# SpeedTracer: A Web usage mining and analysis tool

by K.-L. Wu P. S. Yu A. Ballman

SpeedTracer, a World Wide Web usage mining and analysis tool, was developed to understand user surfing behavior by exploring the Web server log files with data mining techniques. As the popularity of the Web has exploded, there is a strong desire to understand user surfing behavior. However, it is difficult to perform useroriented data mining and analysis directly on the server log files because they tend to be ambiguous and incomplete. With innovative algorithms, SpeedTracer first identifies user sessions by reconstructing user traversal paths. It does not require "cookies" or user registration for session identification. User privacy is protected. Once user sessions are identified, data mining algorithms are then applied to discover the most common traversal paths and groups of pages frequently visited together. Important user browsing patterns are manifested through the frequent traversal paths and page groups, helping the understanding of user surfing behavior. Three types of reports are prepared: user-based reports, path-based reports and group-based reports. In this paper, we describe the design of SpeedTracer and demonstrate some of its features with a few sample reports.

Popularity of the World Wide Web (www) on the Internet has exploded recently. Many organizations have invested a tremendous amount of capital to operate sites on the Web. These Web sites provide communications and services to their employees, customers, and suppliers. With money invested in these sites, there is a strong desire to understand the effectiveness of such investments and to find ways to realize the potential opportunities

provided by the Internet. As a result, it has become important to understand user surfing behavior.

To understand how visitors navigate a Web site, the Web server log files are analyzed. However, it is generally difficult to perform user-oriented data mining or analysis directly on the server log files because they tend to be ambiguous and incomplete. Typical server log files contain the following information about a request: client host Internet Protocol (IP) address, time stamp, method, URL<sup>1</sup> (uniform resource locator) address of the requested document, HTTP<sup>2</sup> (HyperText Transfer Protocol) version, return code (status of the request, i.e., success or error codes), bytes transferred, referrer page URL, and agent (browser and client operating system). The user identifier is usually not available in the log file. Due to the use of proxy servers by Internet Service Providers (ISPs) and firewalls by commercial corporate gateways, true client IP addresses are not available to the Web server. Instead of various distinct client IPs, the same proxy server or firewall IP will be recorded in the server log files, representing requests of different users who come to the Web site through the same proxy server or firewall. This situation creates ambiguity in the log records. Furthermore, some Web pages are generally cached by local clients or various proxy servers, or both, in order to reduce net-

©Copyright 1998 by International Business Machines Corporation. Copying in printed form for private use is permitted without payment of royalty provided that (1) each reproduction is done without alteration and (2) the Journal reference and IBM copyright notice are included on the first page. The title and abstract, but no other portions, of this paper may be copied or distributed royalty free without further permission by computer-based and other information-service systems. Permission to republish any other portion of this paper must be obtained from the Editor.

#### Table 1 Sample log entries from an NCSA HTTPd

```
peo-ill-21.ix.netcom.com - - [24/Feb/1997:00:00:21 +0000]
"GET /images/nudge.gif HTTP/1.0" 200 37 "http://www.internet.ibm.com/"
"Mozilla/2.0 (compatible; MSIE 3.01; Windows NT)"
slip166-72-149-200.wv.us.ibm.net - - [24/Feb/1997:00:00:21 +0000]
"GET / HTTP/1.0" 200 9185 "http://www.ibm.com/Products/"
"Mozilla/2.0 (Win95; I)"
ss5-08.inre.asu.edu - - [24/Feb/1997:00:00:21 +0000]
"GET /commercepoint/html3/purchasing/3_a.html HTTP/1.0" 200 6277
"http://www.internet.ibm.com/commercepoint/html3/purchasing/3.html"
"Mozilla/3.0 (Win95; I)"
peo-il1-21.ix.netcom.com - - [24/Feb/1997:00:00:24 +0000]
"GET /images/isbutton.gif HTTP/1.0" 200 1333 "http://www.internet.ibm.com/"
"Mozilla/2.0 (compatible; MSIE 3.01; Windows NT)"
ss5-08.inre.asu.edu - - [24/Feb/1997:00:00:25 +0000]
"GET /commercepoint/html3/purchasing/images/fea_a.gif HTTP/1.0" 200 1338
"http://www.internet.ibm.com/commercepoint/html3/purchasing/3_a.html"
"Mozilla/3.0 (Win95; I)"
```

work traffic. As a result, log records will be missing for the corresponding accesses to the cached Web pages, resulting in an incomplete log. A more complete discussion of the difficulties in obtaining reliable usage data on the Web can be found in Reference 3.

For example, Table 1 shows a few sample entries of an access log in the combined log format from a National Center for Supercomputing Applications (NCSA)<sup>4</sup> HTTPd.<sup>5</sup> The first entry in Table 1 represents a GET request from a user going through peo-il1-21.ix.netcom.com for file /images/nudge.gif following HTTP/1.0 protocol. The user may or may not be physically logged in on the machine peo-il1-21.ix.netcom.com. He or she may be just using the machine as a gateway to the Internet. The file size of nudge.gif is 37 bytes, and it was successfully transferred. The agent used to view page nudge.gif is MSIE\*\* 3.01 (Microsoft Internet Explorer\*\* 3.01) running on Windows NT\*\*. Finally, the user was referred to the "gif" file from http: //www.internet.ibm.com/. Namely, either file nudge.gif is on the home page of http://www.internet.ibm.com/ or there is a hyperlink to it from the home page.

To solve the problem of proxy servers or firewalls masking user IPs, it generally requires either user registrations or log-ins or the employment of "cookies' between the Web server and client browsers. With log-ins or cookies, a Web server can identify distinct requests made by individual users through a token carried between the user's browser and the server.

But the desire by many, if not the majority of, users to have privacy and remain as anonymous as possible may force many Web servers not to ask for registration or not to use cookies. As a result, there is a strong need for a tool that can analyze user-oriented behavior from the regular server log files without requiring cookies or registrations. SpeedTracer,<sup>6</sup> a Web usage mining and analysis tool, has been developed for such a purpose.

Several Web server log analysis tools have been implemented. Some of these tools are very simple and do not attempt to identify individual user sessions. These packages are simply mechanisms through which a Web master can view the raw Web server statistics, such as hit counts and distributions based on geographic regions. Examples of this type of tool include www.ics.uci.edu/pub/websoft/ wwwstat) and Analog (http://www.statslab.cam.ac. uk/~sret1/analog).

