

Adding the Temporal Dimension to Search - A Case Study in Publication Search

Philip S. Yu
IBM T.J. Watson Research Center
19 Skyline Drive
Hawthorne, NY 10532
psyu@us.ibm.com

Xin Li, Bing Liu
Department of Computer Science
University of Illinois at Chicago
851 South Morgan Street,
Chicago, IL 60607-7053
{xli3, liub}@cs.uic.edu

Abstract

The most well known search techniques are perhaps the PageRank and HITS algorithms. In this paper we argue that these algorithms miss an important dimension, the temporal dimension. Quality pages in the past may not be quality pages now or in the future. These techniques favor older pages because these pages have many in-links accumulated over time. New pages, which may be of high quality, have few or no in-links and are left behind. Research publication search has the same problem. If we use the PageRank or HITS algorithm, those older or classic papers will be ranked high due to the large number of citations that they received in the past. This paper studies the temporal dimension of search in the context of research publication. A number of methods are proposed to deal with the problem based on analyzing the behavior history and the source of each publication. These methods are evaluated empirically. Our results show that they are highly effective.

1. Introduction

Much of the impact of the Web to society is due to the success of Web search engines. The objective of search engines is to find the most relevant pages given a user query. PageRank [6] and HITS [14] are motivated by the observation that a *hyperlink* from a Web page to another is an implicit conveyance of authority to the target page. One can use these algorithms to find important Web pages.

However, an important factor that is not considered by these techniques is the timeliness of search results. The Web is a dynamic environment. Quality pages in the past may not be quality pages now or in the future. In this paper, we study search from the temporal dimension, which is important due to the following reasons:

1. Users are often interested in the latest information. Except for the Web pages that contain well-established facts and classics, most contents on the Web change constantly. New contents are added; ideally, outdated contents are deleted. However, in practice many outdated links are not deleted. This fact prevents the ranking algorithms from retrieving the updated results.
2. Existing Web page evaluation techniques basically favor pages that have many in-links. Thus, older pages are favored because they tend to accumulate more in-links due to longer existence. In contrast, new pages that are of high quality will not be ranked high.

We believe that dealing with the problems related to the temporal dimension of search is of great importance to future developments of search technology. In this paper, we take a step towards this direction.

To understand the issues in greater detail, we coarsely classify Web pages into two types, old pages and new pages. Obviously, there are pages in between. Let us ignore them for simplicity of explanation.

Old pages: These are the pages that have appeared on the Web for a long time. We can also classify these pages into *quality pages* and *common pages*. Quality pages usually have a large number of in-links, and common pages do not have many in-links. Old quality pages can be further classified from the temporal dimension:

1. *Old quality pages that are up-to-date:* As time goes by, the authors of the pages update the contents to reflect the latest developments. Such pages often stay as quality pages.
2. *Old quality pages that are not up-to-date:* The authors of these pages do not update their contents over time. These pages become outdated, and receive fewer and fewer new in-links over time. However, if many Web users do not clean up hyperlinks to these pages, they may still maintain a

sizeable in-links, and would be ranked high in spite of their low value.

Regarding old common pages, we can also classify them into two types from the temporal dimension:

1. *Old common pages that remain common pages:* Most pages on the Web are such pages. As time goes by, they are still common pages, as they do not receive many in-links.
2. *Old common pages that have become important:* These pages were not important in the past, but as time goes by they become valuable pages. This transition may be due to a number of reasons, such as fashion change or quality contents being added.

New pages: These are pages that appeared on the Web recently. New pages can also be grouped into categories:

1. **New quality pages:** These pages are new and are of high quality. However, they received few or no in-links because they are new.
2. **New common pages:** These pages are new and common.

Unlike an older page, a new page receives few or no in-links. It is thus difficult to judge if it is a quality page.

In summary, for a search engine to consider the temporal dimension and the dynamics of the Web, two new problems need to be dealt with in page evaluation:

1. How to assign lower importance scores to those old quality pages that are not up-to-date.
2. How to assign higher importance values to those new quality pages.

