# An Efficient Approach for Ranking of Citations Using Graph Based Model

Deepti Kapila

*Research Scholor,Computer science department*
*RIMT, Mandi Gobindgarh*


Prof. Charanjit Singh

*Assistant Professor,Computer science department*
*RIMT, Mandi Gobindgarh*

**Abstract-   The main essential of the ranking of research papers is relevancy and performance of the papers. Number of ranking algorithms are being used for providing the performance in the network, but all the algorithms face the problem of better output and relevance accordingly time factor. A new algorithm has been developed that can completely remove the problem of research paper's performance with less time. Citation count and content count algorithm is very less implemented together. CCA (Content Citation Author) count algorithm can be defined as self-organized ranking algorithm with citation, content and author rank support. Content Citation Author Count (CCA) ranking is proposed with a view to provide a relevancy order to the search results. This algorithm is a combination of two basic ranking methods i.e. Citation Count ranking and Content based ranking. Also, the ranking system aggregates different ranking scores to produce a new one, which reflects the contribution of the individual scores all together i.e. CCA ranking, Page Rank, time based citation count ranking. In this ranking system, Static and Dynamic ranking is employed to calculate the final rank of the research paper. Static rank is a constant rank and calculated offline while the dynamic rank depends upon the query fired and calculated online. From the study, it has been found that developing a CCA algorithm increases the performance of the papers. On comparing both the networks citation count and content count was found to have better performance than citation count in all aspect.**

**Keywords – Page rank, Citation, CCA ranking, Content**

## I. INTRODUCTION

Numbers of research papers are published every year and these papers span various fields of research. For a new researcher, it becomes a very difficult task to go through the entire repository of research papers in order to determine the important ones. There can be several ways of determining whether a research paper is important depending on the field of work, conference of publication, etc.

An efficient ranking algorithm is important in any information retrieval system. In spite of advances in search engine technologies, there still occur situations where the user is presented with non-relevant search results. For example, when a user inputs a query for some scientific literature, book to a search engine such as Google, it returns a long list of search results consisting of tutorials, news, articles, blogs etc. This happens due to limited crawling by the search engines. As user wants to get the results in short span, it is necessary to rank the pages which are relevant according to the user's input query. To overcome this problem, a new algorithm has been introduced to make retrieval mechanism more effective and relevant for researchers or users.

## II. ORGANIZATION OF PAPER

 The organization of paper is divided into sections as follows: we first introduce about some related work regarding Page Ranking approaches in section III. In section IV, we present our proposed system in which we describe the overall architecture of the system and a framework that computes ranking of research papers on the basis of citations, author rank and content rank at the search engine or by query. Section V and VI, includes the comparison of citation rank algorithm with the proposed algorithm in context of content availability and response time. Finally in section VII and VIII, we conclude the paper and discuss some future directions for the system.

## III. RELATED WORK

This section discusses about various ranking algorithms for ranking the research papers. Research papers have many features based on which different rankings could be performed. These features are citations to the publication, content, authors, publication year and journal of the publication etc. As per research, we concluded that different digital libraries rank their results on the basis of different factors .for example: IEEE Explore offers a ranking on the basis of title; ACM Digital Library gives the choice to select the ranking based upon publication year, citation counts, alphabetically by title or journal and relevancy. All these digital libraries use different ranking algorithms to rank their papers. A brief description of various existing ranking algorithms [3,4] is given below:-

*PageRank Algorithm:*- Surgey Brin and Larry Page [9,10] proposed ranking algorithm in which outgoing links from the paper are also considered along with the incoming links. The weightage of the incoming link is higher if the link is coming from an important paper. To calculate the rank of the paper by this algorithm, a formula is given below-

$$PR\ u = 1 - d + d\ PR(v)N_v \qquad (1)$$

Where u represents a paper, B(u) is the set of papers that point to u, PR(u) and PR(v) are rank scores of papers u and v respectively, $N_v$ denotes the number of outgoing links of paper v, and d is a normalization factor which lies between 0 and 1.

*Page Content Rank Algorithm*: - Jaroslav Pokorny et al.[5] gave a ranking method of page relevance ranking employing Web Content Mining(WCM) technique, called Page Content Rank (PCR). This method combines a number of heuristics that seem to be important for analyzing the content of web pages. Here, page importance is determined on the basis of the importance of terms contained in the page; while the importance of a term is specified with respect to a given query q. PCR uses a neural network as its inner classification structure. In PCR, for a given query q, resultant papers are in turn classified according to their importance. Here a page is represented in a similar way as in the vector model and frequencies of terms in the page are used.

