

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
2 March 2006 (02.03.2006)

PCT

(10) International Publication Number
WO 2006/023765 A2

- (51) International Patent Classification:
G06F 17/30 (2006.01)
- (21) International Application Number:
PCT/US2005/029615
- (22) International Filing Date: 19 August 2005 (19.08.2005)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data:

60/603,140	19 August 2004 (19.08.2004)	US
11/015,583	17 December 2004 (17.12.2004)	US
60/637,684	20 December 2004 (20.12.2004)	US
60/660,798	11 March 2005 (11.03.2005)	US
60/662,680	17 March 2005 (17.03.2005)	US

(US). LINDSEY, Benjamin [US/US]; 48 Haight Street, #20, Menlo Park, CA 94025 (US). VANDEVELDE, Scott [US/US]; 849 College Avenue, Menlo Park, CA 94025 (US). FAIRCHILD, Jason [US/US]; 3 Golden Spar Place, Rolling Hills, CA 90274 (US). PINGEL, Adam, William [US/US]; 1646 Greenfield Ave. #202, Los Angeles, CA 90025 (US).

(74) Agents: BEFFEL, Ernest, J., Jr. et al.; Haynes, Beffel & Wolfeld LLP, P.O. Box 366, Half Moon Bay, CA 94019 (US).

(81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BW, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KM, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NA, NG, NI, NO, NZ, OM, PG, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, SM, SY, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW.

(71) Applicants (for all designated States except US): CLARIA, INC. [US/US]; 555 Broadway Street, Redwood City, California 94063 (US). BENNETT, Dominic [US/US]; 48 Woods Lane, Los Altos, CA 94024 (US).

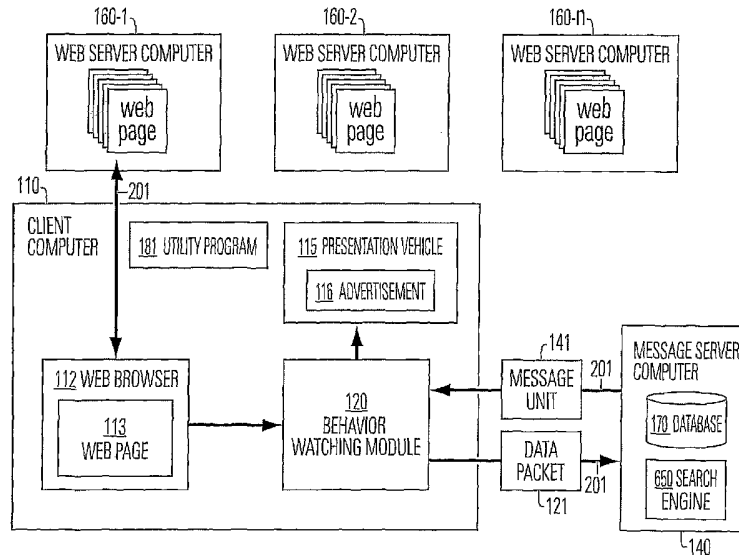
(72) Inventors; and

(75) Inventors/Applicants (for US only): PACZKOWSKI, Remigiusz, K. [PL/US]; 200 Davey Glen Road, #517, Belmont, CA 94002 (US). WESTOVER, Matthew, E. [US/US]; 4441 Edgewood Avenue, Oakland, CA 94602

(84) Designated States (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LS, MW, MZ, NA, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI,

[Continued on next page]

(54) Title: METHOD AND APPARATUS FOR RESPONDING TO END-USER REQUEST FOR INFORMATION



(57) Abstract: The present invention relates to an advanced search engine. The advanced search engine may include a client component for monitoring an end-user's browsing activity, a remote server (may comprise one or more computers) for storing and processing data received from the client component, and a module that process web pages and serves search results to end-users. The advanced search engine may collect web pages for keywords of proven interest, fetch web pages requested by end-users, generate snippets or abstracts of the web pages, eliminate duplicate web pages, rank the importance of the web pages, and provide relevant web pages or links to web pages in response to an end-user search request for information regarding one or more keywords, for example. Particular aspects of the present invention are described in the claims, specification and drawings.

WO 2006/023765 A2



FR, GB, GR, HU, IE, IS, IT, LT, LU, LV, MC, NL, PL, PT,
RO, SE, SI, SK, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA,
GN, GQ, GW, ML, MR, NE, SN, TD, TG).

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

Published:

— *without international search report and to be republished upon receipt of that report*

CLAR 1006-3

METHOD AND APPARATUS FOR RESPONDING TO END-USER REQUEST FOR INFORMATION

Inventors: Remigiusz K. Paczkowski,

Matthew E. Westover,

Dominic Bennett,

Benjamin Lindsey,

Scott Vandavelde,

Jason Fairchild, and

Adam William Pingel

RELATED APPLICATIONS

[0001] This application claims the benefit of US provisional Patent Application Nos. 60/603,140 entitled "Method and Apparatus for Responding to End-User Request for Information" by Westover et al. filed on August 19, 2004; 60/637,684 entitled "Method and Device Publishing Cross-Network User Behavioral Data" by Wohlers et al. filed on December 20, 2004; 60/662,680 entitled "Method and Device for Publishing Behavioral Observations to Customers" by Eagle et al. filed on March 17, 2005; and 60/660,798 entitled "Method and Apparatus for Responding to End-User Requests for Information" by Westover et al. filed on March 11, 2005. This application continues-in-part the disclosure of US Patent Application No. 11/015,583 entitled "Search Engine for a Computer Network" by Anthony G. Martin filed on December 17, 2004. These identified applications are incorporated by reference for all purposes.

FIELD OF THE INVENTION

[0002] The present invention relates to an advanced search engine. The advanced search engine may include a client component for monitoring an end-user's browsing activity, a remote server (may comprise one or more computers) for storing and processing data received from the client component, and a module that process web pages and serves search results to end-users. The advanced search engine may collect web pages for keywords of proven interest, fetch web pages requested by end-users, generate snippets or abstracts of the web pages, eliminate duplicate web pages, rank the importance of the web pages, and provide

CLAR 1006-3

relevant web pages or links to web pages in response to an end-user search request for information regarding one or more keywords, for example. Technical problems solved, measures used and results obtained are discussed below.

BACKGROUND

5 [0003] One approach to search engines, taken by Google, is to organize the world's information and make it universally accessible and useful. Another approach, once taken by Dogpile, is to have a meta-search engine aggregate the results of other search engines. These approaches create a great haystack of results. For instance, the keyword "cheap travel" returns about 18,000,000 results from Google, about 85,800,000 from Yahoo and 68,377,619
10 from MSN, as of summer 2005!

[0004] Much work has been done to float the "needles" to the top of the results haystack, to devise methods of ranking links returned in response to a query. Google's published patent applications propose, in their titles, to use local inter-connectivity, article information, location awareness and other factors to decide on the position of results. Yahoo,
15 Overture and Microsoft also have worked to refine their presentation of results.

[0005] In any set of information, a search term is sometimes not enough to determine what results are sought. In one sense, the search term may be ambiguous, as extensively discussed in Bharat et al., "Generating User Information for Use in Targeted Advertising", US 205/0131762 A1 published June 16, 2005 and in Carrasco et al., "Disambiguation of Search
20 Phrases Using Interpretation Clusters", US 2005/0015366 A1 published January 20, 2005. The term "jaguar" might refer to cars, animals, a football team, or an operating system. Even if a term were unambiguous, different users might prefer to access different information. For instance, teenage travelers, business travelers and luxury travelers look for different travel arrangements and accommodations, potentially using similar search terms.

25 [0006] It is desirable to return the most relevant results, whether in response to a search or, more generally, on an information feed. The growing number of documents published on web sites (and of documents accessible on private servers) invites development of alternative or improved technology to quickly return relevant results responsive to users' queries. In effect, to find the 50 or 100 most relevant web sites for a particular user whose
30 keyword is "cheap travel" and effectively summarize them for the user. This further invites

CLAR 1006-3

development of technologies that personalize the information returned, whether content, sponsored content or advertising, based on the interests of the user.

BRIEF DESCRIPTION OF THE DRAWINGS

- [0007] FIG. 1 is a schematic diagram of an example computer that may be used in
5 various embodiments.
- [0008] FIG. 2 shows a schematic diagram of a computing environment in which this
technology may operate.
- [0009] FIG. 3 schematically illustrates the tracking of a conversion process.
- [0010] FIG. 4 depicts a high level block diagram of components that can function
10 together to produce a search engine that personalizes sponsored and non-sponsored
information delivery.
- [0011] FIG. 5 is an example of a cookie or data record that reports a URL clicked.
- [0012] FIG. 6 is an example of a cookie that stores the most recent commercial
keyword of interest.
- 15 [0013] FIG. 7 is an example of summary data maintained by the behavior watcher
that may be written to persistent memory.
- [0014] FIG. 8 illustrates a bit string further summarized by category.
- [0015] FIG. 9 illustrates buckets that may be applied to recency of using search
words, clicking through, entering URL destinations or other indications of user interest.
- 20 [0016] FIG. 10 illustrates buckets that can be applied to frequency of visits.
- [0017] FIG. 11 illustrates how a category ID can be associated with both the category
name and a revenue priority.
- [0018] FIG. 12 compares the insight that a behavior watching module on a user's
machine gives regarding web surfing with the limited view that a server can achieve.
- 25 [0019] FIG. 13 depicts a web page highly relevant to the keyword "cars".
- [0020] FIG. 14 is a counter example of a page weakly relevant to "cars".
- [0021] FIG. 15 juxtaposes the results for "cars" from Yahoo!, msn Search and Google
with an embodiment of the search engine described herein.
- [0022] FIG. 16 illustrates user differentiation using results from a behavior watching
30 module.

CLAR 1006-3

[0023] FIG. 17 depicts traffic to a particular domain, from a particular search engine, by keyword.

[0024] FIG. 18 depicts an architecture for following links and retrieving content.

[0025] FIG. 19 illustrates the importance of recent behavior in response to banner ads, measured by click through rates.

[0026] FIG. 20 depicts how the window scoring routine begins with the first word of text within the document.

DETAILED DESCRIPTION

[0027] The following detailed description is made with reference to the figures. Preferred embodiments are described to illustrate, not to limit the scope of the claims. Those of ordinary skill in the art will recognize a variety of equivalent variations on the description that follows. Persons of ordinary skill in the art will recognize, however, that the embodiments described can be practiced without one or more of the specific details. In other instances, well-known details are not shown or described to avoid obscuring aspects of the embodiment.

[0028] Being computer-related, it can be appreciated that the components disclosed herein may be implemented in hardware, software, or a combination of hardware and software (e.g., firmware). Software components may be in the form of computer readable program code stored in a computer-readable storage medium, such as memory, mass storage device, or removable storage device. For example, a computer readable medium may comprise computer-readable program code for performing the function of a particular component. Likewise, computer memory may be configured to include one or more components, which may then be executed by a processor. Components may be implemented separately in multiple modules or together in a single module.

[0029] Embodiments and aspects of embodiments described below can be applied to solve various technical problems. One problem applies to a large network: how to monitor and usefully aggregate patterns of communication among users, search engines and documents accessed. In one scenario, the users are surfing the Internet at widely disbursed personal computers; the search engines include Baidu, Google, Yahoo! and MSN; the documents are pages posted on web sites around the world. In another, employees use an

CLAR 1006-3

enterprise intranet with an enterprise search engine to locate reference documents exposed on workgroup servers. Addressing these problems may involve positioning a behavior observing module where it can monitor the communication channels in use and report observations to a server, preferably without disrupting the users' routines. In one embodiment, the behavior
5 observing module may run on the user's personal computer (such as a desktop, laptop or handheld computer or media center device). The module can observe both communications and the status of the user's machine when the communications take place. For example, what search terms and results did a user follow to a particular web site? The module can achieve distributed processing and substantially reduce the resources required to aggregate
10 communication behavior by filtering reports and categorically encoding activity. It may report observations to the server as resources are available or when a user browses to an affiliated domain. As part of the ordinary uploading of cookies to domains accessed, observation cookies can be transmitted to the server. The server can aggregate reported patterns of communication. One technical result is receiving reports from communication
15 monitors positioned across a large network and aggregating patterns, including the status or state of individual computers when communications took place. In turn, the aggregated patterns of communications can be used to pre-organize information for retrieval or publication, in anticipation of a query or contact.

[0030] A related technical problem is how to organize over-abundant electronic
20 records based on the current and recent status of a particular terminal connected to a network, to enhance the relevance of the first records presented to a user. The over-abundant electronic records may come from web sites world wide, such as the web sites for "jaguar." Or, they may be documents stored on workgroup servers. They are over-abundant in the sense that they are too numerous to display on a user's screen without repeatedly pressing "page down"
25 or the like. The current and recent status of the terminal, preferably associated with a particular user, may include web sites visited within the last 24 hours (or some other period) before a search query was submitted. Or, it may include a history of documents retrieved from workgroup servers. Either type of activity may be preprocessed and categorically classified. The period for reporting terminal status may precede a query or request for a
30 personalized electronic journal that selects and filters the information based on the current and recent status of the terminal. The technical method again involves a behavior watching

CLAR 1006-3

module running on the user's computer. In this embodiment, the module summarizes terminal status, publishes the summary to an electronic record (e.g., a cookie) and communicates the record to a search engine or other site that the user contacts. Reporting categorical summaries of status or activity distributes processing and reduces the need for server-based resources. The search engine uses the information, for instance, to determine what sense of "jaguar" is of interest? A highly involved auto category user who searches for "jaguar" would get Jaguar auto related links, while a person with no relevant category involvement would get a mix of auto, animal, etc. links. Category profiles may be developed to further categorize a user's interests. A new parent in the suburbs might be recognized from content accessed and be assigned to a different category for home accessories or cars than a single gen-X 20-something person. Life change events also might be recognized, such as marriage, home buying or parenthood. Like the "jaguar" example, for a "travel Italy" inquiry, the search engine might determine the style of travel that interests the user and organize the over-abundant electronic records accordingly. The technical result is respond to a query or contact based on an electronic report of the status or recent activity of a terminal, selecting from over-abundant electronic records a particular set of records that are most likely relevant to the current and recent status of the terminal.

[0031] Web-wide behavioral targeting differs substantially from site-side behavioral targeting. Practicing site-side behavioral targeting, a group of affiliated sites attempt to identify commercial behaviors. The sites typically serve ads, such as portals and news sites. Messages are displayed on the affiliated sites, responsive to behavior that is recognized from visits to the affiliated sites. The results of site-side behavioral targeting are better than non-behaviorally targeted campaigns, but depend on insight that can be gained from a narrow portion of user's behavior, as illustrated by **FIG. 12**, which is described below. Web-wide behavioral targeting, in contrast to site-side, can identify commercial surfing behavior and, optionally, non-commercial surfing behavior of interest, across all sites, even sites to do not accept advertising. This is significant because 90% of commercial Web activity takes place at sites or on pages that do not display advertising. Therefore, portals see less than 10% of user's total commercial activity, as measured by page views. Non-portal publishers see less than 3% of their users' total commercial activity. A behavior watching module may simplify the relationship with web site operators and allow messages to be displayed across a larger

CLAR 1006-3

network of publishers than participate in site-side behavioral targeting. Web-wide behavioral targeting produces better results than site-side behavioral targeting.

[0032] Referring now to **FIG. 1**, there is shown a schematic diagram of an example computer that may be used in various embodiments. Depending on its configuration, the computer shown in the example of **FIG. 1** may be employed as a client computer, a server computer, or other data processing device. The computer of **FIG. 1** may have less or more components to meet the needs of a particular application. As shown in **FIG. 1**, the computer may include a processor 101, such as those from the Intel Corporation, Advanced Micro Devices or IBM, for example. The computer may have one or more buses 103 coupling its various components. The computer may include one or more input devices 102 (e.g., keyboard, mouse), a computer-readable storage medium (CRSM) 105 (e.g., floppy disk, CD-ROM), a CRSM reader 104 (e.g., floppy drive, CD-ROM drive), a display monitor 109 (e.g., cathode ray tube, flat panel display), a communications interface 106 (e.g., network adapter, modem) for coupling to a network, one or more data storage devices 107 (e.g., hard disk drive, optical drive, FLASH memory), and a main memory 108 (e.g., RAM). Software embodiments may be stored in a computer-readable storage medium 105 forwarding into a data storage device 107 or main memory 108. Software embodiments in main memory 108 may be executed by processor 101. In the example of **FIG. 1**, main memory 108 is shown as having a behavior watching and message delivery program 120, which is further described below. Behavior watching and message delivery program 120 and other programs in main memory 108 may be loaded from a computer-readable storage medium 105, a data storage device 107, or over the Internet by way of communications interface 106, for example. Behavior watching and message delivery program 120 and other programs in main memory 108 may be executed by processor 101.

