# Incorporating Synonyms Into Snippet Based Query Recommendation System

Megha R. Sisode and Ujwala M. Patil

Department of Computer Engineering, R. C. Patel Institute of Technology,
Shirpur, Maharashtra, India
`Megha_sisode@yahoo.com`
`patil_ujwala2003@rediffmail.com`

## ABSTRACT

*Recently, growth of internet has been increased for information retrieval though it is difficult to extract the relevant information in less time. Search engine sometime fails to understand user search intend. Query recommendation can be used to help user to state exactly their information need. Search engine can return appropriate result to meet users' information needs. There are various methods based on history of users and snippets to retrieve the information. But these methods fail to satisfy users need. Therefore in addition of history and snippets with synonyms will do better. Moreover user preferences can be used to build the user profile which will help in effective recommendations. Here for given query recommendation the synonyms are extracted on line. Synonym based method ranks the clicked URLs at the top of the result based on user profile. The performance of the system shows that the synonym based approach give better and effective recommendation for all queries as compared to previous methods.*

## KEYWORDS

*query recommendation, synonym, snippets, information retrieval, user profile*

## 1. INTRODUCTION

As the growth of World Wide Web is increased with increase of size and popularity and the assembly of large scale volumes of web data, thus it is difficult to extract the relevant information that have been used in wide range of application. Many novice users face the difficulty to get the desired information although they use most efficient search engines such as yahoo, google.

The search engine has gain more success and the growth of internet resources is increasing as the web is a repository of large scale updated information. Web search engine is the major platform to extract the needed information to user by posing a query. The web search engine helps user to exploit the required information based on user query. For this purpose search engine provide platform to the users to specify their information need in the form of queries simply as list of keywords. This keyword based user interface causes lots of troubles in search process.

User queries are the most important factors as they are only interface for users to access web pages that affect the performance of search engines. Although users' information needs are

complicated, their queries are usually simple, short and possibly ambiguous. Queries are simple because users are unable to organize complicated queries which can describe their information needs more exactly.

This causes a major challenge in current Web search techniques, which is the understanding of user's information need behind queries. Sometimes it becomes difficult for search engines to understand information need from only queries, so that click through behavior data can be utilized. Thus the query recommendation technique is proposed to present users with a list of possible choices whose information needs are relatively clear to search engines. By this means, users can exactly state their information need by clicking recommendation query links instead of inputting new queries [1].

With the analysis into altavista search engine's query logs it is found that the average length of user queries is 2.35 terms and mostly the user queries are short including around two terms per query, on the other hand the ambiguity of language play essential role, as the users often fail to organize appropriate terms for their search query, so as the search engine returns mismatch results to the same topic and also faces the problem for synonyms and polysemous words that exist in language. Thus query recommendation function helps user to recognize their short formed and possibly ambiguous queries and return appropriate result to satisfy the user information need [1].

Recently the query recommendation has been widely used by the users to satisfy their information needs. According to the survey, it has been observed that approximately 78% users will change their queries with search engine recommendation function if they cannot obtain satisfactory results for their query. So that it is mandatory for search engine to provide good quality recommendations which can express users actual information needs more exactly.

According to research their has been lot of work done for improving result of search engine based on users previous query log data and click behavior so that search engine can locate popular queries which are similar to current query either in content or in click context  [3].

These methods end with suggesting user to adopt a similar and/or frequently adopted query also fails to exactly understand users' information need and also does not consider current users search intent into account as they believe that current user shares similar interest as other user with the same query. These methods also produce improper recommendation for low frequency query as not much candidate queries for them.

In order to solve these problems and give better recommendation results which can satisfy users information need, we have to understand the way to express the users' information need. For better recommendation the query must be formulate properly and well organized manner with more exact meaning. If we observe the way user search on web, then we will come to know that when user clicks certain search result returned by search engine, it does not always mean that user is interested in the resultant document as because she/he has not yet viewed the resultant document. Instead the assumption is the user must be interested in the snippets of the corresponding document because it is the snippets that are actually shown to and read by user.

The synonym based method follows that users' information needs are described in the interaction with search engine, specifically, in the snippets which they have been ever clicked for the results.

Thus on the basis of these assumption, a synonym based query recommendation framework with snippet click model is presented, which include global scale and local scale snippets. With these models keywords are extracted from clicked snippets to make effective recommendations with using synonyms of query word along with snippets and location information. Differently

synonyms based query recommendation methods gives effective and more accurate results as compared to the history and snippet based system.