Both these cases present difficulties to the PageRank or HITS algorithm. Here we attempt to deal with these problems by taking time into consideration in evaluating the page quality. In this work, we investigate these problems in the context of research publication search because of following three reasons:

1. Results in the research publication domain can be objectively evaluated as we can count the number of citations received by a paper in the “future” (from test data) to see if our evaluation is accurate. Future citation count of a paper is a commonly used indicator for its quality. Given a collection of papers, all the citation data is available. In contrast, on the Web it is hard to know when a particular hyperlink is installed.
2. Concepts and entities in both domains are similar. For example, a research paper corresponds to a Web page, and a citation of a research paper corresponds to a hyperlink in a Web page.
3. Publication search is important and useful in its own right. With growing popularity of digital libraries on the Web, searching for relevant and updated publications is becoming increasingly valuable.

Of course, there are differences between Web pages and research publications. For example, on the Web a

page may be deleted, but a published paper cannot be deleted. Hyperlinks can also be added and deleted from a Web page, while for a research paper, once published no reference or citation can be modified. However, the main issues are essentially similar in both domains.

In this paper, we perform a study of the citation based evaluation of research papers, which corresponds to the hyperlink based evaluation of Web pages. We present a number of methods to mine the temporal behavior of each publication. The techniques are evaluated experimentally; the results show that the proposed methods are highly effective.

2. Related Work

Since the PageRank [6] and HITS [14] algorithms were published, a large number of papers on improvements, variations, and speed-up of the algorithms have appeared in literature [1][3][4][5][6][10][11][13][15]. Many applications of the algorithms have also been reported, in both Web search and research publication search, e.g., Web resource discovery [6][3], search considering both hyperlinks and page contents [8], and research paper search [15]. These works are still within the framework of the original algorithms, and do not consider the temporal aspect. [9][16] study the evolution of the Web and identify the same problem as we discussed above. [2] also identifies the problem and makes a limited attempt to tackle the problem with no evaluation.

[15] describes the CiteSeer system, a Web digital library for research publications. It also uses PageRank and HITS algorithm to rank papers by either “hub” or “authority” score. [15] mentions that the temporal aspect should be considered in publication search. However, the topic was not further investigated.

3. The Proposed Techniques

There are many factors that contribute to the ranking of search results. Broadly speaking, we can group them into content based factors and reputation based factors.

Content based factors: These factors are related to the content of documents. They determine the degree of relevance of a document to the user query.

Reputation based factors: Reputation of documents helps to determine the ranking of the relevant documents retrieved. In the context of publication search, reputation factors include the citation count of the paper, the reputation of its authors and journals.

This paper focuses on reputation based factors and studies how the temporal dimension may be included in the evaluation of a research paper’s reputation.

There are two main timing factors for a research paper.

1. The publication date, and
2. The dates that the paper is cited by other papers.

The major algorithm that evaluates the reputation of a paper is PageRank, which is based on citations (or hyperlinks in the Web context) of a paper. To consider time, an obvious approach is to include time in the algorithm. We also describe a linear regression based method to deal with the problem later.

3.1 TimedPageRank

Before describing the TimedPageRank technique, we first introduce the original PageRank algorithm[6]. It is also applicable to research papers. The PageRank (PR) score of a page/paper A is:

$$PR(A) = (1 - d) + d \times \left(\frac{PR(p_1)}{C(p_1)} + \dots + \frac{PR(p_n)}{C(p_n)} \right) \quad (1)$$

where

- $PR(A)$ is the PageRank score of paper A ,
- $PR(p_i)$ is the PageRank score of paper p_i that links to paper A ,
- $C(p_i)$ is the number of outbound links of paper p_i , and
- d is a damping factor, ranging between 0 and 1.

In this work, we still use the damping factor of 0.85, which was used in the original PageRank paper [6]. Initially, the PageRank score for each paper is set to 1. The calculation is done in an iterative fashion until the results finally converge.

