

*Department of Computer Science*

*Publications*

---

*Southern Illinois University Carbondale*

*Year 2006*

---

# Extracting Web User Profiles Using a Modified CARD Algorithm

Shahram Rahimi\*

Lisa Gandy<sup>†</sup>

Bidyut Gupta<sup>‡</sup>

\*Southern Illinois University Carbondale, rahimi@cs.siu.edu

<sup>†</sup>Southern Illinois University Carbondale

<sup>‡</sup>Southern Illinois University Carbondale

Published in Rahimi, S., Gandy, L., Gupta, B. (2006). Extracting Web user profiles using a modified CARD algorithm. 2006 IEEE International Conference on Fuzzy Systems, 582-589. doi: <http://dx.doi.org/10.1109/FUZZY.2006.1681770> 2006 IEEE. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE. This material is presented to ensure timely dissemination of scholarly and technical work. Copyright and all rights therein are retained by authors or by other copyright holders. All persons copying this information are expected to adhere to the terms and constraints invoked by each author's copyright. In most cases, these works may not be reposted without the explicit permission of the copyright holder.

This paper is posted at OpenSIUC.

<http://opensiuc.lib.siu.edu/cs-pubs/21>

# Extracting Web User Profiles Using a Modified CARD Algorithm

Shahram Rahimi, Lisa Gandy, Bidyut Gupta

*Department of Computer Science  
 Southern Illinois University Carbondale  
 Carbondale, IL 62901, USA  
 [rahimi, lgandy, bidyut]@cs.siu.edu*

**Abstract** - Clustering algorithms are widely used methods for organizing data into useful information. The Competitive Agglomeration for Relational Data (CARD) Algorithm is one such clustering algorithm that is designed to organize user sessions into profiles, where each profile would highlight a particular type of user. The CARD algorithm is a viable candidate for web clustering; however, it does have limitations such as an extended execution time. In addition, the methods that prepare the input data for the CARD algorithm's use employs concepts which seem to be incomplete. These limitations of the CARD algorithm are explored and modifications are introduced to yield a more practical and efficient algorithm.

## I. INTRODUCTION

As the World Wide Web has grown in size, an immense amount of data has been produced. If this data is not retrieved and organized correctly then the information that the data can provide is essentially wasted. Users often become frustrated because they are certain that the information that they are looking for is on a particular website, but the website is so large that it is almost impossible to find the information. In this case, an adaptive website can become very helpful. This website would follow the *clickstream* of the user. That is, the website will keep account of web pages that the user is looking at, and recommend other pages that previous users found useful.

The website will implement this adaptability by having a database of typical user profiles. A user profile is defined as an "abstract model that summarizes the relevance of each URL on a site relative to a group of users sharing a similar interest [2]." The user profile will be characterized by previous user sessions. A user session is a set of web pages that a user examined within a specified time period. To organize the user sessions into profiles, the website administrator could examine each user session and manually place them into profiles. However, this is unrealistic for many reasons, for instance, many web sites have millions of users accessing it daily and it is impossible to find patterns by mere manual examination. Therefore, we can easily see that an unsupervised clustering algorithm is an attractive choice for clustering user sessions into profiles of several types of typical users since it relies on user access patterns and is capable of examining large amounts of data in a fairly reasonable amount of time [1].

Although the CARD algorithm has found useful user profiles in prior publications, there are several assumptions made during the implementation of the algorithm that are questionable [1]. For example, when the relational data is being formed, URL paths of HTML documents are compared. We propose comparing web pages with an algorithm that will scan the title of the page, and then scan for keywords. The keywords of the two documents will then be compared and a similarity in the range of [0, 1] will be found. A comparison will then be implemented to find whether the new clusters that result from the modified data are more valid.

In addition, the original algorithm takes a large amount of time to run. From our initial experiments with 1800 web sessions, the algorithm took approximately 6 hours to run. We will analyze methods to decrease the execution time of the algorithm when run serially, and also implement the algorithm in a parallel computing environment. The algorithm easily lends itself to parallel computation since it is primarily based on matrix mathematics.