Nowadays, google is a world most leading search engine with different language interfaces. There exist some limitations with the keyword based searching. One of the web search key issue is that user tend to insert very general queries. That leads huge amount of information to be returned for given query. There are various ways to deal with a huge amount of retrieved web pages for arranging with the proper meaning. Synonyms or word sense disambiguation can be used along with snippets. The synonym based query recommendation approach uses WordNet[*] for discovering the synonyms.

The rest of the paper is organized as follows. Section 2 gives literature survey on work done in query recommendation processing methods; section 3 explains the motivation for synonym based method; then section 3 gives the working of synonym based recommendation method. Section 4 describes experimental setup for synonym based method. Finally conclusion and the future work is explained in section 5.

## 2. LITERATURE REVIEW

Recent researchers have proposed various recommendation systems for online information retrieval using various approaches. A literature survey is done to examine different approaches in order to mine essential features from query log data of search engine.

Ricardo Baeza-Yates et al. had proposed a method for suggesting list of related queries to user based on a query clustering process. This method not only discovers the related queries, but also ranks them according to a relevance criterion. This notion of query similarity has several advantages that it is simple and easy to compute. Moreover, it captures semantic relationships among queries by relating queries that are worded differently but stem from the same topic [2].

Silviu Cucerzan et al. had presented a method to suggest queries based on mining into post-query browsing behaviors referred as "search trails". They utilized user landing pages i.e. the ending pages of search trails to generate query suggestions. For each landing page of a user submitted query they identify queries from query logs that have these landing pages as one of their top 10 results and these queries are used for suggestions [3].

Shen Xiaoyan et al. proposed an effective approach for query suggestions. This approach accepts Chinese web query as input and the approach not only identify related queries already existed in the log of previously submitted queries of search engine but also use synonyms that are extracted from web based corpuses to construct new related queries. Also rank the queries according with degree of relatedness, freshness and effectiveness. This approach proves its effectiveness in recommending related queries for high frequency queries than that of low frequency query [4].

Qi He et al. proposed a novel sequential query prediction approach to grasp a users' search intent based on her/his past query sequence and its resemblance to historical query sequence models mined from massive search engine log data. Differently from previous work done where only single preceding query is used for prediction, this work considers variable number of preceding query and effectively captures more complex context information for recommendation. The Results shows that the sequence-wise approaches significantly outperform the conventional pair-wise ones in terms of prediction accuracy[5].

*\* http://wordnet.princeton.edu/*

Thus the work has one fundamental difference from all previous session-based approaches. As all previous work focuses on pair-wise query relations and uses only a single preceding query for query prediction, proposed method consider variable number of preceding queries and effectively capture more complex context information for query recommendation. Moreover, this approach can automatically determine the optimal context length to be used for query prediction [5].

Hamada Zahera et al. proposed a method for suggesting a list of queries that are related to the user input query based on previously issued queries by the users. Their method was based on clustering process in which groups of semantically similar queries posed by user are detected in order direct them toward their required information need. This method not only discovered the related queries but also rank the query according to a similarity measure [6].

C. Sumathi et al. proposed a session based approach where the proposed method is based on the users' navigational patterns and provide recommendations to fulfill the current users information need. This method had classified and matched an online user based on her/his browsing interests [7].

Poonam Goyal et al. had proposed a method to facilitate users with query recommendations in which the concepts related to the users information need are suggested to the users to satisfy their exact information need. In that they extracted the concepts from the web snippets and have used two weight functions to measure the relevance between query and concept. Related concepts with different meaning are selected and recommended as query suggestions to the users [8].

Ji-Rong Wen et al. had proposed an approach to cluster similar queries to recommend URLs for frequently asked queries of a search engine by using four notions according to: first, the context of the query; second, common clicked URL's between queries; third, string matching of keywords, and fourth is, the distance of the clicked documents in some pre-defined hierarchy. But result of this method generates very sparse distance matrices and this sparsity is diminished using large query logs. Thus string matching features are used to locate similar queries [9].

Osmar Zaiane et al. had used content similarity to recommend similar queries using Query Memory, a data structure that holds the collective query trace and also extra information pertaining to the queries that would help in measuring similarities between queries. Query trace is a log containing previously submitted queries. The major advantage of this method is that it suggests the queries when user is not satisfied by current search result but sometimes this method produces irrelevant result and leaves the choice up to user [10].

Eugene Agichtein et al. shows that incorporating user behaviour data can significantly improve ordering of top results in real web search setting. Also alternatives for incorporating feedback into the ranking process has been examined and explored the contributions of user feedback compared to other common web search features[12].