We now describe the TimedPageRank technique. Since we are interested in the importance of a paper now, a citation occurred a few months ago is clearly more important than one occurred a few years ago. We modify the PageRank technique by weighting each citation. The system calculates the time-weighted PageRank (PR^T) value for each paper as follows. Equation (2) is a modified

$$PR^T(A) = (1 - d) + d \times \left(\frac{w_1 \times PR^T(p_1)}{C(p_1)} + \dots + \frac{w_n \times PR^T(p_n)}{C(p_n)} \right) \quad (2)$$

version of equation (1). In this equation, w_i is the time based weight for each citation. Its value depends on the citation date from paper p_i to A , which is also the publication date of p_i . The earlier the citation occurred, the smaller the weight is. Since exponential average is extensively used in time-series prediction, we choose to decay the weights exponentially according to time,

$$w_i = DecayRate^{(y-t_i)/12}$$

where y is the current time, t_i is the publication time of paper p_i and $(y-t_i)$ is the time gap in months. $DecayRate$ is a parameter. While the effect of $DecayRate$ on the prediction results is further studied in section 4.4, we use 0.5 in the following example to illustrate the concept. For instance, in our training data, the newest papers are

published in December 1999. The citations occurred in December 1999 and December 1998 have the weights of 1 and 0.5 respectively. Note that if $DecayRate$ is 1, the time-weighted PageRank algorithm will be the same as the original PageRank algorithm. Therefore, the $DecayRate$ parameter could be tuned according to the nature of a dataset/the user. When its value is close to 1, the weight decreases slowly with time. It is more suitable for static domains or users that are new to the domain.

Weighting each citation considers recent citations more important. It thus only assesses the importance of a paper from the past. We are also interested in the potential importance of the paper in the future. To evaluate this, we introduce the other parameter called *trend factor*.

Continuing our previous example, for a paper A , $PR^T(A)$ already captures the importance at the end of 1999. How does the importance change through the future year? We assume that this is reflected by the citation change at the end of 1999. Therefore, we mine the past behavior of a paper A to compute its *trend factor* $Trend(A)$:

1. Data preprocessing: We filter out those papers whose citations are lower than one per month because they are not likely to last as time goes by. Due to the same reason, we assign the minimum *trend factor* to them. To make a reliable prediction, we also smooth out the noise in the citation data by using “moving average” of monthly citation data. The moving average citation of a paper in a given month is calculated by averaging its citations in that month and the previous month.
2. Compute the *trend factor* from citation change at the end of the most recent year (i.e. 1999). For a paper A , if its citation count for the 3rd quarter of 1999 is n_t , and its citation count for the 4th quarter is n_f , then we say the trend ratio for A , $r(A) = n_f / n_t$.

If a paper’s age is less than 3 months, there is no sufficient data available to compute its trend ratio. However, we can solve this problem with our source evaluation method, which will be covered later.

After computing all the $r(P)$ for all the papers, we normalize them, so the normalized values are between *minimum trend factor* and 1. The *minimum trend factor* is set to 0.5; the reason is that for any paper, the weight for each previous citation will reduce by half through one year, as in equation (2). The normalized value of $r(P)$ is the *trend factor* of paper P , $Trend(P)$.

Paper A ’s final TimedPageRank (TPR) is:

$$TPR(A) = Trend(A) * PR^T(A) \quad (3)$$

where $PR^T(A)$ is computed using equation (2).

3.2 Source Evaluation: Author and Journal

Although TimedPageRank considers time, it is not

sufficient because the technique is not useful for new papers that have few or no citation from other papers. To assess the potential importance of a new paper, two pieces of source information are useful, the reputations of its authors and its journal. We can make use of time-weighted PageRank to evaluate these two reputations.

Author evaluation: The reputation of an author is based on the research papers that he/she published in the past. We compute author evaluation by averaging the time-weighted PageRank values of all his/her past papers. Let the papers published by author a_j in the past be p_1, p_2, \dots, p_m , the author score $Author(a_j)$ is:

$$Author(a_j) = \frac{\sum_{i=1}^m PR^T(p_i)}{m}$$

where $PR^T(p_i)$ is the time-weighted PageRank score of paper p_i . Here $PR^T(p_i)$ is used rather $PR(p_i)$ as we consider that more recent citations are more representative for the current reputation of the author.

Journal evaluation: A score, $JournalEval(j)$, is also assigned to each journal j by considering papers published by the journal in the past.

Using the author and journal evaluations, we can evaluate paper importance in a variety of approaches.