To provide user-oriented Web usage analysis, user sessions must first be identified. More sophisticated analysis packages identify user sessions with some or all of the following three mechanisms. First, if the Web server provides cookies, it is a trivial task to formulate the session. Every access to the Web server with the same cookie value makes a single session. Second, if the server does not provide cookies, it may require a log-in ID for each browser. The analysis tool can use the log-in IDs to identify sessions. In Reference 7, a data mining tool was developed on the

assumption that log-in IDs are available. Without log-in IDs, the data mining tool cannot perform its intended functions. In fact, most log records do not contain log-in IDs. Lastly, if the Web server does not provide either cookies or user IDs, the analyzer identifies sessions with host addresses. All accesses to the Web server from a given host address are considered to be a session until a predefined amount of time has passed between accesses. As mentioned previously, the use of a proxy and firewall causes all browsers from a given proxy or firewall to be considered a single user. As a result, an identified session may in fact contain many independent user sessions. Several Web analyzers only use the host address to identify sessions, such as SurfReport\*\* (http://software.bienlogic.com/SurfReport) and NetTracker\*\* (http://www.sane.com/products/ NetTracker). Other tools use a combination of methods to identify sessions, such as the Usage Analyst\*\* by Microsoft Corporation (previously Interse http://www.interse.com) and WebTrends\*\* (http: //www. webtrends.com).

In contrast, SpeedTracer uses the referrer page and the URL of the requested page as a traversal step and reconstructs the user traversal paths for session identification. No "cookies" or user registration are required. In Reference 8, an alternative approach to reconstructing user traversal paths was proposed. Instead of using a referrer page, information about the topology (i.e., hyperlink structure) of a Web site (together with other heuristics) was used to identify legitimate traversals. A software agent was first used to perform an exhaustive breadth-first traversal of pages within the Web site in order to construct the topology. However, the topology is not really needed if referrer information is available.

Once user sessions are identified, statistics related to user behavior can be obtained. Interesting user-based statistics include the top N referrers to a Web site, the top N pages most frequently visited by users, the top N pages from or into which users most frequently exit or enter a Web site, the top N browsers most frequently used, the top N IP hosts from which most users come, the demographics (by organization or by country) of users, the distribution of user session durations, the distribution of numbers of pages visited during a user session, and the distribution of depth or breadth of a user session.

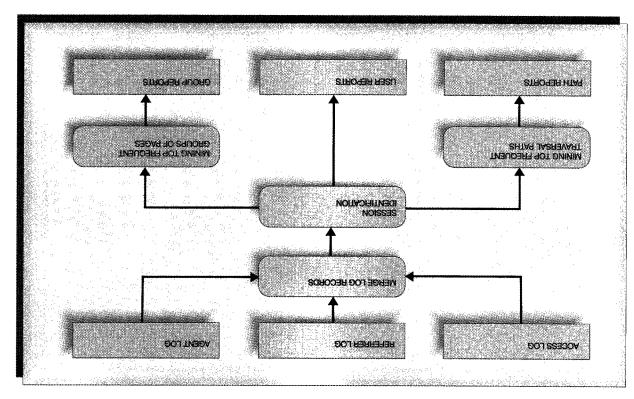
With user sessions, data mining techniques can be applied to obtain interesting user browsing patterns. Data mining has recently been used to discover cus-

tomer buying patterns by many retailers and other service corporations. One of the most important data mining problems concerns mining association rules. 9-12 Given a set of transactions, where each transaction is a set of items, an association rule is an expression of the form  $X \Rightarrow Y$ , where X and Y are sets of items. An example of an association rule is: "30 percent of transactions that contain bread and butter also contain milk; 2 percent of all transactions contain both of these items." Here 30 percent is called the confidence of the rule, and 2 percent the support of the rule. The thrust of mining association rules is to find all association rules that satisfy user-specified minimum support and minimum confidence constraints. In mining association rules, the most important problem is to generate all combinations of items that have the minimal support. These combinations of items are called large itemsets.

In SpeedTracer, we mapped each identified user session into a transaction and then applied data mining techniques to discover the top N most frequented user traversal paths and the top N groups of pages most frequently visited together. These problems are to some extent similar to finding the top N large itemsets for traversal paths and groups of pages. But specific differences exist. A traversal path is a collection of consecutive URL pages in a Web presentation, where one URL is referred to by the immediately preceding URL. The URLs in a traversal path are connected in the Web presentation graph. In contrast, the pages in a group are not necessarily connected among themselves. A frequently visited group of pages may contain two or more disjoint traversal paths. By examining the traversal paths and groups of pages, valuable user browsing patterns can be obtained to improve the organization and linkage of the Web presentation.

Note that finding frequent traversal paths is also to some extent similar to the problem of mining sequential patterns. <sup>13</sup> However, the results from the sequential-pattern-mining method in Reference 13 may contain sequences that do not represent a traversal path in the Web presentation graph. The reason is there may be many backward traversal steps involved in a user session, and pages on two different paths may be recognized as part of a sequential pattern.

In this paper we focus on mining the frequent traversal paths and groups of pages visited together. However, other types of data mining techniques, such as clustering pages, clustering users, or classification,



rogether. the top N groups of pages most frequently visited frequented user traversal paths. Finally, we discover data mining techniques to discover the top N most computed based on these sessions. Next we apply gorithms. Interesting user-based statistics are then sessions are identified using advanced inference alplementation. After log records are processed, user Figure 1 shows the flow diagram of SpeedTracer imin this paper we use the IBM ICS log as an example. NCSA log, the IBM ICS log, and others. For simplicity, to process different kinds of log formats, such as the are stored in separate files. SpeedTracer is designed responding referrer and agent log records even if they tion Server) logs, every access log record has cor-NCSA logs. However, in IBM's ICS (Internet Connecrer or agent log records, or both, may be missing from

User session identification. Given the information available in the server log, a possible approach to grouping various accesses into user sessions is to use both time stamps and agent information. For example, a user session could include all accesses within

could be applied to the derived user sessions. These and other techniques will be explored in future work.

The next section describes the design and implementation of SpeedTracer. The implementation details of session identification, mining of frequent traversal paths, and mining of groups of pages most frequently visited together are presented. The section following the next one shows a few sample reports from SpeedTracer. User reports, path reports, and group reports are highlighted. Finally, we conclude with a summary.

## Design of SpeedTracer

In SpeedTracer, we first process the access, referrer, and agent logs (or just access log, if it is in combined log format) to identify user sessions. If referrer and agent information are stored in separate files, delicate synchronization procedures may be needed to relate an access log record to the corresponding referrer and agent records because log entries may be missing in some of the files. For example, refer-

a predetermined interval if their agents are the same. Unfortunately, this approach cannot distinguish two different clients with the same agent coming from the same proxy within the specified time interval. For instance, in Table 1 the two accesses from ss5-08. inre.asu.edu both use Netscape Navigator\*\* 3.0 (Mozilla/3.0) and come within four seconds. The two accesses can be viewed as from the same user session. But they may be from two different user sessions going through the same proxy server. As the markets for both browsers and desktop operating systems become ever more consolidated, it is highly likely that multiple accesses from different users will have the same agent. For instance, most home users may use the same version of Netscape Navigator browser running on a Windows 95\*\* desktop. Thus, time stamps together with agent information are not sufficient to identify user sessions from the server log

In SpeedTracer we use five key pieces of information from a log record to identify user sessions. They are IP, Timestamp, URL (the requested page), Referral, and Agent. Different IPs or agents obviously indicate different user sessions. If the time stamps indicate that two accesses are separated by more than a prespecified period of time, the accesses are also considered to belong to different sessions. In addition to these obvious rules, SpeedTracer uses the referral page to help more accurately identify user sessions. For each log record, we use the referral page and the requested page URL to form a hyperlink access pair, representing a step in a user traversal path. Each access pair is then used to reconstruct a user traversal path in the Web presentation. The basic idea is that access pairs constitute a connected traversal path during a user session. Note that the traversal path can be forward or backward. Session identification becomes the partitioning of log records into groups so that the access pairs within a group form a connected traversal path. However, because browsers and proxy servers generally use caching to reduce network traffic and improve performance, there are no corresponding log entries for those accesses to the cached pages. As a result, missing access pairs might be in the log files, and these missing access pairs need to be added back during session identification.