Yiqun Liu et al. had presented an approach to focus how to detect users actual information need by extracting the snippets. The snippet based approach considers that users information needs are better described in their interaction with search engine more specifically, in the snippets of the results which ever clicked by users. The key idea of that system follows that if user click certain result from list then it shows that the user has read that particular snippet and interested in that snippet and not in the result. But this system does not consider the location and synonyms also, which may be useful for improving the results of search [1].

## 3. MOTIVATION

Many query recommendation methods used to suggest related queries by extracting information from clicked documents because these documents are expected to contain users' preference and relevance judgments. Different methods use previous query data, history of snippets. But this method does not consider current user search intent exactly, also fail to recommend for low frequency query. In addition to the snippets of clicked documents, the synonyms extracted from the online synonym services can be integrate to improve the performance of search engine. Also user preference can be integrate to form current user profile to specify the current users search intent. So the synonym based recommendation system is expected to give more accurate recommendation based on input.

## 4. SYNONYMS BASED QUERY RECOMMENDATION SYSTEM

The synonym based query recommendation method is based on the assumption that users' information needs are described more specifically in snippets of the results which they ever clicked, this is because when user clicks a certain search result, it does not necessarily mean that she/he is interested with the result because she/he has not yet viewed the resultant document. It is probably mean that she/he is interested in the snippets of the corresponding resultant document because these snippets are actually shown to and read by users. According to this assumption, the synonyms based query recommendation uses clicked snippets also the synonyms extracted from the online synonyms service for considering the synopsis.

In addition with this synonym based recommendation system also considers location information. As in previous method if the user is at some specified location then the system does not consider the location of user and recommends without considering location but in synonym based recommendation method it improves the accuracy in result with considering location information along with synonyms extraction. Figure 1 shows the working of synonym based query recommendation method.

The synonym based recommendation is based on a snippet click models which tries to extract keywords appearing in users' clicked snippets as recommendation. According to Baeza-Yates et al. query recommendation is the method which is used to suggest alternative queries to users in order to help them to specify alternative related queries in their search process [2]. But believe is that the users not only specify alternative related queries but also try to express their information need in the form of query recommendations. Therefore, search engine should recommend queries which are most likely to represent users' information needs.

The synonym based query recommendation task try to rank snippets which are related to the original proposed query on the basis of user profile. Users are interested in the content of snippet because it contains keywords that are related to their information needs and these are actually shown to and read by user. Therefore, the major idea of synonym based query recommendation framework is to locate keywords that appear in snippets clicked by users and can describe users' information need. Differently with previous recommendation methods, it relies on information extracted from users' result click-through process instead of the historical queries fired by other users.
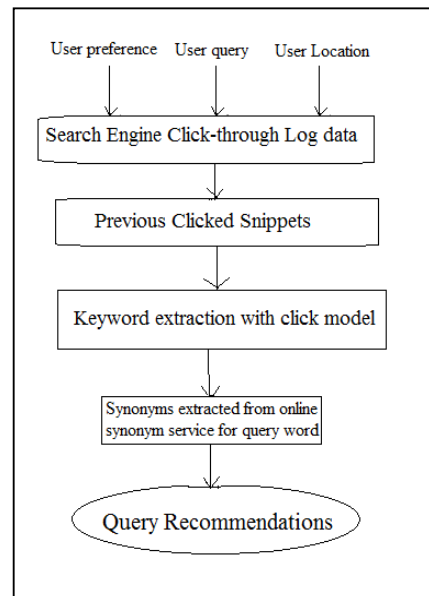
Figure 1.  Synonyms based query recommendation system with considering user preferences and location information

## 4.1. Snippet click mosels incorporating with synonyms

Note that users click a certain resultant document because she/he actually views its corresponding snippet and also expects this document to meet her/his information need. Therefore, the probability of clicking a certain document is decided by both whether user views the snippet and whether user is interested in it. Because user is only able to view the snippet of the document before she/he actually click on the result, then the probability of clicking is decided by whether user is interested in the snippet of the result document; in other words, by whether this snippet meet user's actual information need or not. For synonym based recommendation both a global scale and a local scale snippet click models are used.

These local scale and global scale snippets are adopted to finish the task of query recommendation.

- **Global scale snippet model using synonyms**

For the global scale model, all the clicked snippets for a certain query are treated as a whole ''snippet document''. Therefore, for all clicked snippets, it shows that if user has clicked certain snippet then user must be interested in it and that satisfies the users' information needs. Therefore, a simple TF-based model is used to extract keyword lists from the snippets. For each keyword in the snippets, the recommendation candidates are those with the largest term frequency value, where for a query word W, TF is defined as sum of all appearances of W in all related snippet.

