

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
24 June 2004 (24.06.2004)

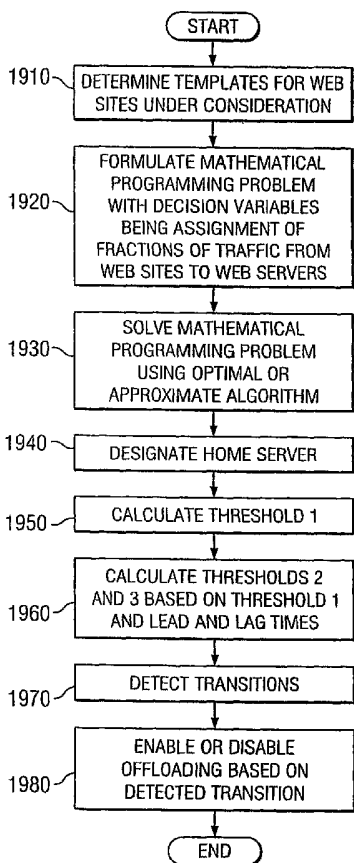
PCT

(10) International Publication Number
WO 2004/054197 A1

- (51) International Patent Classification⁷: H04L 29/06, 12/24
- (71) Applicant (for MC only): COMPAGNIE IBM FRANCE [FR/FR]; Tour Descartes, 2, Avenue Gambetta, F-92400 Courbevoie (FR).
- (21) International Application Number: PCT/EP2003/015016
- (72) Inventors: DIAS, Daniel, Manuel; 3380 Sunny Court, Mohegan Lake, NY 10547 (US). LIU, Zhen; 37 Round-abend Road, Tarrytown, NY 10591 (US). SQUILLANTE, Mark, Steven; 21 Scofield Road, Pound Ridge, NY 10576 (US). XIA, Honghui; 26-L Scenic Drive, Croton on Hudson, NY 10520 (US). YU, Shung-Zheng; Building 635-306, Zhongshan University, Guangzhou, Guangdong 510275 (CN). ZHANG, Li; 1387 Hanover St, Yortown Heights, NY 10598 (US). KING, Richard, Pervin; 117 Old Army Road, Scarsdale, NY 10583 (US).
- (22) International Filing Date: 14 November 2003 (14.11.2003)
- (25) Filing Language: English
- (26) Publication Language: English
- (30) Priority Data: 10/315,335 10 December 2002 (10.12.2002) US
- (74) Agent: DE PENA, Alain; Compagnie IBM France, Direction de la Propriété Intellectuelle, F-06610 La Gaude (FR).
- (71) Applicant (for all designated States except MC): INTERNATIONAL BUSINESS MACHINES CORPORATION [US/US]; New Orchard Road, Armonk, NJ 10504 (US).
- (81) Designated States (national): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH,

[Continued on next page]

(54) Title: APPARATUS AND METHODS FOR CO-LOCATION AND OFFLOADING OF WEB SITE TRAFFIC BASED ON TRAFFIC PATTERN RECOGNITION



(57) Abstract: Apparatus and methods for identifying traffic patterns to web sites based on templates that characterize the arrival of traffic to the web sites are provided. Based on these templates, determinations are made as to which web sites should be co-located so as to optimize resource allocation. Specifically, web sites whose templates are complementary, i.e. a first web site having a peak in arrival traffic at time t1 and a second web site that has a trough in arrival traffic at time t1, are designated as being candidates for co-location. In addition, the present invention uses the templates identified for the traffic patterns of web sites to determine thresholds for offloading traffic to other servers. These thresholds include a first threshold at which offloading should be performed, a second threshold that takes into consideration the lead time needed to begin offloading, and a third threshold that takes into consideration a lag time needed to stop all offloading of traffic to the other servers.

WO 2004/054197 A1



GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, OM, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, UZ, VC, VN, YU, ZA, ZM, ZW.

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

Published:

— with international search report

(84) Designated States (regional): ARIPO patent (BW, GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE,

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.

**APPARATUS AND METHODS FOR CO-LOCATION AND OFFLOADING OF WEB
SITE TRAFFIC BASED ON TRAFFIC PATTERN RECOGNITION**

BACKGROUND OF THE INVENTION

Technical Field:

5 The present invention is directed to apparatus and
methods for co-location and off-loading of web site traffic
based on traffic pattern recognition. More specifically, the
present invention is directed to apparatus and methods for
profiling web sites, determining which web sites should be
10 co-located and when offloading of web traffic to other servers
should be performed.

Description of Related Art:

15 With the increasing popularity of the Internet and its
applications, the allocation of resources in order to provide
a Quality of Service (QoS) has become more difficult. The
complexity of the problem of resource allocation is
exacerbated by the heterogeneity of the Internet
infrastructure and applications as well as the user behaviors.

20 One approach to the resource allocation problem is to
share resources, such as web servers, among multiple entities,
such as web sites, so that peak load conditions for any given
entity can be handled by borrowing resources from other
entities. Such a mechanism is often referred to as
25 co-location. An important problem in this paradigm is
concerned with the clustering of the web sites for the
resource sharing. That is, it is often difficult and
imprecise manual process to identify which web sites should
share resources.

30 Another approach, which is complementary to co-location,
consists of offloading work for an entity to exogenous

resources in a dynamic way. With offloading, web site traffic that is destined for a particular web server that is currently in an overloaded state is redirected to another server to handle the processing of the traffic. A key problem with
5 offloading is determining where and when to offload the work. The known mechanisms for determining where and when to offload work typically fall into the area of load balancing where current state information is used to determine if the work
10 load should be balanced by sending some of the work to other servers. Such work load balancing mechanisms are reactionary and do not make use of known patterns of traffic to begin offloading prior to the servers becoming overloaded.

Thus, it would be beneficial to have an improved apparatus and method for determining which web sites should be
15 co-located and when traffic to web sites should be offloaded to other web servers.

SUMMARY OF THE INVENTION

The present invention provides apparatus and methods for identifying traffic patterns to web sites based on templates that characterize the arrival of traffic to the web sites.

5 Based on these templates, determinations are made as to which web sites should be co-located so as to optimize resource allocation. Specifically, web sites whose templates are complimentary, i.e. a first web site having a peak in arrival traffic at time t_1 and a second web site that has a trough in
10 arrival traffic at time t_1 , are designated as being candidates for co-location.

In addition, the present invention uses the templates identified for the traffic patterns of web sites to determine thresholds for offloading traffic to other servers. These
15 thresholds include a first threshold at which offloading should be performed, a second threshold that takes into consideration the lead time needed to begin offloading, and a third threshold that takes into consideration a lag time needed to stop all offloading of traffic to the other servers.

20

These and other features and advantages of the present invention will be described in, or will become apparent to those of ordinary skill in the art in view of, the following detailed description of the preferred embodiments.

BRIEF DESCRIPTION OF THE DRAWINGS

Figure 1 is an exemplary diagram of a distributed data processing system in which the present invention may be implemented;

5 **Figure 2** is an exemplary block diagram of a server computing device according to the present invention;

Figure 3 is an exemplary block diagram of a client computing device according to the present invention;

10 **Figures 4A-4C** are exemplary diagrams of hourly hits over a period of one week for exemplary commercial web sites;

Figures 5A-5C are exemplary diagrams of daily time-series plots of the exemplary commercial web sites in **Figures 4A-4C**;

15 **Figure 6** is an exemplary diagram of daily coefficient of variation as a function of the daily average for different measures encountered at exemplary commercial web sites;

Figure 7 is an exemplary diagram of daily peak to mean ratio as a function of the daily average for different measures encountered at exemplary commercial web sites;

20 **Figure 8** is an exemplary diagram illustrating requests per hour over one day collected from exemplary commercial web sites for use in clustering the web sites;

Figure 9 is an exemplary diagram illustrating the patterns for four classes of request patterns into which the web sites of **Figure 8** are clustered;

25 **Figure 10A** is an exemplary diagram illustrating the four templates for the four classes of request patterns of **Figure 9**;

30 **Figure 10B** is an exemplary diagram illustrating the day of the week patterns generated by applying the clustering and profiling of the present invention to the empirical data for the web sites of **Figures 4A-4C**;

Figure 11 is an exemplary block diagram of a web site

classification device according to the present invention;

Figure 12 is a flowchart outlining an exemplary operation of the present invention;

Figure 13 is a flowchart outlining an exemplary operation
5 of the present invention for clustering web sites;

Figure 14 is a flowchart outlining an exemplary operation of the present invention for clustering web sites;

Figure 15 is a flowchart outlining an exemplary operation of the present invention for classifying web sites;

10 **Figure 16** is a flowchart outlining an exemplary operation of the present invention for determining which web sites are candidates for co-location;

Figure 17 is a diagram illustrating a template for a web site with thresholds for offloading according to the present
15 invention illustrated;

Figure 18 is a flowchart outlining an exemplary operation of the present invention for offloading traffic according to the present invention;

Figure 19 is a flowchart outlining an exemplary operation
20 of the present invention for a combination of both co-location and offloading; and

Figure 20 is a block diagram of a resource allocation determination system.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

The preferred embodiments of the present invention are implemented in a distributed data processing environment in which traffic data is compiled and used to profile, cluster and categorize web sites. Since the present invention is implemented in a distributed data processing environment, a brief description of this environment will first be provided in order to provide a context in which the present invention operates.

With reference now to the figures, **Figure 1** depicts a pictorial representation of a network of data processing systems in which the present invention may be implemented. Network data processing system **100** is a network of computers in which the present invention may be implemented. Network data processing system **100** contains a network **102**, which is the medium used to provide communications links between various devices and computers connected together within network data processing system **100**. Network **102** may include connections, such as wire, wireless communication links, or fiber optic cables.

In the depicted example, server **104** is connected to network **102** along with storage unit **106**. In addition, clients **108**, **110**, and **112** are connected to network **102**. These clients **108**, **110**, and **112** may be, for example, personal computers or network computers. In the depicted example, server **104** provides data, such as boot files, operating system images, and applications to clients **108-112**. Clients **108**, **110**, and **112** are clients to server **104**. Network data processing system **100** may include additional servers, clients, and other devices not shown.

In the depicted example, network data processing system **100** is the Internet with network **102** representing a worldwide collection of networks and gateways that use the Transmission

Control Protocol/Internet Protocol (TCP/IP) suite of protocols to communicate with one another. At the heart of the Internet is a backbone of high-speed data communication lines between major nodes or host computers, consisting of thousands of commercial, government, educational and other computer systems that route data and messages. Of course, network data processing system **100** also may be implemented as a number of different types of networks, such as for example, an intranet, a local area network (LAN), or a wide area network (WAN).

Figure 1 is intended as an example, and not as an architectural limitation for the present invention.

Referring to **Figure 2**, a block diagram of a data processing system that may be implemented as a server, such as server **104** in **Figure 1**, is depicted in accordance with a preferred embodiment of the present invention. Data processing system **200** may be a symmetric multiprocessor (SMP) system including a plurality of processors **202** and **204** connected to system bus **206**. Alternatively, a single processor system may be employed. Also connected to system bus **206** is memory controller/cache **208**, which provides an interface to local memory **209**. I/O bus bridge **210** is connected to system bus **206** and provides an interface to I/O bus **212**. Memory controller/cache **208** and I/O bus bridge **210** may be integrated as depicted.

Peripheral component interconnect (PCI) bus bridge **214** connected to I/O bus **212** provides an interface to PCI local bus **216**. A number of modems may be connected to PCI local bus **216**. Typical PCI bus implementations will support four PCI expansion slots or add-in connectors. Communications links to clients **108-112** in **Figure 1** may be provided through modem **218** and network adapter **220** connected to PCI local bus **216** through add-in boards.

Additional PCI bus bridges **222** and **224** provide interfaces

for additional PCI local buses **226** and **228**, from which additional modems or network adapters may be supported. In this manner, data processing system **200** allows connections to multiple network computers. A memory-mapped graphics adapter **230** and hard disk **232** may also be connected to I/O bus **212** as depicted, either directly or indirectly.

Those of ordinary skill in the art will appreciate that the hardware depicted in **Figure 2** may vary. For example, other peripheral devices, such as optical disk drives and the like, also may be used in addition to or in place of the hardware depicted. The depicted example is not meant to imply architectural limitations with respect to the present invention.

The data processing system depicted in **Figure 2** may be, for example, an IBM eServer pSeries system, a product of International Business Machines Corporation in Armonk, New York, running the Advanced Interactive Executive (AIX) operating system or LINUX operating system.

With reference now to **Figure 3**, a block diagram illustrating a data processing system is depicted in which the present invention may be implemented. Data processing system **300** is an example of a client computer. Data processing system **300** employs a peripheral component interconnect (PCI) local bus architecture. Although the depicted example employs a PCI bus, other bus architectures such as Accelerated Graphics Port (AGP) and Industry Standard Architecture (ISA) may be used.

Processor **302** and main memory **304** are connected to PCI local bus **306** through PCI bridge **308**. PCI bridge **308** also may include an integrated memory controller and cache memory for processor **302**. Additional connections to PCI local bus **306** may be made through direct component interconnection or through add-in boards. In the depicted example, local area network (LAN) adapter **310**, SCSI host bus adapter **312**, and expansion bus interface **314** are connected to PCI local bus **306** by direct

component connection. In contrast, audio adapter **316**, graphics adapter **318**, and audio/video adapter **319** are connected to PCI local bus **306** by add-in boards inserted into expansion slots. Expansion bus interface **314** provides a connection for a
5 keyboard and mouse adapter **320**, modem **322**, and additional memory **324**. Small computer system interface (SCSI) host bus adapter **312** provides a connection for hard disk drive **326**, tape drive **328**, and CD-ROM drive **330**. Typical PCI local bus implementations will support three or four PCI expansion slots
10 or add-in connectors.

An operating system runs on processor **302** and is used to coordinate and provide control of various components within data processing system **300** in **Figure 3**. The operating system may be a commercially available operating system, such as
15 Windows XP, which is available from Microsoft Corporation. An object oriented programming system such as Java may run in conjunction with the operating system and provide calls to the operating system from Java programs or applications executing on data processing system **300**. "Java" is a trademark of Sun
20 Microsystems, Inc. Instructions for the operating system, the object-oriented operating system, and applications or programs are located on storage devices, such as hard disk drive **326**, and may be loaded into main memory **304** for execution by processor **302**.