[0033] **FIG. 2** shows a schematic diagram of a computing environment in which this technology may operate. In the example of **FIG. 2**, the computing environment includes one or more web server computers 160 (i.e., 160-1, 160-2, ...) , one or more client computers 110, one or more message server computers 140, and other computers not specifically shown. In the example of **FIG. 2**, a client computer 110 communicates with server computers (e.g., a web server computer or a message server computer) over the Internet. As such, arrows 201 denote Internet connections in this example. Intermediate nodes such as gateways, routers,

CLAR 1006-3

bridges, Internet service provider networks, public-switched telephone networks, proxy servers, firewalls, and other network components are not shown for clarity.

[0034] A client computer 110 is typically, but not necessarily, a personal computer such as those running the Microsoft Windows™ operating system, for example. A consumer
5 may employ a suitably equipped client computer 110 to get on the Internet and access computers coupled thereto. For example, a client computer 110 may be used to access web pages from a web server computer 160.

[0035] A web server computer 160 may be a server computer hosting a website, which comprises web pages designed to attract consumers surfing on the Internet. A web
10 server computer 160 may include web pages supporting advertisements, downloadable computer programs, products available for online purchase, and so on. As can be appreciated, a website may be on one or more server computers.

[0036] A message server computer 140 may include the functionalities of a web server computer 160. In one embodiment, a message server computer 140 further includes a
15 database 171. Database 171 may be a commercially available database, such as those available from the Oracle Corporation. Database 171 may store client data received from behavior watching and message delivery programs 120 running in client computers 110. The client data may be transmitted from a client computer 110 to message server computer 140 in a data packet 121. The client data may include navigation and behavioral data obtained by a
20 behavior watching and message delivery program 120 by monitoring a consumer's on-line activities. In the example of **FIG. 2**, message server computer 140 is shown as communicating with one client computer 110 for clarity of illustration. In practice, message server computer 140 receives data packets 121 containing client data from a plurality of client computers 110, each of which has a behavior watching and message delivery program 120. A
25 message server computer 140 may also include downloadable computer programs and files for supporting, updating, and maintaining software components on a client computer 110. The message server computer 140 may include or be connected with a search engine 650.

[0037] Web server computers 160 and message server computers 140 are typically, but not necessarily, server computers such as those available from Sun Microsystems,
30 Hewlett-Packard, or International Business Machines. A client computer 110 may communicate with a web server computer 160 or a message server computer 140 using client-

CLAR 1006-3

server protocol. It is to be noted that client-server computing is well known in the art and will not be further described here.

[0038] As shown in FIG. 2, a client computer 110 may include a web browser 112 and a behavior watching and message delivery program 120. Web browser 112 may be a commercially available web browser or web client. In one embodiment, web browser 112 comprises the Microsoft Internet Explorer™ web browser. A web browser allows a consumer on a client computer to access a web page. In the example of FIG. 2, web browser 112 is depicted as displaying a web page 313 from a web server computer 160. A web page, such as web page 113, has a corresponding address referred to as a "URL" (Uniform Resource Locator), Web browser 112 is pointed to the URL of a web page to receive that web page in client computer 110. Web browser 112 may be pointed to a URL by entering the URL at an address window of web browser 112, or by clicking a link pointed to that URL, for example.

[0039] In one embodiment, behavior watching and message delivery program 120 is downloadable from a message server computer 140 or a web server computer 160. Behavior watching and message delivery program 120 may be downloaded to a client computer 110 in conjunction with the downloading of another computer program. For example, behavior watching and message delivery program 120 may be downloaded to client computer 110 along with a utility program 181 that is provided free of charge or at a reduced cost. Utility program 181 may be a wallet or calendar program, for example. Utility program 181 may be provided to a consumer in exchange for the right to deliver advertisements to that consumer's client computer 110 via behavior watching and message delivery program 120. In essence, revenue from advertisements delivered to the consumer helps defray the cost of creating and maintaining the utility program. Behavior watching and message delivery program 120 may also be provided to the consumer along with free or reduced cost access to an online service, for example.

[0040] Behavior watching and message delivery program 120 is a client-side program in that it is stored and run in a client computer 110. Behavior watching and message delivery program 120 may comprise computer readable program code for displaying advertisements in a client computer 110 and for monitoring the online activity of a consumer on the client computer 110. It is to be noted that the mechanics of monitoring a consumer's online activity, such as determining where a consumer is navigating to the URL of web pages received in

CLAR 1006-3

client computer 110, the domain names of websites visited by the consumer, what the consumer is typing on a web page, what keyword the consumer is providing to a search engine, whether the consumer clicked on a link or an advertisement, when the consumer activates a mouse or keyboard, and the like, is, in general, known in the art and not a further
5 described here. For example, behavior watching and message delivery program 120 may learn of consumer online activities by receiving event notifications from web browser 112.

[0041] Behavior watching and message delivery program 120 may record the consumer's online activity for reporting to message server computer 140. The recorded consumer online activity is also referred to as "client data," and provided to message server
10 computer 140 using data packets 121. Message server computer 140 may use the client data to provide targeted advertisements to the consumer. Message server computer 140 may include the advertisement or data for displaying the advertisement in a message unit 141. In the example of **FIG. 2**, the targeted advertisement is labeled as advertisement 116 and displayed in a presentation vehicle 115. Presentation vehicle 115 may be a pop-under, pop-up, separate
15 browser window, custom browser window, or other means for displaying an advertisement on a computer screen. Techniques for delivering advertisements to client computers using a client-side program are also disclosed in commonly-owned U.S. Application No. 10/152,204, entitled "Method and Apparatus for Displaying Messages in Computer Systems," filed on
20 May 21, 2002 by Scott G. Eagle, David L. Goulden, Anthony G. Martin, and Eugene A. Veteska, which is incorporated herein by reference in its entirety.

[0042] As will be more apparent below, behavior watching and message delivery programs are primarily used to obtain client data for building a search engine index, not necessarily to display presentation vehicles in a client computer 110. That is, a behavior watching and message delivery program does not necessarily have to display advertisements
25 in a client computer 110. This is advantageous in that consumers may be allowed to obtain a free or reduced cost utility program 181 (or other benefits) without having to see advertisements from the provider or sponsor of the utility program.

[0043] **FIG. 3** schematically illustrates the tracking of a conversion process. A "conversion" occurs when a consumer makes an online purchase or registration in response to
30 an advertisement presented to her. As will be more apparent below, the embodiments may use conversion as a factor in ranking a keyword-link combination. In the example of **FIG. 3**, web

CLAR 1006-3

pages 202 (i.e., 202-1, 202-2,...) may be sequentially displayed in the same or separate windows of web browser 112. Each web page 202 includes a page identifier 21 13 (i.e., 210-1, 210-2, ...), which may be a URL. Behavior watching and message delivery program 120 records the 20 URLs of web pages 202 viewed by the consumer as well the amount of time the consumer spent with each web page as navigation data. For purposes of the present disclosure, the web pages "viewed" by consumers are those web pages received in client computers. In the example of **FIG. 3**, navigation data 627 comprises log entries 117 (i.e., 117-1, 117-2, ...). Each log entry 117 includes a machine ID anonymously identifying the client computer 110 (or the consumer), a page identifier, and a time stamp indicating when the log entry 417 was made. The time stamps between log entries 117 provide an estimate of the amount of time the consumer spent viewing the indicated web page. A log entry 117 may be created for each web page 202 viewed by the consumer. For example, a log entry 117-1 may be created when the consumer clicks on a link 501 to receive web page 202-2. In client computer 110, a log entry 117-2 may be created when the consumer receives web page 202-3 in client computer 110, and so on.

[0044] In the example of **FIG. 3**, web page 202-2 is also referred to as a "landing page" because it is the web page directly pointed to by an advertisement having a link in web page 202-1. A website that sells products on line may also have a "confirmation page" 202-5. A confirmation page is a web page provided to the consumer to confirm a just completed online purchase. Alternatively, a confirmation page may involve confirmed submission of registration information. A website may have "intermediate pages" 202-3, 202-4, and so on between a landing page 202-2 and a confirmation page 202-5. An intermediate page may be an online product catalog, shopping cart, and other types of web pages. The page identifiers of landing and confirmation pages of popular or partner websites may be stored in database 171 of message server computer 140, the page identifiers of the landing and confirmation pages may be compared to those in navigation data 627 of client computers 110 to determine the number of consumers making online purchases from a particular web page. As can be appreciated, web pages from which consumers make purchases may be highly relevant, and may thus be deemed a consumer preferred page. Techniques for monitoring consumer purchase behavior are also disclosed in commonly-owned U. S. Application No. 10/464,419, entitled "Generation of Statistical Information in a Computer Network," filed by David L.

CLAR 1006-3

Goulden and Dominic Bennett on June 17, 2003, which is incorporated herein by reference in its entirety.

Process Flow

[0045] FIG. 4 is a high level block diagram of components that can function together to produce a search engine that personalizes sponsored and non-sponsored information delivery. As background, behavior watching modules currently capture behavioral web-navigation information on over 40 million machines in the U.S. and abroad. The information collected will expand and be adapted as needed. This navigational behavior information is currently used to power web-wide GMT ad vehicles (e.g., pop-ups and slide-unders) as well as power many data-rich applications that give unique insights into consumer behavior across the web. The behavior watching module has the ability to see the entire search life-cycle, from observing the initial keyword entered, to the resultant links displayed, to post-click metrics including dwell times, pages viewed, repeat visits, and post-click purchases or registrations. By aggregating a user's behavior, category profiles and life change events can be determined. User sophistication with search engine usage and degree of category involvement can be determined over time. Using this information, one can select and re-rank the search results, filtered by humans who requested the searches, not by bots.

Data Collection

[0046] Web usage statistics are collected (401) using behavior watching modules (120) for users searching on selected search engines across the Internet or an enterprise intranet. The modules send back information related to their searches and how effective their searches were on each particular search engine for each particular keyword. The information is captured in a database, which is loaded daily or at some other concurrent frequency. Information available includes what the user saw and how they responded. URLs seen by the user may be displayed in algorithmic or natural sections the results. Pertinent information available for a single search term includes machine id, keyword, search engine where search was performed, resultant URLs, algorithmic URLs, bidded URLs, paid inclusion URLs, whether or not a URL was clicked, number of pages viewed, dwell time, repeat visits and user metrics such as category involvement and search engine sophistication.

CLAR 1006-3

[0047] Some statistics from US users can bring to life the analytical reach of data collection and ranking. The one million most frequently searched keywords presently account for just more than half of the searches conducted on the major search engines, 53% of the searches. The 45,000 most frequently searched keywords account for 38%. The distribution of keywords entered as searches can be represented by a Zipf distribution, which plots as a straight line on a graph with logarithmic scale on both axes. See, e.g., Jacob Nielson, "Diversity is Power for Specialized Sites", Alertbox (June 16, 2003) accessed August 13, 2005 at <http://www.useit.com/alertbox/20030616.html>. It is estimated that 2,000 behavior watching modules will generate enough data to rank the 10,000 most frequent searches. A larger group of 20,000 behavior watching modules will cover the 45,000 most frequently searched keywords. A base of 40 million behavior watching modules may capture 30 billion rows of data per month, filtered from 150 billion page views observed. Among the page views observed, on the order of 8 billion commercial events per month are noted, including more than 2.5 million purchases. These numbers and even a passing familiarity with statistics should excite the reader to aggregating the intelligence and behavior of a multiplicity of searchers, and presenting by popular acclaim the most significant web sites on the first page of results from a search engine. A search engine that uses aggregated consumer behavior is well-positioned to rate the authority and popularity of pages/documents as responsive to a search request.

[0048] FIGS. 5-11 are examples of data records or cookies that can communicate information from the behavior watching module to a ranking and link-following server and to a search engine that produces personalized results. FIG. 5 is an example of a cookie or data record that reports a URL clicked. A keyword 501, time stamp 502 and list of one or more URLs selected by the user 503 are conveyed, for instance in a cookie when a user reaches a web site affiliated with one of the servers or in a data record when the behavior watching module reports activity, independently of Internet surfing behavior. Additional data regarding conversion status, return conversion status, rank in a list that a user saw, etc. can be included in such a cookie or data record.

[0049] FIG. 6 is an example of a cookie from that stores the most recent commercial keyword of interest. This may be the most recent commercial keyword of interest from a search or a keyword associated with a URL or click-through. More than one recent

CLAR 1006-3

commercial keyword of interest may be stored in one or more cookies, following the example of FIG. 6. One attribute or field of this cookie is the version of the crumb 602. Another is the time written 604, which may be expressed in seconds from 1970 to the time the cookie is written. A commercial flag 606 may be used to categorize the keyword type, for instance
5 between key words that are of current commercial interest and key words that are for sale. The keyword may be included 608, either in text or as a shorthand code. A character set 610 such as UFT-8 may be included to assist with decoding of the keyword 608. This may be important in international settings.

[0050] FIG. 7 is an example of summary data maintained by the behavior watcher
10 120 that may be written to persistent memory. In this example, a category ID 802 is a shorthand for a particular category. Bins are allocated for a current period 702, such as a day. The periods may be of equal length for different lengths. In one embodiment, a bit string implements the bins. A history of 60 days may be implemented by 60 bits plus storage for a category ID. Each bin corresponds to a particular time segment. A bit flag indicates whether
15 there was activity related to the category during the particular time segment. At a regular time corresponding to the shortest period represented by any bin, the bit string is updated. Data is shifted into older bins, making room in the most recent bin. Data regarding recent activity is summarized into the most recent bin, labeled zero in the figure. If bins represent unequal
20 amounts of time, special provisions are made at the boundaries between shorter and longer time bins, either algorithmically or by allocating additional storage. For instance, activity during any one-hour segment during a preceding day may result in setting the bit flag in a day-long time segment. The day-long time segment is flag "0" only if the hour segment flag at midnight is "0" when the midnight update takes place. This is an example of an
algorithmic provision at the boundary between a one-hour segment and a day-long segment.
25 Alternatively, storage may be allocated for bins to record number of visits in a particular time segment. This may significantly increase the amount of storage required, without necessarily changing how a reported-behavior aware server 132 responds. Compact storage is preferred when there are hundreds of categories and the practical limit for the size of a cookie is 2K to 4K bytes. Cookies can be sent to a server that archives user activity data. The user activity
30 data can be used for analysis or targeting of ads. A bit string can be further summarized by category, shown in FIGS. 8-10.

CLAR 1006-3

[0051] **FIG. 9** illustrates buckets that may be applied to recency of using search words, clicking through, entering URL destinations or other indications of user interest. The most recent visit is recorded by membership in a particular bucket. One of 15 or 16 buckets can be identified with four bits. In the figure, the buckets are of unequal length. A single, most recent visit record such as illustrated by **FIG. 6** or an entry kept in memory can be used to assign a category to a particular bucket.

[0052] **FIG. 19** illustrates the importance of recent behavior in response to banner ads, measured by click through rates. A user who had behavior within the last hour in an advertising category was 30 times as likely to respond to a banner ad as one with no behavior over the last 7 days in the category. Behavior in the last hour, two hours or three hours has been observed to generate a huge increase in responsiveness to banner ads. Even behavior within the last seven days correlates strongly with a user's interest. Preferably, the selection of recency bucket sizes reflects the observed power of recent behavior illustrated by the figure.

[0053] **FIG. 10** illustrates buckets that can be applied to frequency of visits. When the bit string of **FIG. 7** is organized by days, activity in a category can be assigned to buckets of different visit frequencies, such as 0, 1, 2, 3, 4 or 5, 6 to 10, 11 to 30, or 31 to 60 days out of the last 60 days in which the user indicated interest. The number of bits set in the bit string is counted and summarized using this categorization.