In session identification, we process the log records one at a time. Each access pair is added to an active session, if possible. If  $(x_i \rightarrow y_i)$  represents an access pair, then the traversal path of a session S of size n can be expressed as follows:

$$S: (x_1 \rightarrow y_1), (x_2 \rightarrow y_2), \cdots, (x_n \rightarrow y_n)$$

where  $x_{i+1} = y_i$ ,  $1 \le i < n$ . A new access pair  $(x_j \rightarrow y_j)$  can be appended to an active session S, if  $x_i = y_k$ ,  $1 \le k \le n$ , or  $x_1 = x_i$ . However, unless  $x_i = y_n$ , a backward access path  $(y_n \to x_n), \dots,$  $(y_{k+1} \rightarrow x_{k+1})$  must first be added back to S to maintain a connected traversal path. For example, if  $(b \rightarrow d)$  were to be appended to a session  $S_i$ :  $(a \rightarrow b)$ ,  $(b \rightarrow c)$ , a backward traversal pair  $(c \rightarrow b)$  has to be first appended to  $S_i$ . Thus, the new  $S_i$  becomes  $(a \to b)$ ,  $(b \to c)$ ,  $(c \to b)$ ,  $(b \rightarrow d)$ .

Apparently, there can be multiple candidate sessions to which a new access pair can be appended. Different criteria or combinations of them can be used to choose one candidate session. For example, one criterion can be the number of backward access pairs needed to be added. Another criterion can be the time-stamp difference between the access pair and a session. The time stamp of a session is the time stamp of its latest appended access pair. Combinations of these two criteria can also be used. For example, one can choose the session with the smallest time-stamp difference with backward access pairs no more than m, or the one with the smallest number of backward access pairs with time-stamp difference no more than q minutes. Advanced inference algorithms are developed for this purpose.

As an example, Table 2 shows the key information of eight example log records used by SpeedTracer for session identification. These eight log records represent requests coming from a gateway for the IBM Watson Research Center. From Table 2, the hyperlink access pairs for the eight log records are  $(-\rightarrow a), (e\rightarrow b), (b\rightarrow c), (-\rightarrow b), (b\rightarrow c),$  $(-\rightarrow f)$ ,  $(a\rightarrow b)$ ,  $(b\rightarrow g)$ , respectively. Here, "-" means that no referral page is available for this access. Using these access pairs and the agent information, we can identify four user sessions as follows:  $S_1: (- \rightarrow a), (a \rightarrow b), (b \rightarrow g)$  from log records 1, 7, and 8;  $S_2$ :  $(e \rightarrow b)$ ,  $(b \rightarrow c)$  from log records 2 and 3;  $S_3$ :  $(- \rightarrow b)$ ,  $(b \rightarrow c)$  from log records 4 and 5; and  $S_4$ :  $(- \rightarrow f)$  from log record 6. Note that if we were to use only time stamp and agent for session identification, we would have grouped log records 1, 2, 3, 7, and 8 as a user session and log records 4, 5, and 6 as another user session. However, from the referral information of both log records 1 and 2, it is obvious that these two are from different user sessions. The access to page b in log record 2 must

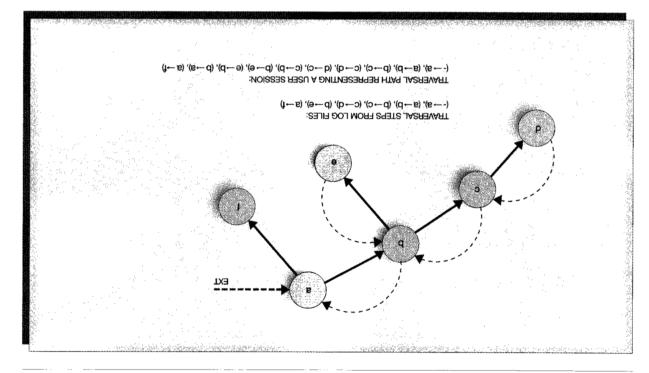


Table 2 Key information in log records used for session identification

jnegA.	IsneteR	חצר	Timestamp	dl	Record
4.1.4 XIA :0.2\tasllisoM	-	B	00:0£:80	moz.mdi.noslsw.wg-15m3fni	Ţ
4.1.4 XIA ;0.2\alpha ilizoM	ð	q	10:05:80	moə.mdi.noətsw.wg-tənrətni	7
A.1.4 XIA ;0.2\7a8lfizoM	q	9	10:06:80	internet-gw.watson.ibm.com	٤
29 niW ;0.2\nsilizoM	<u>-</u>	q	10:06:80	moo.mdl.nosisw.wg-15n15ini	<b>,</b>
29 aiW :0.2\astlixoM	q	Э,	20:05:80	moo.mdi.nostsw.wg-tənrətni	ç
29 niW ;0.2\zsllixoM	-	J	£0:0£:80	moo.mdi.nostsw.wg-təntəini	. 9
A.1.4 XIA :0.2\rankersellizoM	۳.	q	<del>\$0</del> :06:80	moo.mdi.nosisw.wg-15m5ini	L
4.1.4 XIA ;0.2\nsillaoM	q.	8	\$0:06:80	moo.mdi.nosisw.wg-15n15ini	8

ple, in Figure 2,  $(d \to c)$ ,  $(c \to b)$ ,  $(e \to b)$ , and may be missing from the server log files. For examcently visited pages, some of the actual traversal steps  $(n \to f)$ . However, since browsers usually cache re- $(a \leftarrow b), (b \leftarrow c), (c \leftarrow c), ($ versal path in Figure 2 can be described as follows: x. The user session represented by the connected traprowser to page y after the viewer has looked at page x or by clicking on the backward button on the Web

user session since it did not contain a referral. the access to page a may be the deginning of a new be following a previous access to page e, whereas

pe the result of clicking a hyperlink to page y on page ing session. Note that any access pair  $(x \to y)$  can shows a traversal pattern by a Web user during a surfbe illustrated with the following example. Figure 2 sion identification due to proxy or client caching can The need for inferring backward access pairs in ses $(b \rightarrow a)$  may be missing. These missing traversal steps may need to be inferred in order to identify traversal paths and user sessions.