25 Those of ordinary skill in the art will appreciate that the hardware in **Figure 3** may vary depending on the implementation. Other internal hardware or peripheral devices, such as flash read-only memory (ROM), equivalent nonvolatile memory, or optical disk drives and the like, may
30 be used in addition to or in place of the hardware depicted in **Figure 3**. Also, the processes of the present invention may be applied to a multiprocessor data processing system.

As another example, data processing system **300** may be a stand-alone system configured to be bootable without relying

on some type of network communication interfaces. As a further example, data processing system **300** may be a personal digital assistant (PDA) device, which is configured with ROM and/or flash ROM in order to provide non-volatile memory for storing operating system files and/or user-generated data.

The depicted example in **Figure 3** and above-described examples are not meant to imply architectural limitations. For example, data processing system **300** also may be a notebook computer or hand held computer in addition to taking the form of a PDA. Data processing system **300** also may be a kiosk or a Web appliance.

As mentioned previously, the present invention provides a mechanism for categorizing web sites. Such categorization is a support functionality for use in workload characterization, performance modeling, workload and performance forecasting, capacity planning, and the like. Basically, each of these various functions are directed to optimizing resource utilization and making sure that there are enough resources available to handle the traffic experienced by the web site in order to give a determined level of service to client devices. Web site categorization according to the present invention may be a principle support function for ensuring accurate modeling of the web site for use in these other functions.

A fundamental part of the present invention is the discovery that web sites have repeated patterns of traffic characteristics that may be exploited to help solve the problems of workload characterization, performance modeling, workload and performance forecasting, and capacity planning. These patterns may exist for various different measures including number of hits, bytes, page views, visits, hits per visit, page views per visit, seconds per page view, seconds per visit, and the like. Moreover, these patterns may exist for various time scales including monthly, weekly, daily, hourly, and the like. To illustrate the repeated patterns of

traffic characteristics, the number of hits per hour over a week time interval for a plurality of exemplary commercial web sites will be considered.

Figures 4A-4C are exemplary diagrams of hourly hits over
5 a period of one week for exemplary commercial web sites. In the plots shown in **Figures 4A-4C**, the measure of number of hits per hour is used to characterize the incoming request patterns from client devices over each day of the week.

As can be seen from **Figures 4A-4C**, a large number of web
10 sites exhibit daily access patterns for which there is a considerable drop in the request rate (both hits and pages) over Saturdays and Sundays relative to the other days of the week. Some of the web sites across different industries often have clear weekend patterns, as illustrated in the first plot
15 **410**, the third plot **420**, the ninth plot **430** and the eleventh plot **440** in **Figures 4A** and **4B**, where there is a significant drop in the request rate over the weekend. Other web sites, such as **450** have weaker yet still prevalent weekend patterns with less significant drops in the request rate over the
20 weekend. Still other web sites, such as **460**, do not exhibit any weekend patterns.

A similar pattern is prevalent in traffic data for different time scales, e.g., daily, weekly, monthly, and for other measures of request patterns, e.g., number of bytes,
25 page views, visits, hits per visit, etc. **Figures 5A-5C** are exemplary diagrams of daily time-series plots of some of the exemplary commercial web sites in **Figures 4A-4C** for various measures. From **Figures 5A-5C** it can be seen that weekend effects of varying degrees for the web sites are present in
30 these other measures. Certain web sites with weekend patterns also exhibit very consistent week-to-week behavior in which request measures do not change much from one week to the next.

In short, from the empirical data shown in **Figures 4A-4C**

and **5A-5C**, it is clear that many web sites experience patterns in their traffic. In addition to the above empirical characteristics, various statistical measures can be used to identify and examine some of the most complex characteristics of the user request patterns, at different time scales, in traffic for web sites. Such characteristics include, for example, the traffic variability and the peak to mean ratio, which are less visible and can only be computed from the time-series data.

10 Some useful aspects of the variability of a request time series can be examined by the coefficient of variation (CV) of the overall request volume, i.e. the ratio of the standard deviation to the mean. For example, **Figure 6** is an exemplary diagram of the daily coefficient of variation of different
15 measures as a function of the daily average for these measures that are encountered at exemplary commercial web sites. From **Figure 6** it can be seen that under some measures (e.g., daily visits plot **610**), different sites exhibit quite different variabilities but similar daily averages while under some
20 other measures (e.g., daily pageview per visit **620**), different sites exhibit similar variabilities but different daily averages.

One rough measure of burstiness in the request patterns is the ratio of the peak to mean request rate over a certain
25 interval of time. This measure quantifies the peak request volume relative to the average request volume. **Figure 7** is an exemplary diagram of the daily peak to mean ratio for different measures as a function of the daily average for these measures that are encountered at exemplary commercial
30 web sites. From **Figure 7** we observe similar behavior as in **Figure 6**, namely different statistical properties of various measures can provide very different results.

Thus, while the results of the above empirical and statistical data analyses clearly illustrate that web sites

experience patterns in their traffic, most of these results by themselves are not easily exploitable as the basis for our clustering, profiling and classification purposes. The present invention provides a mechanism for exploiting these traffic
5 patterns to generate accurate models of the web sites for use in workload characterization, performance modeling, workload and performance forecasting, and capacity planning.

The present invention may be broken into three primary components: clustering, profiling and characterizing web
10 sites. The first step of characterizing web sites is to generate clusters of web sites based on traffic data obtained for these websites. Once the clusters are identified, each cluster, or class, is profiled to obtain a template for the class. Thereafter, as new traffic data is obtained for a web
15 site, the traffic data may be compared against established templates for the classes in order to categorize the web site into one of the known classes. Alternatively, if the comparison results in the web site being sufficiently different from all of the known classes, a new class may be
20 generated using the traffic data for the web site. This classification may then be used to perform functions such as workload characterization, performance modeling, workload and performance forecasting, and capacity planning, in order to best optimize the available resources for the web site.

As mentioned above, the first step in the operation of
25 the present invention is to obtain traffic data from a plurality of web sites and cluster the web sites based on their traffic data. Clustering involves selecting a measure of traffic data to be used to cluster web sites and then
30 identifying templates of the traffic data with regard to this selected measurement for each of the web sites. The templates are then clustered using a clustering algorithm which identifies groups of templates that are most similar to one another within a given tolerance. These groups, or classes,
35 are the clusters of web sites that will be used to perform

profiling and classification.

As mentioned above, the particular measure of traffic data used to perform the clustering must be selected prior to performing the clustering. Depending on the subset of measurement data used, different clustering results can be obtained. For example, web sites can be clustered according to the load/request patterns, user navigation patterns, site hypertext structures, etc. Each of these clusterings are different and can be considered orthogonal to one another.

5 The particular clustering performed with the present invention may be selected based on the particular implementation of the present invention and the measures that are most important to a user of the present invention.

In an exemplary embodiment of the present invention, the clustering is performed according to load patterns. For this purpose, the number of hits per hour is used as representative traffic data for characterizing the incoming request traffic patterns. **Figure 8** is an exemplary diagram illustrating requests per hour over one day collected from exemplary commercial web sites for use in clustering the web sites according to this exemplary embodiment.

15 Once the measurement of traffic data that is to be used to perform clustering is selected, the traffic data is then analyzed to identify templates, or typical shapes, in the traffic data with regard to this selected measurement. This template is essentially the time series data values for the selected measurement, or a function of these time series data values for the selected measurement.

20 In the exemplary embodiment of the present invention, the weighted average request pattern $l_i(h)$ representing the weighted mean of the hourly request pattern profile that occurs on a web server is utilized as a template for the web site. The weighted average request pattern $l_i(h)$ is obtained using the following equation:

$$l_i(h) = \sum_{d=1}^D (a_i(d)/D) (x_i(h, d)) \quad (1)$$

5 where $x_i(h, d)$ denotes the number of requests from the empirical data that the i th web server receives in the h th hour of day d , $h=0, 1, \dots, n$, and where the weights $a_i(d)$ are the weight for day d of site i so that the workloads of different days are normalized to the same mean. Moreover,
 10 with the exemplary embodiment, the peak load regimes are determined to be the focus of the clustering since they have a more significant impact on web server performance. Thus, each weighted average request pattern is normalized by its maximum value and its peak hour traffic pattern is defined as follows:

$$15 \quad g_i(h) = \begin{cases} l_i(h)/m_i, & l_i(h)/m_i > 0.5, \\ 0.5, & l_i(h)/m_i \leq 0.5 \end{cases} \quad (2)$$

where $m_i = \max_h \{l_i(h)\}$.

A dissimilarity measure between the peak hour patterns is defined as:

$$20 \quad d_{i,j} = \min_{h^d} \max_h \{ |g_i(h) - g_j(h+h^d \bmod 24)| \} \quad (3)$$

where h^d is used as the hourly shift needed when comparing two traffic patterns from different web servers with, for
 25 example, differences in time zones. This dissimilarity measure is the minimum of the maximum difference between the normalized weighted average request pattern for web site i and the normalized weighted average request pattern for web site j , shifted to compensate for the differences in time zones, if

any. This dissimilarity measure is used to identify the normalized weighted average request patterns that are most similar to one another in order to cluster the patterns into classes of web sites.

5 In the exemplary embodiment, a complete linkage, or furthest neighbor, algorithm is used to cluster the normalized weighted average request patterns based on the dissimilarity measure. That is, in a first step, each pattern represents its own cluster and the distances between these patterns are
10 defined by the dissimilarity measure given in equation 3 above. Then, the two patterns with the smallest distances are linked together. The distances between this new cluster and the other clusters (or individual patterns) are defined by the greatest distance between any two patterns in the respective
15 clusters, i.e. by the furthest neighbors. As a result, the algorithm proceeds in subsequent steps to link more and more patterns together and to aggregate larger and larger clusters within a predetermined threshold.

In the exemplary embodiment, using a threshold of in the
20 complete linkage algorithm, the request patterns shown in **Figure 8** are clustered into four distinct classes. **Figure 9** is an exemplary diagram illustrating the patterns for these four distinct classes of request patterns into which the web sites of **Figure 8** are clustered. From **Figure 9**, the following
25 observations about the various classes can be made. For class 1 patterns, the request traffic load increases to the peak level by noon and then goes down significantly in the afternoon. This suggests that users' interests for such web sites are more instantaneous, just like checking the weather
30 report every day. The web sites are popular primarily in their local areas.

For class 2 patterns, the request traffic load increases to the peak level at noon and continues to remain high in the afternoon, but becomes very low in the evening. This suggests that users show their interests for these web sites primarily

5 during working hours. For class 3 patterns, the request traffic load increases to the peak level somewhat after noon and remains high throughout much of the afternoon. In the evening users continue to show some interest for these web sites.

10 For class 4 patterns, the request traffic load remains at a high level over a long period of time, from before noon well into the evening. This suggests that most of the users visit these web sites either during working hours or in their spare time and that the users are probably spread over the country

15 and even the world. By making such observations regarding the characteristics of web sites falling into each of these different classes, it is possible to predict the usage of a web site that is later classified into one of these classes. Thus, from such a

20 prediction, various measures can be employed to handle the traffic that the web site should expect to experience. Thus, according to a preferred embodiment of the present invention, clustering of web sites involves obtaining traffic data for a plurality of web sites, determining a measure of

25 the traffic data to use as a basis for the clustering, identifying a pattern of the traffic data in accordance with the selected measure, defining a dissimilarity or similarity relationship for the traffic data, and then using a clustering algorithm to cluster the web sites based on this dissimilarity

30 or similarity relationship. Once these clusters are identified, profiling of the clusters is performed to identify a template for the cluster that may be used with later classification of web sites.

Profiling involves first determining a metric upon which

35 the traffic profiles will be based. In an exemplary

embodiment, the weighted average load $l_i(h)$ is used to generate the profiles for the identified classes of web sites.

In generating the

profile for a class, a template for the class is identified

5 that is defined as a request pattern that is most similar to all of the members of the class.

One approach to finding the template for a class is to simply average all the members of the class. While this may be done to obtain a template for the class, the result will

10 typically not be a good choice for a template for the class

since an outer member of the class may be far from this

template but close to the templates of other classes. Another

underlying consideration is that the templates defined for

different classes should be far from each other. Therefore,

15 the template for a class is defined so that it minimizes the maximum difference to all members in the class.

Assume that G_k denotes the set of request patterns that belong to class k . Though similar in shape, the members, i.e. the web sites, within class k may actually be located in

20 different time zones, or include any other sources of shifted behavior. To define the template, it is therefore necessary

to first identify the correct shift biases $\{h_i^d, icG_k\}$ so that

upon the shift, all members have the closest shapes. Such

shift biases can be solved via the following mathematical

25 expression:

$$\text{Min}_{\{h_i^d, h_j^d, i, j \in G_k\}} \{ \max_{\{i, j \in G_k\}} [|g_i(h+h_i^d \bmod 24) - g_j(h+h_j^d \bmod 24)|] \}$$

... (4)

That is, upon the shift $\{h_i^d, icG_k\}$, the maximum difference
30 between any two members of the class should be minimized. The particular shift, in an exemplary embodiment, may be identified using a shifting algorithm such as that set forth

below. Other algorithms for identifying the particular shift may be used without departing from the spirit and scope of the present invention.

In an exemplary embodiment, the shift algorithm involves
 5 a first step in which the set $\{h_i^d, icG_k\}$ is to be arbitrary integers between 0 and 23 (possible shift in hours). Then for each icG_k the shift value h_i^d is updated so that:

$$H_i^d = \arg \min_{h_i^d} \{ \max_{\{i, j \in G_k\}} [|g_i(h+h_i^d \bmod 24) - g_j(h+h_j^d \bmod 24)|] \}$$

10 ... (5)

This step is then repeated until $\{h_i^d, icG_k\}$ converges to a local optimum.

Thus, for each icG_k h_i^d is chosen so that upon the shift, member i is close to all other members of the class. The
 15 procedure is then repeated iteratively until no further improvement can be obtained.

Once the optimal shift biases $\{h_i^d, icG_k\}$ are identified, the template of class k is defined to be:

$$T_k(h) = 1/2 \{ \max \{ g_j(h+h_j^d \bmod 24) \} + \min \{ g_j'(h+h_j^d \bmod 24) \} \}$$

20 ... (6)

for $h = 0, 1, \dots, 23$. That is, the template is the most similar pattern to the patterns for all of the members of the class. **Figure 10A** is an exemplary diagram illustrating the
 25 four templates for the four classes of request patterns of **Figure 9** generated using the clustering and profiling of the present invention. Applying this same clustering and profiling to empirical data, different templates may be generated for different traffic effects, such as day of the
 30 week, week of the month, month of the year, etc. For example, as shown in **Figure 10B**, the day of the week patterns are shown

which are generated by applying the clustering and profiling of the present invention to the empirical data for the web sites of **Figures 4A-4C**.