## II. COMPETITIVE AGGLOMERATION FOR RELATIONAL DATA ALGORITHM

The CARD algorithm is characterized by equation (1):

$$J(\mathbf{UB}; \mathbf{X}) = \sum_{i=1}^C \sum_{j=1}^N (u_{ij})^2 d^2(x_j, \beta_i) - \alpha \sum_{i=1}^C \left[ \sum_{j=1}^N u_{ij} \right]^2 \quad (1)$$

subject to  $\sum_{i=1}^C u_{ij} = 1$ , for  $j \in \{1, \dots, N\}$ .

where  $X = \{x_j \mid j = 1, \dots, N\}$  is a set of  $N$  vectors and  $B = (\beta_1, \dots, \beta_C)$  corresponds to  $C$ -tuple prototypes, where each prototype characterizes a cluster. In (1),  $d^2(x_j, \beta_i)$  represents the distance from the feature vector  $x_j$  to the cluster  $\beta_i$  and  $u_{ij}$  represents the membership of the point  $x_j$  to the cluster  $\beta_i$ . The first term given in the equation controls the size and the shape of the clusters. The second term controls the number of clusters [1].

To prepare data for the CARD algorithm, clustering user sessions requires the use of relational data, which can be represented by a matrix of similarities between each session to all other sessions. To define this similarity matrix of size  $N_s \times N_s$  (where  $N_s$  is the number of user sessions) we first assume that each session is defined by a

vector known as  $s$  which has size  $N_{url}$  (where  $N_{url}$  is the number of URLs). Each entry contained in the vector can either contain a 1 or 0 depending on whether  $i^{th}$  URL was visited during the particular user session. Thus, the relation between two user sessions can be found by applying equation (2), which is essentially then number of identical URLs between the two sessions related to the number of URLs used in both sessions [1].

$$S_{1,kl} = \frac{\sum_{i=1}^{N_{url}} s_i^{(k)} s_i^{(l)}}{\sqrt{\sum_{i=1}^{N_{url}} s_i^{(k)}} \sqrt{\sum_{i=1}^{N_{url}} s_i^{(l)}}} \quad (2)$$

for  $k=1..N_s; l=k..N_s$

Equation (2) relies on the fact that only identical URLs between sessions are considered. However, two web pages might not have the same URL but may have quite similar content or usage, to take this into account, syntactic similarity can be used. This syntactic similarity is known as  $S_u$ , and is displayed in equation (3) [1].

$$S_u(i, j) = \min \left( 1, \frac{|p_i \cap p_j|}{\max(|p_i|, |p_j|) - 1} \right) \quad (3)$$

This similarity will lie between [0,1]. For instance the web page [www.cs.siu.edu/~cs491-3/index.html](http://www.cs.siu.edu/~cs491-3/index.html) and the web page [www.cs.siu.edu/materials.html](http://www.cs.siu.edu/materials.html) would have a syntactic similarity of 0.5. This syntactic similarity is then incorporated into the intersession similarity equation to form the updated equation (4).

$$S_{2,kl} = \frac{\sum_{i=1}^{N_{url}} \sum_{j=1}^{N_{url}} s_i^{(k)} s_j^{(l)} S_u(i, j)}{\sum_{i=1}^{N_{url}} s_i^{(k)} \sum_{j=1}^{N_{url}} s_j^{(l)}} \quad (4)$$

However, it can be shown that when the syntactic similarities are low  $S_1$  is a better approximation of session similarity and then the syntactic similarities are high  $S_2$  is a better approximation. Therefore to get the optimal similarity measure past research suggests that the maximal result of the two equations should be taken as  $S_{kl}$ , as shown in equation (5).

$$S_{kl} = \max(S_{1,kl}, S_{2,kl}) \quad (5)$$

For further information concerning the CARD algorithm please refer to [1].

### III. IMPROVING CLUSTERING: THE FIRST MODIFICATION

#### A. Theory

As stated previously, finding  $S_u$  by comparing URL paths can be problematic, in the sense that simply because two web pages are in the same directory structure does not

guarantee that they will have similar content. I proposed comparing each document to all other documents and finding a keyword similarity that will be added to the similarity found in (3). The proposed keyword similarity is given in (6)

$$S_{key}(i, j) = \frac{\sum_{k=1}^{Uniqkey} \min(key_k^{(i)}, key_k^{(j)})}{\sqrt{\sum_{k=1}^{Nkey(i)} key_k^{(i)}} \sqrt{\sum_{k=1}^{Nkey(j)} key_k^{(j)}}} \quad (6)$$

*UniqKey* identifies the number of matching keywords between two documents  $i$  and  $j$ . The identifier  $key_k^{(i)}$  represents the number of occurrences of keyword  $k$  in document  $i$ . Therefore in the numerator of (6), a running total of the minimum occurrences of matching keywords in the compared documents is taken. The denominator of (6) is simply the number of occurrences of all keywords occurring in documents  $i$  and  $j$ . This equation has the attractive property that  $S_{key}(i, i) = 1$  and  $S_{key}(i, j) = S_{key}(j, i)$  and  $S_{key}(i, j) = [0..1]$ .