One way to evaluate paper is based on its authors. Let the authors of the paper p be a_1, a_2, \dots, a_k . The paper's score can be estimated based on the weighted average of author scores:

$$AuthorEval(p) = \frac{\sum_{i=1}^k (Author(a_i))^2}{\sum_{i=1}^k Author(a_i)} \quad (4)$$

Another way is to combine author evaluation and journal evaluation to score each paper. Assume that paper p is published in journal j . We compute a simple average of the author evaluation and journal evaluation scores:

$$AJEval(p) = (JournalEval(j) + AuthorEval(p)) / 2$$

Of course, there are many other ways to compute the combined score. One of alternatives is to compute a weighted average of $JournalEval(j)$ and $AuthorEval(p)$:

$$AJEval(p) = \frac{(JournalEval(j))^2 + (AuthorEval(p))^2}{JournalEval(j) + AuthorEval(p)}$$

We will compare these two variations later on.

It should be noted that after a paper has been published for a while, it is more effective to use TimedPageRank to score the paper than author/journal evaluation. The reason is that author/journal evaluation is only an averaged result of all the papers from the authors/journal.

3.3 Linear Regression

Another simple technique that can be used to score a

paper is linear regression. Specifically, one can use citation counts of the paper received in the latest time period to perform a linear regression to predict the citation count. This predicted citation count can be used as the score of the paper.

If a paper is published only recently, it may have a very few citation. Then, linear regression will not be accurate. In this case, we again used author and journal evaluation to score the new papers. Both author and journal evaluations can be done by actual citation counts of all the papers from the author or the journal. Let the papers published by an author a_j be p_1, p_2, \dots, p_m . Author score is computed as follows:

$$Author(a_j) = \frac{\sum_{i=1}^m count(p_i)}{m}$$

where $count(p_i)$ is the citation count of paper p_i . The score of a paper based on author evaluation is again given by Equation (4).

Journal evaluation can be done similarly. After they are computed, the same method used in section 3.2 can be employed to combine the journal and author scores.

4. Empirical Evaluation

In this section, we evaluate the proposed techniques and compare them with the PageRank algorithm. We use the KDD CUP 2003 research publication data, which is an archive of High Energy Particle Physics publications catalogued by Stanford Linear Accelerator Center.

4.1 Experimental Settings

Our experiments use the standard search paradigm. That is, given a collection of research papers and a user query, the system ranks the relevant papers. Our research focuses on investigating the effect of time on the citation based ranking, so we are not interested in content based factors such as keyword locations, their distances in the paper, etc. We simply assume that a paper is relevant to a query if it contains all the query words

Evaluation method: To evaluate the proposed techniques, we do not compare their rankings directly, which is harder to quantify. Instead, we compare the number of citations that the top ranking papers receive in the following year. i.e., one year after the user performs the search. It is reasonable because the citation count reflects the importance of a paper. If those highly cited papers in the future are ranked high by an algorithm, it indicates that the algorithm is effective in retrieving high quality papers.

Table 1. Comparison results of different methods using all papers

1	2	3	4	5	6	7	8	9	10	11	12	13	14
No. of top papers	Original PageRank		TPR		TPR(AJEval - simple avg)		TPR(AJEval - weighted avg)		LR		LR(AJEval)		Best Citation Counts
10	2516	44%	4236	75%	4312	76%	4382	77%	4219	75%	4136	73%	5661
20	3406	46%	5702	78%	5739	78%	5756	78%	5371	73%	5447	74%	7345
30	4024	48%	6816	81%	6845	81%	6788	81%	6385	76%	6519	78%	8406

4.2 Experimental Results with All Papers

In this set of experiments, we use all the papers from 1992-1999 to perform various evaluations for the proposed methods. We issue 25 queries of frequent physics terms and rank the relevant papers for each query. The data of year 2000 tests various ranking methods.

Table 1 presents the experiment results. Only the results for the top 30 papers are given. The reason for using only top 30 ranked papers is that users seldom have the patience to look at more than even 20 papers.

The experimental results are presented in 3 rows. Each row gives the total citation counts of different methods for a group of papers. The first row is for the top 10 papers (we also call it a group of papers), where the citation count is the summation of the citation counts of all the top ten papers over the 25 queries. Similarly, the second row is for the top 20 papers, and so on.

Below, we explain the results column by column.

Column 1: It lists each group of top ranked papers.