The corresponding global scale snippet algorithm using synonym based recommendation (Algorithm 1) is as follows:

Algorithm 1. Query recommendation based on global scale snippet click model using synonym

QueryRecommendation (Original query Q, Users Click through pattern CLKPAT)

1. Find all the documents clicked for Q in CLKPAT and form a set of document called D;
2. Extract all the snippets of D for query Q by using search engine interface and form a snippet set called S;
3. For snippet set S, extract N keywords by using TF or other keyword extraction algorithms;
4. Extracts the synonyms for all the N keywords from online synonym service.
5. Return these N keywords as recommendation words with considering synonyms.

This algorithm generates a list of keywords for recommendation of query Q. Note that sometimes these keywords may not be directly used for recommendations because they should be combined with the original query to form complete information need. For example, keyword ''free download'' may be returned for query ''Yahoo messenger'', it should be combined with the original query word to form a complete query recommendation word such as ''Yahoo messenger free download''. However, even these keywords cannot be directly adopted as recommendations, these are supposed to meet users' information needs.

- **Local scale snippet model using synonyms**

Differently in a local scale snippet click model each snippet is considered to be treated separately. With the bag-of-words model, a certain clicked snippet can be represented by a set of keywords each with having different TF values. As this consider each snippet separately, and not all keywords appears in the each clicked snippet. So many keywords will have term frequency value as zero, thus it may generate the sparsity problem. The smoothing technique can be used to avoid data sparsity problem and to estimate exactly the information need of user. After this we can consider the probability for each keyword is consider to describe users' information need and the keyword having high probability to satisfy users information need are used to suggest for query recommendation.

The synonym based query recommendation algorithm for local scale snippet is as follows:
Algorithm 2. Query recommendation based on local scale snippet click model using synonyms

QueryRecommendation (Original query Q, Users Click through pattern CLKPAT, Users search interest SI)

1. Find all documents clicked for Q in CLKPAT and form a document set called D;
2. Extract all snippets of D for query Q by using search engine interface and form a snippet set called S;
3. Extracts users search interest SI by combining users profile P and location information L;
4. Recommendation candidate set CANDIDATE = { };
5. For each snippet $S_i$ in snippet set S, if P( Click $_i$ ) is greater than threshold T , then put all words into CANDIDATE set;
6. For all words in CANDIDATE set, extracts the keywords from snippets and also after smoothing task, and form the equation E.
7. Solve E according to Gaussian elimination[*] or other methods.
8. Select N keywords with the largest probability values which indicate greater possibility of describing users' information need;
9. Extracts the synonyms for all N keywords from online synonym service.
10. Return these N keywords as recommendation words with considering the synonyms.

Here also, as similar to Algorithm 1 these N keywords should be combined with the original query to form complete query recommendations. Except the keywords which only appear in snippets with having probability of clicking P(Click) values lower than that of threshold T.

## 5. EXPERIMENTAL SETUP AND RESULTS

For evaluation of the performance of the synonym based recommendation system, the synonym based query recommendation system is run on configuration having Windows 7 with 4GB RAM. The synonyms based recommendation method is implemented with java on android platform. For the system android works at front end and SQLite works at back end to store the database of application. Database is stored in the android device itself with the help of SQLite database. For this the Eclipse software development kit (SDK) is used, which includes the Java development tools to develop an android application, where Eclipse is an integrated development environment (IDE). An android emulator has been created with having query recommendation as an application on it working in contribution with SQLite database.

To evaluate performance of the synonyms based query recommendation framework, practical search engines database has been used and compared performance of synonym based recommendation with current search engine's query recommendation performances. The evaluation is different from most previous researches where the performance of query recommendation is evaluated by how many percentages of users actually clicked these recommendations in practical environment. The synonym based query recommendation method adopts human-annotation based precision-recall metrics for evaluation.

Precision and recall are the basic measures used in evaluating search strategies. Recall is the ratio of the number of relevant records retrieved to the total number of relevant records in the database. It is usually ex pressed as a percentage. Precision is the ratio of the number of relevant records retrieved to the total number of irrelevant and relevant records retrieved. It is usually expressed as a percentage. Recall and precision are inversely related. As the recall increases the value of precision get decreases, and vice-versa.

Click-through rate and user profile are adopted as metrics to evaluate the performance of recommendation algorithms. User profile is built from click patterns, users location patterns and user interest, all of which is available by user once the application is used by that user. For extracting the users location the latitude and longitude has been used. For the synonym based recommendation, click-through data of an application is tracked by the application itself and stored by using SQLite database.