The use of *min* in the numerator deserves some thoughtful consideration. The idea is that if a keyword occurs very seldom in one document and very frequently in another document then the two documents must not be very similar in general. In this case it is better to take the smaller value of keyword frequency to reflect this dissimilarity. If the average was taken then the fact that one document might have a small frequency for one keyword and the compared document has a large frequency might be masked.

Once the keyword similarity is found it is added to the similarity found in equation (5) to find a new similarity measure given below in equation (7). When the combined values of  $S$  and  $S_{key}$  exceed 1, we simply choose 1 as the maximum value. Notice that  $S_{key}$  has been given a weight of two thirds and  $S$  has been given a weight of one thirds. In previous trials  $S_{key}$  was simply added to  $S$ , however at these times, the clusters formed were too large, because all pages had very high similarity measures. At the same instance, however, we want keyword similarity to have more weight than the URL similarity so that when comparing using only URL weights and URLs and keywords weights combined, we see a noticeable difference in cluster size and shape.

$$S_{ukey} = \min((S + S_{key} * 2) / 3, 1) \quad (7)$$

#### B. Implementation

During keyword parsing the Porter Stemming Algorithm is used. The purpose of this algorithm is to remove suffixes that can cause keywords to seem different when they actually have the same meaning. By using this algorithm the performance of the keyword system will improve, and will also be less complex, because there will be less keywords per document [7]. The algorithm used is taken from the URL

<http://www.tartarus.org/~martin/PorterStemmer/java.txt>. After all keywords are found, keyword similarity is generated and added to the URL similarity mentioned previously.

We should keep in mind that the keyword algorithm that was used is not the most state of the art keyword processing algorithm that can be found. It simply considers words between certain html tags (such as <p>, <a>, <title>, <b> and <i>), gleans out words such as possessives, articles and commonly found scripting terms and then stems these words. No document term frequency is found. More advanced keyword processing was not possible due to time constraints. However with further research enhanced keyword processing should give enhanced results.

### C. Results

1) *Performance*: One of the key points of concern when adding keyword similarity to the generation of similarity data using URL syntactic similarity is that this process will add too much execution time to the overall algorithm. Although this computation did incur additional execution time, this time was negligible. For instance before intersession similarity can be found using keywords, the actual keywords must be parsed from the html documents that have been referenced by the users. With 1116 user sessions and 2151 html documents it took approximately 40 seconds to find all keywords. After all keywords have been found, then the keywords can be used to find intersession similarity. The time between finding intersession similarity when using just URLs and when using URLs and keywords to generate intersession similarity is also low. For instance with 1116 user sessions and 2151 URLs it took approximately 58 seconds to find similarity using URLs and approximately 66 seconds to find similarity using URLs and keywords. Fig. 1 illustrates the time differences between using URLs and URLs and keywords for different numbers of user sessions.

2) *Cluster Validity*: Before results of using URL similarity and combined URL and keyword similarity can be examined, the proper validity measures must be chosen. For this paper, the equation that finds clusters is given in equation (8).

$$X_i = \{s^{(k)} \in S \mid d_{ik} < d_{jk} \quad \forall j \neq i\}, \quad \text{for } 1 \leq i \leq C \quad (8)$$

Once we find the clusters using (8), then we find the intra-cluster similarity. This is a cluster validity measure that gives the approximate distance between all pairs of sessions in a cluster. The formula used to find the intra-cluster distance is given in equation (9).

$$\overline{D_{wi}} = \sum_{s^{(k)} \in X_i} \sum_{s^{(l)} \in X_i, l \neq k} d_{kl}^2 / |X_i|(|X_i| - 1) \quad (9)$$

Another useful validity measure for clusters is that of URL probability per profile, given in equation (10) [1].

$$P_{ij} = p(s_j^{(k)} = 1 \mid s_j^{(k)} \in X_i) = \frac{|X_{ij}|}{|X_i|}, \quad (10)$$

$$\text{where } X_{ij} = \{s_j^{(k)} \in X_i \mid s_j^{(k)} > 0\}$$

$P_{ij}$  represents the probability that URL  $j$  will occur in cluster  $i$ , where cluster  $i$  is represented by  $X_i$ . For cluster validity, we tend to favor clusters where the intra-cluster distance is low and the majority of URL probabilities per cluster are high. If the URL probabilities are high then there is a frequent pattern of URL access. A “bad” cluster will generally have high intra-cluster distance and low URL probabilities.