Since a "gif" or "jpg" file typically does not expand a traversal path, we eliminate all log records whose URL contains these graphical files in our session identification. SpeedTracer also takes care of special cases caused by a user's clicking on the "reload" button and "bookmarking" his or her hot links. On a reload, the repeated access pair is discarded since it does not expand a traversal path. If an access is the result of a user accessing the Web page through his or her bookmark or directly typing in the URL, no referrer information is available on the log record. In SpeedTracer, we view this as the beginning of a new session. Once sessions are identified, user-oriented statistics can be obtained.

Interesting user-based statistics are provided by SpeedTracer, including the most frequent N external referrers to a site, the most frequent N visited pages by users, the most frequent N pages that users most often come into and exit from a site, the top N hosts from which most users come to visit a site, the distribution of user session durations, and the number of pages visited in a session. Sample reports and their applications will be presented in the next section. N can be specified by a user of SpeedTracer before it analyzes the log files and prepares the reports.

Mining frequent traversal paths. Once user sessions are identified, the problem of mining frequent traversal paths becomes a matter of discovering the most frequent subpaths common among all the sessions. In finding user traversal patterns, we are only interested in forward traversal subpaths. As a result, SpeedTracer first finds all maximum forward paths in each user session, and then discovers all common subpaths among all the maximum forward paths of user sessions. A maximum forward path is a sequence of maximum connected pages in a Web presentation where no page is previously visited.14 For example, in Figure 2, three maximum forward paths are in this session: (1)  $(a \rightarrow b)$   $(b \rightarrow c)$   $(c \rightarrow d)$ ; (2)  $(a \rightarrow b)$   $(b \rightarrow e)$ ; and (3)  $(a \rightarrow f)$ .

Figure 3 shows the algorithm for finding all maximum forward paths from a user session. Assume  $\{x_1, \dots, x_m\}$  represents a user session and  $\{y_1, \dots, y_m\}$  $y_{i-1}$  represents a string holding a potential maximum forward path. The idea is to examine each page  $x_i$  in the session one at a time and try to expand the

Figure 3 Algorithm for finding all maximum forward paths from a session

```
y_1 = x_1; j = 2; i = 2; flag = YES;
while (i \le m) {
   if (x_i == y_k) for some 1 \le k < j {
    if (flag == YES)
       output \{y_1, \dots, y_{i-1}\}\ as a maximum forward path;
    j = k + 1; i = i + 1; flag = NO:
   else {
     y_j = x_i; j = j + 1; i = i + 1;
      flag = YES;
if (flag == YES) {
  output \{y_1, \dots, y_{i-1}\} as the final maximum forward path;
```

potential maximum forward path by copying  $x_i$  to  $y_i$ , if  $x_i$  is not equal to any  $y_k$  for every  $1 \le k < j$ . Namely, the pages in the potential maximum forward path are all distinct pages, and we are going in the forward direction in the user traversal path. We use a flag to indicate that we are currently moving in the forward direction in path traversal. In contrast, if  $x_i$ is equal to some  $y_k$ ,  $1 \le k < j$ , then we are going backward in the traversal path, and subpath  $\{y_1, \dots, y_n\}$  $y_{i-1}$  can be a maximum forward path if the flag indicates that we have been going in the forward direction before this step. After discovering matched page  $y_k$ , we eliminate pages  $\{y_{k+1}, \dots, y_{i-1}\}$  from the potential maximum forward path by moving j backward to k + 1 for the next iteration, and set the flag to indicate backward direction. At the end, if the flag indicates forward direction, the final subpath is the final maximum forward path for the ses-

As an example, Table 3 shows the values of subpath  $\{y_1, \dots, y_{i-1}\}$  and the flag at the end of each execution step of finding maximum forward traversal paths. We use the user session in Figure 2 as our input. If we represent the session as a sequence of pages visited, this session is  $\{a, b, c, d, c, b, e, b,$ a, f, and the three maximum forward paths identified are  $\{a, b, c, d\}$ ,  $\{a, b, e\}$ , and  $\{a, f\}$ . For the first four steps, we are going in the forward direction and expanding the potential maximum forward path. In step 5, page c is found in subpath  $\{a,$ b, c, d, so a maximum forward path is found, and the traversing direction is reversed. Such a reversal

Figure 4 Algorithm for discovering large traversal path set LP

```
for each F_i {
    for each \{x_1, x_2, \dots, x_m\} in F_i {
        if (m≥k) {
          for (j = 1; j < m - k + 1; j++) {
              if (\{x_j, \dots, x_{j+k-1}\}) is already in LP_k
               increase its corresponding count;
              else if ((support of \{x_j, \dots, x_{j+k-2}\} \ge s_{k-1}) and
                     (support of \{x_{i+1}, \dots, x_{i+k-1}\} \ge s_{k-1}))
               insert \{x_i, \dots, x_{i+k-1}\} into LP_k;
```

Table 3 Example execution steps of finding maximum forward paths

lep	×	Subpath (y <sub>1</sub> ,····, y <sub>j-1</sub> ) :	Flag	Maximus forward path
1	a	{a}	YES	i de la companya de l
2	b	$\{a,b\}$	YES	
3	<b>c</b>	$\{a,b,c\}$	YES	
4	d	$\{a,b,c,d\}$	YES	
5	c	$\{a,b,c,d\}$	NO	$\{a,b,c,a$
6	ь	{a, b}	NO	
7	e	{a, b, e}	YES	
8	b	$\{a,b\}$	NO	$\{a,b,e\}$
9	a	{ <b>a</b> }	NO	
10	f	{a, f}	YES	
11		$\{a,f\}$	YES	$\{a,f\}$

lasts for two steps until step 7 when page e is expanded again to form  $\{a, b, e\}$ . In step 8, page b forces  $\{a, b, e\}$  out as another maximum forward path and reverses the traversal direction. At the end,  $\{a, f\}$  is found as a maximum forward path since the flag indicates forward direction.

Once the maximum forward paths are constructed for each session, we then map the problem of finding the top N frequent traversal paths into the one of finding frequently occurring consecutive subsequences among the maximal forward paths of all user sessions. A large traversal path is a sequence of consecutive pages that appeared in the maximal forward paths of a sufficient number of sessions. The number of sessions in which a large traversal path appears is called its *support*. A large traversal path of size k contains k pages. In this paper, we denote the set of top M large traversal paths of size k as  $LP_k$ .

Note that a significant difference exists between discovering large itemsets in mining association rules and discovering large traversal paths in mining traversal patterns. In a large traversal path, the pages must form a consecutive sequence in a maximal forward path, whereas a large itemset in mining association rules is just a set of items in a transaction.