Previously for snippet based query recommendation the experimental results show that the keywords generated by snippet query recommendation method are more preferred by users than the others. About one third of the recommendation is provided by Baidu and Sogou search engines match snippet based query recommendation algorithm results. In snippet based query recommendation method for all recommendations generated by search engines, some match recommendation keywords generated by snippet system while others not. The comparative results of click-through rate and average amount of user clicks are shown in figure 2 and figure 3[1].
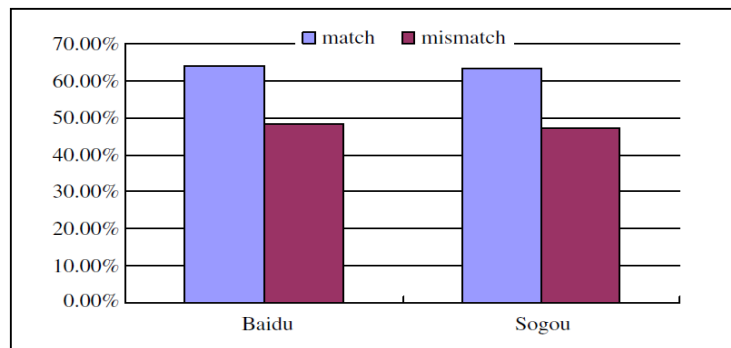
Figure 2. Comparison of click through rate between the recommendations that matches and does not match
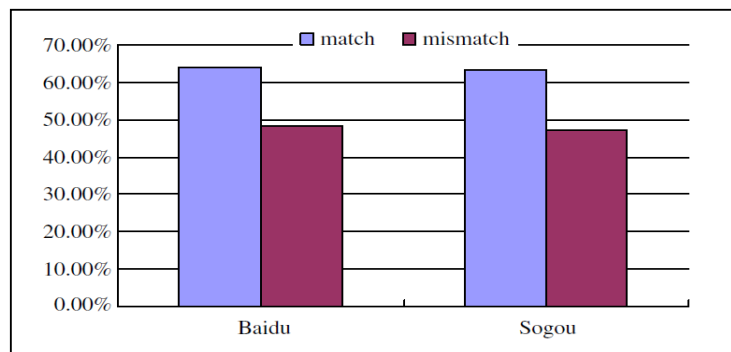


Figure 3. Comparison of average amount of user click between the recommendations that matches and does not match

The synonym based recommendation method have used google search engine database[*], where global scale snippet is done by google itself. For global scale model all clicked snippets for certain query word are treated as a whole. In this users information need is supposed to be related with the snippets. For local scale model each snippet for a query word is treated separately. In this the system can use Algorithm 2 to estimate the probability of keyword in representing users' information need accurately.

The synonym based recommendation method is also incorporated with users' preferences. The method is also integrated with synonyms matching with the help of online synonym service that gives us the requested synopsis for the given words.

The results of the synonym based recommendation system are shown in the form of URLs. It gives URLs because it represents the links from where the snippets are fetched, and then for the result user have to click on the URLs to get the snippets. In case if the user clicks on a URL then the URL appears again in the result then the URL would be ranked at the top of the list. This concept of snippet ranking is based on users' profile.

For snippet based query recommendation the past experimental results show that the keywords generated by snippet query recommendation method are more preferred by users than the others. About one third of the recommendation is provided by Baidu and Sogou search engines match snippet based query recommendation algorithm results. In snippet based query recommendation method for all recommendations generated by search engines, some match recommendation keywords generated by snippet system while others not.

*https://developers.google.com/web-search/docs*
*https://developers.google.com/cloud-sql*

In order to measure the performance of synonym based query recommendation the experimental setup is made with the current well known google search engine database with user profile, user preferences and click through data as metrics. For evaluation the same sample query has been fired to synonym based recommendation system and measured the performance by posing the same query in google search engine. The results are compared on the basis of results returned and delay time.

For example different sample queries are run, and compared with the returned results from google with the same sample query, and the precision and rank of returned results are measured for synonym based query recommendation system. Among these there are variation in time required. Also the synonym based system have more relevant results at the top as this is considering the synonyms returned by online synonym service for query word. Figure shows the representation of the result for sample query run on synonym based query recommendation.

From the observation it has been noticed that as compared with the history and snippet based methods; synonym based recommendation gives better results with synopsis. Also it ranks the link at the top on the basis of users preference and location. The synonym based recommendation method able to satisfy users' information need in less time. For the evaluation of the synonyms based query recommendation method different metrics have been applied such as first percent precision is calculated for differed sample queries and another parameter is used as number of related queries.