Columns 2 and 3: Column 2 gives the result for each group of top papers based on rankings using the original PageRank algorithm. Each result here is the total citation count of each group of top ranked papers for the 25 queries. Each count is obtained from citations that the paper received in year 2000. Column 3 gives the ratio of the total citation count for this method and the total citation count of the ideal ranking (called *best citation count* in Column 14), expressed as a percentage. The ideal ranking is one that ranks relevant papers based on the actual number of citations received by each paper in the following year.

Columns 4 and 5: Column 4 gives the results (citation counts) of the TimedPageRank (TPR) method.

Column 5 gives the ratio of the total citation count for the TimedPageRank method and the total citation count of the ideal ranking, expressed as a percentage.

From Column 4 and Column 5, we observe that TimedPageRank's results are significantly better than those of the original PageRank algorithm.

Columns 6-7, 8-9: Similar to the previous columns, Columns 6-7, and 8-9 give the corresponding results of TPR (AJEval - simple average) and TPR (AJEval - weighted average). These two methods combined TimedPageRank (TPR) with both author and journal evaluation for new papers, but with different

combination functions. Papers are regarded as new if they were published less than 3 months ago. Both methods performed better than TimedPageRank alone. The reason is that TimedPageRank cannot handle new papers well. We also observe that the weighted average approach performs slightly better than the simple average approach for the top 10 papers, which are of more importance for our search. The reason for the difference between the two methods will be analyzed in details in Section 4.3.

Columns 10-11, 12-13 give the corresponding results of linear regression (denoted as LR) based methods. We can see these two methods perform well also, but the results are not as good as TimedPageRank based methods.

Column 14: It gives the best citation count for each group of paper based on the ideal ranking, i.e., ranking relevant papers based on the actual number of citations received in year 2000.

To summarize, both TimedPageRank and Linear Regression based algorithms perform significantly better than the original PageRank algorithm. Moreover, the author and journal evaluation helps improve the prediction results in most cases. Because the fraction of new papers is small (around 3%), the effect of author and journal evaluation is diluted. We will show the experiment results involving only new papers in Section 4.3. Among all the four methods, the TimedPageRank with author and journal evaluation gives the best performance.

To give some indication on the effectiveness of the proposed ranking methods, we find the top 10 most cited papers in 2000. We then use the proposed methods to rank all the papers appeared from 1992 to 1999. Table 2 shows the ranking results.

Column 1: It shows the ranking of the top 10 papers in 2000.

Column 2: It gives the paper IDs of these papers.

Column 3: It gives the rank of each paper using the original PageRank algorithm.

Column 4: It gives the rank of each paper based on TimedPageRank. The new papers are ranked using the combined author and journal evaluation (Section 3.2).

Column 5: It gives the rank of each paper based on Linear Regression. The new papers are ranked using the combined author and journal evaluation (Section 3.3).

Table 2 clearly shows that the ranking results of the

original PageRank algorithm are quite poor. In contrast, our proposed methods perform remarkably well. All the top 10 papers are ranked very high (within 20).

Table 2. Ranks of the top 10 papers

Rank	Paper ID	Original PageRank	TPR (AJ Eval.)	LR (AJ Eval.)
1	9711200	19	1	1
2	9908142	742	8	5
3	9906064	613	6	10
4	9802150	39	2	2
5	9802109	46	4	3
6	9711162	323	11	7
7	9905111	576	9	4
8	9711165	620	20	14
9	9610043	17	12	19
10	9510017	7	13	8

4.3 Results on New Papers Only

In this set of experiments, we use only the new papers. That is, we only use those papers that are less than three months old at the query time. The purpose here is to assess the effectiveness of author and journal evaluations.

This set of experiments does not use TimedPageRank and Linear Regression because these papers have few citations. Note that we did not use query here because each query returns only a few results, as the number of new papers is small. We use the proposed methods to rank all the new papers. Table 3 lists the results.

Column 1: It lists each group of top ranked papers.

Columns 2 and 3: Column 2 gives the total citation count of each group of top papers based on random ranking.

Random ranking is a reasonable method because the new papers are hardly cited by any other papers. Column 3 gives the ratio of the total citation count for this method and the total citation count of the ideal ranking (Column 18), expressed in percentage. The results of random ranking turn out to be very poor.

