methodology by utilizing some different properties, for example, selectiveness and to make a framework competent to catch a progression of inquiries. In the Existing System, it has a high cost as far as processing and communication.

The existing system has three framework architectures. In these three segments i.e. client, server and proxy has been utilized. Client data's are imparted to the proxy. In the proposed framework, data's has selectiveness. It can't be imparted to the protection. At the point when the sought data's are summed up and after that just data's are put away in the history. Just hidid data's are put away into the history. String Similarity Match Algorithm (SSM Algorithm) is superior to the covetous calculation. It accomplishes more exactness in indexed lists.

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 5, May 2015

## IV. PROPOSED FRAMEWORK

In order to execute the queries that focus on content as well as location information, content mining and location based concept utilizes the location and content preferences. Location data was extracted from the web archives or history pages, which was converted into location based data i.e. latitude-longitude pairs. At the point when a client enters a query combine with a latitude-longitude pair, a search circle was made by the framework which was focused at the predefined scope longitude match and recovers data containing area information inside the search circle.

A framework for a protecting privacy for user in personalized web search has been proposed, which process each query according to user customized protection necessities using hashing.

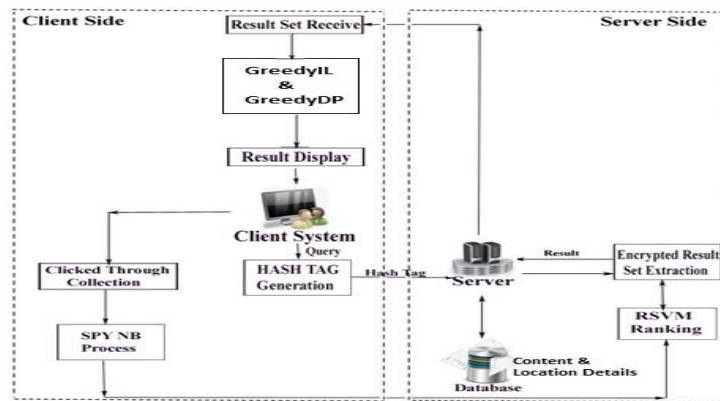


Fig.1. Architecture of proposed PWS

When a user issues a query on the client for web search, it is encrypted using hash tag generated by Hashing technique, as illustrated in Fig 1. The encrypted query, which is in the form of a metadata, is passed to the server. The server then retrieves the data relevant to the query from the database after it has been processed by Taxonomy Management. The taxonomy management consists of content ontology and location ontology. The content ontology and the location ontology together with click through data are used to create feature vectors containing the user location preferences. They will then be changed into an area weight vector to rank the indexed lists as indicated by the client's location preference.

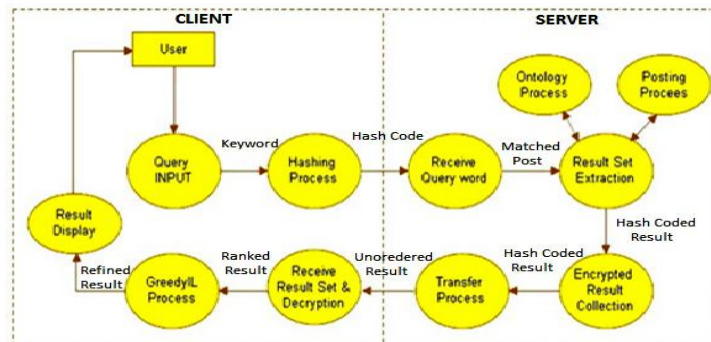


Fig.2. Dataflow Diagram of Personalized Web Search

In the client side, when the user issues input or query to the system for searching any information in the web, the query is encrypted. This is done using Hashing technique, as shown in Fig 4b. Then, the encrypted query is in the form of a metadata. This meta data is passed to the result set extraction. During the process of result set extraction, by using the posting and ontology details present in the server's database, the server processes the query provided by the user. The ontology and posting process are managed by the admin of the system i.e. the server. The server transfers the result set to the client side in the form of encrypted data. The encrypted results are then decrypted. After decryption, GreedyIL and GreedyDP algorithms ensures that all the possible results are provided to the user thereby minimizing information loss and maximizing discriminating power respectively.