[0054] With the categorical coding of recency and frequency in mind, we turn to **FIG. 8**. A single cookie may record data for many category ID's 802. For each category ID, three or more metrics of user interest may be recorded. One or more recency categories 804 can be recorded. The recording of recency can be applied to key words submitted to search engines or to other kinds of user behavior, described above. One or more frequency categories 806 can be recorded, applied to key words submitted or other kinds of user behavior. Banner clicks 808 can be separately categorized by recency, frequency or both. The comments in data column 810 are editorial comments to explain the data in the table and are not part of the data. Category ID 98409 experienced heavy and recent user interest. Category ID 65625 experienced heavy user interest last month. Category ID 68530 experienced light, recent user interest. Category 147374 experienced light activity last month, which probably also should be reflected by a nonzero value under frequency 806.

CLAR 1006-3

[0055] The behavior watcher module 120 preferably sorts the category history information in **FIG. 8** so that the most important categories are listed first. Alternatively, the system could leave the category sorting to the reported-behavior responsive server. As illustrated in **FIG. 11**, a category ID 1102 can be associated with both the category name
5 1104 and a revenue priority 1106. It may be useful to sort category history information by revenue priority 1106, especially if only a portion of the client-sided user behavior data is transmitted to a reported-behavior responsive server.

[0056] The statistics returned regarding web usage may or may not distinguish
10 between commercial and non-commercial keyword searches. This distinction is made at the time of loading into the database where keywords are checked against the ever-changing list of commercial terms, for instance, terms bidded by Overture. If a keyword is determined to be commercial, then it is assigned a keyword id, which may be compatible with the Overture keyword id list. If the keyword is not deemed commercial, then another id is assigned in the data loading process.

15 [0057] These two types of data (commercial and non-commercial) are loaded into separate sections of a data warehouse. At the time of a new search engine loading run, keyword data is extracted from both the commercial and non-commercial tables. The results are joined and unique keyword ids are assigned from a master table. A keyword can sometimes be found in both the commercial and non-commercial contexts. For instance, if
20 the commercial nature of a keyword is tied to whether it's being bidded upon, a keyword which is not commercial today may become a commercial tomorrow, if it receives a bid. This duality of keywords creates non-unique keyword ids as the same keyword may have both a keyword id created by a bidding process and a second one created through the load process. To eliminate this, it is preferred to create and maintain a single unique keyword id for every
25 keyword. This list is updated when new keywords are discovered, and assigns existing ids to keywords already in the system.

[0058] Some preprocessing may be performed by the behavior watching module to simplify the URLs reported. URLs are unwrapped and cleaned in a separate process. URLs are often wrapped by search engines to enable the serving search engine to track clicks on
30 served URLs. There are many different forms of URL wrapping. For example, a wrapped URL from yahoo.com might be: `http://rds.yahoo.com/S=2766679/K=bmw/v=2/SID=w/l=`

CLAR 1006-3

WS1/R=2/SS=100040736/H=1/SHE=0/*-http://www.bmwmotorcycles.com/. Unwrapping the URL produces http://www.bmwmotorcycles.com

[0059] From the server's perspective, the process begins with receipt (402) of behavioral information. The server uses whatever information it receives. From the user's perspective, the behavior watching module will report its observations and the user will receive search results ordered aggregating the user's information with others and/or will receive personalized to the user's recent behavior.

[0060] De-duping may also occur at the URL level. Information for two URLs which are identical is aggregated into one single URL. If two URLs differ even slightly however, (e.g., by a slash) then the two distinct versions are kept and another attempt at de-duplication is made as described below, for instance, using a combination of title and generated snippet.

[0061] Conversion data can be associated to a specific search by an algorithm that ties a search URL click to a specific conversion event, which occurred within a predetermined window. Usage and conversion data are matched for advertiser domains (URLs) that have clicks at the machine id, query time, advertiser domain level. For a particular machine with a click on a particular advertiser domain, if a conversion stat is observed within a predetermined window, then the conversion is attributed to that search click. If the conversion falls outside of the predetermined window, then the search click is not attributed.

[0062] For machine id-advertiser domain pairs that have a conversion stat attributed, subsequent future conversions are attributed as repeat conversions. These attributes also may be carried along and are available for use by a ranking algorithm (404). Metrics included with search data include number of visits, time spent (dwell time) and pages viewed.

[0063] Domain event data is joined to user data (with conversion metrics) at the machine id-advertiser domain level across sources, for combining search behavior for U.S. machines at google.com, msn.com, and yahoo.com. These results are put into time series order within machine id and advertiser domain. Domain events which occur within a predetermined time period following a search click are assigned as post-click metrics for that search click on that particular advertiser domain. If post-click metrics cannot be assigned to a particular search-click, the record is thrown out.

[0064] Several of the domain event data elements are subject to inaccuracies manifested in the client-sent stats. Both time spent and pages viewed are occasionally

CLAR 1006-3

misreported, and at other times, accurately reported, but in need of logically driven limits to be imposed.

[0065] Time spent can be misrepresented by machines having bad or inaccurate clocks. It can be accurately represented but in need of caps in such a case when a machine is left on a particular domain for an extended period of time. In order to cap outliers and to maintain a reasonable threshold for time spent on a site post click (403), a time limit of 30 minutes has been employed. A cap for pages viewed has also been implemented and set at 5. Other time and pages viewed caps may be substituted. These caps can be implemented on the behavior watching client side or after data is received at a host. In addition, user activity can be monitored by the behavior watching module so that extended periods of inactivity are not counted as dwell time.

[0066] When data is joined and aggregated at the machine-id, keyword, and URL level, the resulting data structure may include: machine-id, keyword, keyword_id, URL, URL_ID, domainid (corresponds to the domain of the URL), clicks, dwell_per_click, pages_per_click, conversions_per_click, rank_position (from search results list viewed by the user). Optionally, only keywords of predetermined interest (402) may be processed. For instance, keywords having commercial interest, such as bidded keywords, may be processed.

[0067] Outliers optionally may be removed (403), to avoid scoring anomalies. One example of an outlier is a link that is returned only once by a search engine and followed with enthusiasm by the user when was returned. A single strong sample point can give a link an unbeatable average score. A link that appeared just before a keyword was rescored could potentially be ranked in the top position for that keyword on the basis of a single sample! Accordingly, one example of optionally removing outliers is to not rank links unless they have been followed a predetermined number of times. For instance, if a link has not been selected by users and followed at least 10 or 100 times, it might remain unranked until its activity level reached the predetermined level or threshold.

Ranking Algorithm

[0068] A URL ranking algorithm (404) has been developed to identify and rank links for any given keyword. Many variations on combining the aggregated observations have

CLAR 1006-3

merit, as discussed below. One combination uses URL click rate and dwell time metrics (time spent at the domain and/or number of pages viewed), to select and rank URLs.

[0069] Optionally, user responses can be normalized for the position of a URL in a search result set. Position produces an inherent bias in URL click data for search results, which may be desirable or not. URLs occupying higher ranks garner higher clicks. In order to account for this bias, a normalization algorithm was developed to put clicks on links in disparate positions on equal footing.

[0070] Data is aggregated for each position and average click through rates, average time spent, and average pages viewed are calculated. For each rank position (1 – n), there are at least three average aggregate measures of that position’s importance: click through rate, time spent and pages viewed. Normalization of these measures can be expressed as:

$$\text{Pages Viewed Inflation Factor for Position } n = \frac{\text{Avg Pages Viewed (Position 1)}}{\text{Avg Pages Viewed (Position } n)}$$

$$\text{Dwell Time Inflation Factor for Position } n = \frac{\text{Avg Dwell Time (Position 1)}}{\text{Avg Dwell Time (Position } n)}$$

$$\text{Pages Viewed Inflation Factor for Position } n = \frac{\text{Avg Pages Viewed (Position 1)}}{\text{Avg Pages Viewed (Position } n)}$$

15 This embodiment calculates and ranks top URLs for any keyword (404) based upon observed user metrics.

[0071] This embodiment may re-rank results based upon a time share metric, which corresponds to an individual machine’s percentage vote. The algorithm takes into account user web surfing patterns and effectively places users on equal voting for relevant links. The premise is that a user has a certain amount of time which is spent on a site post a search click. These times are totaled to form the individual user’s total time value which was spent viewing sites post search clicks. Percentages of the total time are then calculated for each URL click made by that particular user, resulting in a time fraction vote.

[0072] In using this methodology, users who in general spend less time surfing the Web have the same voting power as users who tend to spend longer amounts of time. This

CLAR 1006-3

evens the playing field across all categories. Alternatively, other embodiments may take into account other factors which make up an individual user's profile. Users identified as category experts may have a higher vote. For example, a user highly involved in the electronics category may have his vote count more for links clicked than a user who is new to the category. The voting blocks may take place within a particular category, and not across all categories as a whole.

[0073] The following detailed computational example applies to a specific keyword-URL pair through the aggregation and cleansing process.

[0074] Links in higher positions garner higher clicks and hence possess higher click through rates. It follows that users also tend to spend longer amounts of time and view more pages at URLs occupying these higher ranks. In order to account for this bias, a normalization process is applied.

[0075] Average clicks, average dwell time, and average pages viewed are calculated for each position regardless of keyword or URL combinations. These numbers are shown below as Position Averages. Average fraction clicks, dwell time and page views are also shown for this keyword, URL combination. These average fractions correspond to the percentage of each metric devoted by all machines to each keyword, URL, position grouping.

Kwd	URL	POS	Position Average Share of Clicks	Position Average Share of Dwell Time	Position Average Share of Page Views	Average Fraction Clicks	Average Fraction Dwell Time	Average Fraction Page Views
01	101	1	0.15	0.14	0.18	.25	.14	.08
01	101	3	0.13	0.12	0.15	.21	.17	.09
01	101	7	0.09	0.11	0.1	.32	.18	.11

[0076] Inflation factors are calculated for each position and applied to the appropriate observed metrics, normalizing them for position. The normalizing equation:

$$\text{Inflation Factor For Position } n = \frac{\text{Average (Position 1)}}{\text{Average (Position } n)}$$

Clicks, dwell time, and page views in lower positions are factored up by the appropriate inflation factor observed for that particular position. In the example below, the average fraction dwell time for position 7 pre-normalization was .32. After the application of the

CLAR 1006-3

inflation factor, the average fraction dwell time was .53. This number is now a normalized dwell time.

Kwd	URL	POS	Clicks Inflation Factor For Position <i>n</i>	Dwell Time Inflation Factor For Position <i>n</i>	Page Views Inflation Factor For Position <i>n</i>	Normalized Clicks	Normalized Dwell Time	Normal Page V
01	101	1	1.00	1.00	1.00	0.25	0.14	0.08
01	101	3	1.15	1.17	1.20	0.24	0.20	0.11
01	101	7	1.67	1.27	1.80	0.53	0.23	0.20

[0077] Weighted averages are calculated for each metric at each position, and totals are calculated across the positions for clicks, and all of the weighted average fraction measurements; clicks, time and page views.

[0078] Over time, with new reports of web usage from behavior watching modules, URLs for a particular keyword will adjust their positions. A URL in position 1 this week may be in position 7 the next. Adjusting for position is therefore a cleanup and adjustment process.

Kwd	URL	POS	Absolute Clicks	Weighted Average Fraction Clicks	Weighted Average Fraction Time	Weighted Average Fraction Page Views
01	101	1	155	38.75	21.70	12.40
01	101	3	124	30.05	24.59	13.39
01	101	7	98	52.27	22.45	19.40
01	101	ALL	377	121.06	68.74	45.20

[0079] The final result is the Weighted Average of Normalized Totals for Keyword 01, URL 101, at any position. Computed for every keyword, URL combination a single score is calculated for each metric needed for the v17.1 algorithm. These metrics are now normalized for position, and for multiple rank occurrences.

KWD	URL	POS	Absolute Clicks	Weighted Average Fraction Clicks	Weighted Average Fraction Time	Weighted Average Fraction Page Views
01	101	ANY	377	0.32	0.18	0.12

[0080] Ranks for scoring can be based on time that a user spends viewing pages on the domain. The information received from the behavior watching module may limit the maximum amount of time that will be assigned for any viewing session or it may track the

CLAR 1006-3

user's behavior, such as window navigation between programs, mouse clicks or mouse movement, and disregard periods of inactivity when calculating dwell time.

[0081] Ranking may ignore links that were selected by users less than a predetermined number of times, which may be predetermined as a fixed number or a function of traffic for the keyword or category. Ignoring outlier links (403) may avoid giving a high ranking to a link that was rarely presented by the search engines and followed only once or twice.

[0082] Ranks for scoring also can be based on a combination of click through rate, dwell time and the number of pages or documents viewed after following the link. Combining these factors, in some instances one of the factors will dominate: all or more than two thirds of the ranking weight may be assigned to just one of click through rate, dwell time or number of pages or documents viewed. Alternatively, they may be equally weighted, plus or minus 10%, or the factors may be assigned weighting ratios of approximately 2-1-1, plus or -10%, so that one factor is given approximately half of the combined weighting.

[0083] Another factor that can be used in ranking is return visits. If the user returns to the domain within a predetermined time after leaving it or within a predetermined number of navigation events, the user's return to the site can be assigned significance. Return to the site may reflect a favorable impression after considering other sites.

[0084] Conversion from browser to buyer or registered lead can be considered to be particularly worthwhile as a factor. Again, conversion may include both a purchase in the domain and a registration. In some instances, such as car or home purchases, registration may be more realistic measure, because the purchase may be impractical or infrequently completed at a web site. Return conversion also may be taken into account.

[0085] Results may be segregated for analysis by search engine and ranks scored. Then, the separate rank scores may be combined into an overall ranking.

[0086] Statistical or other analysis can be applied within categories or keywords to determine which combination of ranking factors best attracts users to follow a link responsive to a search. It is anticipated that ranking information will be used differently among categories of keywords. Time spent will be important in the auto category. Conversions will be much more important music downloads category.

CLAR 1006-3

[0087] Optionally, click segmentation bands may be applied. These bands give precedence to URLs with high numbers of clicks. Employing these bands may improve the resultant links on selected algorithms.

[0088] Segmentation bands are identified based upon total clicks received by a particular URL. For instance:

Tier 1: URLs with 100+ clicks

Tier2: URLs with between 50 and 99 clicks

Tier3: URLs with between 10 and 49 clicks

Tier4: URLs with less than 10 clicks

10

[0089] URLs for a particular keyword are first put into the appropriate segmentation band. Once the band is identified, these URLs are set in descending order by rank score.

[0090] A predetermined number of links, such as the top 15 links (405), may be selected for data collection, to be followed by a spider engine (406).

15 [0091] Three tables are generated as output from the rank process:

Keyword table (KEYWORDID, KEYWORD)

URL table (DOMAINID, URL_ID, URL)

Rank Table (keyword, URL, rank)

Following Links

20 [0092] Traditional crawling programs at other search engines (ex: Slurp at Yahoo!, Googlebot at Google, MSNBot at MSN) crawl the entire web in search of relevant pages to index to be used in determining the rank order of links to display for a given keyword. The embodiment disclosed here, in contrast, is given a succinct number of URLs to crawl, which may optionally be selected (405) from links reported by the behavior watching module (401).

25 These links are pre-ranked, hence this information retrieval process needs not determine the relative importance of a given URL from its connections to others, but rather to obtain the best possible descriptive information from the URL.

[0093] This embodiment takes a specific set of URLs and performs several specific tasks: It strips out all HTML tags and returns first 1,000k or another predetermined chunk of the text on the page to a file. It takes and stores a mapping from the text object's value into a

30

CLAR 1006-3

uniform scalar space to be used as a text signature or text fingerprint. It calculates an MD5 or other fingerprint of the document (with or without html tags). It calculates a summary count of the characters within the text extracted from the document.

[0094] This method may be implemented by a Java application which operates in a
 5 Linux environment as illustrated by **FIG. 18**. The application (1811) is a brokered one which allows for a large and scalable number of threads to run simultaneously across multiple Linux servers. In the current configuration, each Linux server (1821, 1822) can handle up to 10 spider instances (1831-33, 1834-36). Each spider in turn runs 200 worker threads (1841-43, 1844-46). At capacity, this configuration is able to spider 800 links a second.