*Citation Count ranking algorithm:* - One of the most frequent used ranking algorithms for measuring a scientist's reputation, named Citation Count was proposed by Joeran Beel et al. [6]. In this algorithm, the importance of the paper is based upon the number of citations to it. More the number of citations to the paper, higher would be its rank. This is the most commonly used ranking algorithm in digital libraries. Citation count ranking is defined as:

$$CC_i = |I_i| \qquad (2)$$

Where CCi is the citation count of the paper i, |Ii| is the number of citations to paper i.

*Graph based model:* - In graph based model, the pages are sorted according to the importance of citations and author journal. It is based on link analysis, but both citation count and graph based model is not solving any clear purpose to rank or score the papers.

*Popularity Weighted Ranking:-algorithm:* Yang Sun and C. Lee Giles [11] proposed a new Page Rank Algorithm with improved performance, named as *Popularity Weighted Ranking algorithm.* It came with the concept of popularity of the venue of publication. It means this algorithm considers the importance of the venue along with the weighted of incoming links to the paper.

Although many ranking algorithms have been proposed in the literature, there still exist many problems which need attention. Some of the existing algorithms rely on the link structure of the publications (web structure mining), whereas others look for the content in the publications (web content mining), while some use a combination of both. As per the review conducted, following problems have been discovered.

- Existing citation count algorithm does not take into account the importance of citing paper.

- Next problem is that Graph based model relies only on relationships between directed graphs and their nodes, but it totally ignores important pages related to contents. It takes extra calculations to find the author ranking and the time impact of citations.

- Importance of the paper/article to be ranked is totally ignored in this method. Importance of the paper is depends upon mainly three factors: (i) how many papers referred the paper/article, (ii) how many browsed the article, and (iii) in which journal or conference, the paper is published.

<div align="center">IV. PROPOSED ALGORITHM</div>

For better search experience in digital libraries, a novel ranking method for ordering the research publications is proposed. In this section, the detailed description of the proposed ranking mechanism and various modules used in the system is given.

### A. CONTENT CITATION AUTHOR COUNT RANKING ALGORITHM

This algorithm is a combination of two basic ranking methods i.e. Citation Count ranking and Content based ranking. Citation Count is a frequently used ranking algorithm for measuring a scientist's reputation. This method uses the citation graph of the web to determine the ranking of scientific work. This method states that if a publication has more number of citations (incoming links) to it than publication become important. But some authors misuse it to have their publications highly ranked during search. To overcome this problem, the proposed algorithm uses the content of the paper which cited the publication along with the number of citations. In this algorithm, the relevancy score between the publication and the paper which cited the publication is computed on the basis of their content. To check whether the papers are related or not, it uses the summaries instead of comparing the whole content of the papers.

*Objectives of proposed method:*

The major objectives of the proposed method could be summed up as follows:

- To study various citation counts using ranking algorithms.

- To develop optimized algorithm using citation count and graph based model.

- To reduce retrieving time and citation limitation of back links of web pages.

- Compare and analysis of algorithm with citation algorithm.

### B. METHODOLOGY FOR PROPOSED   ALGORITHM

The ranking algorithm methodology is divided into 5 phases to achieve the desired goal:

Phase 1: Collect the required information in this phase, and implement the layout of the project and create a database that will be queried for static and dynamic implementation.

Phase 2: Further, implementation of citation count will be included and implement citation count algorithm, then all the citations will be assigned to all pages in the database according to their links.

Phase 3: In this phase graph based i.e. author score will be implemented, pages will be searched according to query and based on the author and its score and then rank will be assigned to the pages.

Phase 4: Citations and author ranking will be combined to create a hybrid approach and that will give us final results.