5 The templates for the classes generated using the clustering and profiling described above may be used to recognize incoming request patterns. First, the templates are used to recognize the existing samples and to check if the results of this recognition are consistent with the clustering. Thereafter, the request pattern is matched to a
10 class of requests based on the template. From this matching, certain characteristics of the web site traffic may be discerned based on the characteristics of the other web sites that are part of the class.

Preferably, in order to perform the classification,
15 equations 1 and 2 above are used to compute the normalized peak hour request pattern for newly collected data. Thereafter, equation three is preferably used to compute the distance measures between the incoming request pattern data and the data for the four templates. Based on these distance
20 measures, a closest matching template may be identified and the incoming request pattern classified into the corresponding class.

Since the predetermined templates may not cover all possible request patterns, a new type of request pattern or an
25 extraordinary request pattern may be far from all of the predetermined templates, i.e. the minimum distance is greater than a threshold amount. In such cases, the new request pattern may be added as a new template for a new class of request patterns.

30 By clustering, profiling and classifying web sites according to the present invention, characteristics about web site traffic may be identified based on the web sites falling into the same class. This classification may be used with many different types of applications including traffic
35 prediction, capacity planning, hot-spot detection, dynamic

off-loading, web site co-location, and the like.

With regard to traffic prediction, one key issue in capacity planning is the prediction of workload behavior. The prediction mechanism needs to capture the characteristics of
5 long-term trends, periodicity, dependency and variability. It is difficult to use a single technique to capture all of these factors. Therefore, a more accurate approach would be to use a hybrid technique pertaining to both macro and micro level statistics. While long-term trends may be measured using
10 linear regression methods, the periodicity at different scales (e.g., monthly , weekly, daily, etc.) may be handled using the clustering and profiling technique of the present invention.

The clustering approach of the present invention can
15 greatly simplify the capacity planning task. With the present invention, the templates (or profiles) of different clusters (or classes) may be used to analyze the capacity demand for each individual profile, the impact of the scaling factor, and the mixtures of the profiles for servers in a web server farm,
20 cluster, or the like. When a new customer comes along, the classification technique of the present invention may be used to determine the cluster to which the new customer's traffic belongs so as to adjust the capacity requirements, if necessary. In the same way, short-term capacity planning
25 decisions may also be easily adjusted if some web sites cause the clusters to change because of special events, web site redesign, etc.

With special events, e.g., holiday sales for e-commerce web sites, some web sites can be heavily loaded and thus,
30 require certain additional operations in order to fulfill the needed quality of service. Examples of such operations include offloading and adding new resources. Such special events represent "hot-spots." With hot-spot detection, the goal is to detect the hot-spots so that appropriate mechanisms
35 for handling the hot-spots may be triggered. With the

profiling approach of the present invention, these hot-spots can be detected once it is observed that the current workload is deviating significantly from typical behavior which is described by the templates.

5 As mentioned above, sometimes dynamic offloading operations are needed in order to alleviate server overload, such as when a hot-spot is encountered. One way to achieve this dynamic offloading is to create new (or use different) versions of the web pages with references of offloadable
10 objects (such as images) to the server onto which some of the extra load can be offloaded (such as Akamai servers). With proper use of the templates of the present invention, it is easy to determine what is the threshold beyond which traffic should be offloaded. Moreover, the offloading scheme may be
15 started before the server is saturated in order to account for the lag time of the offloading scheme. This new threshold, i.e. the threshold accounting for lag time, may be determined from the template together with the lag time.

With web site co-location, the goal is to share resources
20 among multiple web sites so that peak load conditions for any given web site can be handled by borrowing resources from the other entities. An important problem in this paradigm concerns the clustering of the web sites for resource sharing. Based on the observed traffic templates, one can easily
25 identify the shapes of the peak regimes for the different sites, as well as the different traffic peaks and valleys at different times and at different geographical time zones. Optimization tools may then be used to achieve load balancing across a number of web sites in order to obtain the smoothest
30 possible overall peak loads. The optimization problem can be considered as a general bin packing problem where the items are the templates. This also can be formulated as an integer programming problem.

It should be noted that while the above embodiments of
35 the present invention have been described with regard to

request patterns, the present invention is not limited to such. Rather, any measurement data for web sites may be used to perform the clustering, profiling and classification of the present invention. For example, the present invention may
5 operate on server utilization data, bandwidth consumption data, or the like.

Figure 11 is an exemplary block diagram of a web site classification device according to the present invention. The elements shown in **Figure 11** may be implemented in hardware,
10 software, or any combination of hardware and software without departing from the spirit and scope of the present invention. In a preferred embodiment, the elements of the web site classification device are implemented as software instructions executed by one or more processors.

As shown in **Figure 11**, the web site classification device includes a controller **1110**, an input/output interface **1120**, a web site data storage device **1130**, a clustering engine **1140**, a profiling engine **1150**, and a classification engine **1160**. The elements **1110-1160** are in communication with one another via
20 the control/data signal bus **1170**.

The controller **1110** controls the overall operation of the web site classification device and orchestrates the operation of the other elements **1120-1160**. The controller **1110** receives web site traffic data from web sites via the input/output
25 interface **1120** and stores this web site traffic data in the web site data storage device **1130**. The controller **1110** then instructs the clustering engine **1140** to cluster the web sites for which data is stored in the web site data storage device **1130**.

30 Once the clustering is performed, and the web sites are assigned to particular classes of web sites, the controller **1110** instructs the profiling engine **1150** to generate a template, or profile, for each of the classes. The clustering and profiling may be performed on a periodic basis such that

the clusters maintained by the web site classification device, and their corresponding templates or profiles, are updated as new web site traffic data becomes available. Alternatively, the clustering and profiling may be updated each time traffic data is received from a new web site so as to include this new web site into the clusters and templates maintained by the web site classification device.

When a traffic data is received from a new web site, the traffic data is compared to the templates for the various classes maintained by the web site classification device. From this comparison, a determination may be made as to whether the traffic data for the new web site fits the template for one of the classes or is sufficiently dissimilar to all of the templates for the classes so as to warrant the creation of a new class and template based on the traffic data for the new web site. From this classification of the traffic data for the new web site, the new web site is assigned to a particular class of web sites or is used to generate its own class of web sites. As noted above, this classification may then be used by different processes to perform resource management functions such as workload characterization, performance modeling, workload and performance forecasting, capacity planning, and the like.

Figure 12 is a flowchart outlining an exemplary operation of the present invention. As shown in **Figure 12**, the operation of the present invention starts with receiving web site traffic data from a plurality of web sites (block **1210**). The web site traffic data is then clustered using the clustering technique discussed above (block **1220**). For each cluster, or class, a template or profile is generated (block **1230**).

The operation then waits for traffic data from a new web site (block **1240**). A determination is then made as to whether traffic data for a new web site is received (block **1250**). If not, the operation returns to block **1240** and continues to wait

for traffic data from a new web site. If traffic data for a new web site is received, the traffic data for the new web site is classified based on the existing classes of web sites (block 1260). The operation then ends.

5 **Figure 13** is a flowchart outlining an exemplary operation of the present invention for clustering web sites. As shown in **Figure 13**, the clustering operation of the present invention starts by selecting a traffic data measurement by which to perform the clustering (block 1310). Templates of
10 the traffic data for the web sites with regard to the selected traffic data measurement are then identified (block 1320). The templates are then clustered using a clustering algorithm (block 1330).

Figure 14 is a flowchart outlining an exemplary operation
15 of the present invention for profiling web sites. As shown in **Figure 14**, the operation starts with determining a metric upon which the traffic profiles will be based (block 1410). Templates for each web site's traffic data in the class are generated based on this selected metric (block 1420). The
20 templates are then shifted by a shift bias amount, if necessary, to compensate for factors such as different time zones (block 1430). A profile is then generated by selecting a profile that is most similar to all of the templates for the web sites in the class (block 1440). This operation may be
25 performed for each class of web sites.

Figure 15 is a flowchart outlining an exemplary operation of the present invention for classifying web sites. As shown in **Figure 15**, the operation for classifying web sites starts with the receipt of traffic data for a new web site (block
30 1510). A template for the traffic data for the new web site is generated (block 1520) and compared to the templates for the existing classes of web sites (block 1530). A determination is made as to whether a matching template is identified (block 1540). If so, the web site is classified

into the class associated with the matching template (block 1550). If not, a new class is generated using the traffic data for the new web site as a basis for generating a template for the new class (block 1560). The operation then ends.

5 Thus, the present invention provides mechanisms for the clustering, profiling and classification of web sites based on their traffic data. With the present invention, similarities between web sites with regard to their experienced traffic may be identified for use in predicting and planning for workloads
10 that are most likely to be experienced in the future. Thus, the present invention provides a tool through which resource management for web servers may be performed with regard to the web sites they host.

 The clustering, profiling and classification of web sites
15 in the manner discussed above may be applied to resource allocation to obtain an optimized resource allocation for a group of web sites. With the present invention, once the templates for web sites in a group of webs sites are obtained through the mechanisms discussed above, these templates may be
20 used to determine which web sites in the group are candidates for co-location and when traffic should be offloaded to other web servers. These functions of co-location and offloading may further be combined to obtain greater optimization of the resource allocation to provide a required Quality of Service
25 (QoS) to client devices that access those web sites.

 With co-location, the goal is to co-locate web sites which have peaks and troughs in traffic at different times based on the templates associated with the web sites. That is, the goal is to have a single server (or a suite of
30 servers) that handles traffic for two or more web sites wherein the templates compliment each other to obtain a fairly consistent resource utilization or at least to provide the required Quality of Service at all time periods.

 There are several variants of the co-location problem can
35 be solved. One example is that the web server must have a

given capacity for traffic. This is the typical case in which a web server provider wants to keep a certain capacity at all time epochs to handle any non-normal traffic volumes that may be experienced.

5 This type of co-location problem is solved by the present invention by first determining the template of the traffic for each web site under consideration for co-location. This template is determined from the traffic logs of the web site. The template may be obtained through the clustering,
10 profiling, and/or classification mechanisms described previously.

Once the templates for each of the web sites under consideration are obtained, an integer programming (IP) problem is formulated. This integer programming problem has
15 decision variables that are the assignment of the web sites to web servers. The constraints of the integer programming problem are the capacity limitation for any server at any time epochs. The objective function is any increasing function of the minimum distances between the server capacity and the peak
20 load at different servers.

After the integer programming problem is formulated, it is solved using optimal or approximate algorithms. Standard methods and algorithms to solve integer programming problems include branch and bound, cutting plane algorithms, LP
25 relaxation methods, and the like.

For example, assume that there are I web sites that must be served with a total of J servers, each having a capacity C_j , where $j = 1$ to J . Suppose that the traffic of Web site i can be characterized by template $T_i(t)$, where $i = 1$ to I . The
30 binary decision variables are x_{ij} , where $i = 1$ to I and $j = 1$ to J , such that $x_{ij} = 1$ if site i is assigned to server j , otherwise $x_{ij} = 0$.

The co-location problem can then be formulated as the following binary integer programming problem, where b is the
35 target utilization:

min b

s.t. $\sum_{i=1}^I x_{ij} T_i(t) \leq b C_j, j=1, \dots, J, t \geq t_0.$

5

$\sum_{j=1}^J x_{ij} = 1, i=1, \dots, I.$
 $b[1$

where x_{ij} 's are binary integers.

10 In the above approach, the objective function can also be another load balancing criteria, such as the empirical variances of the loads, for example. Some slackness in the constraints may also be introduced such that the peak loads are below the server capacities by a certain amount. In this way, the mean response time of the requests will be upper-
 15 bounded.

Another example of the co-location problem is to minimize the total cost of the web servers, given that the different web servers have different costs and the cost is incurred if
 20 at least one web site is assigned to it. In this case, the co-location problem may be solved in the following manner.

As with the previous variation of the co-location problem, the first step of the process is to identify the template of the traffic for each web site under consideration
 25 for co-location. Once the templates are identified, the integer programming problem is formulated in which the decision variables are again the assignment of the web sites to the web servers. The constraints are the capacity limitation for any server at any time epochs. The objective
 30 function, however, is the sum of the costs of all web servers to which at least one web site is assigned. This integer programming problem is solved using optimal or approximate

algorithms.

As an example of this variation of the co-location problem, consider the case described previously but where the objective is to minimize the total cost of the web servers given that using web server j will cost P_j , where $j=1$ to J . In such a case, the objective function of the previous example is replaced with the following:

$$\min \sum_{j=1}^J P_j \sum_{i=1}^I 1(Sx_{ij}>0),$$

where the function $1(a)$ equals 1 if statement a is true, and 0 otherwise. Thus, the binary integer programming problem is formulated as follows:

$$\begin{aligned} \min \sum_{j=1}^J P_j \sum_{i=1}^I 1(Sx_{ij}>0), \\ \text{s.t. } \sum_{i=1}^I x_{ij} T_i(t) \leq bC_j, \quad j=1, \dots, J, \quad t \geq 0. \\ \sum_{j=1}^J x_{ij} = 1, \quad i=1, \dots, I. \end{aligned}$$

where x_{ij} 's are binary integers.

A simpler version of the above problem is that the server capacities are identical and so are the server costs. In this case, the problem reduces to a two dimensional bin packing problem where the bins are the servers and the items are the traffic templates. Thus, the problem can be solved in the following manner.

First, as with all of the co-location problems, the templates for the web sites under consideration for co-location are identified. Then the bin packing problem is

formulated with the objective being to minimize the number of bins used. Thereafter, the two dimensional bin packing problem is solved using fast or approximate algorithms. Such fast and approximate algorithms are known in the art and more information regarding such fast and approximate algorithms may be found, for example, in E.G. Coffman, Jr. et al., "Approximation Algorithms for Bin-Packing -- An Updated Survey," Algorithm Design for Computer System Design, edited by Ausiello, Lucertini, and Serafini, Springer-Verlag, 1984 and David S. Johnson, "Fast Algorithms for Bin Packing," Journal of Computer and System Sciences 8, pages 272-314, 1974.