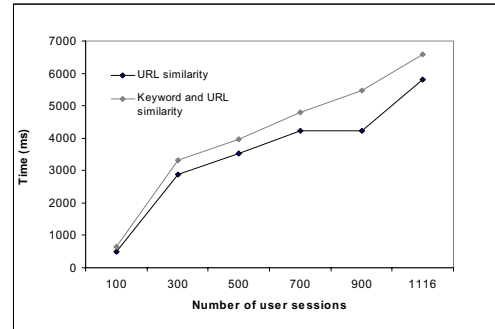


Fig. 1 Time to find intersession similarity using just URLs and URLs and keywords

### 3) Discussion

The settings for the CARD algorithm to obtain the results were following: at the beginning of the algorithm 50 rows of the membership matrix were chosen randomly to be the centers of clusters. The time constant  $\tau$  was set to 10, the decay constant  $n_0$  was set to 0.0002. A cluster was discarded if its cardinality ( $N_i$ ) was less than 22.0.

When using URL similarity only, five clusters were generated (see Table I). In the results obtained, cluster 0 is primarily referring to Dr. McGlinn's pages (a computer science professor at Southern Illinois University), cluster 12 is primarily referring to the CS pages and in some instances the cs200b web site, cluster 13 is referring to a combination of Dr. Wainer's cs485 class and Mr. Fong's cs311 and cs220 classes (Mr. Fong is a computer science instructor).

TABLE I  
RESULTS WHEN USING URL SIMILARITY

$I$	$ X_i $	URL	$P_{ij}$	$D_{\mathcal{W}_i}$
0	310	/~mcglinn/Courses/cs202/index.html	0.50	0.19
		/~mcglinn/Courses/cs202/home.html	0.46	
		/index.html	0.28	
		/materials.html	0.23	
		/Courses/cs202/labsessions/Coords.html	0.22	
		/Courses/cs202/homework/index.html	0.22	
		/Courses/cs202/labs/index.html	0.20	
		/~mcglinn/Courses/courses.html	0.12	
		/~mcglinn/index.html	0.12	
		/~mcglinn/home.html	0.12	
12	695	/index.html	0.41	0.22
		/materials.html	0.24	
		/~cs200b-1/index.html	0.13	
13	5	/~wainer/485F04/485F04OL.html	0.40	0.26
		/materials.html	0.20	
		/index.html	0.20	
		/~wainer/485F04/assign485F04.html	0.20	
		/~wainer/485F04/485HW1F04.html	0.20	
		/~rahimi/flairs05/topics.html	0.20	
		/~rahimi/flairs05/dates.html	0.20	
		/~rahimi/flairs05/cfp.html	0.20	
		/~kfong/index.html	0.20	
		/~kfong/cs311/menu.html	0.20	
		/~kfong/cs311/main.html	0.20	
		/~kfong/cs311/lab_assignments/lab1/index.html	0.20	
		/~kfong/cs311/lab_assignments/index.html	0.20	
		/~kfong/cs311/index.html	0.20	
		/~kfong/cs311/homework/hw3.html	0.20	
		/~kfong/cs220/menu.html	0.20	
		/~kfong/cs220/main.html	0.20	
		/~kfong/cs220/java/index.html	0.20	
		/~kfong/cs220/info.html	0.20	
		/~kfong/cs220/index.html	0.20	
25	47	/computing/unix/tar.html	0.04	0.25
		/~wainer/guide/guideLocal.html	0.02	
34	58	/~kfong/cs220/index.html	0.26	0.22
		/~kfong/cs220/menu.html	0.26	
		/~kfong/cs220/main.html	0.24	