Phase 5: Final result will be validated and will be compared with citation count and graph model i.e. link analysis and author ranking.
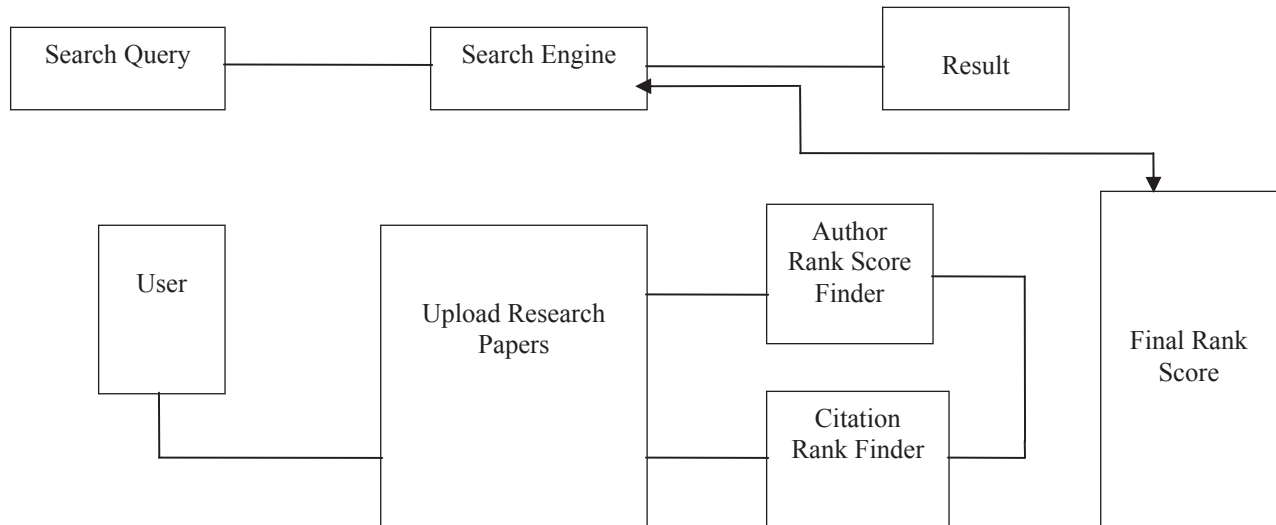The proposed methodology of an algorithm is described below in figure 1.

Figure 1.        Framework of proposed methodology

## V. RESULT ANALYSIS OF CCA ALGORITHM

This section presents the implementation details and the experimental results that have been performed over the proposed ranking system. Various queries from different contexts were submitted on the system to carry out an analysis. A comparison of the results obtained by using the proposed algorithm with the existing algorithms was done and the observation obtained thereof were analyzed to prove the efficiency of the proposed system.

I.   *Home page with upload new papers or upload papers by submitting queries*

Fig 2 (a) displays the home page of the proposed system for digital libraries. This page provides two options:-
- *Upload the new research paper*

- *Search papers by submitting queries*

In the upload section, a new research paper is uploaded in the database from the repository. The upload section is protected through password as this section can be accessed by authentic users or administrators.
Depending upon the success or failure of the uploading action, different outcomes are returned to the user. In the search section, upon submitting the user query, if the research papers related to the query exist in the database then their links are returned to the user otherwise an error page is displayed to the user.
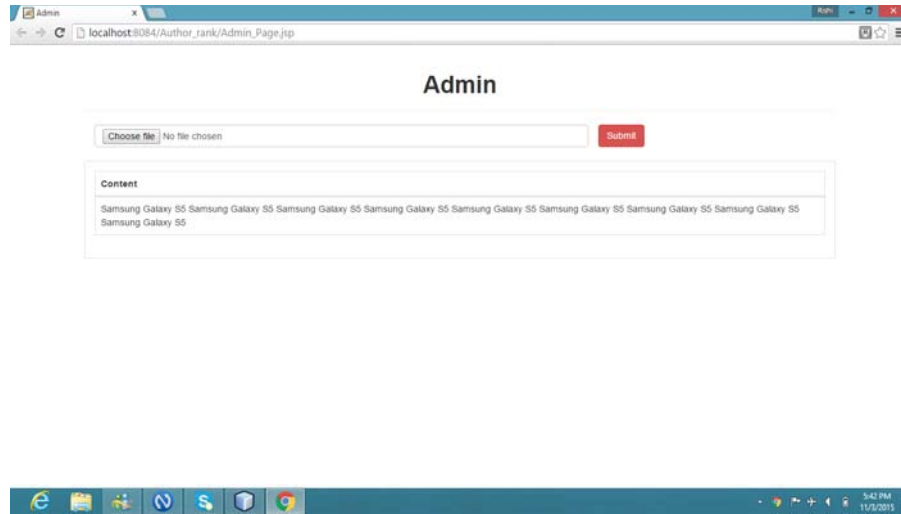
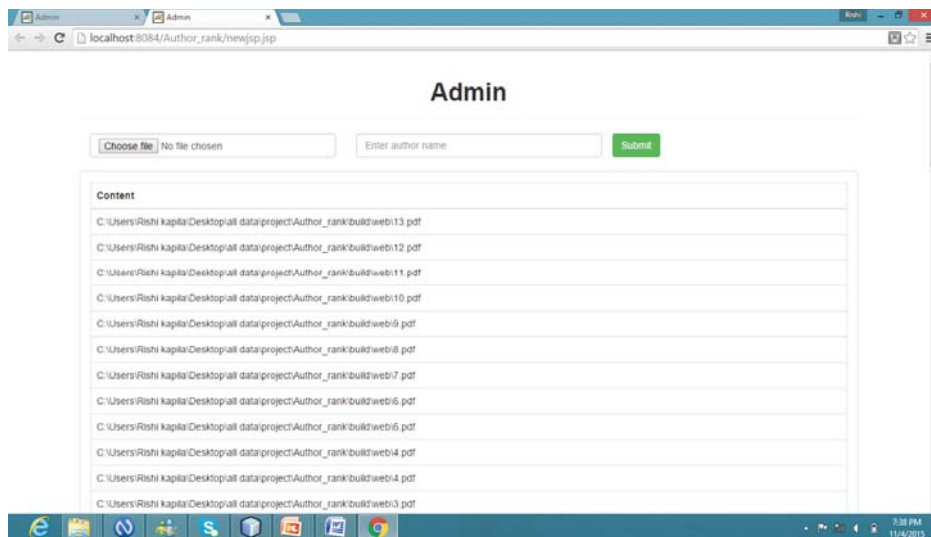Figure 2(a) Home Page of the proposed system