In figure 4, the vertical axis represents the % precision calculated for the query and the horizontal axis represents the number of results matched with the user search intent. Also the performance of the method incorporating synonyms with snippet based recommendation is measured with having percent of matching recommendation and the rank of the required result as a metrics of recommendation for different sample queries.

It can be observed from the results that the percentage of precision increases if the query presented by the user is high frequency query and also the time required to returns the result get increase with low frequency query.
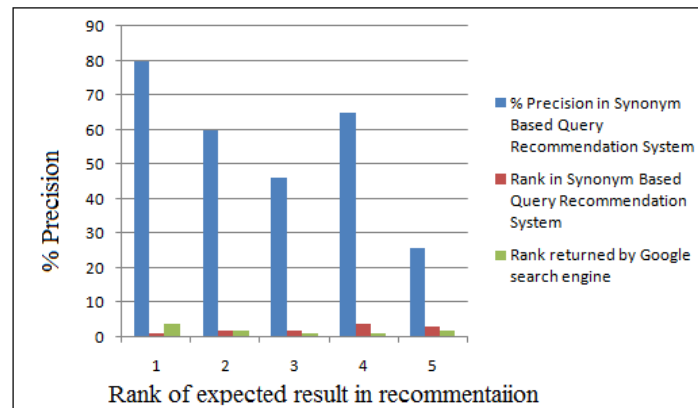


Figure 4.  Result of synonym based query recommendation for few sample queries

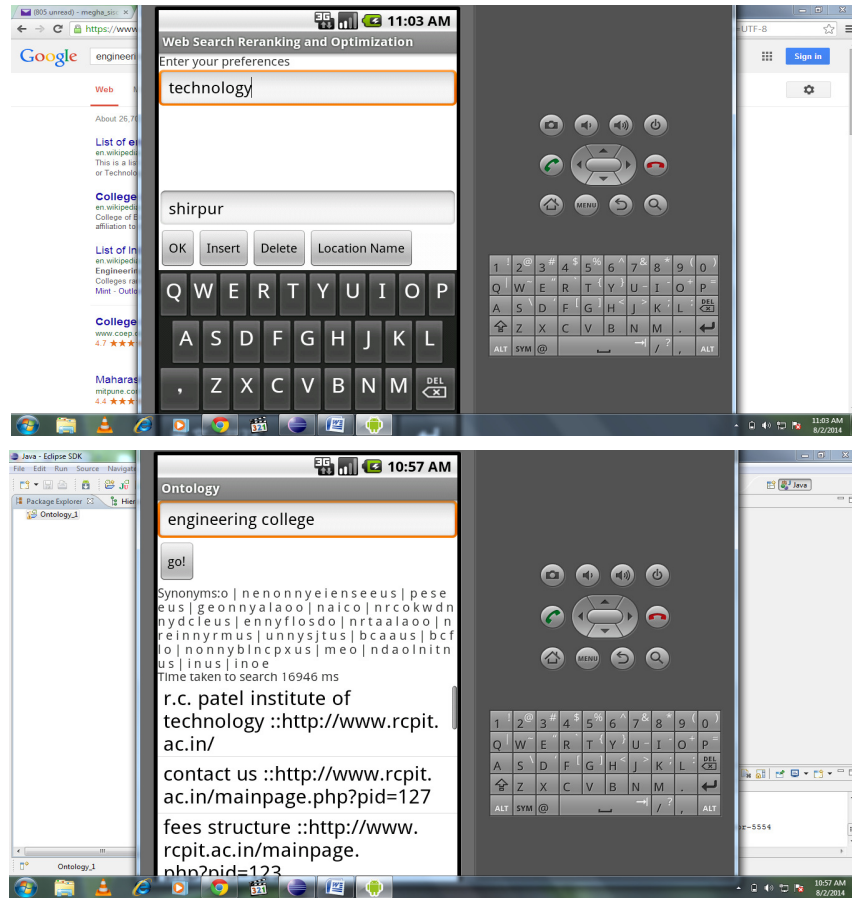- **SNAPSHOTS FOR SYNONYMS BASED QUERY RECOMMENDATION**