TABLE II  
RESULTS WHEN USING URL SIMILARITY AND KEYWORD SIMILARITY

$i$	$ X_i $	URL	$P_{ij}$	$D_{w_i}$
0	36	/index.html /~cs200b-1/index.html /materials.html /~kfong/cs311/index.html /~kfong/cs311/main.html /~kfong/cs311/menu.html /~cs200a-1/index.html /~cs200b-1/Lab/All_Lab_assignments.html	0.28 0.25 0.17 0.17 0.14 0.14 0.11 0.11	0.19
1	75	/~mcglinn/Courses/cs202/index.html /~mcglinn/Courses/cs202/home.html /~mcglinn/Courses/cs202/homework/index.html /~mcglinn/Courses/cs202/labs/index.html /~mcglinn/Courses/cs202/labsessions/Coords.html /~mcglinn/Courses/courses.html /~mcglinn/Courses/courses.html /~mcglinn/Courses/cs202/McGlinnSlides/index.html	0.81 0.71 0.41 0.39 0.28 0.17 0.17 0.12	0.36
6	94	/~cs200b-1/Class_Assignments/assign3_10_13b.html /~cs200b-1/Lab/All_Lab_assignments.html	0.11 0.09	0.12
8	20	/~mcglinn/Courses/cs202/index.html /~mcglinn/Courses/cs202/home.html /index.html /~mcglinn/Courses/cs202/homework/index.html /~mcglinn/Courses/cs202/labs/index.html /materials.html /~mcglinn/Courses/cs202/labsessions/Coords.html /~mcglinn/home.html /~mcglinn/index.html /~mcglinn/Courses/courses.html /faculty.html	0.95 0.95 0.90 0.80 0.60 0.60 0.55 0.40 0.40 0.40 0.35	0.33
9	375	/~kfong/cs220/main.html /~kfong/cs220/index.html /~kfong/cs220/menu.html	0.15 0.15 0.15	0.14
10	504	/index.html /materials.html	0.61 0.39	0.2
12	12	/~mcglinn/Courses/cs202/index.html /~mcglinn/Courses/cs202/home.html /~mcglinn/index.html /~mcglinn/home.html /~mcglinn/Courses/courses.html /~mcglinn/Courses/cs202/labs/index.html /index.html /~mcglinn/Courses/cs202/labsessions/Coords.html /materials.html /~mcglinn/Courses/cs202/JavaInformation/index.html /~mcglinn/Courses/cs202/homework/index.html	1.00 1.00 0.67 0.67 0.58 0.42 0.42 0.33 0.33 0.25 0.25	0.32

In addition cluster 25 is referring to a UNIX description page and cluster 34 is referring to primarily the cs220 pages. All of the clusters are valid clusters, in the sense that their intra-cluster similarity ( $D_{in}$ ) is low and that inside the clusters the general URL access pattern is high.

The only cluster that seems questionable is cluster 25, where the URL access patterns are very low. Because of the unusually low URL access pattern this cluster could be characterized as a collection of unrelated pages, more like the leftovers from the other clusters.

When using URL and keyword similarity together (see Table II), seven clusters are generated. It seems that the clusters from just using URL similarity have been broken down into more identifiable clusters when incorporating keyword similarity. For example when using URL similarity alone, cluster 0 is primarily composed of Dr. McGlinn's pages and all other clusters do not reference these pages. However when keyword similarity, clusters 1, 8, and 12 all reference Dr. McGlinn's pages.

The URLs that cause these three clusters to deviate, are the following: in cluster 1, the url `/~mcglinn/Courses/cs202/McGlinnSlides/index.html` is present, in cluster 8 the two urls `/~mcglinn/home.html` and `/~mcglinn/index.html` are present, and in cluster 12, the two urls `/~mcglinn/Courses/cs202/JavaInformation/index.html` and `/~mcglinn/Courses/cs202/omework/index.html` are present. Table III shows the URL access probabilities for the five previously mentioned urls in cluster 0 when using URL similarity only and clusters 1, 8 and 12 when incorporating keyword similarity as well.

What appears to be happening is that when using keyword similarity the similarity of the URLs given in Table 1 is increasing, since these pages share a common theme with the other pages in Dr. McGlinn's web site.

This increased similarity gives these URLs a chance to compete when the clusters are being formed and lend enough similarity that they help delineate the one cluster referring to Dr. McGlinn's web site in Table 1 into three separate clusters in Table 2.