10 [0095] Total threads working for a single broker can be arrived at by the following equation: With i number of spiders each having j number of workers (threads):

$$\text{Total threads for } i \text{ \# of Spiders: } \sum_{i=1}^n \left(\sum_{j=1}^n \text{worker } j \right)$$

[0096] The dual-broker model (1821, 1822) can segregate keywords by keyword
 15 velocity. General keywords are funneled through a robust, heavy duty version of the ranking algorithm. Fast moving keywords (e.g., news, current events) can be processed through a nimble, express version of the ranking algorithm, which uses less history. Keyword velocity is a measure of how quickly the popularity of a keyword changes. The highest velocity keywords can be selected by comparing the number of keyword searches in the last 24 hours (day 0) against the 24 hours before that (day-1). A different time span, such as four or eight
 20 hours, can be used, of course. How far the ration day 0/day-1 varies from "1" is the keyword velocity. If the ration is less than 1, the keyword is becoming less popular, "old news." If the ratio is much more than 1, the keyword may relate to a new story. Generally, a predetermined number of relatively high velocity keywords are re-indexed at a predetermined interval or as resources permit. In one embodiment, the top 10,000 keywords are re-indexed each day.
 25 While one metric of keyword velocity or volatility has been described, variations are anticipated.

[0097] The heavy duty version handles the ranked keyword URL pairs. These ranked keywords URL pairs are made available through an Oracle table on a database. The

CLAR 1006-3

URL_TABLE includes: DOMAINID, URLID, URL, LENGTH, SIGNATURE_H,
SIGNATURE_T, SPIDER_DATE and HOST

[0098] The DOMAINID, URLID, and URL fields are populated from a reference
database prior to following the links. After the link-following process for a specific URL, the
5 LENGTH, SIGNATURE_H, SIGNATURE_T, SPIDERDATE, and HOST fields are written
back to the database.

[0099] Brokers use Java Database Connectivity (JDBC) to connect in to the Oracle
database. The broker accesses the URL_TABLE from the ranking process. The broker
makes a request for 1/100th of the total number of domains which are available in the URL
10 table for which SPIDER_DATE is null. All URLs associated with these domains are
extracted by the broker where they are grouped by domain. Individual spider boxes talk to the
Broker via Remote Method Invocation (RMI) requesting URLs for domains 1,000 domains at
a time. Domains are then passed from the spider to a worker who takes all of the URLs
associated with its domain and operates upon those URLs.

15 [00100] URLs are passed to the workers grouped by domain in order to accommodate
generally accepted crawling or link following practices so as not to swamp domains with
thousands of requests simultaneously. It is a generally accepted practice to not access a single
domain with more than one request at a time. The link following process (406) respects this
generally accepted principle by assigning each worker all URLs associated with a given
20 domain.

[00101] The link following process (406) is a robust, scalable, application which
fetches content and calculates statistics from a specific URL. Once a worker receives a
domain and its associated URLs, it accesses that URL using HTTP protocols. If a good
response code is received, a link following worker goes to work on that page. The worker
25 receives a 200 response code (status OK) more than 98% of the time. If the page returns an
HTTP code indicating a redirect (codes 301, 302, 303 and 307), further action must be taken
by the worker or system in order to obtain information about that URL. A worker will follow
up to 5 redirects from an initial URL before abandoning. Once the worker reaches an end
point, the following tasks take place: Acquire HTTP return code from the URL. If a good
30 response code is acquired: Identify title meta tag if available; calculate an MD5 fingerprint of

CLAR 1006-3

the entire document (both HTML and text); parse HTML from the page; and write back first 1,000k of text to disk.

[00102] Once the content is parsed and written back to the disk, a subsequent operator takes over. This operator makes several calculations used for the document fingerprint and writes those and other statistics back to the Oracle database. The system writes back the following fields to Oracle: URL_TABLE, DOMAINID, URLID, URL, LENGTH, SIGNATURE_H, SIGNATURE_T, SPIDER_DATE and HOST.

[00103] LENGTH is a count of characters in the text of the document (first 1,000k). This feature can be used for de-duping URLs later in the process (408). SIGNATURE_H is the MD5 hash code signature. SIGNATURE_T is a CRC32 checksum code of the text (first 1000k). SPIDER_DATE indicates the date and time that the particular URL was accessed. HOST pertains to which spider machine stored the text of the URL.

[00104] The following system may create three different measures designed to aid in document de-duplication (409). This de-duplication process aims at identifying documents that are identical or very similar within a given keyword result set. In a prior step not separately shown, URLs are de-duped at the URL level. Easily identified duplicates such as two occurrences of the exact same URL are eliminated. The system attempts to eliminate URLs that do not appear to point to the same page, but in fact do. In one embodiment, mathematical signatures (fingerprints) are taken for each URL and compared to other URLs within a given keyword result set. Three exemplary signatures are a length signature, an MD5 signature and a CRC32 checksum. Other signatures may be substituted.

[00105] For the length signature, the character length of the text document is calculated. This measure aids in the de-duping process to aid in giving context to a page which has been identified as a duplicate. For instance, if two sites show identical MD5 and CRC32 signatures, but have very disparate URLs, the signature is analyzed. If the length signature is low, meaning the page is small, it is likely that these two URLs share, for instance, a standard warning screen as would be found prior to entering an adult content site.

[00106] An MD5 signature typically is a way to verify data integrity. The MD5 algorithm takes as input a message of arbitrary length and produces as output a 128-bit "fingerprint" or "message digest" of the input. The MD5 algorithm is intended for digital signature applications, where a large file must be "compressed" in a secure manner. The

CLAR 1006-3

system computes an MD5 signature for the entire document, reducing the identity comparison process to a 128-bit comparison, for instance.

[00107] A CRC32 checksum generates the cyclic redundancy checksum polynomial of 32-bit lengths. This is usually used to validate the integrity of data being transmitted. The CRC is a "digital fingerprint" of a file, With CRC32 you can "melt down" a huge 20 MB (or even much bigger) file to have a small, handy reference to it, a single 32-bit number like 7d9c42fb (hexadecimal notation) which reflects the entire contents of this huge file. The system computes a CRC32 signature of the text of the document, giving insights into the text content of the page.

10 [00108] Another signature that can be calculated and used is Rabin's fingerprinting algorithm, for instance Broeder's implementation, which produces a compact checksum.

[00109] Any of the checksums or fingerprints can be applied to the whole document, the whole document less HTML tags stripped away, the selected chunk of the document that is cached, the title and snippets or some other predetermined excerpt from the document.

15 More or less than 1,000k of the document can be used.

[00110] The process completes a run for a particular URL with data being written to an Oracle database and a spider box. The Oracle database receives fingerprint information (length, MD5, crc32), spider date/time, and host location information written to URL_TABLE and a spider date/time stamp written to KEYWORD table. The spider box receives files for data links that it followed: URL, title (if it was obtained during the initial fetch from the URL) and text of the document (first 1,000k) to be used for snippet generation. The text contains elements of the meta description and the body of the document

Snippet Generation

[00111] A snippet generation process generates titles and snippets for display (407).

25 The snippet process takes a keyword phrase and URL combination, comes up with the best title describing that URL, and creates the best snippet (i.e., abstract, description) for that URL outlining in a 200 character space the information contained in the URL that pertains to the keyword. Snippet generation follows the link following process. Snippets are created from the text of the document retrieved from the chosen URL.

30 [00112] The keyword "somec bicycles" produces the following sample text for display:

CLAR 1006-3

Upland Sports Group: Somec Bicycle Frames

... is a mystique to Italian **bicycles**; an aura born of precision, innovative materials and artistry. Societa Meccanica, better known as **Somec**, has contributed to ...
www.uplandsg.com/somec/ - 12k - Aug 17, 2004 - [Cached](#) - [Similar pages](#)

In this example, the title is "Upland Sports ... Frames". The snippet is the two lines following the title. The URL is on the bottom line.

[00113] Titles are usually generated from the title of the page retrieved when a link is followed. Most sites annotate the title of the page for search engines through the use of HTML meta tags. A tag identifying the title is present on over 97% of all URLs.

[00114] In the 3% of URLs for which the HTML tags do not supply a title, the process composes a title. If there is text available for the URL, the process takes the first approximately 70 characters of text (respecting word boundaries) and creates a title. If there was no text generated from the URL, the domain name is stripped from the URL (all information between www and .com) and displayed as the title.

[00115] Snippet generation is a mix of art and science. The process creates snippets leveraging mathematical equations and linguistic science. In one embodiment, snippets can be comprised of 1 single sub-snippet, or up to 3 sub-snippets separated by ellipses (...). A scoring algorithm decides which sub-snippets when combined (or not in the case of a single sub-snippet) produce the best score.

[00116] The snippet scoring algorithm is a multi-step process which scores various portions of the document's text. In four parts, it includes keyword tokenization, window scoring, window trading and final determination.

[00117] Keyword tokenization is applied because keywords are not always single words. Keywords are often multi-word phrases. The process tokenizes or single outs individual words within a phrase. Identifying individual word tokens typically includes searching for word separators such as spaces, periods, commas, or colon/semicolons. Once the tokenization of the keyword phrase is complete, the window scoring routine can commence.

[00118] In one version of window scoring, windows of three different sizes are calculated within the text of the document, for instance, for sub-snippet of lengths 200

CLAR 1006-3

characters, 100 characters and 66 characters. For various display formats and user preferences, a total snippet length of 50, 100, 200, 500 or 1,000 characters may work best. The example here was 200 characters. **FIG. 20** depicts how the window scoring routine begins with the first word of text (2003) within the document (2004). In the first case, the window spans 200 characters. This window receives a score based upon several metrics. Once the score is calculated, the window moves 1 character to the right and calculates another score for the next window. Alternatively, the window may move one word to the right, aligning either the window start or window end with word boundaries. This process iterates (2004, 2005) until a score for the final window is obtained. Intelligence may be built into the algorithm to eliminate the need to compute the score for each entire window, but rather add the additional information to the previous window's score and subtract the information which was lost in the move. Titles (2002) also may be captured, as described above.

[00119] When the process is complete, there may exist:

- i windows of length 200 (where $i = \text{document length} - 200$)
- 15 j windows of length 100 (where $j = \text{document length} - 100$)
- k windows of length 66 (where $k = \text{document length} - 66$)

[00120] Window scoring may be based on one or more metrics, such as the number of unique tokens found within the window, the total number of tokens found within the window, and/or the longest string of tokens found within the window. A formula for each window is computed from a combination of these metrics and assigned to that window.

[00121] In the case where there is one 200 character snippet, the window with the highest score is chosen. The two highest scoring windows of length 100 are chosen for the two sub-snippet model. The three highest scoring windows of length 66 are chosen for the three sub-snippet model.

25 [00122] The best scores are calculated for each model (1, 2, or 3 sub-snippets). A final algorithm may be applied when 2 or 3 windows are eligible for a snippet. If the global window score can be increased by one window giving up characters to another, then that action is seen as a gain and it is taken. If the global window score cannot be raised in this manner, the snippets are used without trading.

30 [00123] The output from snippet generation may include 5 different scores: Score of single sub-snippet model; score of non-traded two sub-snippet model; score of traded two

CLAR 1006-3

sub-snippet model; score of the non-traded three sub-snippet model; and/or score of the traded three sub-snippet model. Of these, the single highest score is chosen and that sub-snippet model is applied to that keyword, URL combination.

[00124] For a sample keyword = "red dog run", the following steps may be followed:

5 [00125] Step1: Tokenize keyword into three tokens:

- Token1 = red
- Token2 = dog
- Token3 = run

[00126] Step2a: Locate instances of the tokens within the text document

10 -----red-----dog-----dog-----
 ----run-----red-----red-----
 -----red-dog-run-----dog-----
 ---run---red-----run---run-----

15 [00127] Step2b: Score the windows and identify the top ones. In this example, the three sub-snippet model, the best 3 windows were calculated.

-----red-----dog-----dog-----
 [-----run-----red--]------red----- (1)
 ---[----red-dog-run--]------dog----- (2)
 --run----[-red----run---]-run----- (3)

20 [00128] Step3: Allow for trading to occur. In this case, if window 1 can give up some of the non-token containing characters within it's left edge to window 3. This allows window 3 to expand and include the final token 'run', increasing the overall global score of the snippet.

-----red-----dog-----dog-----
 25 -----[run-----red--]------red----- (1)
 ---[----red-dog-run--]------dog----- (2)

CLAR 1006-3

---run-----[-red-----run----run]----- (3)

[00129] In an alternative embodiment, the snippet generation process may involve the creation of an approximately 200 byte field used as a descriptor for the associated link.

5 Snippet generation takes place post spidering and is created from the complete text of the document associated with the chosen URL or at least the portion of the document stored.

Personalization and Ranking

[00130] Within the ranking algorithm, there is the ability to select anonymous users who, based upon their behavioral profile, would have their votes for particular categories of links count more than other users.

[00131] Users who are heavy searchers (based upon their observed search behavior) would have their votes count more on links that they click more than the votes of novice searchers on that same link. In this way, the search experts would help produce more relevant ranking results.

15 [00132] Similarly, users who are highly involved in a particular category would have their votes count higher in that category than users who have no involvement in that category. Using behavior watching modules, one can identify users who are highly involved in various categories such as digital cameras, fantasy sports, or automobiles. For example, a user identified as being highly involved in the digital camera category would have his vote count more for links he clicked after a search for 'Canon G3' than a user who is new to the category searching on that same keyword.

[00133] Identification of a user's category involvement status also drives personalization. A user with a high degree of involvement in a particular category would get different results from a user identified as less involved. This personalized results serving would require the presence of a cookie like object available on a particular machine. This lifestyle cookie would provide the search engine with a behavioral profile of the user, obtained from the users category navigational patterns. These category specific navigational patterns would be obtained from information contained in a categorization structure that also can be used for targeted advertising. For commercial purposes, a budget category or likely budget can be inferred from sites visited. Visitors to IKEA and Target are likely to have a

CLAR 1006-3

different budget for apparel than visitors to Sachs Fifth Avenue or Bloomingdale's. Similarly, Hyatt Hotels are in a different budget category than youth hostels.

[00134] Personalization based on observed communications is much more powerful than user-entered customization, because research shows that only 8-14 percent of users manually personalized their content. Personalization highly correlates with pages viewed at a domain: users who personalize have been reported to view 130 percent more pages at the domain than users who do not personalize.

[00135] Sometimes different behavioral profiles can be leveraged to make a difference in search results. Other times, differences between two users' behavioral profiles does not help in the context of a particular search keyword.

[00136] Some examples are helpful. First, an ambiguous search terms example: A highly involved auto category user who searches for "jaguar" would get more Jaguar auto related links than jaguar animal related links as compared to a normal mix of auto and animal related links for someone with no identifiable category involvement. Identification and usage of these behaviorally profiles could slant results, without completely replacing results. In the example above, the auto category involved user could get 100% auto results, or just a larger percentage of auto results than found among popular websites.

[00137] Next, a sub-category identification example: Three users search for the keyword "rental car". Three separate sets of results come up, each personalized for the users. Each user has a particular behavioral profile obtained from their past navigational patterns observed within the travel category. These behaviors are readily identifiable from the observed communications.

[00138] User 1: Frequent business traveler – his rental car results would be slanted toward the business traveler car rental results, possibly more about frequent rental points, etc.

[00139] User 2: Budget traveler – his rental car results would be slanted toward the budget traveler; rent-a-wreck type results, specials on sub-compact cars etc.

[00140] User 3: Luxury Traveler – his rental car results would be slanted toward the high-end luxury traveler; sports car rentals, classic car rentals, etc.

[00141] FIG. 12 compares the insight that a behavior watching module on a user's machine gives regarding web surfing with the limited view that a server can achieve. These statistics were generated to compare what a Yahoo! server knows about user behavior,

CLAR 1006-3

compared to what the behavior watching module sees. The x-axis 1202 identifies categories of web sites. The y-axis is a relative scale, with Yahoo! server observation scaled to "1". In each stack, Yahoo! 1203 has a visibility of 1. The behavior watching module has a relative score 1204 that is higher and varies strongly between categories. This surprising degree of difference strongly distinguishes the on-machine behavior watching module from visibility that even a widely used search engine with many affiliated sites can achieve.