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 5, May 2015

## A. Ontology for Web Personalization

Ontological user profile is the representation in which the unified context model for a user is represented as an instance of a reference domain ontology in which concepts are annotated by interest scores derived and updated implicitly based on the information access behavior. To utilize the user context to personalize search results by re-ranking the results returned from a search engine for a given query is the main objective. The semantic knowledge is an essential part of the user context so domain ontology is used as the fundamental source of semantic knowledge in our framework.

An ontological approach to user profiling has proven to be successful in addressing the cold-start problem in recommender systems where no initial information is available early on upon [19]. At the time when learning client interest, the system performs inadequately until enough data has been gathered for client profiling. Utilizing ontology in terms of the profile allows the initial user behavior to be matched with existing domain ontology concept. In this approach, the purpose of using ontology is to identify topics that might be of interest to a specific Web user. In the Fig.3, the hierarchical relationship among the concepts is taken into consideration for building the ontological user profile as the annotations for existing concepts is updated using spreading activation

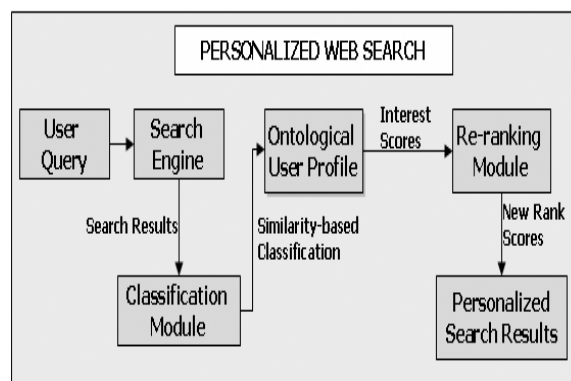


Fig.3. Personalized Web Search based on Ontological User Profiles

Like content ontology, the location ontology combined with click through data from the client is used to create feature vectors which contain the preference of user location. Then it will then be transformed into a location weight vector for the purpose of re-ranking the search results according to the user's preferences.

## B. Hash Tag Generation for Query Terms and Privacy Data

Client query and location are kept protected into the server which is hash coded by MD5 (One way encryption). Message-Digest Algorithm 5 is a broadly used cryptographic capacity with 128-bit hash value. MD5 has been utilized in a wide mixture of security applications, and is additionally ordinarily used to check the trustworthiness of records. A MD5 hash is commonly communicated as a 32-digit hexadecimal number.

The information message is separated into lumps of 512-bit hinders (sixteen 32-bit integers), the message is cushioned so that its length is detachable by 512. The cushioning fills in as takes after: first a solitary bit, 1, is affixed to the end of the message. This is trailed by the same number of zeros as are obliged to bring the length of the message up to 64bits less than 512. The remaining bits are topped off with a 64-bit number speaking to the length of the first message, in bits. The MD5 calculation utilizes 4 state variables, each of which is a 32 bit whole number (an unsigned long on most frameworks). These variables are cut and diced and are (in the end) the message digest.

MD5 comprises of 64 of these operations, gathered in four rounds of 16 operations as shown in Fig.4.  $F$  is a nonlinear capacity; one capacity is utilized as a part of each round  $M_i$  denotes a block of 32-bit of the message input, and  $K_i$  denotes a constant of 32-bit, different for each operation.  $\lll_s$  denotes a left bit rotation by  $s$  places;  $s$  varies for each operation.  $\boxplus$  denotes addition modulo  $2^{32}$ .