Thus, as shown in **Figure 16**, with each variance of the co-location problem, the operation for determining the assignment of web sites to web servers, and thus, the co-location of web sites, starts by identifying the templates for the web sites being considered (block **1610**). The co-location problem is then formulated as an integer programming problem with an objective function that embodies the desired optimization for allocation of the web sites (block **1620**). The integer programming problem is formulated such that the decision variables are the assignment of the web sites to web servers and the constraints are the capacity limitation for any server at any time epoch. Once the integer programming problem is formulated, the integer programming problem is solved with optimal or approximate algorithms (block **1630**). As noted above, the integer programming problem may be reduced to a bin packing problem when certain characteristics, such as capacity and cost, are considered identical for all servers.

As mentioned previously, the present invention may further be applied to offloading of traffic to other web sites. The offloading mechanism involves redirecting part of the incoming requests to backup or other less-loaded servers when a web site becomes heavily loaded and thus, requires

certain additional service capacity in order to fulfill the needed quality of service. When the offloading mechanism is turned on, requests can be redirected by, for example, using HTTP redirection response code 301 or 302 to have the client's browser retry the request at a different web server or placing a load-balancing device, like IBM Network Dispatcher or Cisco Global Director, in front of the web servers and using the ability of those devices to send requests to different servers at different times under control of the present invention.

The key question concerns how to proactively detect or predict such heavily loaded situation thus triggering the offload mechanism at the right time. The offloading according to the present invention will be described with reference to **Figure 17** which illustrates a template for a web site with thresholds for offloading according to the present invention also depicted.

With offloading according to the present invention, the templates of traffic for each web site under consideration are determined from the traffic logs of the web sites in the manner previously described. **Figure 17** illustrates one such template for a web site. In a preferred embodiment, the templates are determined in terms of the number of requests per unit of time.

A first threshold **1710** is calculated with respect to the traffic intensity, above which a fraction of the traffic should be offloaded from the home web server in order to guarantee a Quality of Service (QoS) criteria. This threshold **1710** is determined based on the marginal distribution of the total traffic and the amount of traffic that is offloadable. More precisely, if h_1 is the threshold **1710**, a is a proportion of offloadable traffic, $E[S_o]$ and $E[S_{no}]$ are the expected service requirement per offloadable and non-offloadable request, C is the capacity of the bottleneck resource under consideration, and b is the target resource utilization which

is determined by the QoS that the system needs to guarantee, then the first threshold **1710** may be determined by:

$$h_1 = bC / (a E[S_o] + (1-a)E[S_{no}])$$

Once the first threshold **1710** is computed in the manner
5 described above, a second threshold **1720** is computed with
respect to the traffic intensities in the transition regime
from low to high. A third threshold **1730** is computed with
respect to the traffic intensities in the transition regime
from high to medium. The calculation of the second threshold
10 **1720** takes into account the lead time to meet the threshold
1710 based on the reaction time of the dynamic offloading
mechanisms and the statistical pattern of the traffic.
Similarly, the calculation of the third threshold **1730** takes
into account the lag time based on the reaction time of the
15 dynamic offloading mechanisms and the statistical pattern of
the traffic. Thus, the second threshold **1720** may be
calculated by subtracting a lead time from threshold **1710** and
the third threshold **1730** may be calculated by adding a lag
time to the threshold **1710**.

20 Once the templates are identified and the thresholds are
calculated, online transitions from low to high regimes and
from high to medium regimes are detected. The detection can
be carried out using several statistical techniques. One such
technique is change-point detection which consists of
25 computing the marginal distributions of the incoming traffic.
With this technique, the deviation of this marginal
distribution is detected. More information regarding
change-point detection may be found, for example, in Carlstein
et al., "Change-Point Problems," IMS Lecture Notes - Monograph
30 Series, vol. 23, 1994. A simpler way to detect the
transitions is to compute the moving average of the traffic
intensity. This moving average may be taken at quite coarse

time scales on the order of minutes or tens of minutes. For example, one way of computing the moving average is to compute the weighted sum of the previous estimates together with a number of recent observations. For example:

$$5 \quad \text{Estimate}_{\text{new}} = c_1 T(t) + \dots + c_k T(t-k+1) + \\ (1-c_1 - \dots - c_k) \text{Estimate}_{\text{add}}$$

where $c_1 \dots c_k$ are the nonnegative weights. The requirement is that $c_1 + \dots + c_k = 1$.

With the above approach, the offloading mechanism is
 10 turned on when the traffic intensity exceeds threshold **1720** and is turned off when the traffic intensity falls below threshold **1730**.

Figure 18 is a flowchart outlining an exemplary operation of the present invention for offloading traffic according to
 15 the present invention. As shown in **Figure 18**, the operation starts with identifying the templates for each web site under consideration (block **1810**). The first threshold is computed with respect to traffic intensity based on the marginal distribution of the total traffic and the amount of traffic
 20 that is offloadable (block **1820**). The second and third thresholds are then calculated based on the first threshold and the lead and lag times of the offload mechanism (block **1830**).

Thereafter, the on-line traffic is monitored to determine
 25 if there are transitions from low to high regimes and high to medium regimes (block **1840**). A determination is made as to whether a transition is detected (block **1850**). If not, the operation determines if an end condition occurs, e.g., offloading is disabled (**1860**). If so, the operation ends.
 30 Otherwise, the operation returns to block **1840** and continues to monitor for transitions.

If a transition is detected, a determination is made as

to whether it is a low to high transition (block **1870**). If so, the traffic offloading mechanism is enabled and traffic is offloaded to another server (block **1880**). If the transition is not low to high, then the transition is a high to medium transition and the offloading mechanism is disabled (block **1890**).

Thus, as described above, the present invention of clustering, profiling and classifying web sites based on their traffic logs to thereby generate templates of the traffic for the web sites, may be applied to resource allocation mechanisms such as co-location and offloading. As mentioned previously, the present invention may further be applied to a combination of co-location and offloading. The combination problem of co-location and offloading will often arise in a situation where a web server hosts several web sites and when the load of the server exceeds its capacity, it must offload the offloadable traffic to other web servers. With such a situation, present invention may be implemented in the following manner to perform offloading of co-located web sites.

As with all of the previous mechanisms, the templates for each of the web sites under consideration must first be determined. A mathematical programming problem is then formulated in which the decision variables are the assignment of fractions of traffic from web sites to web servers. The constraints are the capacity limitation for any server at any time epochs. The objective function is the sum of the costs of all web servers to which at least one web site is assigned.

As an example, assume similar conditions as described in the previous co-location problem. Note that the decision variables x_{ij} 's no longer need to be binary integers. Instead, x_{ij} denotes the fraction of traffic from web site i that can be assigned to server j , where $i= 1$ to I and $j= 1$ to J . Thus, $0 \leq$

$x_{ij} \leq 1$, and the mathematical programming problem is formulated as follows:

$$\begin{aligned}
 & \min S \sum_{j=1}^J \sum_{i=1}^I P_j 1(Sx_{ij} > 0), \\
 & \text{s.t. } \sum_{i=1}^I x_{ij} T_i(t) \leq bC_j, \quad j=1, \dots, J, \quad t \geq 0. \\
 & \sum_{j=1}^J x_{ij} = 1, \quad i=1, \dots, I. \\
 & 0 \leq x_{ij} \leq 1, \quad i=1, \dots, I \text{ and } j=1, \dots, J.
 \end{aligned}$$

The mathematical programming problem, once formulated, is solved using optimal or approximation algorithms. The mathematical programming problem has linear constraints and becomes a linear programming problem if the server costs are assumed to be proportional to their utilization.

For each web site, a web server is designated as its home server among those to which some portion of the traffic is assigned. The fraction of the traffic assigned to the home server will also be referred to as the threshold 1 of the web site at its home server.

For each web site, a threshold 2 is computed with respect to the traffic intensities in the transition regime from low to high, and a threshold 3 is computed with respect to traffic intensities in the transition regime from high to medium. These thresholds take into consideration the lead and lag times required by the offloading mechanism. These lead and lag times are a combination of the reaction time of the dynamic offloading mechanisms, the statistical pattern of the traffic, and the web site's share of the capacity at the home server among all other web sites using this server as a home server.

For each web site, the on-line transitions from low to

high regimes and from high to low regimes are detected. As previously discussed, the detection can be carried out using several statistical techniques such as the change point detection based on marginal distribution or the moving average technique.

For each web site, the offloading mechanism is either turned on when the traffic intensity exceeds threshold 2 or turned off when the traffic intensity falls below threshold 3. The destinations to which the traffic is offloaded and the fractions of traffic to send to these destinations are defined by the solution of the linear programming problem solved previously.

Figure 19 is a flowchart outlining an exemplary operation of the present invention for a combination of both co-location and offloading. As shown in Figure 19, the operation starts with the templates for each of the web sites under consideration being determined (block **1910**). A mathematical programming problem is then formulated in which the decision variables are the assignment of fractions of traffic from web sites to web servers (block **1920**). The constraints of this mathematical programming problem are the capacity limitation for any server at any time epochs. The objective function is the sum of the costs of all web servers to which at least one web site is assigned.

The mathematical programming problem, once formulated, is solved using optimal or approximation algorithms (block **1930**). Then, for each web site, a web server is designated as its home server among those to which some portion of the traffic is assigned (block **1940**) and a threshold 1 of the web site is defined as its portion of the traffic assigned to its home server (block **1950**).

For each web site, a threshold 2 is computed with respect to the traffic intensities in the transition regime from low to high, and a threshold 3 is computed with respect to traffic

intensities in the transition regime from high to medium (block **1960**). These thresholds take into consideration the lead and lag times required by the offloading mechanism. For each web site, the on-line transitions from low to high regimes and from high to low regimes are detected (**1970**). The offloading mechanism is either turned on when the traffic intensity exceeds threshold 2 or turned off when the traffic intensity falls below threshold 3 (**1980**).

Figure 20 is a block diagram of a resource allocation determination system. As shown in **Figure 20**, the resource allocation determination system includes a controller **2010**, an interface **2020** to the web site classification device of **Figure 11**, an input/output interface **2030** for sending and receiving resource allocation messages to and from web servers and offloading mechanisms, a co-location determination device **2040**, and an offloading enablement device **2050**. The elements **2010-2050** are in communication with one another via the control/data signal bus **2060**.

The elements shown in **Figure 20** may be implemented in hardware, software, or any combination of hardware and software. In a preferred embodiment, the elements **2010-2050** are implemented as software instructions executed by one or more processors.

The controller **2010** controls the overall operation of the resource allocation determination system and orchestrates the operation of the other elements **2020-2050**. The controller **2010** receives request for determination of co-location, offloading, or a combination of co-location and offloading, via the interface **2030** and instructs the elements **2040-2050** to determine co-location and/or offloading using the templates for the web sites as obtained from the web site classification device via interface **2020**. The elements **2040** and **2050** perform operations for determining co-location of web sites and/or

offloading in the manners previously described and return results to the controller **2010**. The controller **2010** may then transmit messages to web servers and/or offloading mechanisms for performing operations to co-location web sites and/or
5 offload traffic to other web servers.

Thus, the present invention provides apparatus and methods for determining the co-location of web sites based on templates of traffic patterns identified using clustering, profiling, and/or classification of web site traffic data
10 obtained from traffic logs of the web sites. Moreover, the present invention provides apparatus and methods for determining offloading of traffic from a web server to other web servers based on such templates of traffic patterns for web sites. With the present invention, dynamic determination of
15 optimal co-location of web sites and offloading may be performed to obtain a required guaranteed quality of service.

It is important to note that while the present invention has been described in the context of a fully functioning data processing system, those of ordinary skill in the art will
20 appreciate that the processes of the present invention are capable of being distributed in the form of a computer readable medium of instructions and a variety of forms and that the present invention applies equally regardless of the particular type of signal bearing media actually used to carry
25 out the distribution. Examples of computer readable media include recordable-type media such a floppy disc, a hard disk drive, a RAM, and CD-ROMs and transmission-type media such as digital and analog communications links.

The description of the present invention has been
30 presented for purposes of illustration and description, but is not intended to be exhaustive or limited to the invention in the form disclosed. Many modifications and variations will be apparent to those of ordinary skill in the art. The embodiment was chosen and described in order to best explain
35 the principles of the invention, the practical application,

and to enable others of ordinary skill in the art to understand the invention for various embodiments with various modifications as are suited to the particular use contemplated.

CLAIMS:

1. A computer program product in a computer readable medium for allocating resources to a plurality of web sites, comprising:

5 first instructions for identifying a traffic pattern for each web site in the plurality of web sites;

second instructions for identifying a template for each web site in the plurality of web sites based on the traffic pattern; and

10 third instructions for allocating resources to web sites in the plurality of web sites based on the identified templates for each web site in the plurality of web sites.

2. The computer program product of claim 1, wherein the third instructions for allocating resources to web sites

15 include:

instructions for identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites; and

20 instructions for allocating resources to the two or more web sites based on the identification of the two or more web sites being candidates for co-location.

3. The computer program product of claim 1, wherein the third instructions for allocating resources to web sites include:

25 instructions for calculating, for each web site in the plurality of web sites, a first threshold based on the template for the web site;

30 instructions for calculating a second threshold based on the first threshold, wherein offloading of traffic is enabled when a traffic intensity for the web site meets or exceeds the second threshold; and

instructions for calculating a third threshold based on the first threshold, wherein offloading of traffic is disabled when a traffic intensity of the web site meets or falls below the third threshold.

- 5 4. The computer program product of claim 2, wherein the third instructions for allocating resources to web sites includes:

instructions for calculating, for each web site in the plurality of web sites, a first threshold based on the
10 template for the web site;

instructions for calculating a second threshold based on the first threshold, wherein offloading of traffic is enabled when a traffic intensity for the web site meets or exceeds the second threshold; and

- 15 instructions for calculating a third threshold based on the first threshold, wherein offloading of traffic is disabled when a traffic intensity of the web site meets or falls below the third threshold.

5. The computer program product of claim 3, further
20 comprising:

fourth instructions for monitoring traffic of a web site of the plurality of web sites, on a web server to determine if the traffic exceeds the second threshold; and

- 25 fifth instructions for offloading at least a portion of the traffic to another web server if the traffic of the web site exceeds the second threshold.

6. The computer program product of claim 5, further comprising:

sixth instructions for monitoring traffic of a web site
30 of the plurality of web sites, on a web server to determine if the traffic falls below the third threshold; and

seventh instructions for disabling the offloading of

traffic to the another web server if the traffic of the web site falls below the third threshold.