Assume that  $P_{k,M}$  is the Mth largest traversal path in  $LP_k$  (the support of  $P_{k,M}$  is the Mth largest);  $s_k$  is the support of  $P_{k,M}$ ,  $F_k$  is the set of maximum forward paths for session  $S_i$ ; and  $\{x_1, x_2, \dots, x_m\}$  is the sequence of pages representing a maximum forward path in  $F_i$ . The algorithm for constructing  $LP_k$ , k > 1 can be described as in Figure 4. After each  $LP_k$  is constructed, the top N traversal paths are then reported. In SpeedTracer, we set M to be greater than N in computing each  $LP_k$ . Namely, we computed more large traversal paths for each  $LP_k$  than we reported.  $LP_1$  is the set of all single pages, and  $s_1$  is the number of user sessions that referenced the Mth hottest page.

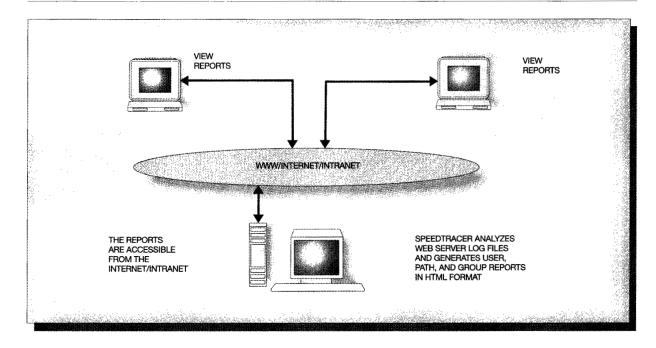
The idea of constructing  $LP_k$  is to find a candidate path of size k,  $\{x_i, \dots, x_{j+k-1}\}$ , from a maximum forward path and then compute its occurrence count among the maximum forward paths of all user sessions. The condition is that the two consecutive subsequences of size k-1,  $\{x_i, \dots, x_{i+k-2}\}$  and  $\{x_{j+1}, \dots, x_{j+k-1}\}\$ , are among the top M largest traversal paths in  $LP_{k-1}$ . For example, assume session D, E and  $\{G, H\}$ . To consider candidate large paths of size 3 for inclusion in  $LP_3$ , we would test three candidate large paths:  $\{A, B, C\}, \{B, C, D\},\$ and  $\{C, D, E\}$ . If both  $\{A, B\}$  and  $\{B, C\}$  are among the top M large paths in  $LP_2$ , then  $\{A, B,$ C} is a candidate for  $LP_3$ . Similarly, if both  $\{B, C\}$ and  $\{C, D\}$  are in  $LP_2$ , then  $\{B, C, D\}$  can be included in  $LP_3$ .

In Figure 4, for each candidate large path of size k,  $\{x_j, \dots, x_{j+k-1}\}\$ , from the maximum forward paths of a user session, we increment its occurrence count if it already is in  $LP_k$ . There are (m - k + 1) total

Figure 5 Algorithm for generating candidate groups CG,

```
Sort the groups in LG_{k-1} in lexicographical order;
for each group \{x_1, \dots, x_{k-1}\} in LG_{k-1} {
  for each group \{y_1, \dots, y_{k-1}\} in LG_{k-1} such that x_2 = y_1, \dots, x_{k-1} = y_{k-2} {
     construct a new group G = \{x_1, \dots, x_{k-1}, y_{k-1}\};
     test all other combinations of subgroups of G with size (k-1);
     if (all such subgroups are among the top M groups in LG_{k-1})
       add G into CGk;
```

Figure 6 Reports can be viewed from the Internet or an intranet

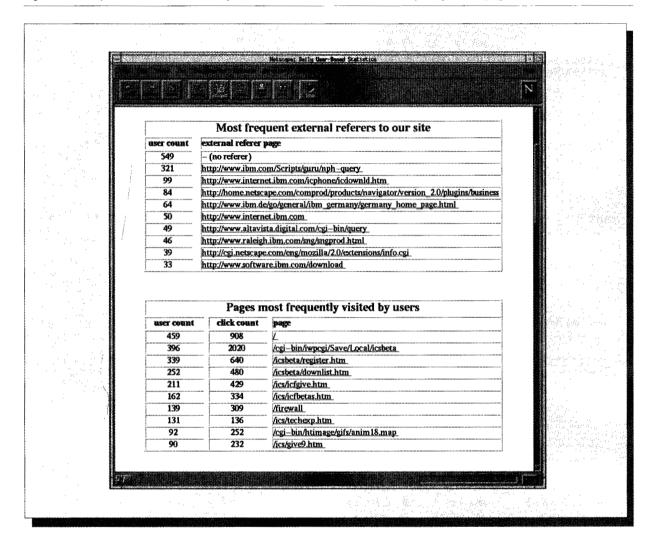


consecutive subsequences of size k from  $\{x_1, x_2, \cdots, x_n\}$  $x_m$ . For example, there are three such candidate consecutive subsequences of size 3 ( $\{A, B, C\}$ ,  $\{B, C\}$ ) (C, D), and (C, D, E)) from a maximum forward path of size 5 ( $\{A, B, C, D, E\}$ ) as we just illustrated above. We examine each one of them. If the subsequence of size k has not already been included in  $LP_k$ , we test to see if it can be. If yes, we add  $\{x_j, \dots, x_{j+k-1}\}\$  to  $LP_k$ ; otherwise we do nothing. The conditions here are based on the fact that if a traversal path of size k is among the top M largest

in  $LP_k$ , then its two consecutive subsequences of size k-1 must be also among the top M largest in  $LP_{k-1}$ . Obviously, if m < k, nothing needs to be done for this maximum forward path. For instance, nothing needs to be done for  $\{G, H\}$  for  $LP_3$ . Note that for each k, all the maximum forward paths of all user sessions are scanned only once.

Mining groups of pages most frequently visited. Frequent traversal paths identify pages that are on the same forward path in a Web presentation. These

Figure 7 Sample statistics from user reports: external referrers and most frequently visited pages



pages represent consecutive subsequences in the maximum forward paths of user sessions. However, there might be groups of pages not on the same traversal path but frequently visited together by users. By examining both frequented traversal paths and frequently visited groups, valuable information can be obtained to improve the organization and linkage of the Web presentation. For example, in Figure 2, pages a, b, e, and f may be visited most frequently by users, but these four pages are not on the same path in the Web presentation. Thus, it may be better to provide an HTTP link from page e to page f so that most users would not have to traverse backwards from page e to b, then to page a before they can go to page f.