Figure 2 (b) Uploaded papers for admin

I.    *(a) Result on the basis of author rank*

If the paper does not exist in the database then it goes through various processing modules. First of all, information about the research paper i.e. authors, titles, references etc is extracted and stored in the database. This information is stored in the database as shown in Fig. 3(a). Fig 3(a) shows the information like title, authors, papered of all the uploaded papers stored in the Content Store database.

II.    *(b) Result on the basis of reference details*
Fig. 3(b) shows the reference details of all the uploaded papers. The highlighted portion shows the reference details of the newly uploaded paper with its papered 20.
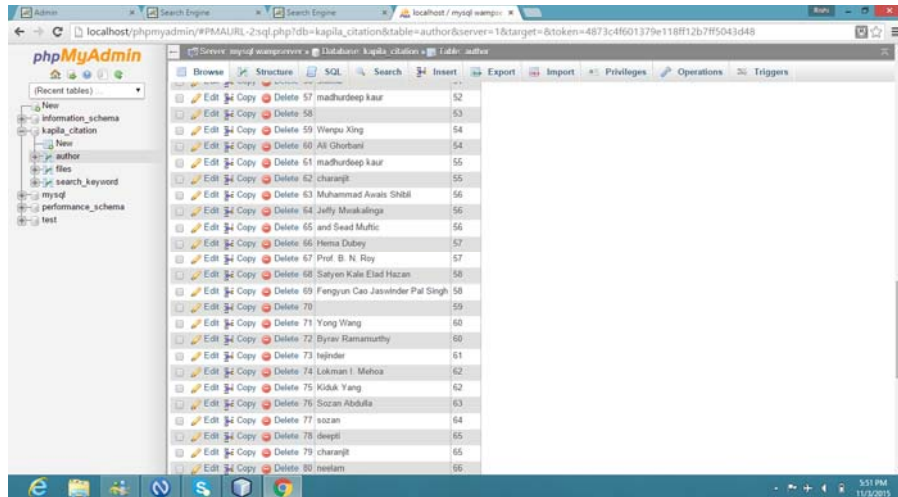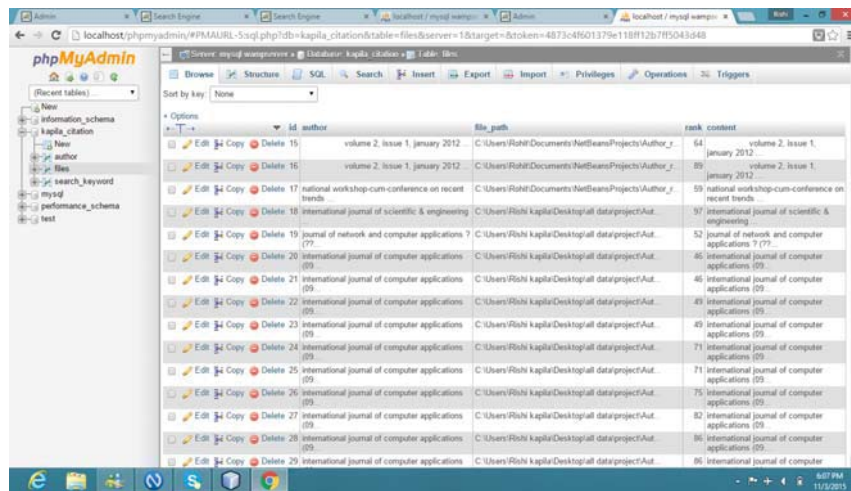
Figure. 3(a) Research papers on the basis of author rank