7. The computer program product of claim 2, wherein the instructions for identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites include:

instructions for identifying a first web site having peaks in traffic at a first set of time periods; and

instructions for identifying a second web site having peaks in traffic at a second set of time periods different from the first set of time periods.

8. The computer program product of claim 7, wherein the second web site has a trough in traffic at approximately a same time as the first web site has a peak in traffic.

9. The computer program product of claim 2, wherein the instructions for identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites include:

instructions for identifying the two or more web sites such that a capacity for traffic of a web server on which the two or more web sites are located remains constant.

10. The computer program product of claim 2, wherein the instructions for identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites include:

instructions for identifying the two or more web sites such that a the total cost of all web servers hosting the plurality of web sites is minimized.

11. The computer program product of claim 2, wherein the

instructions for identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites include:

instructions for formulating an integer programming
5 problem having decision variables that are an assignment of the plurality of web sites to a plurality of web servers, a constraint that is a capacity limitation for any server at any time epoch, and an objective function that is an increasing
10 function of minimum distances between server capacity and peak load at different ones of the plurality of web servers.

12. The computer program product of claim 2, wherein the instructions for identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites include:

15 instructions for formulating an integer programming problem having decision variables that are an assignment of the plurality of web sites to a plurality of web servers, a constraint that is a capacity limitation for any server at any time epoch, and an objective function that is a sum of costs
20 of all web servers to which at least one web site is assigned.

13. A method of allocating resources to a plurality of web sites, comprising:

identifying a traffic pattern for each web site in the plurality of web sites;
25 identifying a template for each web site in the plurality of web sites based on the traffic pattern; and
allocating resources to web sites in the plurality of web sites based on the identified templates for each web site in the plurality of web sites.

30 14. The method of claim 13, wherein allocating resources to web sites includes:

identifying two or more web sites from the plurality of

web sites that are candidates for co-location based on the templates for the two or more web sites; and

allocating resources to the two or more web sites based on the identification of the two or more web sites being
5 candidates for co-location.

15. The method of claim 13, wherein allocating resources to web sites includes:

calculating, for each web site in the plurality of web sites, a first threshold based on the template for the web
10 site;

calculating a second threshold based on the first threshold, wherein offloading of traffic is enabled when a traffic intensity for the web site meets or exceeds the second threshold; and

15 calculating a third threshold based on the first threshold, wherein offloading of traffic is disabled when a traffic intensity of the web site meets or falls below the third threshold.

16. The method of claim 14, wherein allocating resources to
20 web sites includes:

calculating, for each web site in the plurality of web sites, a first threshold based on the template for the web site;

25 calculating a second threshold based on the first threshold, wherein offloading of traffic is enabled when a traffic intensity for the web site meets or exceeds the second threshold; and

30 calculating a third threshold based on the first threshold, wherein offloading of traffic is disabled when a traffic intensity of the web site meets or falls below the third threshold.

17. The method of claim 15, further comprising:

monitoring traffic of a web site of the plurality of web sites, on a web server to determine if the traffic exceeds the second threshold; and

5 offloading at least a portion of the traffic to another web server if the traffic of the web site exceeds the second threshold.

18. The method of claim 17, further comprising:

10 monitoring traffic of a web site of the plurality of web sites, on a web server to determine if the traffic falls below the third threshold; and

disabling the offloading of traffic to the another web server if the traffic of the web site falls below the third threshold.

19. The method of claim 14, wherein identifying two or more 15 web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites includes:

identifying a first web site having peaks in traffic at a first set of time periods; and

20 identifying a second web site having peaks in traffic at a second set of time periods different from the first set of time periods.

20. The method of claim 19, wherein the second web site has a 25 trough in traffic at approximately a same time as the first web site has a peak in traffic.

21. The method of claim 14, wherein identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites includes:

30 identifying the two or more web sites such that a capacity for traffic of a web server on which the two or more

web sites are located remains constant.

22. The method of claim 14, wherein identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites includes:

identifying the two or more web sites such that a the total cost of all web servers hosting the plurality of web sites is minimized.

23. The method of claim 14, wherein identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites includes:

formulating an integer programming problem having decision variables that are an assignment of the plurality of web sites to a plurality of web servers, a constraint that is a capacity limitation for any server at any time epoch, and an objective function that is an increasing function of minimum distances between server capacity and peak load at different ones of the plurality of web servers.

24. The method of claim 14, wherein identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites includes:

formulating an integer programming problem having decision variables that are an assignment of the plurality of web sites to a plurality of web servers, a constraint that is a capacity limitation for any server at any time epoch, and an objective function that is a sum of costs of all web servers to which at least one web site is assigned.

25. An apparatus for allocating resources to a plurality of web sites, comprising:

means for identifying a traffic pattern for each web site

in the plurality of web sites;

means for identifying a template for each web site in the plurality of web sites based on the traffic pattern; and

5 means for allocating resources to web sites in the plurality of web sites based on the identified templates for each web site in the plurality of web sites.

26. The apparatus of claim 25, wherein the means for allocating resources to web sites includes:

10 means for identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites; and

means for allocating resources to the two or more web sites based on the identification of the two or more web sites being candidates for co-location.

15 27. The apparatus of claim 25, wherein the means for allocating resources to web sites includes:

means for calculating, for each web site in the plurality of web sites, a first threshold based on the template for the web site;

20 means for calculating a second threshold based on the first threshold, wherein offloading of traffic is enabled when a traffic intensity for the web site meets or exceeds the second threshold; and

25 means for calculating a third threshold based on the first threshold, wherein offloading of traffic is disabled when a traffic intensity of the web site meets or falls below the third threshold.

28. The apparatus of claim 26, wherein the means for allocating resources to web sites includes:

30 means for calculating, for each web site in the plurality of web sites, a first threshold based on the template for the web site;

means for calculating a second threshold based on the first threshold, wherein offloading of traffic is enabled when a traffic intensity for the web site meets or exceeds the second threshold; and

5 means for calculating a third threshold based on the first threshold, wherein offloading of traffic is disabled when a traffic intensity of the web site meets or falls below the third threshold.

29. The apparatus of claim 27, further comprising:

10 means for monitoring traffic of a web site of the plurality of web sites, on a web server to determine if the traffic exceeds the second threshold; and

means for offloading at least a portion of the traffic to another web server if the traffic of the web site exceeds the
15 second threshold.

30. The apparatus of claim 29, further comprising:

means for monitoring traffic of a web site of the plurality of web sites, on a web server to determine if the traffic falls below the third threshold; and

20 means for disabling the offloading of traffic to the another web server if the traffic of the web site falls below the third threshold.

31. The apparatus of claim 26, wherein the means for identifying two or more web sites from the plurality of web
25 sites that are candidates for co-location based on the templates for the two or more web sites includes:

means for identifying a first web site having peaks in traffic at a first set of time periods; and

30 means for identifying a second web site having peaks in traffic at a second set of time periods different from the first set of time periods.

32. The apparatus of claim 31, wherein the second web site has a trough in traffic at approximately a same time as the first web site has a peak in traffic.

33. The apparatus of claim 26, wherein the means for
5 identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites includes:

means for identifying the two or more web sites such that a capacity for traffic of a web server on which the two or
10 more web sites are located remains constant.

34. The apparatus of claim 26, wherein the means for identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the templates for the two or more web sites includes:

15 means for identifying the two or more web sites such that a the total cost of all web servers hosting the plurality of web sites is minimized.

35. The apparatus of claim 26, wherein the means for identifying two or more web sites from the plurality of web
20 sites that are candidates for co-location based on the templates for the two or more web sites includes:

means for formulating an integer programming problem having decision variables that are an assignment of the plurality of web sites to a plurality of web servers, a
25 constraint that is a capacity limitation for any server at any time epoch, and an objective function that is an increasing function of minimum distances between server capacity and peak load at different ones of the plurality of web servers.

36. The apparatus of claim 26, wherein the means for
30 identifying two or more web sites from the plurality of web sites that are candidates for co-location based on the

templates for the two or more web sites includes:

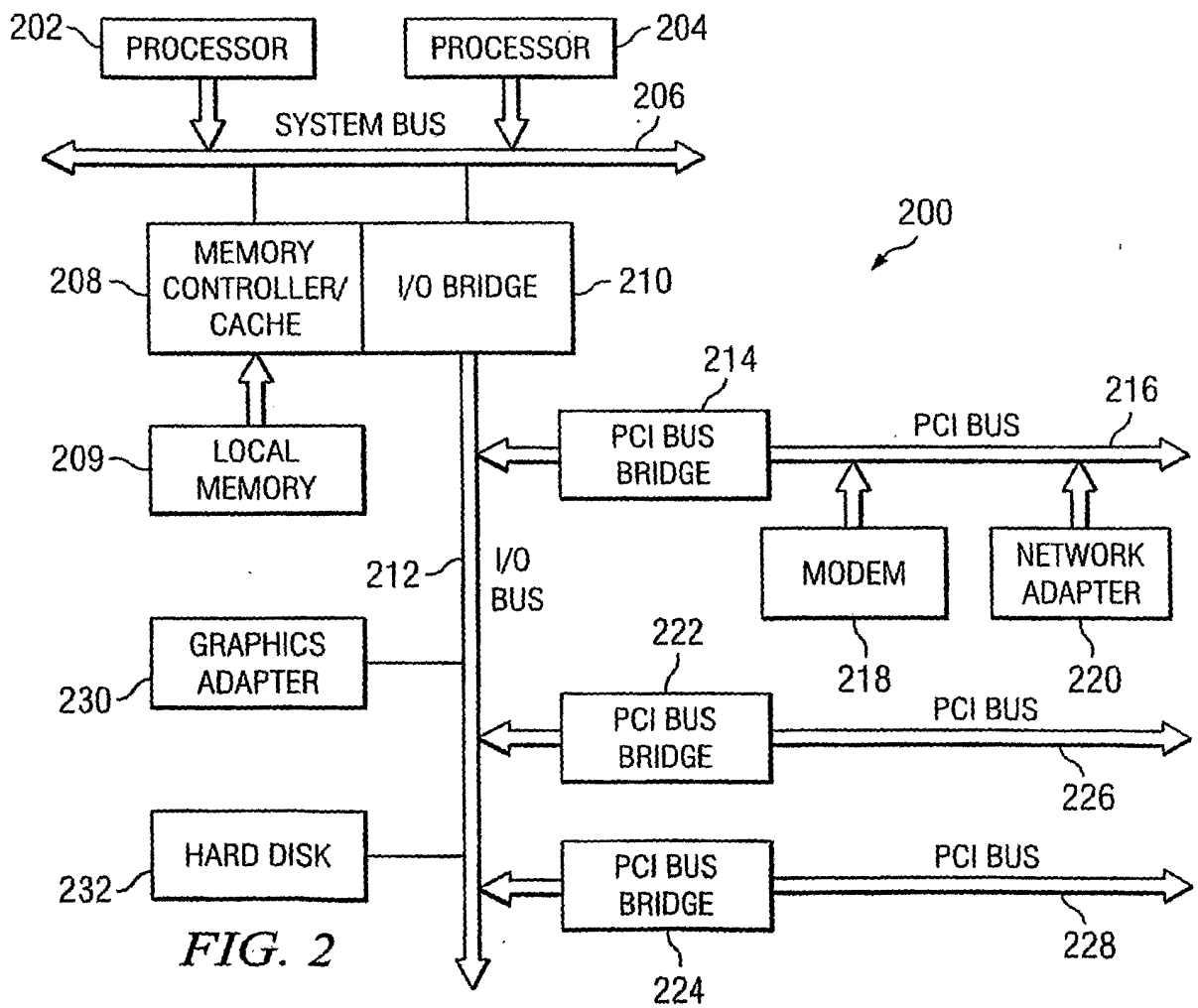
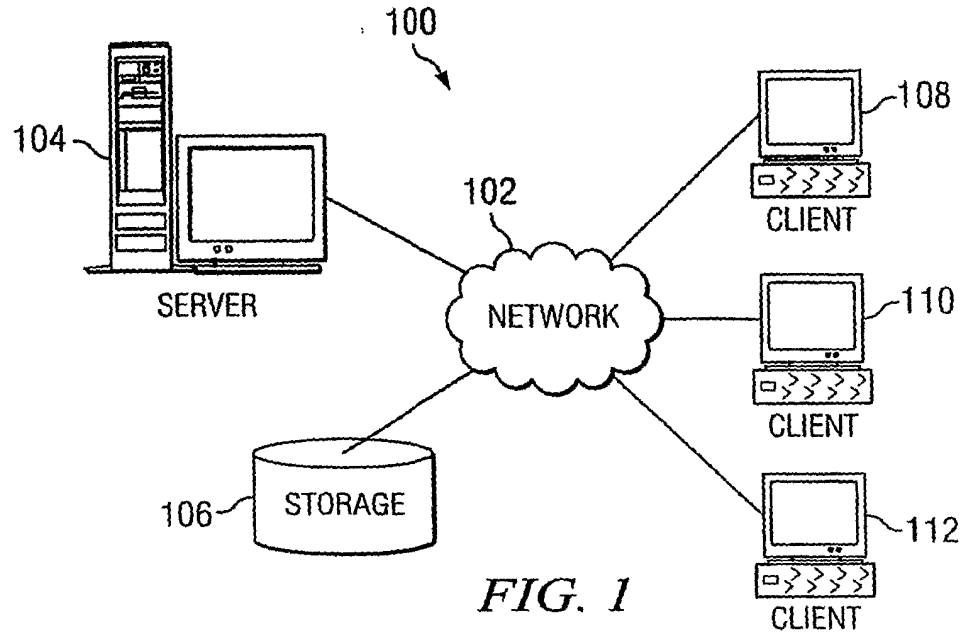
means for formulating an integer programming problem having decision variables that are an assignment of the plurality of web sites to a plurality of web servers, a
5 constraint that is a capacity limitation for any server at any time epoch, and an objective function that is a sum of costs of all web servers to which at least one web site is assigned.