[00142] FIG. 13 depicts a web page highly relevant to the keyword "cars". Applying a ranking algorithm, the CarsDirect home page scores 97 of 100. Factors used to calculate this score may include initial visit metrics, subsequent visit metrics and user level metrics.

10 Observing users' responses to a particular site across search engines responsive to the keyword "cars", the behavior watching module observed a click through rate of 12 percent. An average user viewed five pages at this web site, following the initial link, before they left the site. The average user spent 55 seconds looking at the pages and the conversion rate, here a registration more likely than a purchase, was 80 percent. On subsequent occasions, within a
15 predetermined time after the search or the first visit, an average returnee viewed 15 pages and spent 300 seconds on the site. Returnees registered with a conversion rate of 15 percent. Among users following this link for this keyword, category (e.g., car category) involvement was high, ranked "Expert". Search engine usage was "Maven", indicating that common users were more frequent than computer experts. For this category, among U.S. users, the
20 geographic location of the web site was important. Web sites in the U.S. were more interesting to users (e.g., car buyers) than overseas web sites.

[00143] FIG. 14 is a counter example of a page weakly relevant to "cars". The title of this article matches the keyword, but the click through rate and other metrics indicate that users considered this page much less relevant to searching for "cars" than the CarsDirect
25 home page. The calculated relevancy score for this page is only 11.

[00144] FIG. 15 juxtaposes the results for "cars" from Yahoo!, MSN Search and Google with an embodiment of the search engine described herein. Using a relevancy index, this embodiment produces more relevant search results. The relevancy index is the ranking score, which reflects users' responses to links in search result lists. Sample scores are: MSN
30 598, Google 524 and Yahoo! 487. This significantly improved score for this embodiment

CLAR 1006-3

strongly distinguishes the use of an on-machine behavior watching module and the aggregate behaviors of a multitude of users from the ranking algorithms of the leading search engines.

[00145] **FIG. 16** illustrates user differentiation using results from a behavior watching module. A large sample of user behavior, particularly shopping in the apparel category, is reported in this figure. Users were classified in nine boxes along attributes of amount of browsing in the category (1601) and degree of loyalty (1602) to a particular web domain (alternatively, to a vendor, brand or other category). Amount of browsing was classified as light, medium or heavy. One quarter of users were in each of the heavy and medium usage categories and the remaining half in the light usage category. Loyalty is classified as loyal, a switcher among sources or loyal to competitors. The 3-by-3 arrangement of attributes produces nine boxes. Among loyal shoppers, about one million unique observations were made. The heavy users had a higher click through rate when a particular result was displayed in response to a search, the rate being 7.90%, compared to 5.40% for medium users and just 3.00% for light users. After a click-through event, heavy users converted to purchase or registration at a rate of 8.00%, compared to 7.20% for medium users and 6.50% for light users. Not surprisingly, there are more switchers and competitive loyal unique observations than for the particular domain observed. Personalization of content, such as ads, may depend on which of the nine boxes the user fits. For instance, a loyal, heavy user will get a different message (a la frequent flier programs) than a medium to light switcher. A competitor-loyal heavy user might get a trial offer that would be completely misdirected to a loyal user. The strategy for advertising a heavy, loyal user is to insulate them from competitors. For medium to light loyal users, the strategy may target up-selling and cross-selling. For heavy to medium switchers, the objective is to increase share. Heavy to medium competitive loyal users are targeted for acquisition. Light switchers and light competitive loyal users may be ignored, to focus advertising resources.

[00146] Return visit data from the behavior watching module can assist an advertiser in measuring the effectiveness of a particular ad. User differentiation by box can further be associated with selection of ads and evaluation of ad effectiveness.

[00147] Cross-browsing of users also can be reported. Users can be selected by follow-through, for instance all click-throughs or all users with conversions. The users with a

CLAR 1006-3

conversion at a particular domain (or vendor or brand, for instance) can be rated by the frequency of their visits to competitors' domains (or vendors or brands).

[00148] FIG. 17 depicts traffic to a particular domain, from a particular search engine, by keyword. This information may be useful to a company devising ads and also when
5 bidding on keywords for sponsored content positioning on search engine result pages. Data from the behavior watching module can yield impressions, click, click-through rate, conversion rate and position on the search engine's results page.

Some Particular Embodiments

[00149] The present invention may be practiced as a method or device adapted to
10 practice the method. The same method can be viewed from the perspective of a user at their terminal or personal computer or on the server side, collecting information from users. The invention may be an article of manufacture such as media impressed with logic to carry out computer-assisted method.

[00150] A device embodiment, from the user perspective, may be embodied in a
15 module running on the user's computer and collecting behavioral observations, coupled to a server that responds to the behavioral observations with information personalized to the user.

[00151] While the present invention is disclosed by reference to the preferred
embodiments and examples detailed above, it is understood that these examples are intended in an illustrative rather than in a limiting sense. Computer-assisted processing is implicated in
20 the described embodiments. Accordingly, the present invention may be embodied in methods aggregating of communication patterns, pre-processing links responsive to keyword searches, responding to keyword searches using aggregated communication patterns to rank the responsive links, and responding to keyword searches using recent and current navigation information systems to resolve ambiguities and/or personalize responses based on user
25 characteristics. Other embodiments, as devices, include logic and resources to carry out these methods. As systems, still other embodiments include behavior watching modules on terminals, servers that process or respond to the behavioral data, or both. Other embodiments include media impressed with logic to carry out the methods, data streams impressed with logic to carry out the methods, or computer-accessible services that carry out the methods. It
30 is contemplated that modifications and combinations will readily occur to those skilled in the

CLAR 1006-3

art, which modifications and combinations will be within the spirit of the invention and the scope of the following claims.

[00152] The embodiments disclosed may be practiced as a method or device adapted to practice the method. They may be practiced as a system. The same method can be viewed
5 from the perspective of a server collecting data or a behavior watching module generating data. The invention may include an article of manufacture such as machine readable memory including instructions to carry out the method or a data stream including instructions to carry out the method.

[00153] One embodiment is a method of selectively collecting web pages that may be
10 returned to users in response to search requests. This method includes receiving from behavior watching modules operating on a multiplicity of users' computers, information regarding the users' search engine usage across one or more search engines. This information includes at least keywords (including phrases) submitted by particular users to the search engines, links selected by the particular users from results returned by the search engines, and
15 at least one of dwell time or documents viewed by the particular users when following the selected links. The method further includes using the search engine usage information to choose keywords of interest and to choose links selected by the users corresponding to the keywords of interest.

[00154] The multiplicity of users in this embodiment may exceed 2000 or 100,000
20 users, depending on the desired coverage of frequently submitted search terms. The information received from behavior watching modules may span a plurality of unaffiliated search engines, providing data that no single search engine could collect.

[00155] This embodiment further may include ranking, for the chosen keywords, the corresponding links based on user access rates calculated from the search engine usage
25 information and at least one of the dwell time or the documents viewed. The method further may include limiting the collection of content based on a predetermined application of the ranking, such as following only the top 10 or top 15 links for a particular keyword.

[00156] In some embodiments, the search engine usage information is known to have
30 been qualified before receipt to discount dwell time during periods when a particular user is likely to have been inattentive.

CLAR 1006-3

[00157] This embodiment further may include selecting snippets from the content of a document reached by following one of the chosen links. This snippet selecting may include applying a window of predetermined length to the document, repeatedly shifting the window through the document. The keyword of interest, if it is a phrase, is divided into words. The method further includes repeatedly calculating one or more window scores for the words in the window, including a count of instances of the words in the window and a measure of adjacency of the words in the window. With the window scores, a plurality of non-overlapping window positions are chosen as snippets. The snippets or references to the snippets are stored in a machine readable memory. This method may be refined in many ways. The repeated shifts of the window may align the window so that its start or end matches a word boundary. As the window is repeatedly shifted, the score within the window may be calculated by reducing the score for characters removed from the window and increasing the score for characters added to the window.

[00158] Snippet generation may be refined by trading off lengths of the snippets, approximately maintaining a combined length of all the snippets. As the lengths of the snippets are traded off, the one or more window scores are recalculated. The method further includes choosing a particular trade-off of snippet links using the recalculated window scores, optionally merging the snippets, and storing the chosen snippets or references to the chosen snippets in a machine readable memory.

[00159] Viewed from the perspective of a network of behavior watching modules, the method includes using a multiplicity of behavior watching modules operating on a multiplicity of users' computers, the behavior watching modules collecting information regarding the users' search engine usage across one or more search engines. The information collected includes at least keywords (including phrases) submitted by particular users to the search engines, links selected by the particular users from results returned by the search engines, and at least one of dwell time or documents viewed by the particular users when following the selected links. The method further includes the behavior watching modules electronically reporting the search engine usage information to one or more affiliated servers. As above, the multiplicity of users in this embodiment may exceed 2000 or 100,000 users, depending on the desired coverage of frequently submitted search terms. The information

CLAR 1006-3

received from behavior watching modules may span a plurality of unaffiliated search engines, providing data that no single search engine could collect.

[00160] The behavior watching modules may qualify the search engine usage information before reporting, to discount dwell time during periods when a particular user is likely to have been inattentive.

[00161] The behavior watching module embodiment may further include a particular user invoking a search engine with a keyword and receiving results from the search engine organized to reflect the reported search engine information.

[00162] A system embodiment includes at least one listener module and at least one computer or cluster of computers in communication with the listener module. The listener module includes a network interface that receives from behavior watching modules operating on a multiplicity of users' computers connected by network, information regarding the users' search engine usage across one or more search engines. The information includes at least keywords (including phrases) submitted by particular users to the search engines, links selected by the particular users from results returned by the search engines, and at least one of dwell time or documents viewed by the particular users when following the selected links. The computer or cluster of computers in communication with the listener module operate one or more modules that include logic and resources adapted to process the search engine usage information and choose keywords of interest and the links selected by the users corresponding to the keywords of interest, to follow the chosen links to collect at least part of the content of documents addressed by the chosen links, and to associate the collected content with the corresponding keyword of interest.

[00163] The multiplicity of users in this embodiment may exceed 2000 or 100,000 users, depending on the desired coverage of frequently submitted search terms. The information received from behavior watching modules may span a plurality of unaffiliated search engines, providing data that no single search engine could collect.

[00164] Another embodiment is a method of selecting snippets from content of a document that may be returned to users in response to search requests. The snippets selecting includes applying a window of predetermined length to the document, repeatedly shifting the window through the document. If the keyword of interest is a phrase, is divided into words. The method further includes repeatedly calculating one or more window scores for the words

CLAR 1006-3

in the window, including a count of instances of the words in the window and a measure of adjacency of the words in the window. With the window scores, a plurality of non-overlapping window positions are chosen as snippets. The snippets or references to the snippets are stored in a machine readable memory. This method may be refined in many
5 ways. The repeated shifts of the window may align the window so that its start or end matches a word boundary. As the window is repeatedly shifted, the score within the window may be calculated by reducing the score for characters removed from the window and increasing the score for characters added to the window.

[00165] Snippet generation may be refined by trading off lengths of the snippets,
10 approximately maintaining a combined length of all the snippets. As the lengths of the snippets are traded off, the one or more window scores are recalculated. The method further includes choosing a particular trade-off of snippet links using the recalculated window scores, optionally merging the snippets, and storing the chosen snippets or references to the chosen snippets in a machine readable memory.

15 [00166] Snippet generation may be combined with duplicate elimination. The method may further include evaluating to more documents based on their titles and the calculated snippets. If the titles and snippets or fingerprints or checksums calculated from the titles and snippets match or are very similar, the documents may be declared and handled his duplicates.

[00167] In snippet generation, the count of instances of the words in the window may
20 include both counting how many of the words in the keyword appear in the window and how many times the words appear in the window.

[00168] A further embodiment is a method of ranking including receiving, from
behavior watching modules operating on a multiplicity of users computers, information regarding the users search engine usage across the plurality search engines. Information
25 received includes at least keywords (including phrases) submitted by particular users to search engines, links selected by the particular users from results returned by the search engines, and at least one of dwell time or documents viewed by the particular users when following the selected links. The method further includes ranking links corresponding to a particular carry keyword, using, from the search engine usage information, both (user access rates) and at
30 least one of ((the 12th time) or (the documents viewed)). Then, the ranks are stored in machine readable memory, available in response to search in vote with a particular keyword.

CLAR 1006-3

[00169] Various weightings of user access rates, dwell time and number of pages viewed can be applied. One aspect of the data to which weightings are applied may be that the dwell time is known to have been qualified before receipt to discount periods when a particular user is likely to have been inattentive. Another aspect is that the data may be
5 filtered to remove the outlying observations of search engine usage, eliminating links selected by particular users less than a predetermined number of times. Similarly, dwell times that exceed a predetermined length of time may be discounted after data is received.

[00170] Taking into account user access rates, dwell time and number of pages viewed, one variation has one of the factors dominating the other two. The dominant factor may be
10 the click-through rate, the dwell time or the number of pages viewed at the domain. Another variation gives balanced weighting to the factors, approximately equal plus or -10%. Yet another variation assigns weighting ratios of approximately 2-1-1, plus or -10%, so that one of the three factors is given approximately half the combined weighting.

[00171] Two additional factors that may be considered in ranking our return visits
15 within a predetermined time the link was selected and conversion events. As described above conversion events may be purchases or registrations. Conversions and repeat conversions may separately be given weight.

[00172] For ranking, the combined data from a single user may be given a weight of one or all weight reflecting the amount of time spent searching or the number of links
20 followed. But a single user is given a normalized weight, result biasing by behavior is mitigated by reducing the bots behavior to a single vote.

[00173] The rankings in step of this embodiment may further include groupings like the links into bands and ranking the select links within the bands, as described above. The ranking step also might include separately will ranking links selected by users of particular
25 search engines and then combining rankings from the particular search engines into an overall ranking.

[00174] A further embodiment is a method resolving ambiguity in a keyword search terms submitted by user to a search engine. This method includes receiving, from a behavior watching module operating on the user's computer, information regarding the user's recent
30 navigation behavior. The navigation behavior includes accessing documents that are not all associated with a particular server-side behavior data collection network and further includes

CLAR 1006-3

at least one of a mouse click-through, enter key stroke or other selection action and at least one keyword derived from context of the selection action. The server determines a plurality of interest categories to which the search term plausibly belongs and selects among the plurality of interest categories using the recent navigation information.

5 [00175] One aspect of the behavior watching modules operation is that it may categorize behaviors by interest categories and recency. The documents accessed may be web pages and the server-side behavior data collection network may be an advertising network.

[00176] From this perspective of the behavior watching module, this embodiment may be practiced as a method of supplementing a search request with information that reflects a user's recent navigation behavior among documents accessible via a network. This method includes operating a behavior watching module on the users computer, collecting information regarding the users recent navigation behavior. The recent navigation behavior information includes at least identifying documents electronically accessed by the user; categorizing the documents accessed; and indicating when documents in the categories were accessed. The behavior watching module submits at least a summary of the recent navigation behavior information along with a query to one or more servers including a search engine. The user receives from the search engine a set of results by asked to an interest category corresponding to the submitted recent navigation behavior information. Alternatively, the user may receive from the search engine a set of results by us to a budget category corresponding to the submitted recent navigation behavior information. More generally, any category that can be discerned from the recent navigation behavior can be reflected in a bias of search engine results received in response to the query.

[00177] Another embodiment is a method of personalizing content in response to a keyword search term submitted by a user to a search engine. As with the disambiguation method above, this method includes receiving, from a behavior watching module operating on the user's computer, information regarding the user's recent navigation behavior. The navigation behavior includes accessing documents that are not all associated with a particular server-side behavior data collection network and further includes at least one of a mouse click-through, enter key stroke or other selection action and at least one keyword derived from context of the selection action. The server determines a plurality of likely budget categories

30

CLAR 1006-3

to which results belong and selects among the likely budget categories using the recent navigation information.

[00178] For embodiments using recent navigation behavior, the navigation behavior may reflect recent navigations, including navigations within the last hour and the last day.

5 Alternatively, it may include recent navigations within periods of time determined to correlate to user responsiveness to information. The user navigation information may be summarized into interest categories were likely budget categories when it is received or after it is received. More generally, any category that can be discerned from the recent navigation behavior can be reflected in a bias of search engine results received in response to the query.