Figure 3(b) Research papers on the basis of reference details

III.  *Result on the basis of query at search engine*

After submitting the query, there are finally two possible outcomes:-

- *Search engine page*

- *Rank of papers with citation value*

User has two options either to view the paper or to view the summary of the paper. Firstly, User can select the paper from the list according to his requirement. Then, he clicks the option i.e. view paper and view summary according to his wish. The relevant result through query is shown in figure 3(c).
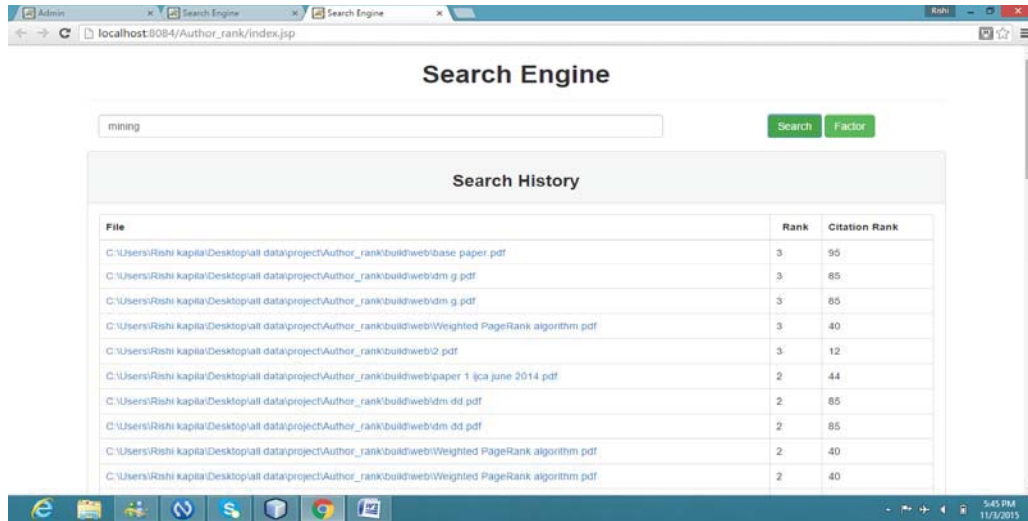
Figure. 3(c) Research papers on the basis of submitting query at search engine

## VI. COMPARISON OF PROPOSED ALGORITHM WITH CITATION COUNT

In this section, the results of the proposed *CCA ranking algorithm* are compared with the existing *Citation Count (CC) ranking algorithm* [19,37] with the help of graph. The figure 4(a) is a graphical representation of response time according to user query. If user searches for mining keyword the research papers of mining keywords are displayed and the value of one research paper of keyword mining is compared to another research paper related to keyword mining**.** If the paper value 1 is compared to paper value 2 and the number of URL are taken 3 then content availability graph displays how two papers are different from their contents to each other. The figure 4(b) shows graph for content availability with paper values.