To mine the frequently visited groups from user sessions, we need the distinct pages in each session. Thus, any duplication of pages caused by backward traversals was first eliminated in each session. Unlike traversal paths where the ordering of pages in a sequence is important, there is no ordering in a group of pages. Similar to mining the traversal paths described above, let us assume that  $LG_k$  is the set of top M largest groups, each consisting of k pages, and the support of the smallest group in  $LG_k$  is  $s_k$ . Unlike mining traversal paths, however,  $LG_{k+1}$  cannot be efficiently constructed directly because of the numerous possible combinations of candidate groups of size k + 1 from each session. For example, a maximum forward path  $\{x_1, \dots, x_m\}$  has (m - k + 1)

Figure 8 Sample statistics from user reports: IP/hosts from which most users come

		Netscape: B	ally User-Based	itatistics			Thal	
			1				N	
	trans							
No.		ndololisti kananasian			everzes rootsiere Pisiololololories keeleni	ilbanik Manada in Makada in Angara da sai	7	
II	IP/Host fron	n which	most use	ers come	to visit us			
l					multi-page user	rs		
	IP/Hest	total users	1-page users	user count	avg. duration	avg. pages	A CONTRACTOR OF THE CONTRACTOR	
	socks4.raleigh.ibm.com	55	36	19	10 min 16 sec	4		
	socks2.raleigh.ibm.com	54	35	19	1 min 41 sec	2		
	content 15a. advantis.com	47	35	12	3 min 31 sec	2		
	socks1.raleigh.ibm.com	40	17	23	0 min 47 sec	3		
	mpngate5.ny.us.ibm.com	35	15	20	2 min 33 sec	3		
ll l	mpngate2.ny.us.ibm.com	26	12	14	2 min 42 sec	2		
H	mpngate6.ny.us.ibm.com	25	15	10	3 min 17 sec	3		
H	webgate.yamato.ibm.co.jp	19	5	14	2 min 29 sec	3		
11	mpngate4.ny.us.ibm.com	16	9	7	1 min 37 sec	3		
II II	mpngate3.ny.us.ibm.com	15	4	11	4 min 52 sec	4		
H	piweba3y-ext.prodigy.com	11	9	2	0 min 34 sec	2		
	medea.castle.riga.lv	11	4	7	1 min 48 sec	4		
	menlo.ge.com	11	3	8	4 min 2 sec	2		
H	150.37.236.28	11	10	1	25 min 48 sec	3		
	www-ba6.proxy.aol.com	11	5	6	2 min 40 sec	3		
	513park-209-A.pols.indiana.edu	10	5	5	4 min 51 sec	2		
	socks3.raleigh.ibm.com	10	7	3	3 min 22 sec	2		
	129.35.251.68	10	5	5	6 min 15 sec	5		
	orange.ge.com	10	2	8	2 min 38 sec	2		
	202.44.144.7	9	7	2	0 min 49 sec	2		
5.78								
Security states				AND PROPERTY OF PERSONS ASSESSMENTS		te vide viditario i colò si the tarial è di i qu		

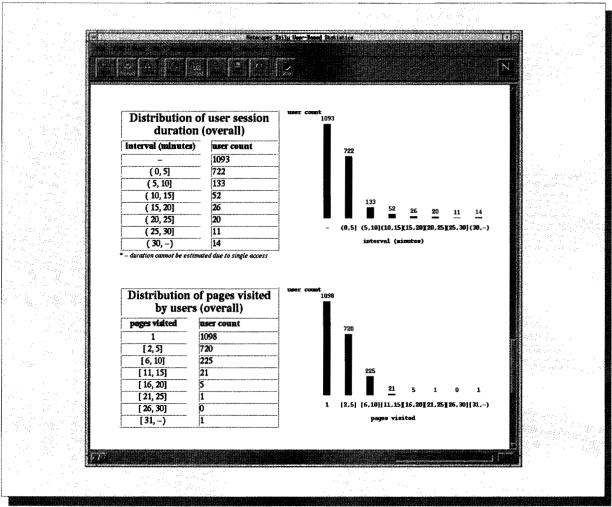
candidate traversal paths of size k. But a session  $\{x_1, \dots, x_m\}$  will have  $C_k^m$  (or m!/k!(m-k)!) candidate groups of size k. As a result, SpeedTracer constructed  $LG_{k+1}$  by (1) generating a set of candidate groups of size k+1, denoted as  $CG_{k+1}$ , from  $LG_k$ , and (2) counting the occurrences of each group in  $CG_{k+1}$  against all sessions. The approach to mining frequently visited groups of pages is thus similar to the discovery of large itemsets in most prior literature of mining association rules.  $^{9,10,12}$ 

Similar to  $LP_1$ ,  $LG_1$  contains the top M hottest pages referenced by user sessions. As pointed out in Reference 12, because of the nature of  $C_2^M$ , the computations of  $CG_2$  and  $LG_2$  can be substantially more demanding than other  $CG_k$  and  $LG_k$  for k > 2. This

is in contrast to the case of mining traversal paths, where the candidate paths of  $LP_2$  cannot be more than the number of links in the Web presentation. Therefore, special treatment is needed. <sup>12</sup>

The task of generating  $CG_k$  (candidate groups set) from  $LG_{k-1}$  can be described as in Figure 5. Note that we generate  $CG_k$  from  $LG_{k-1}$ . To simplify enumerating all possible combinations of groups, we sort the M groups in  $LG_{k-1}$  based on their lexicographical order. The basic idea here is to find all possible groups of size k from  $LG_{k-1}$  based on the fact that for a group of size k to be a candidate group in  $CG_k$ , all its subgroups of size k-1 must be in  $LG_{k-1}$ . So, we first try to construct a group of size k for each  $\{x_1, \dots, x_{k-1}\}$  in  $LG_{k-1}$  by finding all the  $\{y_1, \dots, x_{k-1}\}$  in  $LG_{k-1}$  by finding all the  $\{y_1, \dots, y_k\}$ 

Figure 9 Sample statistics from user reports: distributions of user session duration and number of pages in a session



 $y_{k-1}$ } in  $LG_{k-1}$  such that  $x_2 = y_1, \dots, x_{k-1} = y_{k-2}$ . The new k-sized group is thus an expansion of  $\{x_1, \dots, x_{k-1}\}$  with  $\{y_{k-1}\}$ . In order for such new k-sized group to be included into  $CG_k$ , all the combinations of (k-1)-sized subgroups of it must all be in  $LG_{k-1}$ .

Once  $CG_k$  is generated, we scan all the user sessions one at a time and increment the occurrence count of each candidate group in  $CG_k$  if a session contains all the pages in the group. Upon completion, the top M candidate groups in  $CG_k$  become  $LG_k$ .

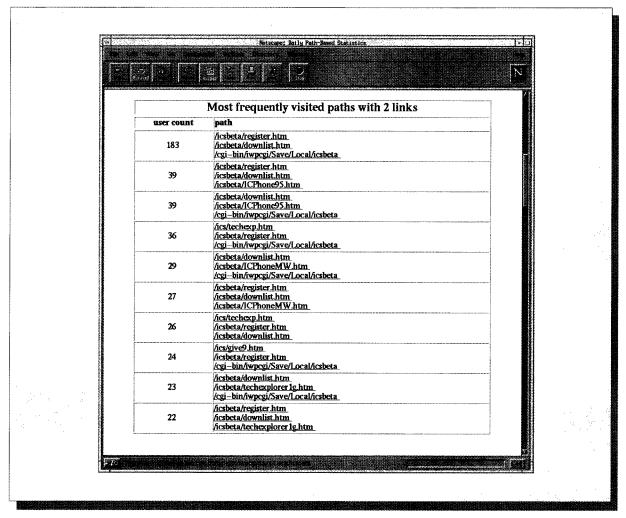
#### Sample reports from SpeedTracer

Now that we have described the implementation of SpeedTracer, we present some sample reports here.