Another interesting development is that when using URL similarity the URLs for cs200b class were simply part of a cluster that primarily referred to several types of computer science pages in cluster 12, however when incorporating keyword similarity the cs200b class has been assigned its own cluster.

However, using keywords has rewards as well as problems. If we examine the results further, we see that when incorporating keyword similarity we often see low URL probabilities. Also, in cluster 0, the url probability for URLs `index.html` and `~cs200b-1/index.html` the URL probability is reasonable, but the other URLs in the cluster do not have very high URL probability. One very possible reason for this is that we were unable to parse keywords from all documents. The session data was obtained in Fall 2004 but the CARD Algorithm was not fully tested until early Spring 2005. Therefore several web pages were already taken off the web, rendering

them impossible to parse. In addition, several pages were password protected and due to time constraints we were not able to gather all of the passwords, therefore we could not obtain the HTML source for these pages for parsing. Therefore, these pages had no added keyword similarity. However, keep in mind that the URL similarity was factored in. Therefore in some instances, course web pages will still be found together, for instance, although Mr. Fong's 220 class was inaccessible for keywords, the pages were still clustered together.

TABLE III  
URL ACCESS PROBABILITIES WHEN USING ONLY URL SIMILARITY AND  
WHEN INCORPORATING KEYWORD SIMILARITY

URLs	Cluster 0 $P_{ij}$	Clusters 1, 8, 12 $P_{ij}$
<code>/~mcglinn/Courses/cs202/Slides/index.html</code>	0.040	0.120 (cluster 1)
<code>/~mcglinn/home.html</code>	0.120	0.400 (cluster 8)
<code>/~mcglinn/index.html</code>	0.120	0.400 (cluster 8)
<code>/~mcglinn/Courses/cs202/JavaInfo/index.html</code>	0.030	0.250 (cluster 12)
<code>/~mcglinn/Courses/cs202/homework/index.html</code>	0.210	0.250 (cluster 12)

#### IV. PERFORMANCE IMPROVEMENT: THE SECOND MODIFICATION

##### A. Parallel Implementation

When clustering user sessions into profiles, the actual CARD algorithm is relatively fast. Most of the time spent in execution is during the early data preparation phase. For example (2) and (4), where the session similarity is computed, take the majority of time to execute. However (2) and (4) can be simply executed by using for loops. This combination of for loops and matrix mathematics make the algorithm easy to implement in parallel.

If we look closely at (2) and (4), we can observe that  $s[i][j] = s[j][i]$ . For example if we find the similarity between session 0 and all other sessions yielding the row  $S[0][0..N_{url}-1]$  then on subsequent iterations we will not need to find any other session's similarity to session 0 because this has already been computed. In general when using (2) and (4) on a serial computation we only need an answer matrix  $S$  which has the following number of elements exhibited in equation (11).

$$numElements = N_s * \sum_{i=0}^{N_s-1} (N_s - i) \quad (11)$$

Therefore when we compute (2) and (4) we only compare session 0 to sessions 1 through  $N_s$ , session 1 to sessions 2 through  $N_s$  and so on. In general with each new session we do one less comparison with the other sessions.

In a parallel computation, where we will denote  $N_{proc}$  as the number of processors available, it would seem most appropriate that each processor should calculate the similarity

of  $N_s / N_{proc}$  number of sessions to all other sessions. For example if there are 1116 sessions and 9 processors then each processor will calculate the intersession similarity for 124 sessions. If the number of sessions does not evenly divide by the number of processors, then we add one session per processor, until all remaining sessions have been distributed.

To begin the parallel computation the root node in the parallel cluster reads in a matrix  $s$  of dimensions  $N_s \times N_{url}$ . It then distributes to each processor a *beginID* and a *endID*. The *beginID* and *endID* are the ids of the sessions for which each processor will compute intersession similarity. For example if there are 1116 sessions and 9 processors, then processor 0 will compute intersession similarity for sessions 0-123, as will processor 1 for sessions 124-147 and so on until processor 8 will compute intersession similarity for sessions 992-1115. Once each processor has received its *beginID* and *endID*, then the root node will distribute the row *beginID* through the  $N_s - 1$  row to each processor.