Columns 4-5, 6-7: They give the corresponding results of AuthorEval method (4-5) and JournalEval methods (6-7). Here both methods used the time-weighted PageRank (PR^T) values of all the papers of the authors/journals (see Section 3.2). AuthorEval and JournalEval have similar performance, and both outperform random ranking.

Columns 8-9, 10-11: They give the corresponding results of the combined method (AJEval - simple average) and (AJEval - weighted average). Both techniques

combine author and journal evaluations, but with different combination functions. The simple average is the mean of the two scores, while the weighted average result is closer to the higher of the two scores as the higher score has more weight. We can see that two combined methods perform significantly better than the method based on individual criterion.

Usually quality papers are from top authors and top journals. However, it is not always the case. For example, a young researcher with potential may only have a lower author evaluation score, but he or she still is able to publish quality papers in top journals. Similarly, a new journal could be underestimated. Therefore, when there is a large gap between the journal evaluation score and the author evaluation score, the weighted average method remedies the problem, and outperforms the simple average method.

Columns 12-13, 14-15 and 16-17 give the results of the three methods using citation count of published papers to evaluate authors or journals. They also perform better than random ranking but significantly worse than the time-weighted PageRank based evaluation of authors and journals.

Column 18: It gives the total citation count for each group of papers based on the ideal ranking, i.e., ranking papers based on the actual number of citations received by each paper in 2000.

To summarize, we observe that the PR^T based evaluation of authors and journals is better than the citation count based method due to timed weight included in the former technique. Our results in Columns 9 and 11 of Table 3 show that the citations of top papers from our prediction account for more than half of that from ideal rank. Considering that there is no citation information regarding the new papers, our results are quite promising.

4.4 Sensitivity Analysis

When introducing the TimedPageRank concept, we pointed out that DecayRate is tunable for a given dataset to reach an optimal result. Our experiments show that TPR(AJEval) is an effective scoring technique. Therefore, we apply a range of DecayRate values in TPR(AJEval - simple average) and study the relation between the scoring effectiveness and DecayRate. A set of DecayRate {0.2, 0.3, 0.5, 0.7, 0.8, 1.0} are experimented, and the results are listed in Table 4.

Column 1: It lists each group of top ranked papers.

Table 3. Comparison results of different methods using only new papers

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
No. of top papers	Random		Author Eval. (PR^T)		Journal Eval. (PR^T)		AJEval (PR^T - simple avg)		AJEval (PR^T - weighted avg)		Author Eval. (count)		Journal Eval (count)		AJEval (count)		Best Citation count
10	73	10%	215	29%	185	25%	366	49%	373	49%	93	12%	134	18%	173	23%	754
20	192	16%	315	26%	414	35%	634	53%	804	67%	212	18%	215	18%	236	20%	1194
30	219	15%	543	37%	692	46%	861	58%	1013	68%	367	25%	282	19%	459	31%	1481

Table 4. Comparison results of the TPR(AJEval - simple average) technique using different DecayRate values

1	2	3	4	5	6	7	8	9	10	11	12	13	14
No. of top papers	DecayRate = 0.2		DecayRate = 0.3		DecayRate = 0.5		DecayRate = 0.7		DecayRate = 0.8		DecayRate = 1.0		Best Citation Counts
10	4257	75%	4287	76%	4312	76%	4243	75%	4341	77%	3308	58%	5661
20	5698	78%	5791	79%	5739	78%	5733	78%	5645	77%	4552	62%	7345
30	6769	81%	6868	82%	6845	81%	6881	82%	6537	78%	5511	66%	8406

Columns 2 and 3: Column 2 gives the results (citation count) of the TPR(AJEval) prediction method with a DecayRate = 0.2. Column 3 gives the ratio of the total citation count for the TPR(AJEval) method (DecayRate = 0.2) and the total citation count of the ideal ranking (column 14).

Columns 4-5, 6-7, 8-9, 10-11, 12-13 have the similar meanings as columns 2-3. The only difference is that DecayRate varies from 0.3 to 1.0 in these experiments.

Columns 14 lists the same data showed in Column 14 of Table 1. It gives the best citation count for each group of papers based on the ideal ranking, i.e., ranking relevant papers based on the actual citation number of each paper in year 2000.