Three types of reports are generated by SpeedTracer: user reports, path reports, and group reports. These reports are generated in HTML format so that one can view the reports through a browser from the Internet or an intranet (see Figure 6). Hot links are also provided so that one can click on them through a browser and go to the original pages to see what they are. Java\*\* applets are used to show various charts in the user reports.

To demonstrate some of the features of Speed-Tracer, we processed the server log files generated at one of the IBM Web sites running IBM ICS. There were 37 984 entries in the access, referrer, and agent log. Note that the three IBM ICS logs contain exactly the same number of entries. However, there might be a big discrepancy among the three log files from

Figure 10 Sample statistics from path reports with two links



other Web servers, such as the NCSA HTTPd. Sometimes log entries are simply dropped by the Web server in one of the files. Special logic in SpeedTracer is designed to synchronize the three files.

Figure 7 shows a snapshot of two statistics from the user reports. One is the top 10 most frequent external referrers to the server site, and the other is the top 10 most frequently visited pages. Note that "external" is with respect to the server site whose log files are being analyzed. The largest user count in the external referrer table is "no referrer." When a user visits a URL from his or her bookmark or by directly typing in the URL, there is no referrer information for such an access. But, the large count is due mostly to the fact that many of the accesses in-

volve CGI (Common Gateway Interface) programs. As a result, these accesses were treated as single-page sessions. The most frequent external referrer table can be used to measure the effectiveness of Web advertisements. It indicates the top external URLs from which most user sessions begin. If one has paid to place an advertisement on a certain site but the user count for this external referral is consistently low, one immediately realizes that the money was not well spent.

In the table for the most frequently visited pages in Figure 7, both click counts and user counts are provided. Differences between click and user counts do exist, and some of them can be substantial. As expected, the most frequently visited page by user count

Figure 11 Sample statistics from group reports consisting of three pages

	Netarinos Daily Group-Resed Statistics	escanii: 20150 (101
	arenis til	[N
Mos	at frequently visited groups consisting of 3 pages	
user count	page groups	
195	/cgi_bin/iwpcgi/Save/Local/icsbeta /icsbeta/downlist.htm /icsbeta/register.htm	1
70	/cgi – bin/iwpcgi/Save/Local/icsbeta /_ /icsbeta/register.htm	
62	/cgi_bin/iwpcgi/Save/Local/icsbeta /ics/icfgive.htm /icsbeta/register.htm	
58	/cgi_bin/iwpcgi/Sav <i>e/Local/icsbeta</i> /icsbeta/downlist.htm /_	
58	ricsbeta/downlist.htm // /icsbeta/register.htm	Anna Maria de Anna de
51	/cgi_bin/iwpcgi/Save/Local/icsbeta /icsbeta/downlist.htm /ics/icfgive.htm	
50	/_ /ics/icfgive.htm /icsbeta/register.htm	
49	/icsbeta/downlist.htm /ics/icfgive.htm /icsbeta/register.htm	
45	/_ /ics/icfgive.htm /cgi –bin/htimage/gifs/anim18.map	
41	/cgi_bin/iwpcgi/Save/Local/icsbeta /icsbeta/downlist.htm /icsbeta/ICPhone95.htm	TOTAL (1.0.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1.1

is the home page (/). However, on other experiments, we found that the home page is not always the hottest page visited by users because some users may have bookmarked some page other than the home page and go directly to it. In fact, this can be verified by the statistics that show the top pages to which users most often enter a Web site as provided by Speed-Tracer (not shown here).

Figure 8 shows other interesting statistics from the user reports: the top 20 IP/host names from which most users come to visit the Web site, the total number of user sessions from each IP/host, the average duration, and the number of pages visited for these user sessions. There are a lot of single-page sessions because of the lack of referrer information. The overall distributions of user session duration and number of pages visited by users are also provided in Figure 9. The majority of the user sessions last less than 10 minutes and are visited by less than 10 pages. Java applets were used to draw charts on the user's browser based automatically on the data in the re-

Figure 10 shows sample statistics from the path reports. SpeedTracer presents the top N most frequently visited paths with different numbers of links. In Figure 10, we only present the top 10 most frequently visited paths with two links. These paths are forward paths, meaning that one can follow the path to visit each page. Figure 11, in contrast, shows the top 10 most frequently visited groups consisting of three pages. These pages may or may not be on the same path. Even if they are, they may be visited at

Figure 12 Sample statistics from path reports with five links

		N
planticide. The shared transport of their activated cited in the cited of the cited	Most frequently visited paths with 5 links	and an individual of the control of
user count	path	kiter cililati in hazori Müstler, haliteler
11	/ics/icfgive.htm /ics/give11.htm /icsbeta/register.htm /icsbeta/downlist.htm /icsbeta/AIX41g.htm /cgi_bin/ivpcgi/Save/Local/icsbeta	
10	/ics/icfgive.htm /ics/give12.htm /icsbeta/register.htm /icsbeta/downlist.htm /icsbeta/NT41g.htm /cgi_bin/iwpcgi/Save/Local/icsbeta	
9	/_ /ics/icfgive.htm /ics/give1.htm /ics/beta/register.htm /ics/beta/downlist.htm /ics/beta/AIX41g.htm	
9	/ /ics/icfgive.htm /ics/give11.htm /ics/beta/register.htm /ics/beta/downlist.htm /cgi_bin/iwpcgi/Save/Local/ics/beta	
7	fics/icfgive.htm fics/give.13.htm ficsbeta/register.htm ficsbeta/downlist.htm ficsbeta/OS241g.htm fcgi_bin/iwpcgi/Save/Local/icsbeta	
7	/_ /ics/icfgive.htm /ics/give13.htm /icsbeta/register.htm	7

various times via other intermediate pages. For example, the top path in Figure 10 and the top group in Figure 11 contain the same three pages. But their user counts are different. The user count for the path is less than that for the group of pages because there are many different ways to visit these three pages.

By comparing Figures 10 and 11, we notice that only the first group and the 10th group in Figure 11 are also among the top 10 frequented paths in Figure 10. However, seven of these eight remaining groups contain pages on the top 10 paths with five links (see Figure 12) except (/, /ics/icfgive.htm, and /cgibin/htimage/gifs/anim18.map). Such findings suggest that many users may have traveled very deep through various paths before they found the commonly desired pages. A simplified design to shorten the depth of the traversal paths might be warranted. Since HTTP links are typically embedded in a very complex way, an examination of both frequently visited paths and groups can help a Web site to better organize its presentation.