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 5, May 2015

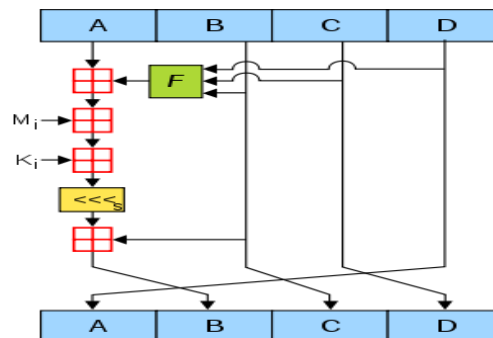


Fig .4. MD-5 operation

Initialized the variables as follows:

A = 0x63252301

B = 0xEDFVAB89

C = 0x98BDESFE

D = 0x10396476.

Now on to the actual meat of the algorithm: the main part of the algorithm uses four functions to thoroughly goober the above state variables. Those functions are as follows:

$F(X,Y,Z) = (X \& Y) | (\sim(X) \& Z)$

$G(X,Y,Z) = (X \& Z) | (Y \& \sim(Z))$

$H(X,Y,Z) = X \wedge Y \wedge Z$

$I(X,Y,Z) = Y \wedge (X | \sim(Z))$

Where  $\&$ ,  $|$ ,  $\wedge$ , and  $\sim$  are the bit-wise AND, OR, XOR, and NOT operators

The transformation of the state variables from their initial state is achieved by using the state variable and input message along with the above functions. Then the message digest is generated. After this, the message digest is stored in the state variables namely A, B, C, D. Example: let us consider after the digest:

A = 0x01234567;

B = 0x89ABCDEF;

C = 0x1337D00D

D = 0xA5501010

Then the message digest would be:

67452301EFCDAB890DD03713010151A5 (required hash value of the input value).

Example:

MD5 ("Java is a programming language")

= 9e107d9d372cc6826bd81d3542a419d6

## C. Personal Behavior Collections using SPY NB Algorithm

Spy NB (Naïve Bayesian) earns user behavior models from preferences extracted from click through data. Assuming that users only click on documents that are of interest to them, Spy NB treats the clicked documents as positive samples, and predicts reliable negative documents from the unlabeled (i.e., un clicked) documents. To do the prediction, the "spy" technique incorporates a novel voting procedure into Naive Bayesian classifier to predict a negative set of documents from the unlabeled document set.

## D. Personalized Searching using Adaptive Ranking Technique

Ranking Support Vector Machine (RSVM) is utilized to take in a customized positioning capacity for rank adjustment of the list items as indicated by the client location and content data. A feature vector can be considered as a point in the feature space because of its ability to represent each document. A set of location and content concepts are



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 5, May 2015

extracted from the search results as document features. The main objective of RSVM is to find a ranking function using the preference pair and holds document preference pairs as many as possible.

## V. PSEUDO CODE

- Step 1: Read the Keyword (“query”) submitted by the user.
- Step 2: Identify the user\_ID, user\_location of the user.
- Step 3: Extract the location details (latitude and longitude) of the user\_location.
- Step 4: Encrypt the keyword using MD5 hashing.  
    `enc = MD5(keyword)`
- Step 5: Find matching post of the keyword.  
    `result= getPostingSearch(keyword,enc)`
- Step 6: Calculate content distance.  
    `avg= (count_list/total)*100;`  
    `avg_list=(avg/2);`  
    `con_avg_list= content_distance (posting, keyword);`
- Step 7: Calculate location distance.  
    `total_dist=total_distance + distance;`  
    `avg= 100-avg;`  
    `avg_list= (avg/2);`  
    `loc_avg_list= location_distance (latitude, longitude);`
- Step 8: Find out the ranking score using SPY-NB and RSVM.  
    `sum= con_avg_list + loc_avg_list;`
- Step 9: Display the result.

## VI. EXPERIMENTAL RESULTS

The proposed framework effectively provided a secured search interface to the user without compromising the personalized search quality. Though this framework focuses on the preference of the user, it returns quality search results to the user. This framework successfully improved the search accuracy by minimizing information loss and maximizing discriminating power.