10 **[00179]** The embodiments and various aspects of the embodiments described above may be practiced as a machine readable memory including instructions to carry out the methods and aspects of methods described or a data stream including the machine-readable instructions. Further, a device may include one or more servers, personal computers or other computer devices having logic and resources adapted to practice the methods and aspects of
15 methods described.

[00180] We claim as follows:

CLAR 1006-3

CLAIMS

- 1 1. A method of selectively collecting web pages that may be returned to users in
2 response to search requests, the method including:
3 receiving, from behavior watching modules operating on a multiplicity of users'
4 computers, information regarding the users' search engine usage across one or more
5 search engines, including at least
6 keywords (including phrases) submitted by particular users to the search engines,
7 links selected by the particular users from results returned by the search engines, and
8 at least one of dwell time or documents viewed by the particular users when
9 following the selected links;
10 using the search engine usage information, choosing keywords of interest and the links
11 selected by the users corresponding to the keywords of interest;
12 following the chosen links to collect at least part of the content of documents addressed
13 by the chosen links; and
14 associating the collected content with the corresponding keyword of interest.
- 1 2. The method of claim 1, wherein the multiplicity of users exceeds 2,000 users and the
2 received search engine usage information spans a plurality of unaffiliated search engines.
- 1 3. The method of claim 1, wherein the multiplicity of users exceeds 20,000 users.
- 1 4. The method of claim 2, further including:
2 ranking, for the chosen keywords, the corresponding links based on user access rates
3 calculated from the search engine usage information and at least one of the dwell time or
4 the documents viewed.
- 1 5. The method of claim 4, further including limiting the collection of content based on a
2 predetermined application of the ranking.

CLAR 1006-3

1 6. The method of claim 1, wherein the search engine usage information is known to have
2 been qualified before receipt to discount dwell time during periods when a particular user is
3 likely to have been inattentive.

1 7. The method of claim 1, further including selecting snippets from the content of a
2 document reached by following one of the chosen links, the snippet selecting method
3 including:

4 applying a window of predetermined length to the document, repeatedly shifting the
5 window through the document;

6 dividing of the keyword of interest (if it is a phrase) into words;

7 repeatedly calculating one or more window scores for the words in the window, including

8 a count of instances of the words in the window, and

9 a measure of adjacency of the words in the window;

10 choosing as snippets a plurality of non-overlapping window positions using the window
11 scores; and

12 storing the snippets or references to the snippets in a machine readable memory.

1 8. The method of claim 7, wherein the repeated shifts of the window align the window
2 start or the window end with a word boundary.

1 9. The method of claim 7, further including:

2 trading off lengths of the snippets, approximately maintaining a combined length of all
3 snippets;

4 recalculating the one or more window scores to the snippets of changed length;

5 choosing a particular trade off of snippet lengths using the recalculated window scores,
6 optionally merging the snippets; and

7 storing the chosen snippet(s) or references to the chosen snippet(s) in a machine readable
8 memory.

CLAR 1006-3

1 10. A machine readable memory including instructions to carry out the method of
2 claim 1.

1 11. A data stream including machine readable instructions to carry out the method of
2 claim 1.

1 12. A device including one or more servers having logic and resources adapted to
2 practice the method of claim 1.

1 13. A method of generating data using a network of behavior watching modules that
2 transmit data to an affiliated system of servers, the method including:

3 using a multiplicity of behavior watching modules operating on a multiplicity of users'
4 computers, the behavior watching modules collecting information regarding the users'
5 search engine usage across one or more search engines, including at least

6 keywords (including phrases) submitted by particular users to the search engines,

7 links selected by the particular users from results returned by the search engines, and

8 at least one of dwell time or documents viewed by the particular users when

9 following the selected links; and

10 the behavior watching modules electronically reporting the search engine usage
11 information to one or more affiliated servers.

1 14. The method of claim 13, wherein the multiplicity of behavior watching modules
2 exceeds 2,000 users and the behavior watching modules collect search engine usage
3 information across a plurality of unaffiliated search engines.

1 15. The method of claim 13, wherein the multiplicity of behavior watching modules
2 exceeds 20,000 users.

1 16. The method of claim 13, wherein the behavior watching modules qualify the
2 search engine usage information before reporting to discount dwell time during periods when
3 a particular user is likely to have been inattentive.

CLAR 1006-3

1 17. The method of claim 13, further including a particular user invoking a search
2 engine with a keyword and receiving results from the search engine organized to reflect the
3 reported search engine information.

1 18. A machine readable memory including instructions to carry out the method of
2 claim 13.

1 19. A data stream including machine readable instructions to carry out the method of
2 claim 13.

1 20. A device including a personal computer having logic and resources adapted to
2 practice the method of claim 13.

1 21. A computer-implemented system that selectively collects web pages that may be
2 returned to users in response to search requests, the system including:

3 at least one listener module with a network interface that receives from behavior
4 watching modules operating on a multiplicity of users' computers connected by a
5 network, information regarding the users' search engine usage across one or more search
6 engines, including at least

7 keywords (including phrases) submitted by particular users to the search engines,

8 links selected by the particular users from results returned by the search engines, and

9 at least one of dwell time or documents viewed by the particular users when
10 following the selected links; and

11 at least one computer or cluster of computers operating one or more modules in
12 communication with the listener module that includes logic and resources adapted to

13 process the search engine usage information and choose keywords of interest and the
14 links selected by the users corresponding to the keywords of interest,

15 follow the chosen links to collect at least part of the content of documents addressed
16 by the chosen links, and

17 associate the collected content with the corresponding keyword of interest.

CLAR 1006-3

1 22. The system of claim 21, further including at least 2,000 behavior watching
2 modules operating on the users' computers, in communication with the listener module.

1 23. A method of selecting snippets from content of a document that may be returned
2 to users in response to search requests, the snippet selecting including:

3 applying a window of predetermined length to the document, repeatedly shifting the
4 window through the document;

5 dividing the keyword of interest (a phrase) into words;

6 repeatedly calculating one or more window scores for the words in the window, including

7 a count of instances of the words in the window, and

8 a measure of adjacency of the words in the window (when the keyword is a phrase);

9 choosing as snippets a plurality of non-overlapping window positions using the window
10 scores; and

11 storing the chosen snippet(s) or references to the chosen snippets in a machine readable
12 memory.

1 24. The method of claim 23, wherein the repeated shifts of the window align the
2 window start or the window end with a word boundary.

1 25. The method of claim 23, further including evaluating whether two or more
2 documents should be considered exact or near-duplicates based on the titles and the snippets
3 from the documents.

1 26. The method of claim 23, wherein the count of instances includes both counting
2 how many of the words appear in the window and how many times the words appear in the
3 window.

1 27. A machine readable memory including instructions to carry out the method of
2 claim 23.

1 28. A data stream including machine readable instructions to carry out the method of
2 claim 23.

CLAR 1006-3

1 29. A device including one or more servers having logic and resources adapted to
2 practice the method of claim 23.

1 30. A method of ranking, the method including:

2 receiving, from behavior watching modules operating on a multiplicity of users'
3 computers, information regarding the users' search engine usage across a plurality of
4 search engines, including at least

5 keywords (including phrases) submitted by particular users to the search engines,

6 links selected by the particular users from results returned by the search engines, and

7 at least one of dwell time or documents viewed by the particular users when
8 following the selected links;

9 ranking links corresponding to a particular keyword, using, from the search engine usage
10 information, both (user access rates) and at least one of ((the dwell time) or (the
11 documents viewed)); and

12 storing the ranks in machine readable memory, available in response to a search invoked
13 with the particular keyword.

1 31. The method of claim 30, wherein the search engine usage information is known to
2 have been qualified before receipt to discount periods when a particular user is likely to have
3 been inattentive.

1 32. The method of claim 30, the ranking step further including filtering outlying
2 observations of the search engine usage before ranking the links, by eliminating links selected
3 by the particular users that were returned by the search engines less than a predetermined
4 number of times.

1 33. The method of claim 30, the ranking step further including filtering outlying
2 observations of the search engine usage before ranking the links, by discounting dwell times
3 that exceeded a predetermined length of time.

CLAR 1006-3

1 34. The method of claim 30, the ranking step further including weighting the dwell
2 time information to dominate over the user access rates and the documents viewed
3 information.

1 35. The method of claim 30, the ranking step further including weighting the user
2 access rates to dominate over the dwell time and the documents viewed information.

1 36. The method of claim 30, the ranking step further including weighting the
2 documents viewed information to dominate over the user access rates and the dwell time
3 information.

1 37. The method of claim 30, the ranking step further including assigning
2 approximately equal weighting, plus or minus 10 percent, to the user access rates, the dwell
3 time and the documents viewed information.

1 38. The method of claim 30, the ranking step further including assigning weighting
2 ratios of approximately 2-1-1, plus or minus 10 percent, so that one of the user access rates,
3 the dwell time or the documents viewed information is given approximately half of a
4 combined weighting.

1 39. The method of claim 30, the receiving step further including information about
2 return visits to a selected link within a predetermined time, the ranking step further including
3 using the return visits information.

1 40. The method of claim 30, the receiving step further including receiving
2 information regarding conversion of accessing the selected links into purchases or
3 registrations and the ranking step further includes using the conversion information.

1 41. The method of claim 30, whereby result-biasing bot behavior is mitigated by
2 further including in the ranking step, weighting the search engine usage information to give
3 the search engine usage by a particular user a normalized weight approximately equal to other
4 users with more or less search engine usage.

1 42. The method of claim 30, the ranking step further including grouping selected
2 links into bands and ranking the selected links within the bands.

CLAR 1006-3

1 43. The method of claim 30, the ranking step further including separately ranking
2 selected links by the particular search engine that returned the links, then combining the
3 separate rankings into an overall ranking.

1 44. A machine readable memory including instructions to carry out the method of
2 claim 30.

1 45. A data stream including machine readable instructions to carry out the method of
2 claim 30.

1 46. A device including one or more servers having logic and resources adapted to
2 practice the method of claim 30.

1 47. A method of resolving ambiguity in a keyword search term submitted by a user to
2 a search engine, the method including:

3 receiving, from a behavior watching module operating the user's computer, information
4 regarding the user's recent navigation behavior which

5 includes accessing documents that are not all associated with a particular server-side
6 behavior data collection network, and

7 further includes at least one of a mouse click-through, enter keystroke or other
8 selection action and at least one keyword derived from context of the selection
9 action;

10 determining a plurality of interest categories to which the search term plausibly belongs;
11 and

12 selecting among the plurality of interest categories using the recent navigation
13 information.

1 48. The method of claim 47, wherein the recent navigation behavior information is
2 categorized by interest categories and by recency of behavior.

1 49. The method of claim 47, wherein the documents are web pages and the server-
2 side behavior data collection network is an advertising network.

CLAR 1006-3

1 50. A machine readable memory including instructions to carry out the method of
2 claim 47.

1 51. A data stream including machine readable instructions to carry out the method of
2 claim 47.

1 52. A device including one or more servers having logic and resources adapted to
2 practice the method of claim 47.

1 53. A method of supplementing a search request with information that reflects a
2 user's recent navigation behavior among documents accessible via a network, the method
3 including:

4 operating a behavior watching module on the user's computer, collecting information
5 regarding the user's recent navigation behavior that includes

6 identifying documents electronically accessed by the user;

7 categorizing the documents accessed;

8 indicating when documents in the categories were accessed ; and

9 submitting at least a summary of the recent navigation behavior information along with a
10 query to one or more servers including a search engine; and

11 receiving from the search engine a set of results biased to an interest category
12 corresponding to the submitted recent navigation behavior information.

1 54. The method of claim 53, further including the behavior watching module filtering
2 the documents accessed and categorizing the indication of when the categories were accessed.

1 55. A machine readable memory including instructions to carry out the method of
2 claim 53.

1 56. A data stream including machine readable instructions to carry out the method of
2 claim 53.

1 57. A device including a personal computer having logic and resources adapted to
2 practice the method of claim 53.

CLAR 1006-3

1 58. A method of personalizing content returned in response to a keyword search term
2 submitted by a user to a search engine, the method including:

3 receiving, from a behavior watching modules operating the user's computer, information
4 regarding the user's navigation behaviors which includes

5 includes accessing documents that are not all associated with a particular server-side
6 behavior data collection network, and

7 further includes at least one of a mouse click-through, enter keystroke or other
8 selection action and at least one keyword derived from context of the selection
9 action;

10 determining a plurality of likely budget categories to which documents responsive to the
11 search term plausibly would be of interest; and

12 selecting at least one among the plurality of likely budget categories using the navigation
13 information.

1 59. The method of claim 58, wherein the user navigation information reflect recent
2 navigations, including navigations at least within the last hour and last day.

1 60. The method of claim 58, wherein the user navigation information is summarized
2 into one or more of the likely budget categories when received.

1 61. The method of claim 58, wherein the user navigation information is summarized
2 into one or more of the likely budget categories after being received.

1 62. The method of claim 58, further including responding to the user with results
2 biased toward high ranking of at least some results in the selected likely budget category.

1 63. The method of claim 58, further including responding to the user by filtering
2 responsive documents favoring results in the selected likely budget category.

1 64. A machine readable memory including instructions to carry out the method of
2 claim 58.

1 65. A data stream including machine readable instructions to carry out the method of
2 claim 58.

CLAR 1006-3

1 66. A device including one or more servers having logic and resources adapted to
2 practice the method of claim 58.

1 67. A method of supplementing a search requests with information that reflects a
2 user's recent navigation behavior among documents accessible via a network, the method
3 including:

4 operating a behavior watching module on the user's computer, collecting information
5 regarding the user's navigation behavior that includes

6 identifying documents electronically accessed by the user;

7 categorizing the documents accessed;

8 indicating when documents in the categories were accessed ; and

9 submitting at least a summary of the navigation behavior information along with an
10 information request to a search engine; and

11 receiving from the search engine a set of results biased to a likely budget corresponding
12 to the submitted navigation behavior information.

1 68. The method of claim 67, further including the behavior watching module filtering
2 the documents accessed and categorizing the indication of when the categories were accessed.

1 69. A machine readable memory including instructions to carry out the method of
2 claim 67.

1 70. A data stream including machine readable instructions to carry out the method of
2 claim 67.

1 71. A device including one or more servers having logic and resources adapted to
2 practice the method of claim 67.