Figure 5. Results of synonym based query recommendation for one of the sample query
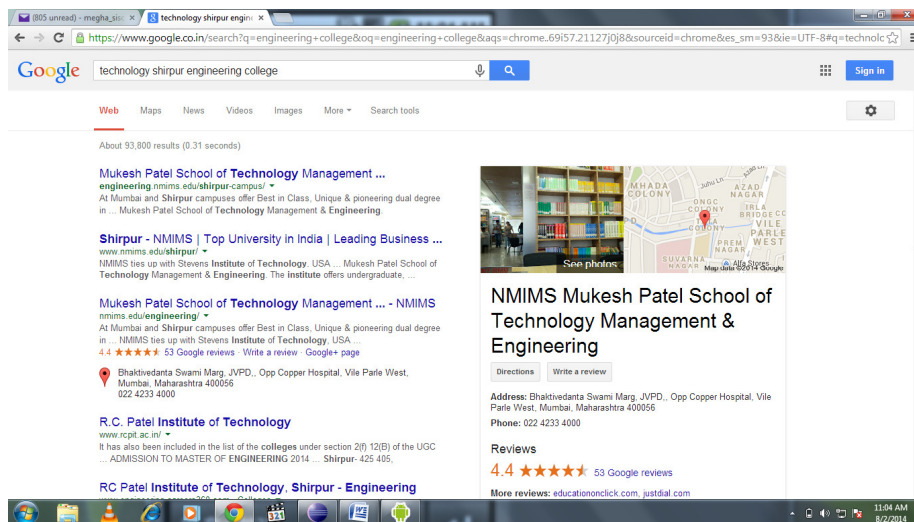


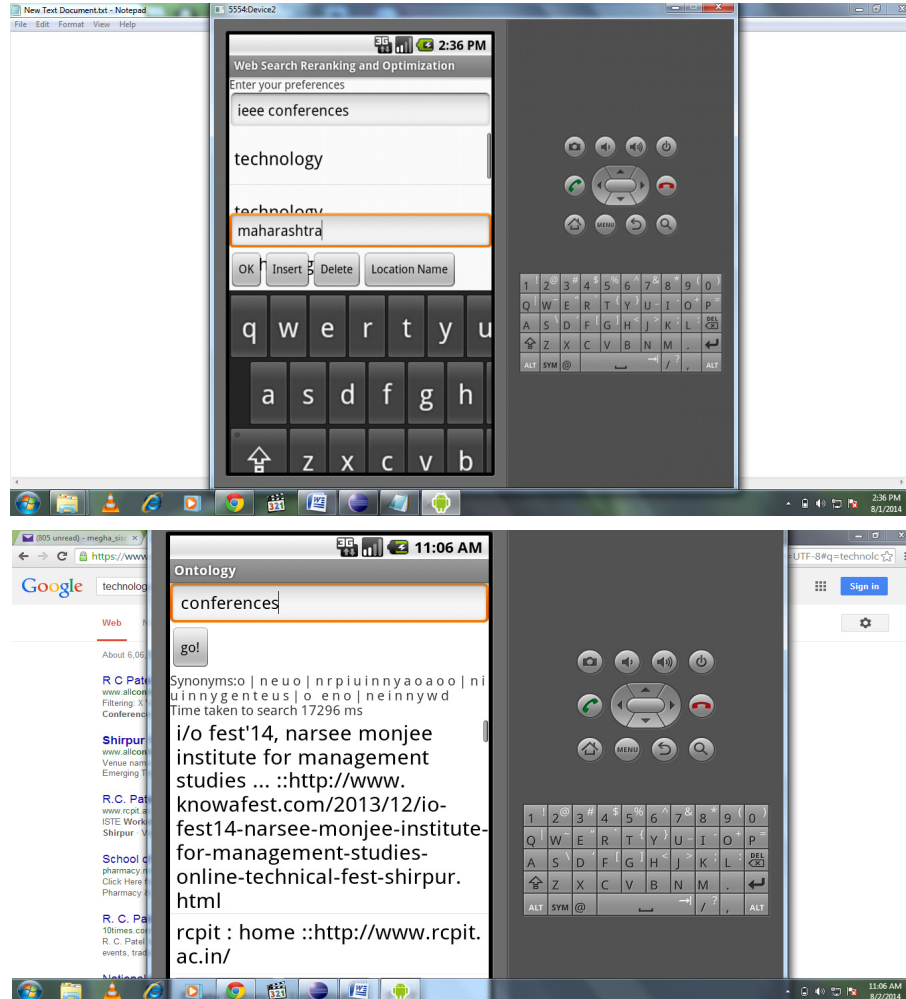Figure 6. Results of Google search engine for the same sample query