Before each processor can compute (2) and (4) it must allocate space in which to hold its intersession similarity. If we define  $N_{rows}$  as the number of sessions per processor then the number of elements for the matrix  $S$  per processor would be found in equation (12).

$$num\_matrixS\_elements = N_{rows} * \sum_{i=0}^{N_{rows}-1} (N_{rows} - i) \quad (12)$$

Once the answer matrix  $S$  has been formed, (2) can be computed. The modified equation, (2) for parallel will simply find intersession similarity for sessions *beginID* through *endID* as related to sessions *beginID*+1 through  $N_s$ . Similarly, a modified equation, (4) can be computed, this equation also finds intersession similarity for sessions *beginID* through *endID* as related to sessions *beginID*+1 through  $N_s$ .

### B. Results and evaluation

The computational complexity of finding the intersession similarity is  $O(n^4)$ . We will consider  $p$  to be the total number processing nodes available for the computation. Therefore, we can define the efficiency ( $E$ ) as the time for the serial algorithm (referred to as  $T_c$ ) divided by the product of the parallel time ( $T_p$ ) and  $p$  [8]:

$$E = \frac{T_c}{pT_p} \quad (13)$$

And the speed up to be [8]:

$$S = \frac{T_c}{T_p} \quad (14)$$

An efficiency of 1 implies that the scalability of the system is unlimited as the problem size grows. Moreover

not considering the main memory and cache limitation effects in a single computing node; speedup is bounded by a threshold value subject to Amdahl's law [8].

For implementation and evaluation of the concepts presented, a dedicated OSCAR cluster [9] of nine nodes was utilized. All compilations were done using the LAM MPI libraries, packaged with the OSCAR middleware. Each of the computational nodes consists of a single Pentium 4 class processor, with 786 megabytes of memory, a clock speed of 2.0 gigahertz and a 8 KB on board cache. The computer that was used for serial computation had a Pentium 4 class processor, with 512 megabytes of memory, a clock speed of 2.0 gigahertz and 512 KB on board cache.

The data that was used in testing ranged from 100 to 1116 user sessions. The serial and parallel runtimes are presented in Fig.3. The speedup of the parallel implementation is illustrated in Fig. 4 and its efficiency is shown in Fig. 5.

The speedup and efficiency seem to reach a peak level at about 500 user sessions and then the speedup levels to about 4. In general we have sped up the algorithm by about 3.5 times. Although this may seem modest, if we look at Fig. 3 we can see that with 1116 user sessions it originally took 4.4 hours to find the session similarity when using serial computation and when using parallel computation we have decreased this time to 1.25 hours.

The suggested use of this parallel computation is to first use the serial algorithm to find the similarity of all URLs to each other, and the occurrence of each URL in each session. This would give us two input files, *urlSimilarity.txt*, which would contain a matrix of size  $N_{url} \times N_{url}$ , and *urlOccurrence.txt*, which would also contain a matrix of size  $N_s \times N_{url}$ . These files could then be sent to the OSCAR cluster and the session similarity could then be computed on the cluster. After the session similarity has been computed, the cluster will output a file called *sessionVsSessionSimilarity.txt* which would be size  $N_s \times N_s$ . This file could then be fed back into the serial computation so that the CARD algorithm can complete.

One relatively simple way to speed up the parallel computation further is to reduce the times that matrices are being sent from the root node to the other nodes. In general in the computation discussed, when the root node sends the similarity matrix to each matrix, this could be stopped, and instead each matrix could read in its portion of the similarity matrix from a file. In this way the time to send the matrix is completely diminished and all that is left is file access time and data retrieval from the file, which is quite fast when using the C programming language. However, this implementation was not performed due to time constraints.

## V. CONCLUSION AND FUTURE WORK

Clustering algorithms have become an invaluable way to form user sessions into profiles. Once these profiles have been formed, then a recommendation engine can use these profiles to give user's suggestions on which URLs will be of interest to them. This can not only be useful, but profitable.



For example a company who is selling items can use the profiles to make suggestions to a user about other items that they are interested in. Improvements in clustering are being made very quickly and these algorithms are becoming faster and more reliable.

In Section 2 we have shown that incorporating keyword similarity into URL similarity creates some clusters which have greater intra-cluster similarity. These clusters also show high URL probability rates. However, when incorporating keywords we have also seen that results are not perfect and that some clusters are formed with low intra-cluster similarity. A probable reason for this discrepancy is HTML documents that are offline and therefore cannot be parsed for keywords. In Section 3 we have shown that there can be a large speedup in overall computation when we find similarity using parallel computing.