1/14

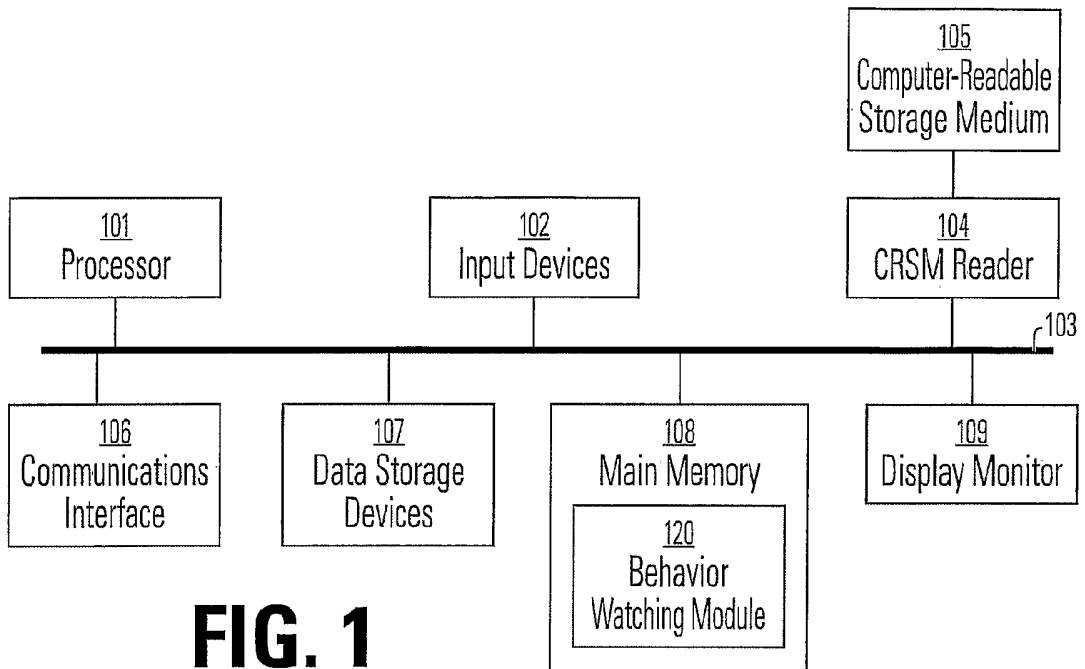


FIG. 1

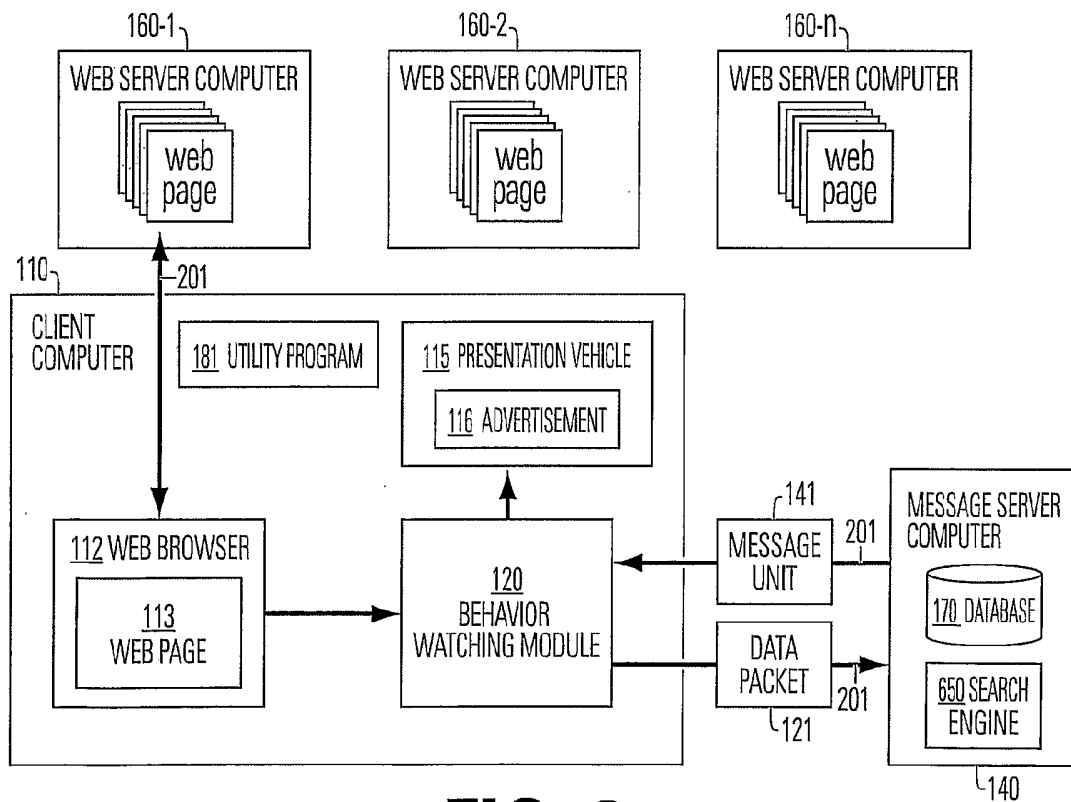


FIG. 2

2/14

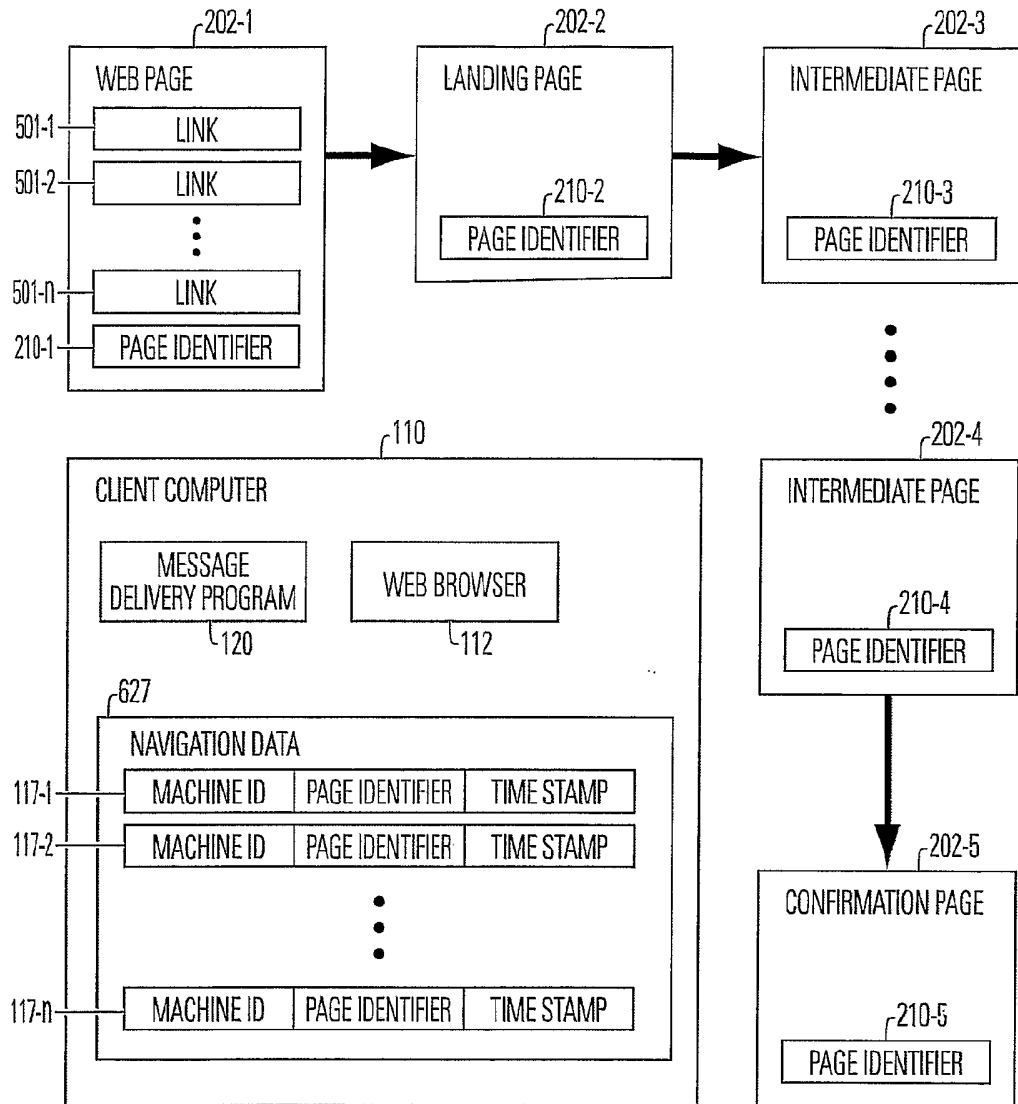


FIG. 3

3/14

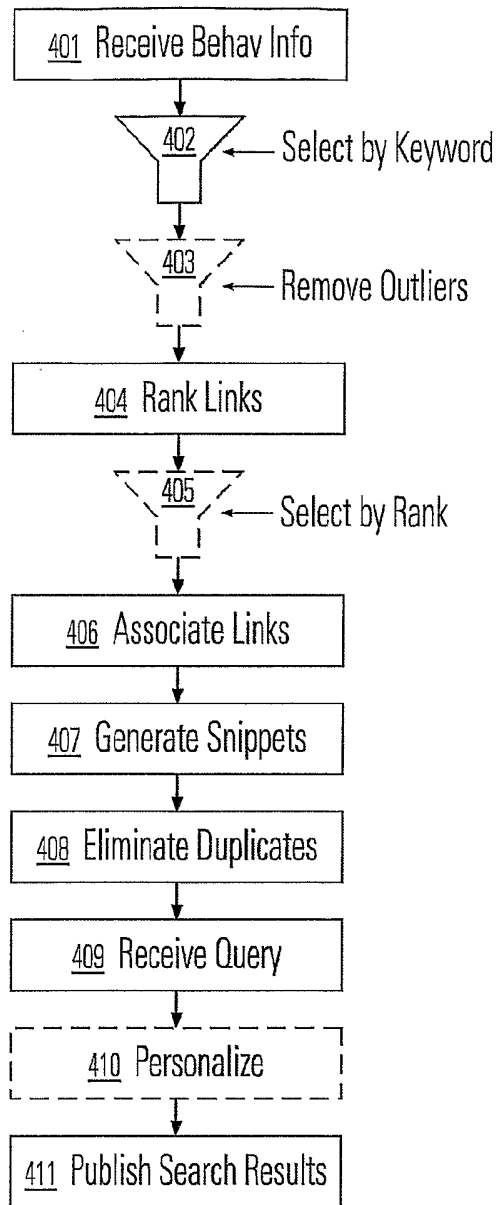


FIG. 4

501	502	503
Keywords	Timestamp	URLs clicked

FIG. 5

4/14

<u>602</u> Version of crumb	<u>604</u> Time Written	<u>606</u> Commercial Flag	<u>608</u> Keyword	<u>610</u> Character Set
3	secs since 1970 hex. eg: 4166FAFC	0 = no, 1 = yes, 2 = unknown (always 1)	2005 Chevrolet Corvette	utf-8

FIG. 6

<u>802</u> CatID	<u>702</u> Day	0	-1	-2	-3
10494	<u>704</u> Hit?	1	0	1	1

FIG. 7

<u>802</u> CategoryID	<u>804</u> Recency	<u>806</u> Frequency	<u>808</u> Banner Clicks	<u>810</u> Comment
10494	3	4	4	
98409	1	6	4	Heavy and recent
65625	14	6	3	Heavy last month
68530	1	0	0	Light but recent
147374	12	0	0	Light last month

FIG. 8

5/14

Bucket #	0	1	2	3	4	5	6	7
Time Slice	N/A	0-5 min	5-15 min	15-30 min	30-60 min	1-2 hrs	2-4 hrs	4-12 hrs
Bucket #	8	9	10	11	12	13	14	15
Time Slice	12-24 hrs	1-3 days	3-7 days	7-14 days	14-21 days	21-30 days	30-45 days	45-60 days

FIG. 9

Bucket #	0	1	2	3	4	5	6	7
# Days w/ a visit	0	1	2	3	4 or 5	6 to 10	11 to 30	31 to 60

FIG. 10

<u>1102</u> CategoryID	<u>1104</u> Category Leaf Name	<u>1106</u> Revenue Priority
9870	Uncategorized	0
10494	Luxury Auto	1
98409	SUVs	2
95625	Payroll Software	3
68530	Florida except Orlando	4

FIG. 11

6/14

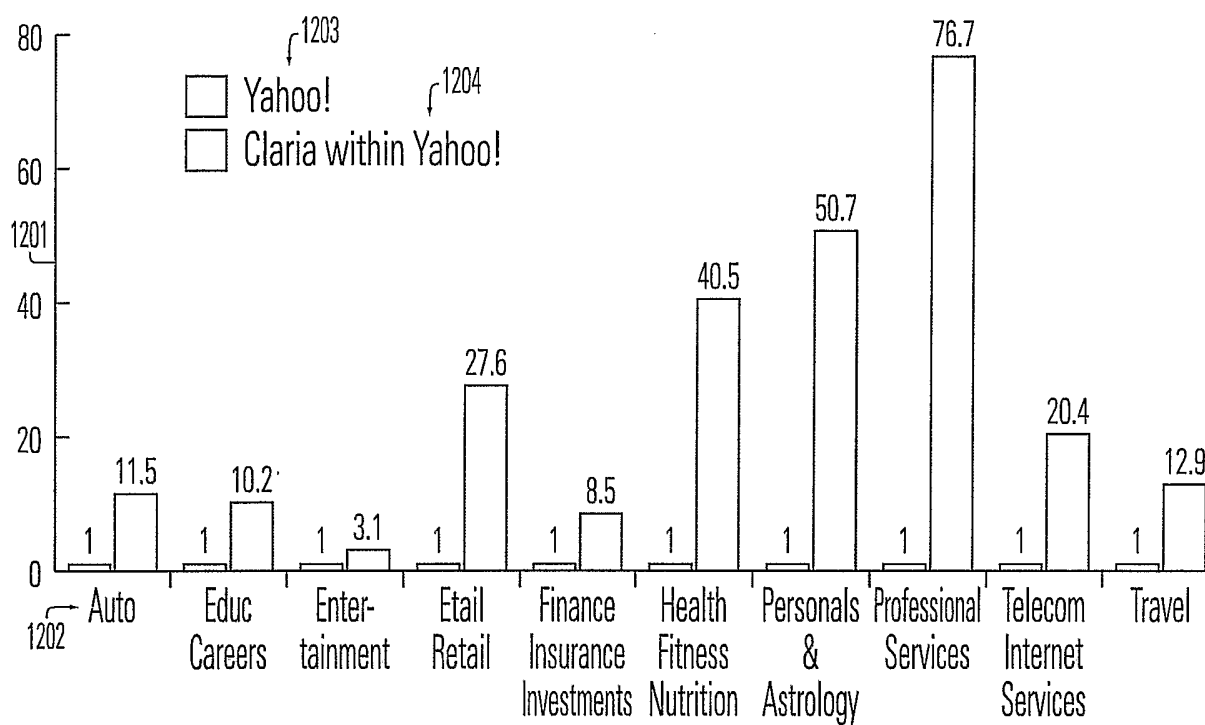


FIG. 12

CarsDirect.com: Research new cars, finance and buy new and used cars or sell - Internet Explorer

File Edit View Favorites Tools Help

Address <http://www.carsdirect.com/home>

1.800.431.2500 | Live Help | Need Help?

Home Research New Cars Used Cars Trade-Ins Sell Finance My CarsDirect

CarsDirect®

America's #1 way to buy cars online!

NEW Save Big with Cash Back!
get a price

more about

USED Find Your Car
500,000+ listings, updated daily

Car Compare
Compare ANY vehicles you choose.

Certified Pre-Owned
Sponsored by Lexus CPO

Rated #1 top internet site by *Forbes*, *Time*, *PC Magazine*, and *Yahoo! Internet Life*, among others... See our awards

Search by Category

- Products & Services
- Great Financing Rates
- Extended Warranties
- Anti-theft products
- LoJack Protection
- Auto Accessories
- Free Credit Report
- Free VIN Check
- Auto Seller's Kit
- Auto Repair Directory

GO →

- Employee Discounts Offered to the Public
- Our Low Price Guarantee

ZIP: GO →

Learn about CPO

Search Listings

ZIP: GO →

Internet

Factors Used To Determine Relevance Score

Initial Visit Metrics—

Click Rate: 12% (high)

Pages Viewed: 5

Dwell Time: 55 sec.

Purchase/Registration: 80% yes

Subsequent Visit Metrics—

Pages Viewed: 15

Dwell Time: 300 sec.

Purchase/Registration: 15% yes

User Level Metrics—

Category Involvement Level: Expert

Search Engine Usage Level: Maven

Geography Relevance Score: High

Relevance Score: 97

FIG. 13

8/14

USATODAY.com - Some hybrids not as reliable as gas-powered models - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Media Print Mail News Links

Address http://www.usatoday.com/money/autos/front.htm Go Links

Advertisement

Get Happy
Get in, get a great deal and get happy.
Drivers wanted!

USA TODAY Marketplace Cars Jobs Franchises Business Opportunities Travel Real Estate Tickets More

Home News Travel Money Sports Life Tech Weather Search

powered by Google

Inside Cars
Cars home
Research a car
Buy a car
Sell a car
Incentives
Recalls
Car Talk
Travel »
Jobs »

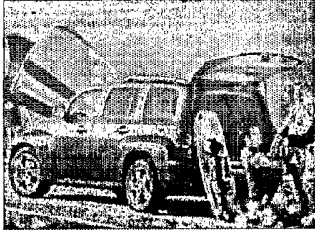
Cars

E-MAIL THIS PRINT THIS SAVE THIS MOST POPULAR SUBSCRIBE

Posted 8/16/2005 8:59 PM Updated 8/16/2005 9:33 PM

Some hybrids not as reliable as Gas-powered models

When Chevrolet developed its 2005 Cobalt small car to replace the aging Cavalier, it also worked up an entirely different vehicle using similar underpinnings. That's the 2006 HHR.