Figure 7.  Results of synonym based query recommendation for another sample query
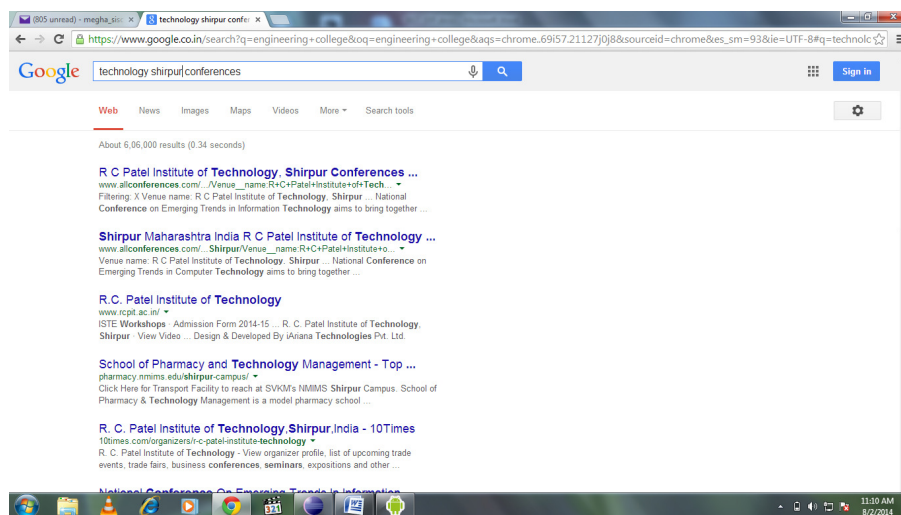


Figure 6.  Results of Google search engine for the same sample query

## 6. CONCLUSION

Mostly many query recommendation systems try to use the previous queries which are similar in either manner with current query. But these methods lack to exactly understand users information need. In order to improve the performance the synonyms and location information is added with snippet information. Global and local scale snippet click models have been used with google search engine log data along with synonyms, which are retrieved for the query by online synonym service in addition with user preferences. By analyzing the results returned by search engine google and compare these results with synonym based query recommendation system. It has been observed that synonym based query recommendation is more efficient. In addition to that synonym based query recommendation performs better for low frequency query. In future, we hope to extend this approach to make use of correctly identified intent for query rewriting by fetching users current location automatically to improve the performance of searching.

### REFERENCES

[1]    Liu, Miao, Zhang, Ma, and L. Ru, (2011)  "How do users describe their information need: Query recommendation based on snippet click model", Expert Systems with Applications, vol. 38(11), pp. 13847-13856.

[2]    Baeza-Yates, Hurtado, and M. Mendoza, 2005 "Query recommendation using query logs in search engines", Current Trends in Database Technology-EDBT 2004 Workshops, Springer Berlin Heidelberg.

[3]    Cucerzan, and R. White, 2007 "Query suggestion based on user landing pages", Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval, ACM.

[4]    Xiaoyan, Bo, Junliang, and M. Xiangwu, 2008 "An effective method for chinese related queries recommendation", In Software Engineering, Artificial Intelligence, Networking, and Parallel/Distributed Computing, 2008, SNPD'08, Ninth ACIS International Conference on, pp. 381-386, IEEE.

[5]    He, Jiang, Liao, Hoi, Chang, Lim, and H. Li, 2009 "Web query recommendation via sequential query prediction", In Data Engineering, 2009, ICDE'09, IEEE 25th International Conference on, pp. 1443-1454, IEEE.

[6]    Zahera, El Hady, and W. El-Wahed, 2010 "Query Recommendation for Improving Search Engine Results", In World Congress on Engineering and Computer Science (WCECS), San Francisco, USA, vol. 1.

[7]    Sumathi, Padmaja Valli, and T. Santhanam, 2010 "Automatic recommendation of web pages in web usage mining", International Journal on Computer science and Engineering (IJCSE), vol. 2, pp. 3046-3052.

[8]    Goyal, and N. Mehala, 2011  "Concept based query recommendation", In Proceedings of the Ninth Australasian Data Mining Conference, vol. 121, pp. 69-78, Australian Computer Society, Inc.

[9]    Wen, Nie, and H. Zhang, 2001 "Clustering user queries of search engine", In Proceeding of the 10th international conference on World Wide Web, 2013, pp. 162-168.

[10]  Zaiane, and A. Strilets, 2002 "Finding similar queries to satisfy searches based on query traces", In Advances in Object-Oriented Information Systems, pp. 207-216, Springer Berlin Heidelberg.

[11]  Baeza-Yates and B. Ribeiro-Neto, 1999 "Modern information retrieval", vol. 463, New York: ACM press.

[12]  Eugene, Brill, and S. Dumais, 2006 "Improving web search ranking by incorporating user behavior information", In Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval, pp. 19-26, ACM.