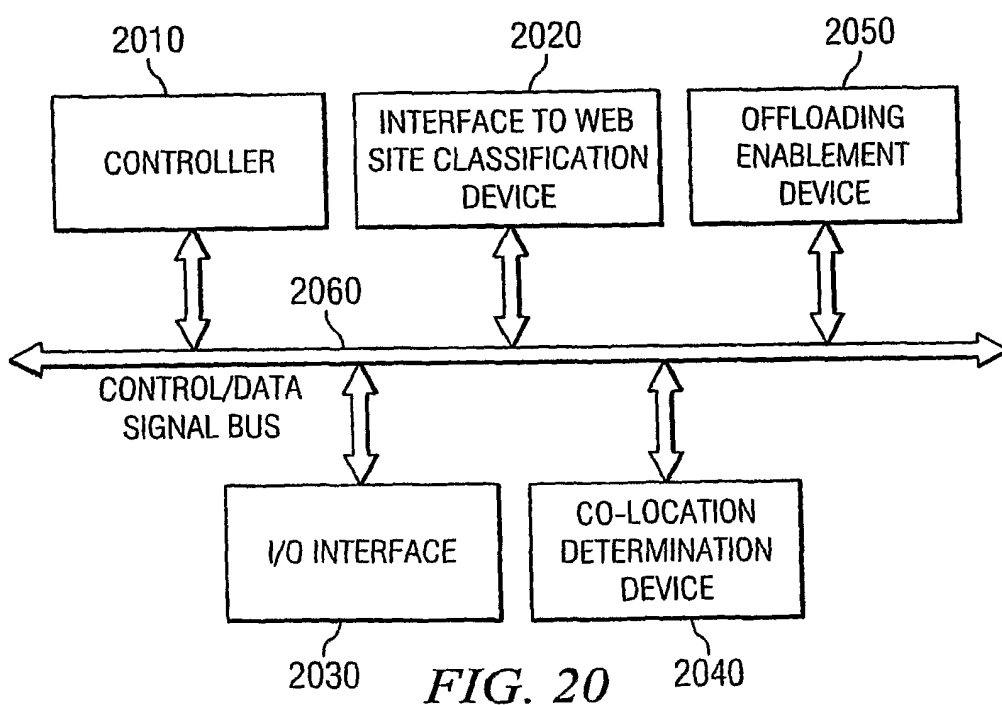
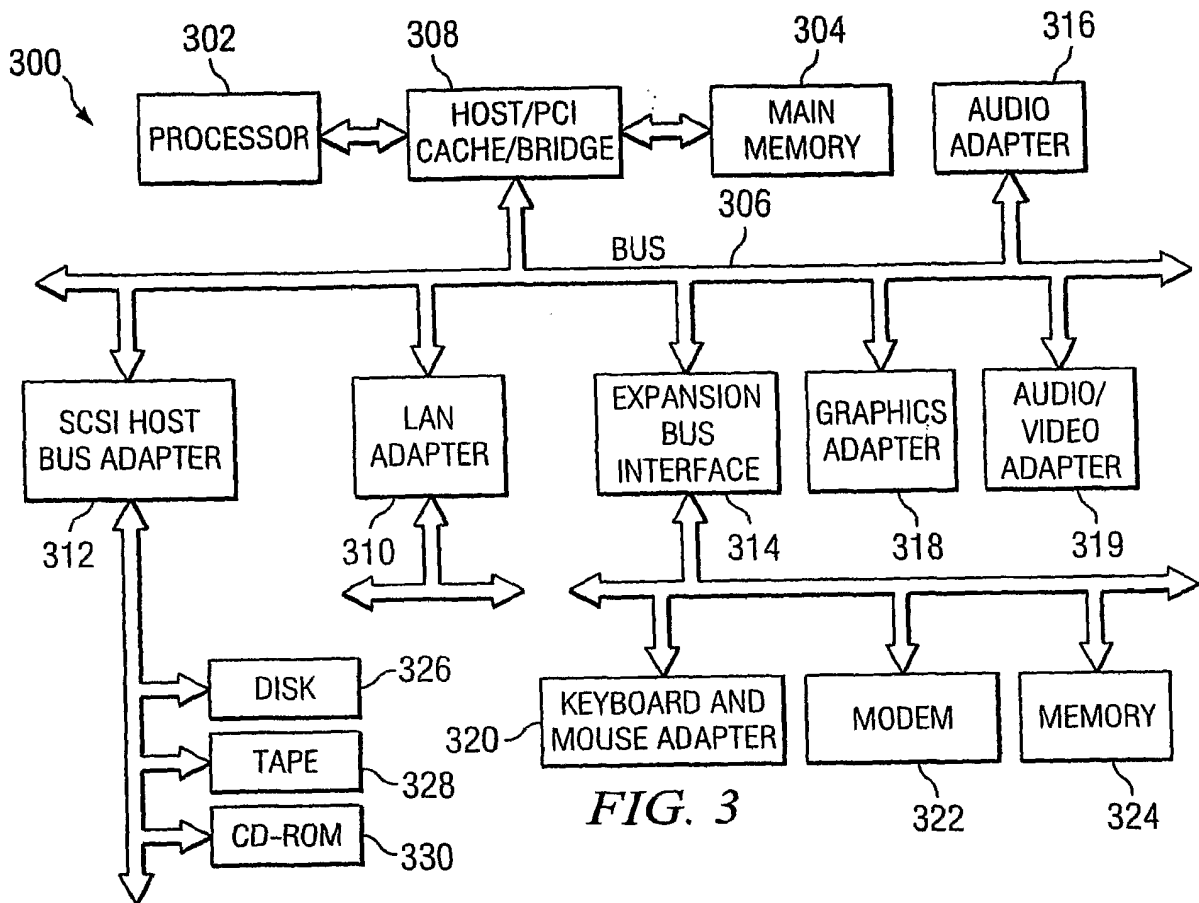
37. A method for deploying computing infrastructure, comprising integrating computer readable code into a computing
10 system, wherein the code in combination with the computing system is capable of performing the following:

identifying a traffic pattern for each web site in the plurality of web sites;

15 identifying a template for each web site in the plurality of web sites based on the traffic pattern; and

allocating resources to web sites in the plurality of web sites based on the identified templates for each web site in the plurality of web sites.