The new HHR is inspired by the 1949 Chevrolet Suburban

HHR stands for heritage high-roof. Heritage refers to the 1949 Chevy Suburban utility vehicle that inspired the front end and bulging fenders of the HHR. High-roof refers to the tall roofline, meant to suggest a roomy vehicle despite its compact outside dimensions.

Think back to the Cavalier. GM used the same underpinnings for another vehicle then, too. It was the Pontiac Sunfire, so similar to Cavalier that the two swiped buyers from each other.

This time the two models that share hardware look different, are different sizes.

RELATED CHARTS

- Microsoft to offer two Xbox
- Archaeologists find treasure
- Wine and beer super-substitute

Advertisement

Get Happy
Get in, get a great deal and get happy.
Drivers wanted!

Marketplace Provided by cars.com

RESEARCH A CAR
Select a make [GO]

BUY A CAR
Select a make [GO]

Want to sell your car?
Click here to create your ad.

Today's Top Money Stories

- Closing arguments begin in Vloxx trial -- 11:05 AM
- AOL worker who stole e-mail list sentenced - 5:07 PM

Internet

FIG. 14

9/14

YAHOO! search	msn Search	Google	Newton
<p>1 Cars.com features access to automol and reports performance da www.cars.com/</p> <p>2 CarsDirect.com research, build and buy a www.carsdirect.com/</p> <p>3 AutoTrader features classified listings and vans including car ca www.autotrader.com/</p> <p>4 Kelley Blue Book official site provides value including invoice and MS www.kbb.com/</p> <p>5 Carsforsale.com offers used cars for sale by www.carsforsale.com/</p> <p>6 Yahoo! Autos Research new car pricing, online, check out used car autos.yahoo.com/</p> <p>7 ClassicCars.com includes chats and forums directory. www.classiccars.com</p> <p>8 Edmunds.com new and used online buyin comparisons. www.edmunds.com/</p> <p>9 Buy Classic Cars features directory of deale Includes links for value g www.buyclassiccars.com/</p> <p>10 Car Talk the official site of Click ur and Ray Magliozzi cartalk.cars.com/</p>	<p>1 Cars.com homepage cars.com is the leading p new and used car pricing</p> <p>2 New car prices, used ...Most Fuel-Efficient No Fuel-Efficient Cars. Vehi www.edmunds.com/</p> <p>3 Used Cars, Sell Your AutoTrader.com - Your C listings and information d inventory.autotrader.com</p> <p>4 CarsDirect.com CarsDirect.com is Ameri Compare, Pricing, and Fi www.carsdirect.com</p> <p>5 SUPERCARS.NET Super Cars. An extensive automobiles ... With over www.supercars.net</p> <p>6 New Cars Used Car cars, used cars: prices and next car or truck online w</p> <p>7 globe-gawheels.com Sale Auto Industry Con Search By Price, Search: www.globe-gawheels.c</p> <p>8 CARS - Vehicle and National Highway Traffi Equipment Information www.nhtsa.dot.gov/ea</p> <p>9 Used Cars for sale in New Cars. Used Cars. En Your Car. More than Car</p> <p>10 KCIC - Cars Ordering Tips. If ordering moving on. For more hel www.pueblo.gsa.gov/o</p>	<p>1 Cars.com homepage: Listings and research cars.com is the leading provider of automotive clasifie new and used car pricings, reviews, and ownership...</p> <p>2 Jaguar Cars Worldwide. www.jaguar.com/</p> <p>3 CarsDirect.com: Research new cars, finance It's easy to research new or used cars or sell your own CarsDirect.com ... CarsDirect.com is the... www.carsdirect.com/home/</p> <p>4 New car prices, used car pricing, auto reviews Find cars for sale in your neighborhood. Select a Make car in front of ... www.edmunds.com/</p> <p>5 SUPERCARS.NET Super Cars. An extensive Motors. With over 32 cars on display, this has to be the show ever assembled for the enginee. JUL ... www.snp</p> <p>6 Lotus Cars USA Official site. Includes dealership listings, photographs, model profile information. www.lotuscars.com/</p> <p>7 Travelocity: Airline Tickets, Hotels, Cars, Vac Find the best prices for airline tickets on the Internet w finder, make reservations for hotels, cruises and rental o</p> <p>8 Welcome to Volvo Cars Choose Your Country ... www.volvocars.com/</p>	<p>1 Cars.com features access to automotive classifieds and dealers, price info, news and reports, performance data... www.cars.com/</p> <p>2 CarsDirect.com research, build and buy a new or used car online. www.carsdirect.com/</p> <p>3 AutoTrader features classified listings on used cars, trucks, and vans including car care, warranties, loans, and more. www.autotrader.com/</p> <p>4 Kelley Blue Book official site provides values on new and used cars, trucks, and vans, including invoice and MSRP. www.kbb.com/</p> <p>5 New car prices, used car pricing, auto reviews by Edmunds ...Most Fuel-Efficient Non-Hybrid Cars. Best-Selling Vehicles. Most Fuel-Efficient Cars. Vehicles for Shorter Drivers ... www.edmunds.com/</p> <p>6 New Cars Used Cars for Sale Autobytel New and New cars, used cars: prices and information and research on how to buy your next car or truck online with our vehicle pricing www.autobytel.com</p> <p>7 SUPERCARS.NET Super Cars. An extensive Motors. With over 32 cars on display, this has to be the most complete show ever assembled for the enginee. JUL ... www.supercars.net/</p> <p>8 Car Place weekly review of vehicles test-driven by auto editor Robert Bowden, plus summaries of prior reviews... www.thecarplace.com</p> <p>9 Autoshows.com provides searchable databasc of auto shows, as well as links to product reviews, dealers, and www.autoshows.com/</p> <p>10 Los Angeles Times cars.com Los Angeles Times cars.com is the leading provider of automotive classified listings and new and used car pricing, www.latimes.com/classified/cars</p>

Search for "cars"
Total Page Relevancy Index:

Yahoo	487
Google	524
MSN	598
Newton	864

FIG. 15

10/14

		1602 ↓		
	←	Loyal	Switcher	Competitive Loyal
Heavy	1601 →	Uniques 333,253 CTR 7.90% Conv. Rate 8.00%	Uniques 500,332 CTR 6.20% Conv. Rate 4.10%	Uniques 1,990,098 CTR 3.82% Conv. Rate 2.15%
Medium		Uniques 467,000 CTR 5.40% Conv. Rate 7.20%	Uniques 898,922 CTR 5.50% Conv. Rate 3.15%	Uniques 1,210,000 CTR 2.20% Conv. Rate 1.05%
Light		Uniques 211,009 CTR 3.00% Conv. Rate 6.50%	Uniques 524,422 CTR 3.10% Conv. Rate 2.80%	Uniques 1,299,878 CTR 1.40% Conv. Rate .55%

FIG. 16

11/14

Advertiser	Short Domain	Traffic Source	Keyword	Impressions	Clicks	Click Through Rate (CTR%)	Conversion Rate	AVG Rank in Listing
orbitz.com		google.com	acapulco	258	0	0.00%	0.00%	2.9
			aer lingus	249	3	1.20%	0.00%	2.4
			aeromexico	307	3	0.98%	0.00%	2.5
			air canada	2	0	0.00%	0.00%	21.5
			air china	211	3	1.42%	0.00%	1.8
			air fare	565	17	3.01%	0.00%	3.9
			air fares	731	44	6.02%	6.82%	3.3
			air ticket	227	19	8.37%	0.00%	1.6
			air tickets	353	31	8.78%	0.00%	1.5
			air tran	16	0	0.00%	0.00%	3.1
			air travel	307	29	9.45%	3.45%	1.9
			airfare	2,268	187	8.25%	3.74%	1.8
			airfares	967	105	10.86%	1.90%	1.8
			airline	726	40	5.51%	2.50%	3.7
			airline fares	493	33	6.69%	6.06%	2.2
			airline flights	549	65	11.84%	15.38%	1.6
			airline ticket	603	61	10.12%	0.00%	1.9
			airline tickets	10,717	800	7.46%	5.75%	2.0
			airlines	4,318	125	2.89%	3.20%	1.5
			airplane tickets	405	11	2.72%	0.00%	2.9

FIG. 17

12/14

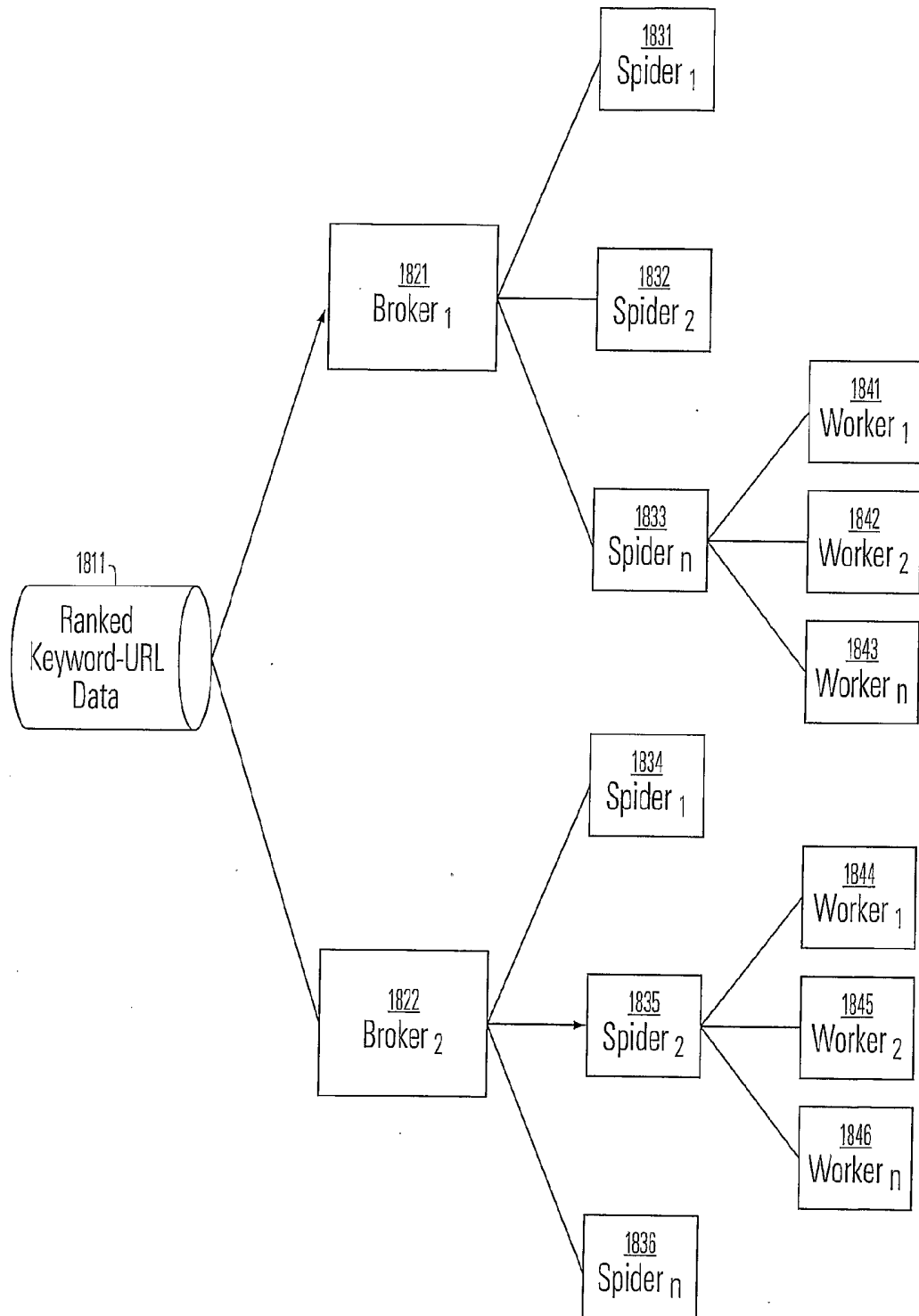


FIG. 18

13/14

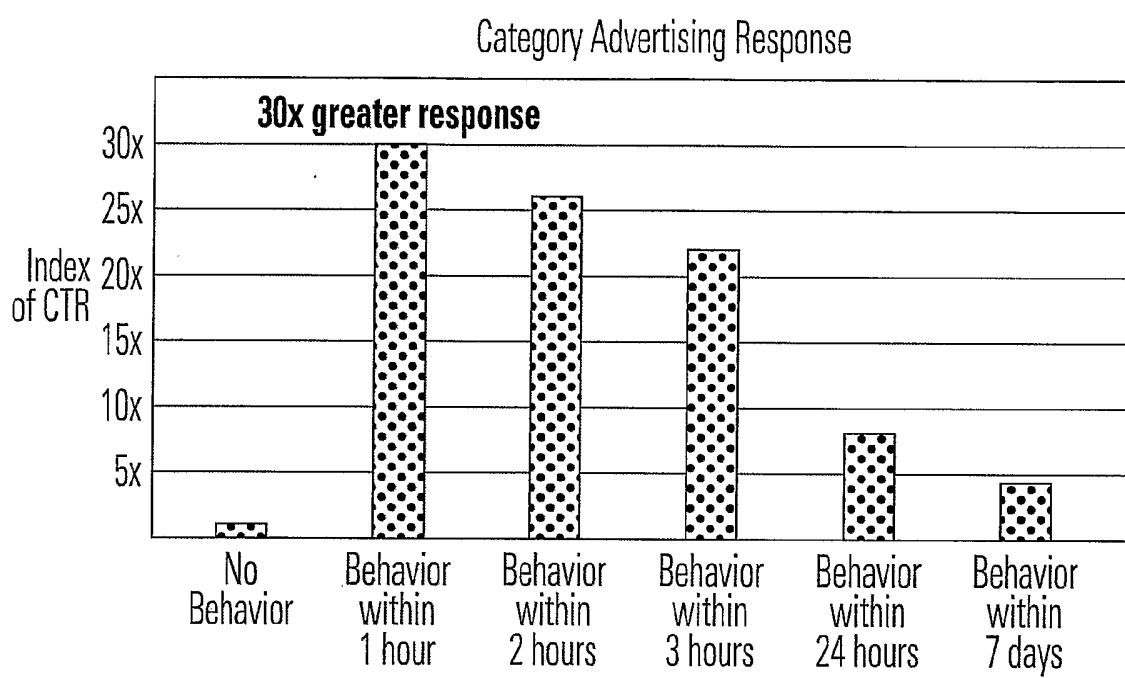


FIG. 19

14/14

2001

 Lorem ipsum nonsense —2002

2003: Lorem ipsum dolor sit amet, consectetur
adipiscing elit, sed do eiusmod tempor
incididunt ut labore et dolore magna
aliqua. Ut enim ad minim veniam, quis
nostrud exercitation ullamco laboris nisi
ut aliquip ex ea commodo consequat.

2004: Duis aute irure dolor in reprehenderit in
voluptate velit esse cillum dolore eu
fugiat nulla pariatur. Excepteur sint
occaecat cupidatat non proident, sunt in
culpa qui officia deserunt mollit anim id
est laborum.

2005: Fusce neque wisi, ornare sed, aliquet sit
amet, ultricies a, est. Ut rhoncus eros
nec mauris. Phasellus vitae felis nec elit
tincidunt auctor. Aliquam congue, metus
ultrices tincidunt blandit, elit pede
dapibus turpis, at accumsan urna diam
sed purus. Pellentesque imperdiet. Nulla
mattis erat eget lorem. Maecenas at
magna. Sed in sem. Suspendisse at orci
a justo ultricies iaculis. Etiam dictum.
Quisque pellentesque hendrerit orci.
Etiam nibh orci, fermentum eget,
placerat sit amet, faucibus id, elit. Nulla
sollicitudin congue leo.

Duis justo sem, sodales nec,
scelerisque id, vestibulum nec, diam.
Phasellus enim. Duis sodales tincidunt
enim. Sed metus. Nulla facilisi. Ut
facilisis risus ut purus. Suspendisse
magna. Cras nec purus in sem posuere
elementum. Aenean ac arcu ut augue
condimentum euismod. Fusce ut metus.
Fusce vel ante varius neque rhoncus
placerat. Ut quis quam vitae purus
facilisis vestibulum.

FIG. 20