The query submitted by the user is encrypted using MD-5 Hashing Technique (Message Digest- 5) as shown in Fig.5. The user’s query and location are kept protected into the server which is hash coded by MD5 (One way encryption). Message-Digest Algorithm 5 is a broadly used cryptographic capacity with 128- bit hash value. MD5 has been utilized in a wide mixture of security applications, and is additionally ordinarily used to check the trustworthiness of records. A MD5 hash is commonly communicated as a 32-digit hexadecimal number. Therefore, this framework successfully added a new dimension of security to the personalized web search.



# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 5, May 2015

s_no	s_query	u_code
6	e6d96502596d7e7887b76646c5f615d9	5
7	1b32aa85a594d41c539d7c90963daf6b	5
8	1b32aa85a594d41c539d7c90963daf6b	5
9	1b32aa85a594d41c539d7c90963daf6b	5
10	e6d96502596d7e7887b76646c5f615d9	5
11	e6d96502596d7e7887b76646c5f615d9	6
12	40203abe6e81ed98cbc97cdd6ec4f144	6
13	40203abe6e81ed98cbc97cdd6ec4f144	6
14	40203abe6e81ed98cbc97cdd6ec4f144	6
15	aa4eceb10187d79aebd416b7d75d2659	6
16	aa4eceb10187d79aebd416b7d75d2659	6
17	1b32aa85a594d41c539d7c90963daf6b	6
18	e6d96502596d7e7887b76646c5f615d9	6
19	e6d96502596d7e7887b76646c5f615d9	6
20	40203abe6e81ed98cbc97cdd6ec4f144	6
21	e6d96502596d7e7887b76646c5f615d9	6
*	(NULL)	(NULL)

Fig.5. Hash Coded Query

s_no	s_query	matched_post	positive	negative	u_code
1	e6d96502596d7e7887b76646c5f615d9	1,2,6,7,9,20	6	1,2,7,9,20	5
2	e6d96502596d7e7887b76646c5f615d9	1,2,6,7,9,20	(NULL)	1,2,6,7,9,20	5
3	e6d96502596d7e7887b76646c5f615d9	1,2,6,7,9,20	(NULL)	1,2,6,7,9,20	5
4	e6d96502596d7e7887b76646c5f615d9	1,2,6,7,9,20	(NULL)	1,2,6,7,9,20	5
5	e6d96502596d7e7887b76646c5f615d9	1,2,6,7,9,20	(NULL)	1,2,6,7,9,20	5
6	e6d96502596d7e7887b76646c5f615d9	1,2,6,7,9,20	(NULL)	1,2,6,7,9,20	5
7	1b32aa85a594d41c539d7c90963daf6b	1	(NULL)	1	5
8	1b32aa85a594d41c539d7c90963daf6b	1	(NULL)	1	5
9	1b32aa85a594d41c539d7c90963daf6b	1	(NULL)	1	5
10	e6d96502596d7e7887b76646c5f615d9	1,2,6,7,9,20	(NULL)	1,2,6,7,9,20	5
11	e6d96502596d7e7887b76646c5f615d9	1,2,6,7,9,20	(NULL)	1,2,6,7,9,20	6
12	40203abe6e81ed98cbc97cdd6ec4f144	26,27	27,27	26	6
13	40203abe6e81ed98cbc97cdd6ec4f144	26,27	27,27,27,27	26	6
14	40203abe6e81ed98cbc97cdd6ec4f144	26,27	(NULL)	26,27	6
15	aa4eceb10187d79aebd416b7d75d2659	24,25	(NULL)	24,25	6
16	aa4eceb10187d79aebd416b7d75d2659	24,25	(NULL)	24,25	6
17	1b32aa85a594d41c539d7c90963daf6b	1	(NULL)	1	6
18	e6d96502596d7e7887b76646c5f615d9	1,2,6,7,9,20	(NULL)	1,2,6,7,9,20	6

Fig.6. SPY-NB Process

Click-through data is a search engine log that records for every query the outcome rundown displayed to the client and also the links clicked on by the client. Click-through data can be recorded with minimal overhead and without trading off the usefulness and handiness of the search engine. Specifically, contrasted with express client criticism, it doesn't include any overhead for the client. The query and the returned positioning can without much of a stretch be recorded at whatever point the subsequent ranking is shown to the client. For recording the snaps, a basic intermediary framework can keep a log document.