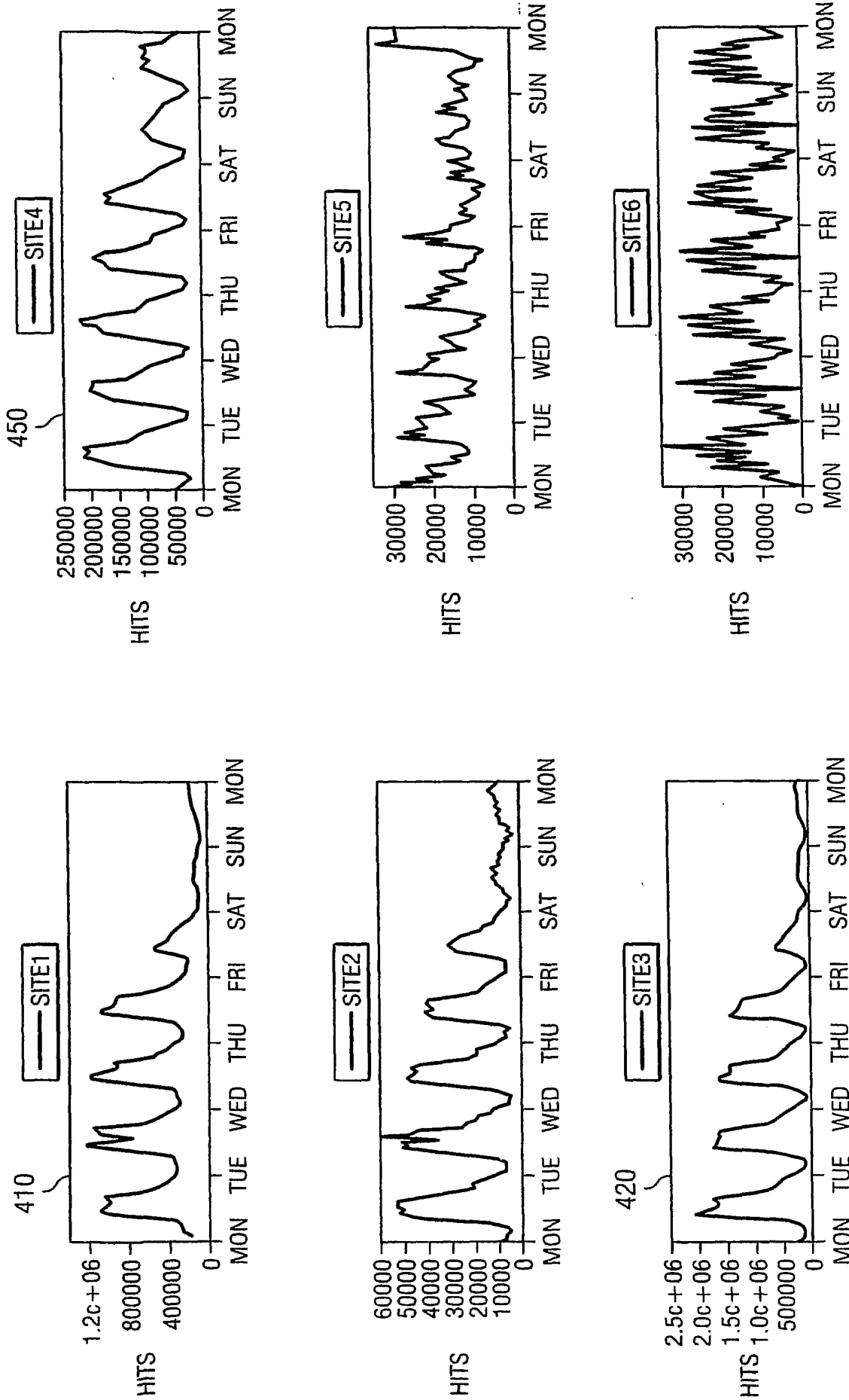


FIG. 4A

4/16

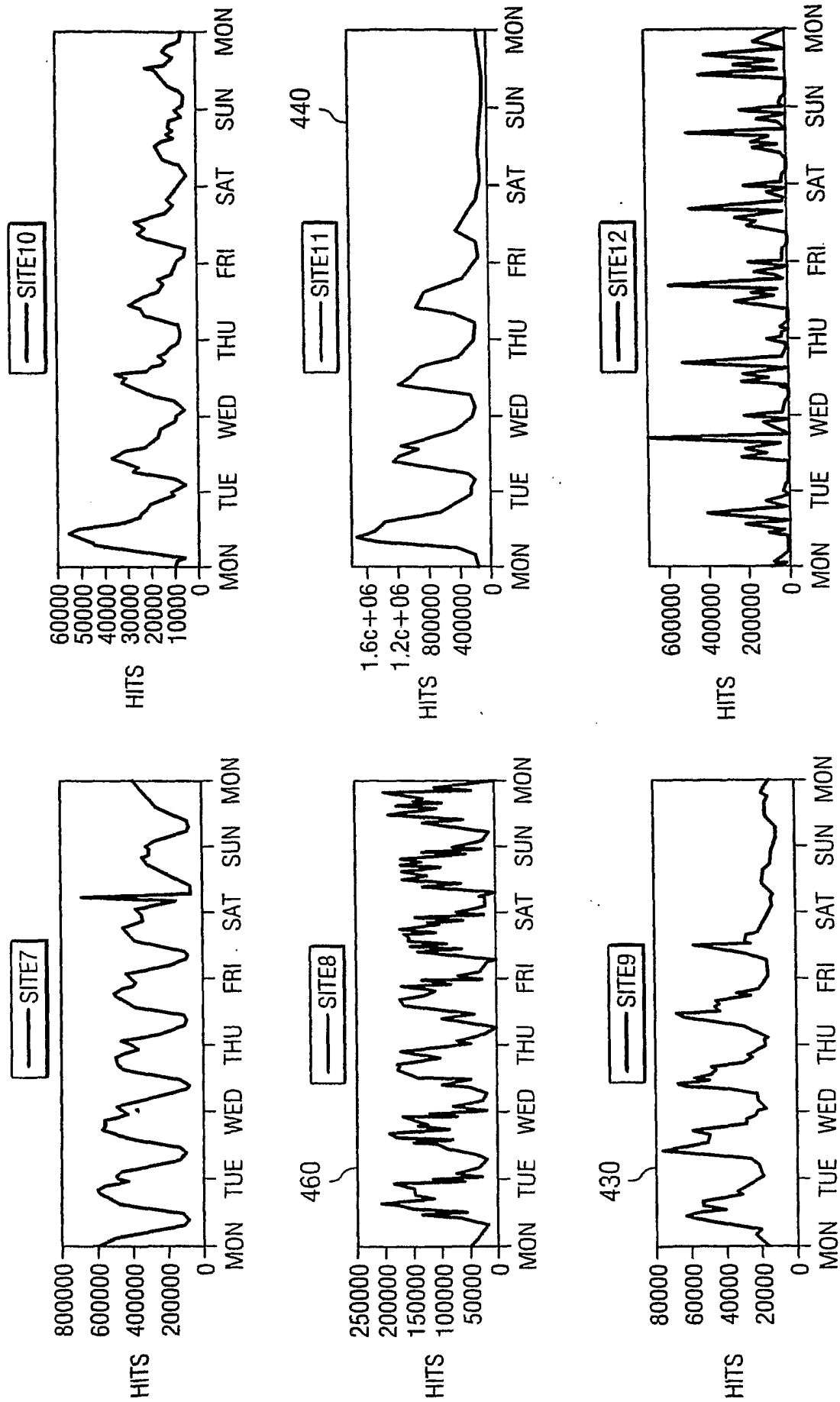


FIG. 4B

5/16

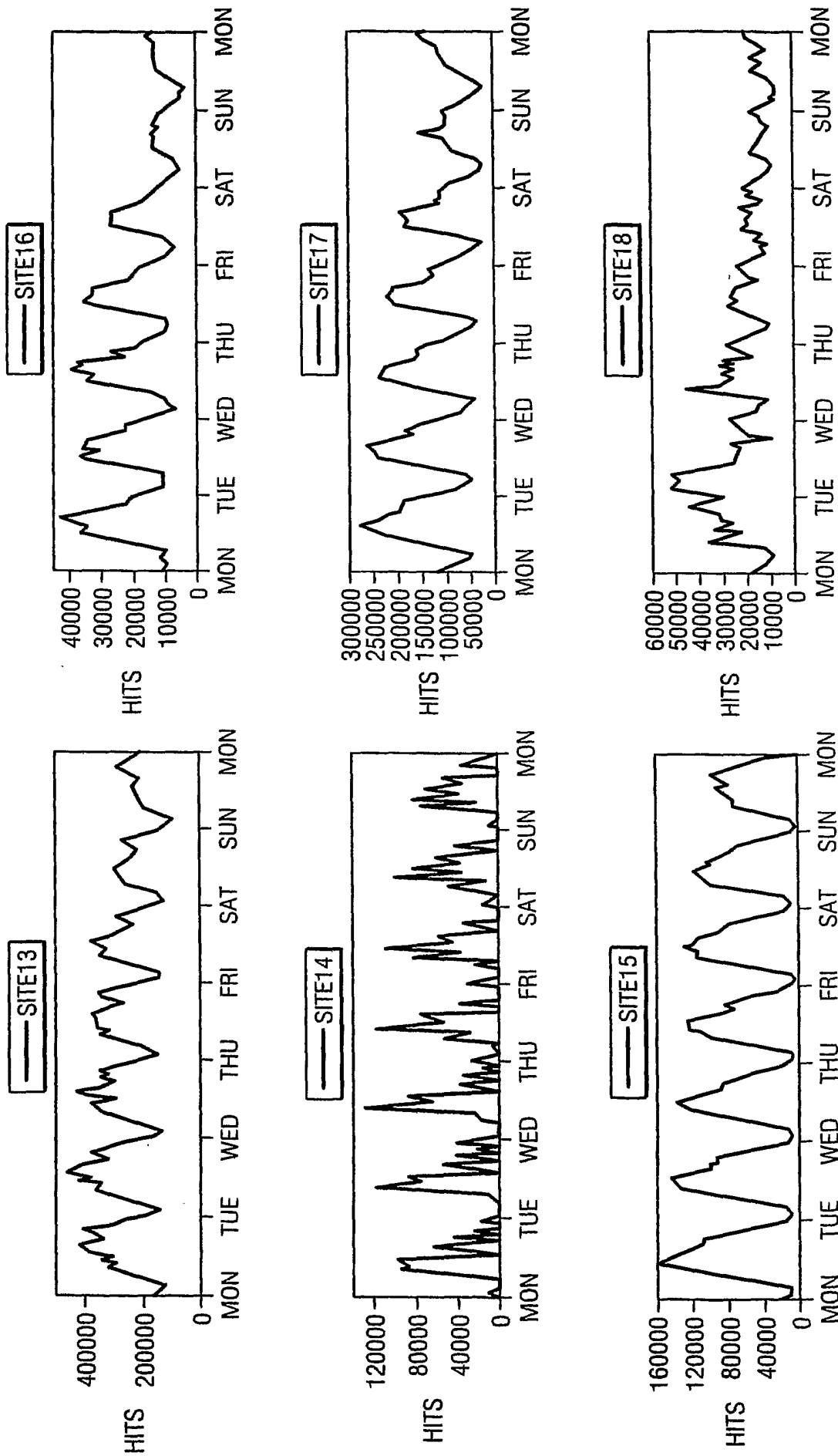
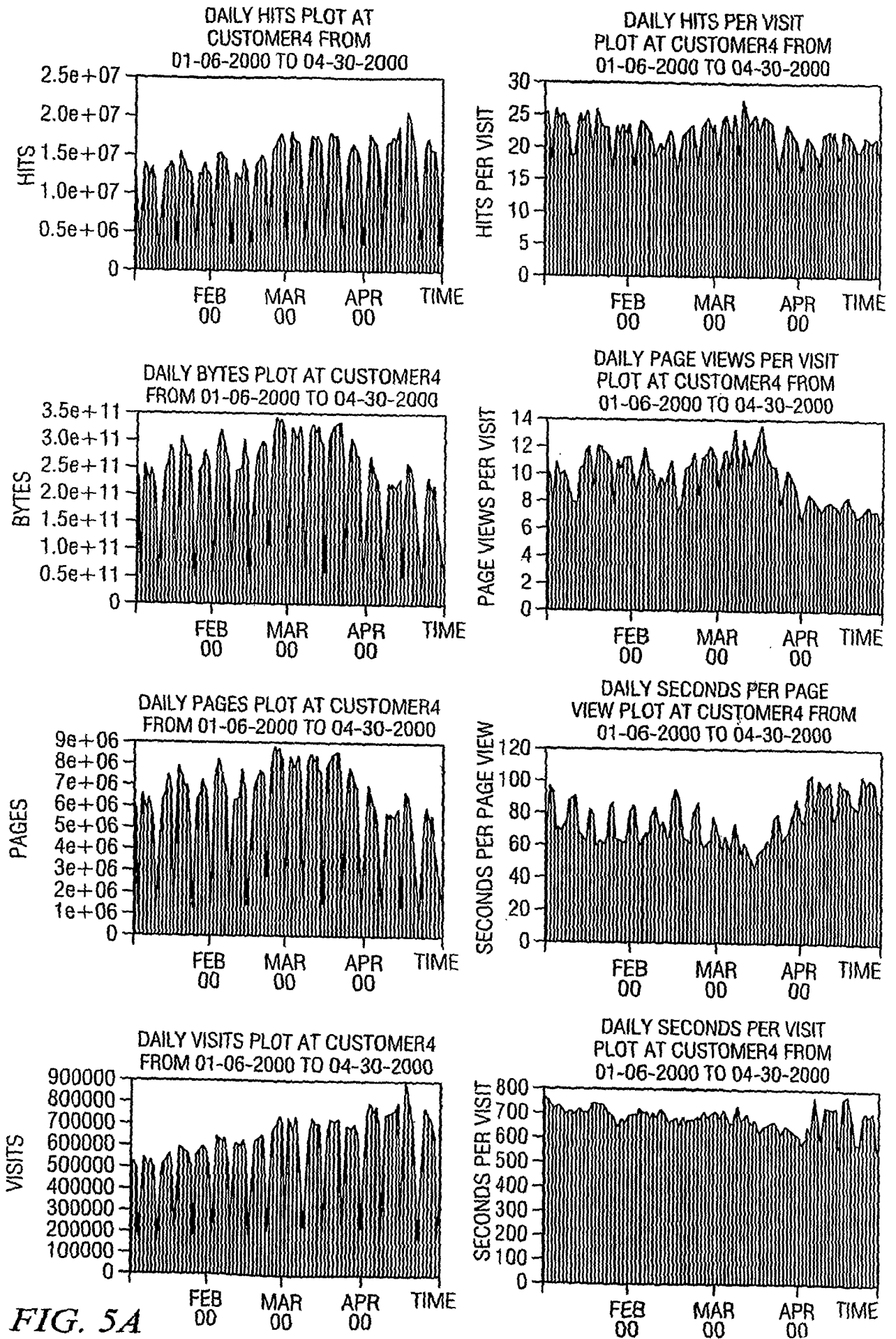


FIG. 4C



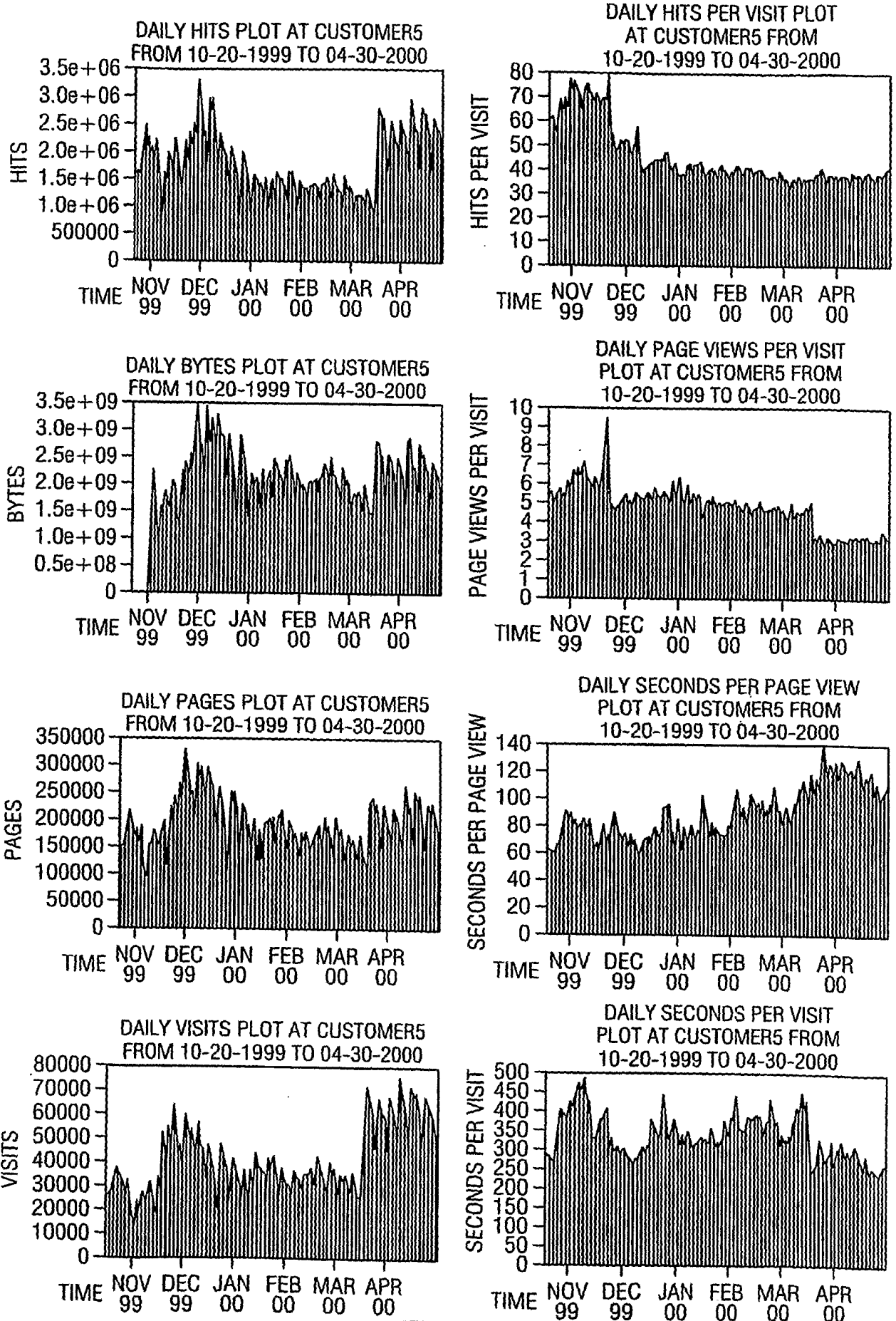


FIG. 5B

8/16

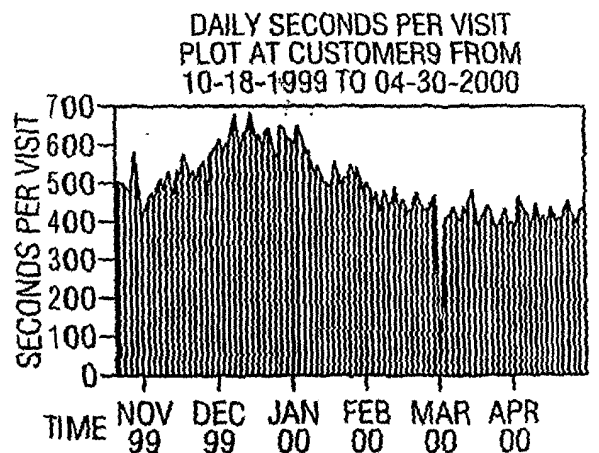
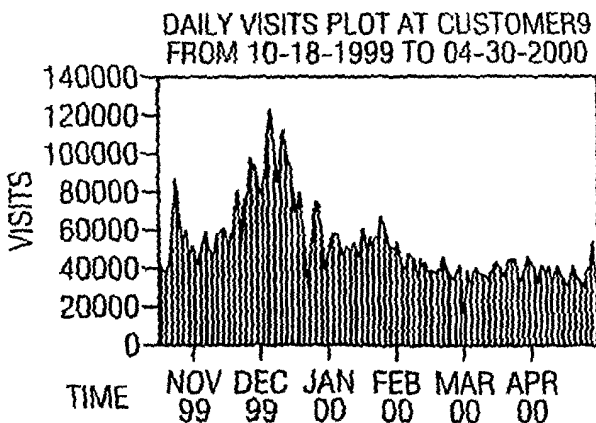
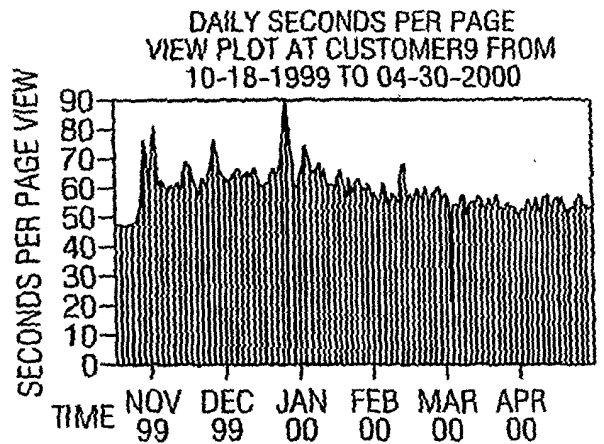
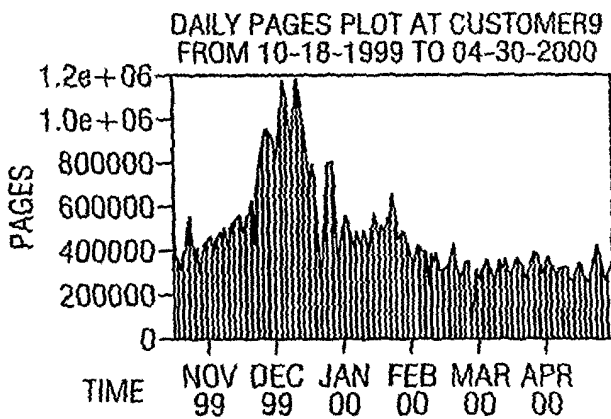
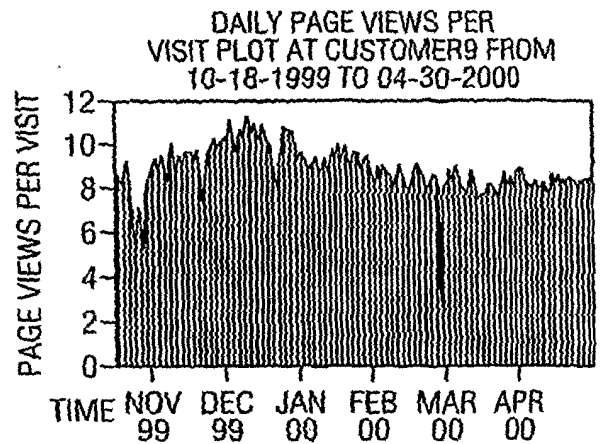
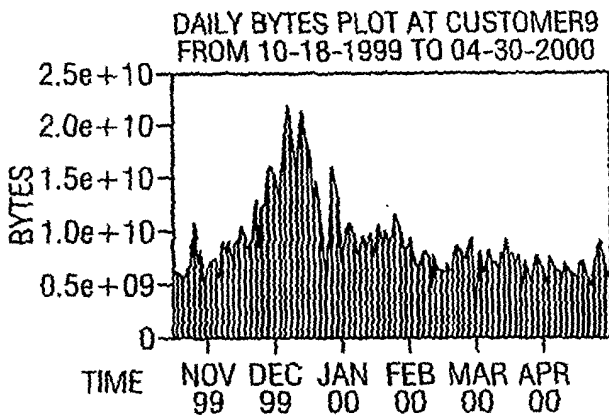
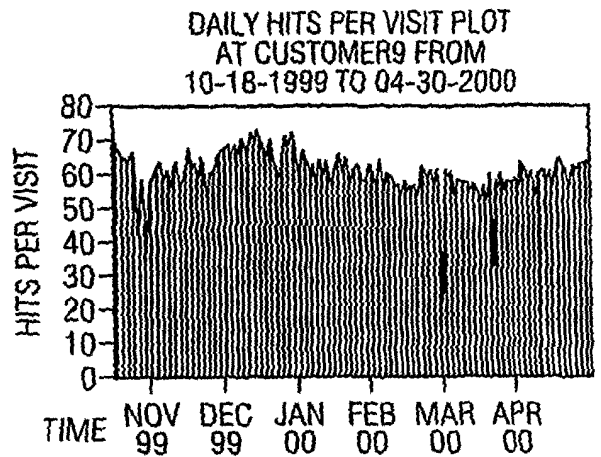
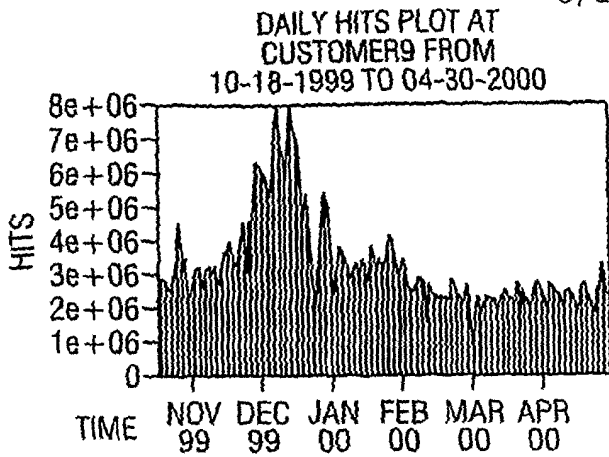