The comparison of ranking algorithms can be done on the basis of number of type of mining, I/P parameters, space required etc. The table 1 shows the comparison between three algorithms named as citation count, content based and CCA (content citation author) rank algorithm. The comparison is done on the basis of different criteria i.e. basic description, type of mining, processing time, space requirement, degree of relevance, input parameters and scanning options.
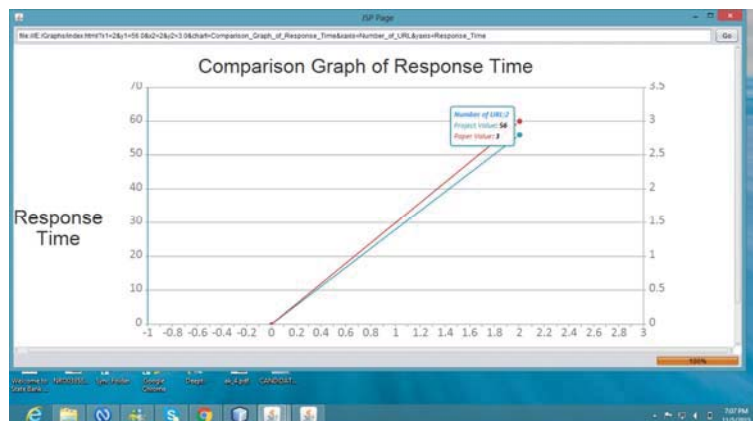


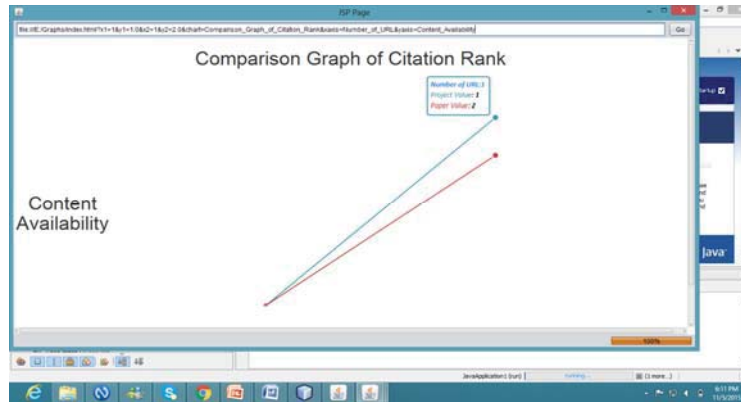Figure. 4(a) Graph Representation of response time with paper value

Figure. 4(b) Graph representation of citation rank with paper value and their content

## VII. CONCLUSION

A new ranking algorithm, Content Citation Author Count (CCA) ranking is proposed with a view to provide a relevancy order to the search results in the digital libraries. This algorithm is a combination of two basic ranking methods i.e. Citation Count ranking and Content based ranking. Also, the ranking system aggregates different ranking scores to produce a new one, which reflects the contribution of the individual scores all together. It may be noted that as compared to existing algorithms like page content rank and citation count rank, the proposed CCA ranking algorithm provides more informative and relevant results about a query. The more significant papers with respect to their content and incoming links are better identified by the CCA ranking algorithm as compared to the existing page content rank and citation count ranking algorithms .In the proposed ranking system, the relevancy of the papers to a given query is better determined, as compared to the existing ranking systems because the dynamic rank measures the similarity between the query and research papers and assign a relevancy score to papers. The existing ranking algorithms, which are commonly used to rank the important papers/articles in digital libraries, rely on links. But the proposed CCA ranking algorithm takes the importance of both the links and relevancy of the papers/articles.

| Criteria for comparison | Citation Count Ranking algorithm | Content Based Ranking algorithm | CCA Ranking Algorithm |
|---|---|---|---|
| Basic Description | Rely only on links | Rely only on content of the paper. | Rely on links as well as content of the paper. |
| Type of Mining | Structure Mining | Content Mining | Combination of Structure and Content Mining |
| Degree of Relevance with Query | Does not check the relevancy | Checks the relevancy | Checks the relevancy |
| Different Scanning Options | No Scanning | Scans the full Paper | Scans only the Summary of the Paper |
| Processing Time | Low | High | Medium |
| Space Requirement | Low | High | Medium |
| I/P Parameters | Backlinks | Paper | Backlinks and Summary |

Table 1 Comparison of citation count and content count with CCA rank

## VIII. FUTURE SCOPE

Although the proposed approach for ranking the papers seems effective, there is still some scope to improve the retrieval of relevant information contained in digital libraries. As part of future work, an efficient page ranking algorithm in terms of time response, accuracy of results, importance of the results and relevancy of results should be developed so that the quality of search results can be improved.

Moreover, the proposed ranking utilizes the content and structural information about the papers to calculate the rank, the users' browsing information stored in historical logs can be utilized to provide more effective and relevant results to the user. Thus, a new ranking algorithm could be developed by considering two or more such factors.