Spy-NB is a compelling intends to produce the positive and negative datasets as indicated in the Fig.6, from which exact preference fragment pairs can be inferred for improving the ranking capacity. This methodology likewise takes care of the issue that a client may skip some applicable connections when he or she outputs down the outcome rundown, prompting the extraction of wrong preference pairs.

SPY-NB created an arrangement of preference that is then fed into the RSVM (Ranking Support Vector Machine) algorithm for advancing the ranking function for the client. At that point, a ranked item is displayed as a pair comprising of a query and a hit archive, only items that have the same query are ranked in one epitome. Everything is represented by a feature vector, which records elements and relating component weights to permit an end client to "tune" a ranking capacity taking into account coupling a query unwinding system with a ranking SVM. A ranking function acknowledges a item feature vector and produces significance scores for the item. The higher the score, the higher this item will be ranked.

# International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 5, May 2015

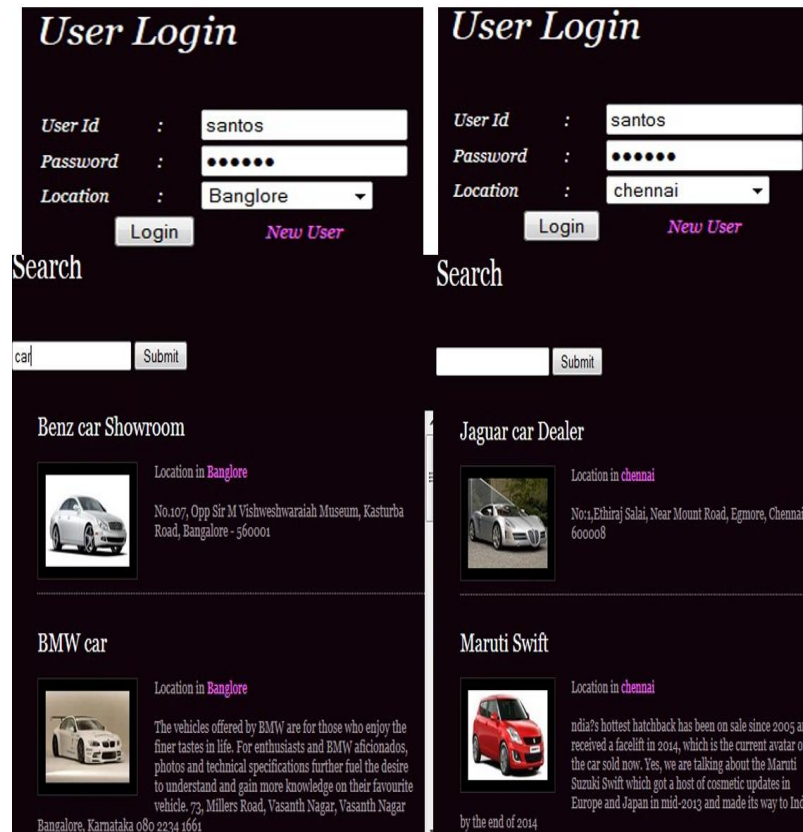


Fig.7. Search Result Comparison based on Location

The search results in Fig.7 are re-ranked based on the location. The user while logging in to the personalized web search interface sets their location (e.g. “Bangalore”). The search engine will return the result based on the location specified by the user. This shows that the location specified is given the main preference while others are kept secondary. Thus, successfully personalized the web search and at the same time, preserved the privacy of the user without providing much performance degradation.