Suggested improvements are to test modifications 1 and 2 on larger data sets. Another improvement would be to change the dataset from a computer science web site to a larger and more diverse website. Also, for the sake of using keywords, one should make sure that the majority of referenced HTML documents are online and if password protected all passwords are known so that the HTML source can be retrieved. In this way, also we think that a more diverse, complete and/or a larger dataset would demonstrate a more marked improvement due to the modifications. Also, when finding updated similarity the keyword similarity is given a weight of 2/3 and the URL similarity is given a weight of 1/3. It would be interesting to modify these weights and find the weights that give the best clustering outcome.

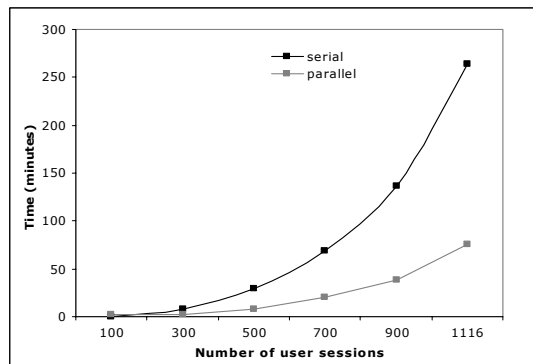


Fig. 3 Serial and parallel run times when finding session similarity

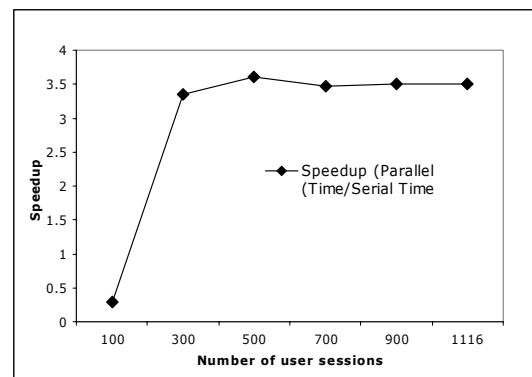


Fig. 4 Speedup Graph

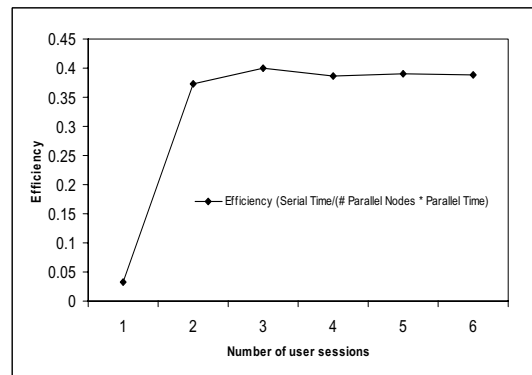


Fig. 5 Efficiency Graph

## REFERENCES

- [1] O. Nasroui, R. Krishnapuram et al. "Extracting web user profiles using relational competitive fuzzy clustering" International Journal on Artificial Intelligence Tools, Vol. 9, No. 4, 2000, pp. 509-526.
- [2] O. Nasraoui, C. Petenes, "An intelligent web recommendation engine based on fuzzy approximate reasoning" 12th IEEE International Conference on Fuzzy Systems, Vol. 2, May 2003, pp. 1116 – 1121.
- [3] R. Krishnapuram, A. Joshi et al, "Low-complexity fuzzy relational clustering algorithms for Web mining," IEEE Transactions on Fuzzy Systems, Vol. 9, No. 4, Aug. 2001, pp. 595-607.
- [4] M. Vazirgiannis, "Data Mining: Concepts and Techniques," Tutorial Paper, Proc. ADBIS 2001, Vilnius, Lithuania.
- [5] F. Heylighen, "Collaborative filtering," website: <http://pespmc1.vub.ac.be/COLLFILT.html>, March 2001.
- [6] G. Linden, B. Smith, and J. York, "Amazon.com Recommendations: item-to-item collaborative filtering", IEEE Internet Computing, Vol. 7, No. 1, Jan/Feb 2003, pp. 76-80
- [7] M. F. Porter, "The Porter Stemming Algorithm," Program, Vol. 14, No. 3, 1980, pp. 130-137.
- [8] A. Grama, Introduction to parallel computing 2<sup>nd</sup> edition, Addison Wesley, 2003.
- [9] "OSCAR (Open source cluster application resources)," website: <http://oscar.openclustergroup.org>.