The results indicate that 0.3-0.7 is the optimal range for DecayRate in this paper collection. When DecayRate is lower than 0.2, the system heavily focuses on very recent citations; papers with less recent citations are absent from the predicted top papers even if they might be important. Failing to include these papers in the results lowers the overall ranking quality. On the contrary, when DecayRate is close to 1.0, the system does not distinguish the timing difference of citations at all. Older quality papers that are not up-to-date are favored because of their longer history. As a result, some new quality papers are excluded from the top rank.

5. Conclusions

This paper studies the temporal dimension of search. So far, limited research work has been done to consider time in either publication search or Web search. Here we attempted to propose a number of techniques to remedy the situation. Our evaluation results show that the proposed techniques are highly effective. Although in this work we used publication search as the testbed, the proposed methods can be adapted to Web search because concepts in the two domains are largely parallel.

6. References

- [1] D. Achlioptas, A. Fiat, A. Karlin, and F. McSherry. "Web search via hub synthesis", *FOCS*, Las Vegas, Nevada, 2001, pp. 500-509.
- [2] R. Baeza-Yates, F. Saint-Jean, C. Castillo, "Web Structure, Dynamics and Page Quality", *SPIRE*, Lisbon, Portugal, 2002, pp. 117-130.
- [3] K. Bharat and M. Henzinger. "Improved algorithms for topic distillation in a hyperlinked environment", *SIGIR*, Melbourne, Australia, 1998, pp. 104-111.
- [4] A. Borodin, J. S. Rosenthal, G. O. Roberts, and P. Tsaparas, "Finding authorities and hubs from link structures on the World Wide Web", *WWW*, Hong Kong, China, 2001, pp. 415- 429.
- [5] A. Broder, R. Kumar, F. Maghoul, P. Raghavan, S. Rajagopalan, R. Stata, A. Tomkins, and J. Wiener, "Graph structure in the Web", *WWW*, Amsterdam, The Netherlands, 2000, pp. 309-320.
- [6] S. Brin, L. Page. "The anatomy of a large-scale hypertextual Web search engine", *WWW*, Brisbane, Australia, 1998, pp. 107 - 117.
- [7] S. Chakrabarti, B. Dom, D. Gibson, J. Kleinberg, P. Raghavan, and S. Rajagopalan, "Automatic resource compilation by analyzing hyperlink structure and associated text", *WWW*, Brisbane, Australia, 1998, pp. 65 - 74.
- [8] S. Chakrabarti, M. van den Berg, and B. Dom, "Focused crawling: A new approach to topic-specific Web resource discovery", *WWW*, Toronto, Canada, 1999, pp. 1623-1640.
- [9] J. Cho and S. Roy, "Impact of Web Search Engines on Page Popularity", *WWW*, New York, USA, 2004. pp. 20 - 29.
- [10] M. Diligenti, M. Gori, and M. Maggini, "Web page scoring systems for horizontal and vertical search", *WWW*, Honolulu, USA, 2002, pp. 508 - 516.
- [11] S. Dill, R. Kumar, K. S. McCurley, S. Rajagopalan, D. Sivakumar, and A. Tomkins. "Self-similarity in the Web", *VLDB*, Roma, Italy, 2001, pp. 69-78.
- [12] R. Fagin, R. Kumar, K. S. McCurley, J. Novak, D. Sivakumar, J. Tomlin, and D. Williamson, "Searching the workplace Web", *WWW*, Budapest, Hungary, 2003, pp. 366 - 375.
- [13] S. D. Kamar, T. Haveliwala, C. D. Manning, and G. H. Golub, "Extrapolation methods for accelerating PageRank computations", *WWW*, Budapest, Hungary, 2003, pp. 261 - 270.
- [14] J. Kleinberg. "Authoritative sources in a hyperlinked environment", *SODA*, San Francisco, USA, 1998, pp. 668 - 677.
- [15] S. Lawrence, K. D. Bollacker, and C. L. Giles, "Indexing and retrieval of scientific literature", *CIKM*, Kansas City, USA, 1999, pp. 139 - 146.
- [16] A. Ntoulas, J. Cho, and C. Olston, "What's New on the Web? The Evolution of the Web from a Search Engine Perspective", *WWW*, New York, USA, 2004, pp. 1 - 12.