## VII. CONCLUSION AND FUTURE WORK

A framework for privacy-preserving in personalized web search using ontology and hashing technique has been implemented successfully. Although personalized search was unclear, whether personalization is constantly effective on different queries for different users and under different search context, the proposed framework has overcome this limitation. In the existing method, client’s data were imparted to the proxy and hence privacy protection was compromised. The implicitly collection of personal data may easily reveal gamut of users’ private life and it’s dangerous. However, the proposed framework provides more secure and innovative idea for privacy protection and personalized web search. To achieve privacy protection, this system provides metadata and user’s query encryption which is achieved using MD-5 hashing technique. The proposed system supports runtime profiling search result are retrieve based on the user, not on the system; which was not present earlier. Re-ranking the search results based on the interest scores and location preferences is effective in presenting the most relevant results to the user and hence SPY-NB (Naïve Bayes) and Ranking Support Vector Machine (RSVM) are employed. To make system more efficient, users’ profiles are maintained at client side and all the possible query results are provided to the user with the help of GreedyIL and GreedyDP, thereby minimizing information loss and maximizing discriminating power respectively. And also taxonomy (Ontology) repository concept is used to achieve personalized web search in this system. In future work, the proposed framework can be integrated into Web browsers like Internet Explorer, Mozilla Firefox, etc. and also be developed as a mobile application (Android/iOS).





# International Journal of Innovative Research in Computer and Communication Engineering

*(An ISO 3297: 2007 Certified Organization)*

**Vol. 3, Issue 5, May 2015**

## REFERENCES

- [1] Lidan Shou, He Bai, Ke Chen, and Gang Chen, "Supporting Privacy Protection in Personalized Web Search," Knowledge and Data Engineering, IEEE Transactions, vol. 26, Issue: 2, pp. 453 - 467 , 2014.
- [2] Y. Xu, K. Wang, B. Zhang, and Z. Chen, "Privacy-Enhancing Personalized Web Search," Proc. 16th Int'l Conf. World Wide Web (WWW), pp. 591-600, 2007.
- [3] Z. Dou, R. Song, and J.-R. Wen, "A Large-Scale Evaluation and Analysis of Personalized Search Strategies," Proc. Int'l Conf. World Wide Web (WWW), pp. 581-590, 2007.
- [4] J. Teevan, S.T. Dumais, and E. Horvitz, "Personalizing Search via Automated Analysis of Interests and Activities," Proc. 28th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR), pp. 449-456, 2005.
- [5] K. Sugiyama, K. Hatano, and M. Yoshikawa, "Adaptive Web Search Based on User Profile Constructed without any Effort from Users," Proc. 13th Int'l Conf. World Wide Web (WWW), 2004.
- [6] X. Shen, B. Tan, and C. Zhai, "Implicit User Modeling for Personalized Search," Proc. 14th ACM Int'l Conf. Information and Knowledge Management (CIKM), 2005.
- [7] F. Qiu and J. Cho, "Automatic Identification of User Interest for Personalized Search," Proc. 15th Int'l Conf. World Wide Web (WWW), pp. 727-736, 2006.
- [8] P.A. Chirita, W. Nejdl, R. Paiu, and C. Kohlschütter, "Using ODP Metadata to Personalize Search," Proc. 28th Ann. Int'l ACM SIGIR Conf. Research and Development Information Retrieval (SIGIR), 2005.
- [9] A. Pretschner and S. Gauch, "Ontology-Based Personalized Search and Browsing," Proc. IEEE 11th Int'l Conf. Tools with Artificial Intelligence (ICTAI '99), 1999.
- [10] Ahu Sieg, Bamshad Mobasher and Robin Burke, "Web Search Personalization with Ontological User Profiles," Proc. 16<sup>th</sup> ACM Conference on information and knowledge management, pp.525-534, 2007.