#### **Summary**

In this paper, we described the design of Speed-Tracer, a Web usage mining and analysis tool. It reconstructs user traversal paths to identify user sessions even if user identities are hidden behind proxy servers or firewalls. No "cookies" or user registration are required for user session identification. Useroriented statistics are provided, such as the most frequent external referrers, the most frequently visited pages, the distributions of user session durations and number of pages visited. In addition, the most frequented traversal paths and the most frequently visited group of pages are also reported by Speed-Tracer. We also presented a few snapshots of sample reports generated with SpeedTracer.

### **Acknowledgment**

The authors would like to express their sincere gratitude to Robert Kreigh, Michael Stokes, Arnold Goldberg, and Ajay Balusu at IBM Raleigh for their helpful discussions during the course of developing SpeedTracer as part of various IBM Lotus Go, Lotus Go Pro, and OS/390\* ICSS version 2.2 product offerings.

\*Trademark or registered trademark of International Business Machines Corporation.

\*\*Trademark or registered trademark of Microsoft Corporation, Bien Logic, Inc., Sane Solutions, LLC, Software, Inc., Netscape Communications Corp., or Sun Microsystems, Inc.

#### Cited references and notes

- The uniform resource locator is used to uniquely identify a resource on the Internet. An example of a URL is "http: //www.ibm.com/," which represents the IBM home page on the Internet.
- HyperText Transfer Protocol is the basic protocol used by the Web to transfer documents between a browser and a Web server
- 3. J. Pitkow, "In Search of Reliable Usage Data on the WWW," Proceedings of Sixth International World Wide Web Conference (1997).
- The National Center for Supercomputing Applications is located in the University of Illinois at Urbana-Champaign, Illinois.
- NCSA HTTPd is an HTTP/1.0-compatible server for making hypertext and other documents available to Web browsers. It is copyrighted by the University of Illinois and owned by the university.
- SpeedTracer is available for download from IBM Alpha-Works<sup>TM</sup> at http://www.alphaworks.ibm.com.
- B. Mobasher et al., Web Mining: Pattern Discovery from World Wide Web Transactions, Technical Report 96-050, Department of Computer Science, University of Minnesota, Minneapolis (September 1996).
- 8. P. Pirolli, R. Rao, and J. Pitkow, "Silk from a Sow's Ear: Extracting Usable Structures from the Web," *Proceedings of 1996 Conference on Human Factors in Computing Systems* (1996), pp. 118–125.
- R. Agrawal, T. Imielinski, and A. Swami, "Mining Association Rules Between Sets of Items in Large Databases," Proceedings of ACM SIGMOD International Conference on Management of Data (1993), pp. 207–216.
- agement of Data (1993), pp. 207–216.

  10. R. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules in Large Databases," *Proceedings of Very Large Data Bases* (1994), pp. 478–499.

- 11. J. Han and Y. Fu, "Discovery of Multiple-Level Association Rules from Large Databases," *Proceedings of the 21st VLDB Conference* (1995), pp. 420–431.
- J.-S. Park, M.-S. Chen, and P. S. Yu, "An Effective Hash Based Algorithm for Mining Association Rules," *Proceedings of ACM SIGMOD International Conference on Management of Data* (1995), pp. 175–186.
- R. Agrawal and R. Srikant, "Mining Sequential Patterns," Proceedings of 11th International Conference on Data Engineering (1995), pp. 3–14.
- M.-S. Chen, J. S. Park, and P. S. Yu, "Data Mining for Path Traversal Patterns in a Web Environment," *Proceedings of International Conference on Distributed Computing Systems* (1996), pp. 385–392.

Accepted for publication August 5, 1997.

Kun-Lung Wu IBM Research Division, Thomas J. Watson Research Center, P.O. Box 704, Yorktown Heights, New York 10598 (electronic mail: klwu@watson.ibm.com). Dr. Wu received a B.S. degree in electrical engineering from the National Taiwan University in 1982 and M.S. and Ph.D. degrees in computer science from the University of Illinois at Urbana-Champaign in 1986 and 1990, respectively. From 1985 to 1989 he was a research assistant at the Center for Reliable and High-Performance Computing, Coordinated Science Laboratory, University of Illinois at Urbana-Champaign. In the summer of 1986 he also worked as a consultant at Texas Instruments. Since March 1990 Dr. Wu has been with the IBM Watson Research Center, where he is currently a research staff member in the Software Tools and Techniques group in the Internet Technology Department. His current research interests include data mining tools for the World Wide Web, Internet applications, interactive information warehousing, database transaction and query processing, multimedia system designs, and network-centric information services. Dr. Wu has served as an organizing and program committee member for various IEEE conferences, and he is a member of the IEEE, ACM, and Phi Kappa Phi.

Philip S. Yu IBM Research Division, Thomas J. Watson Research Center, P.O. Box 704, Yorktown Heights, New York 10598 (electronic mail: psyu@watson.ibm.com). Dr. Yu received a B.S. degree in E.E. from the National Taiwan University in 1972, M.S. and Ph.D. degrees in E.E. from Stanford University in 1976 and 1978, respectively, and an M.B.A. from New York University in 1982. He has been with the IBM Watson Research Center since 1978 and is currently the manager of the Software Tools and Techniques group in the Internet Technology Department, of which one project focuses on developing algorithms and tools for Internet applications, such as the Web usage mining tool. His current research interests include data mining, Internet applications, database systems, multimedia systems, parallel and distributed processing, disk arrays, computer architecture, performance modeling, and workload analysis. Dr. Yu has published more than 210 papers and over 150 research reports and invention disclosures. He holds or has applied for 50 U.S. patents. He is a Fellow of the ACM and the IEEE and was an editor of IEEE Transactions on Knowledge and Data Engineering, also serving as a guest coeditor of a special issue on mining of databases. In addition to serving as program committee member of various conferences, he was the program cochair of the 11th International Conference on Data Engineering and the program chair of the 2nd International Workshop on Research Issues on Data Engineering: Transaction and Query Processing. He will be the general chair of the 14th International Conference on Data Engineering. Dr. Yu has received several honors, including Best Paper Award, and from IBM two Outstanding Innovation Awards, an Outstanding Technical Achievement Award, a Research Division Award, and 19 Invention Achievement Awards.

Allen Ballman IBM Research Division, Thomas J. Watson Research Center, P.O. Box 704, Yorktown Heights, New York 10598. Mr. Ballman is a part-time researcher at the IBM Watson Research Center, exploring Web usage mining tools. Along with his interest in developing tools for the Internet, his research interests are concentrated in the areas of network-based simulation techniques, caching proxy servers, multimedia databases, and memory cache designs. Mr. Ballman received his B.A. in computer science and mathematics from Macalester College in 1992 and is currently pursuing a Ph.D. in computer science at the State University of New York-Stony Brook. He is a member of the ACM and USENIX.

Reprint Order No. G321-5665.