## IX. ACKNOWLEDMENT

REFERENCES

[1]  Lawrence, S. and Gills; Accessibility of information on the web. C.L.1999.

[2]  A. Arasu, J.Cho, H. Garcia-Molina, A. Paepcke and S. Raghavan; Searching the web. ACM transactions on internet technology; 2001.

[3]  Lawrence Page, Sergey Brin, Rajeev Motwani and Terry Winograd; The pagerank citation ranking bringing order to the web. technical report, computer science department, Stanford university; 1998.

[4]  RankDex; The RankDex search engine. available online at http://rankdex.gari.com/.

[5]  Tomasic, A. and Garcia Moline, H. 1993; Performance of inverted indices in shared nothing distributed text documents information retrieval system. In preceeding of 2nd international conference on parallel and distributed information system; jan-1993.

[6]  M. Henginger; Hyperlink analysis on the web. 2003; available online at http://www.cad.eecs.berkeley.edu/tah/170/Notes/170-google.ppt

[7]  Craig Silverstein, Monika henginger, hannes marais and Michael moricz; Analysis of very large altavista query log. tech. Report 1998-014, Digital SRC, 1998.

[8]  J.R. Seeley; The net of reciprocal influence: A problem in treating sociometric data. Canadian Journal of psychology, 1949, pp. 3:234-240.

[9]  Krishna Bharat and George A. Mihaila; When experts agree: using non- affiliated experts to rank popular topics. World Wide Web, 2001, pp. 597-602.

[10]  Naresh Barsagade; Web Usage Mining And Pattern Discovery: A Survey Paper. CSE 8331,2003.

[11]  M. Krishnamurthy; Open access, open source and digital libraries: A current trend in university libraries around the world. A General Review, Emerald Group Publishing Limited.

[12]  R. Cooley, B. Mobasher and J. Srivastava, Web Mining: Information and Pattern Discovery on the World Wide Web. Proceedings of the 9th IEEE International Conference on Tools with Artificial Intelligence, pp. 558- 567, 1997.

[13]  N. Duhan, A.K. Sharma and K.K. Bhatia; Page Ranking Algorithms: A Survey. Proceedings of the IEEE International Conference on Advance Computing, 2009.

[14]  R. Kosala, and H. Blockeel, Web Mining Research: A Survey. SIGKDD Explorations, Newsletter of the ACM Special Interest Group on Knowledge Discovery andData Mining, vol. 2, no. 1, pp. 1-15, 2000.

[15]  A.Gulli, A.Signorini; The Indexable web is More than 11.5 Billion pages. Proceedings of the 14th World Wide web Conference, 2005

[16]  A.Broder; Web Searching Technology Overview. Advanced school and Workshop on Models and Algorithms for the World Wide web, 2002

[17]  B.J.Jansen, A.Spink, J.Bateman, T.Saracevic; Real life Information Retrieval: A study of user queries on the web. ACM SIGIR Forum, 1998.

[18]  Joeran Beel, Bela Gipp; Google Scholar's Ranking Algorithm: The Impact of Citation Counts (An Empirical Study). In Rcis 2009: Proceedings of The IEEE International Conference on Research Challenges In Information Science, 2009.

[19]  L. Marian, M. Rajman; Ranking Scientific Publications Based on Their Citation Graph. Master Thesis, CERNTHESIS, 2009.

[20]  L. Marian, J. Yves LeMeur, M. Rajman, M. Vesely; Citation Graph Based Ranking in Invenio.  ECDL, pp. 236-247, 2010.

[21]  Sergey Brin and Larry Page; The anatomy of a Large scale Hypertextual Web Search Engine". In Proceedings of the Seventh International World Wide Web Conference, 1998.

[22]  A. Sidiropoulos, Y. Manolopoulos; Generalized Comparison of Graph-based Ranking Algorithms for Publications and Authors. Journal of Systems and Software archive, vol. 79, issue 12, pp. 1679-1700, 2006Y.

[23]  Sun, C. L. Giles; Popularity Weighted Ranking for Academic Digital Libraries. In 29th ECIR, pp. 605-612, 2007.

[24]  Kleinberg J.; Authorative Sources in a Hyperlinked Environment.  Proceedings of the 23rd annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 1998.

[25]  Math Explorer's Club; The Mathematics of Web Search. http://www.math.cornell.edu/~mec/Winter2009/RalucaRemus/Lecture4/lecture4.html