FIG. 5C

9/16

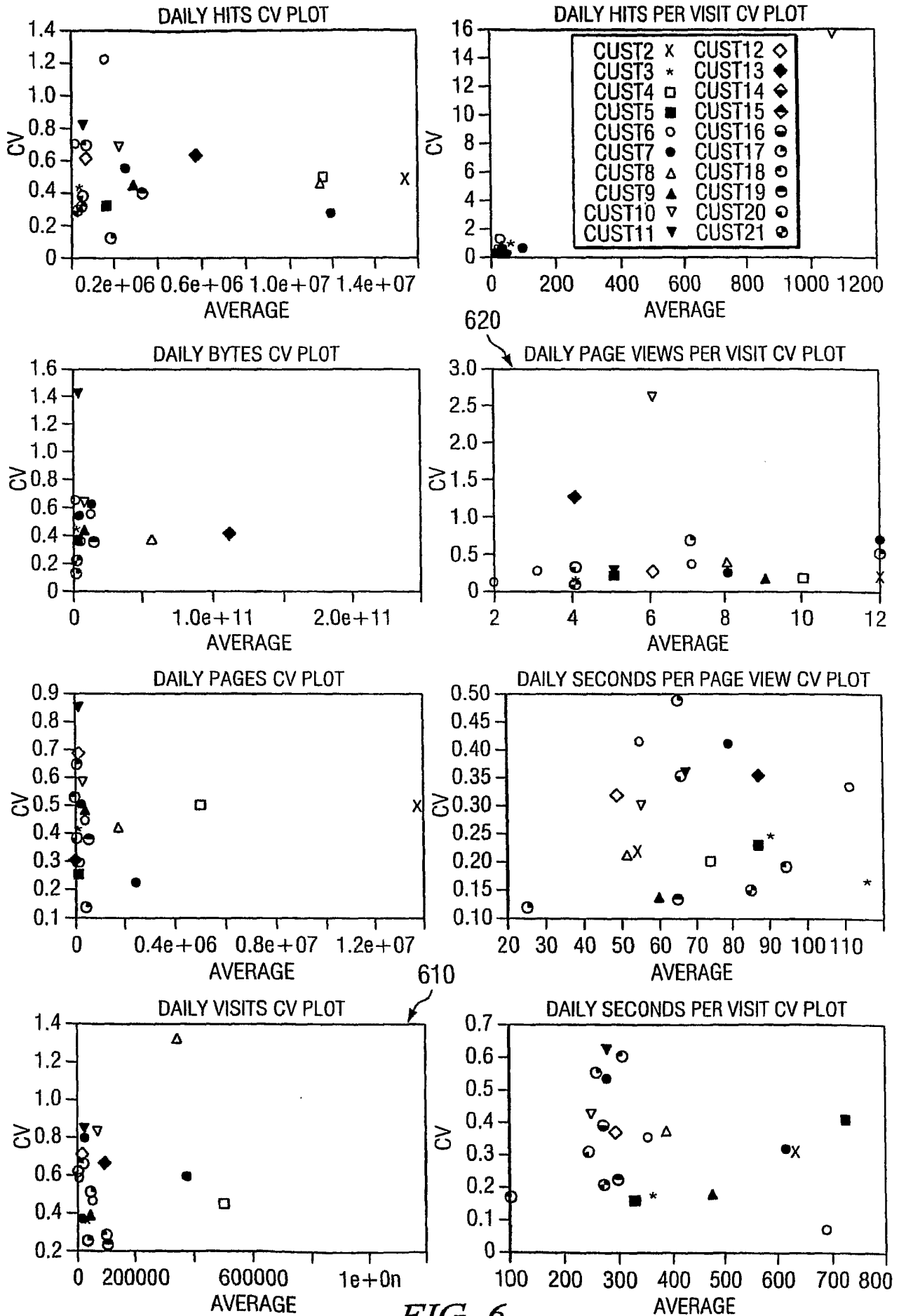


FIG. 6

10/16

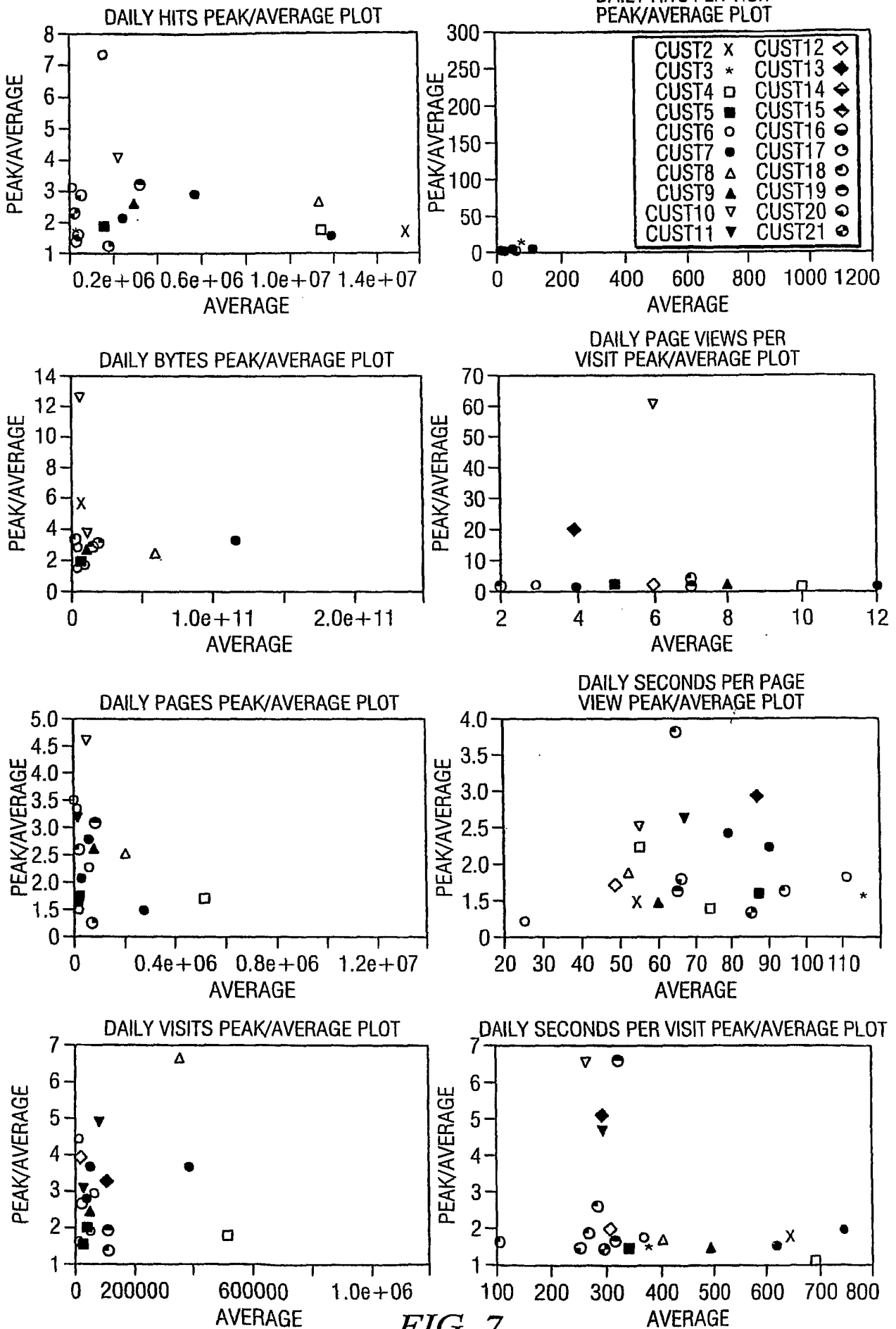


FIG. 7

11/16

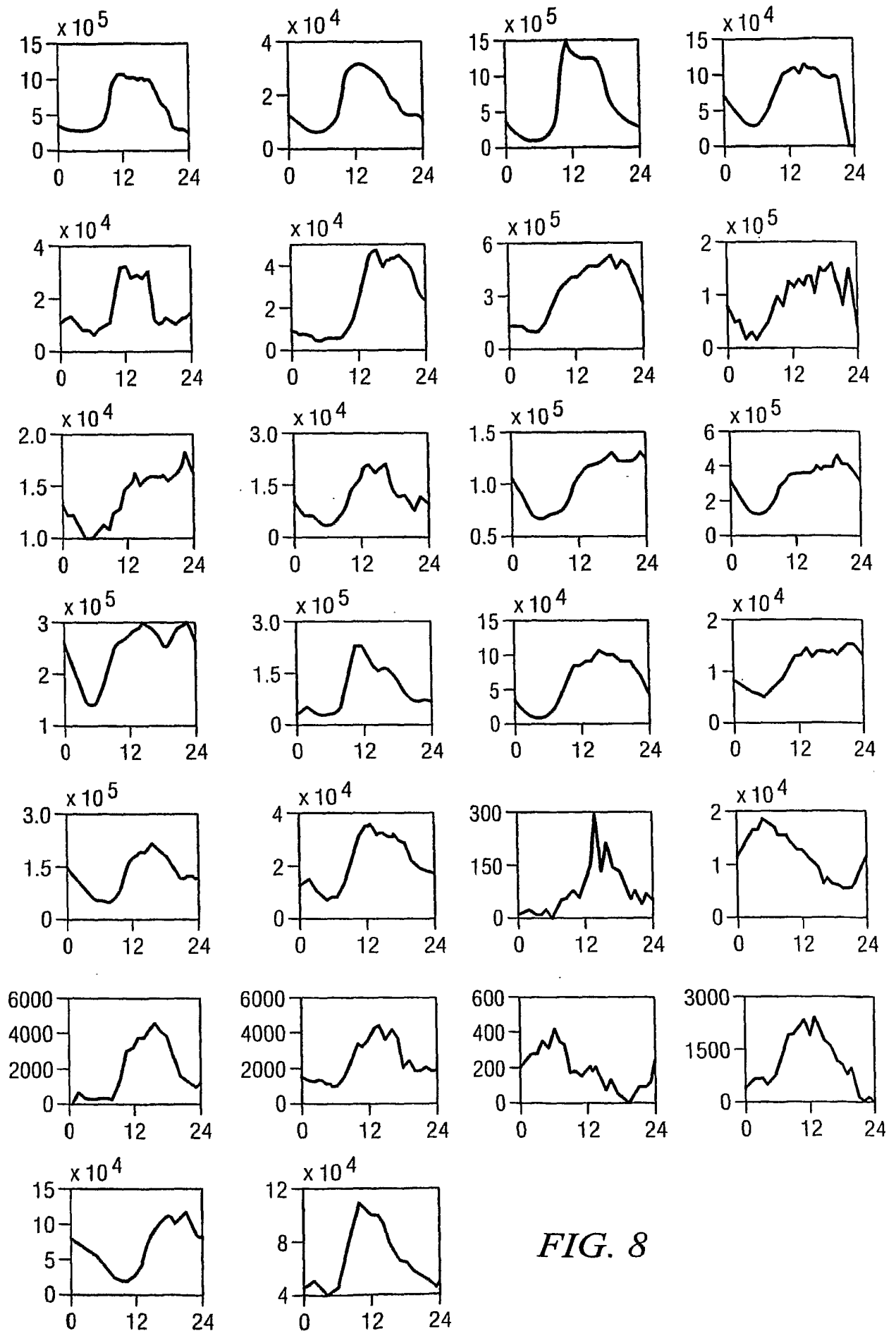


FIG. 8

12/16

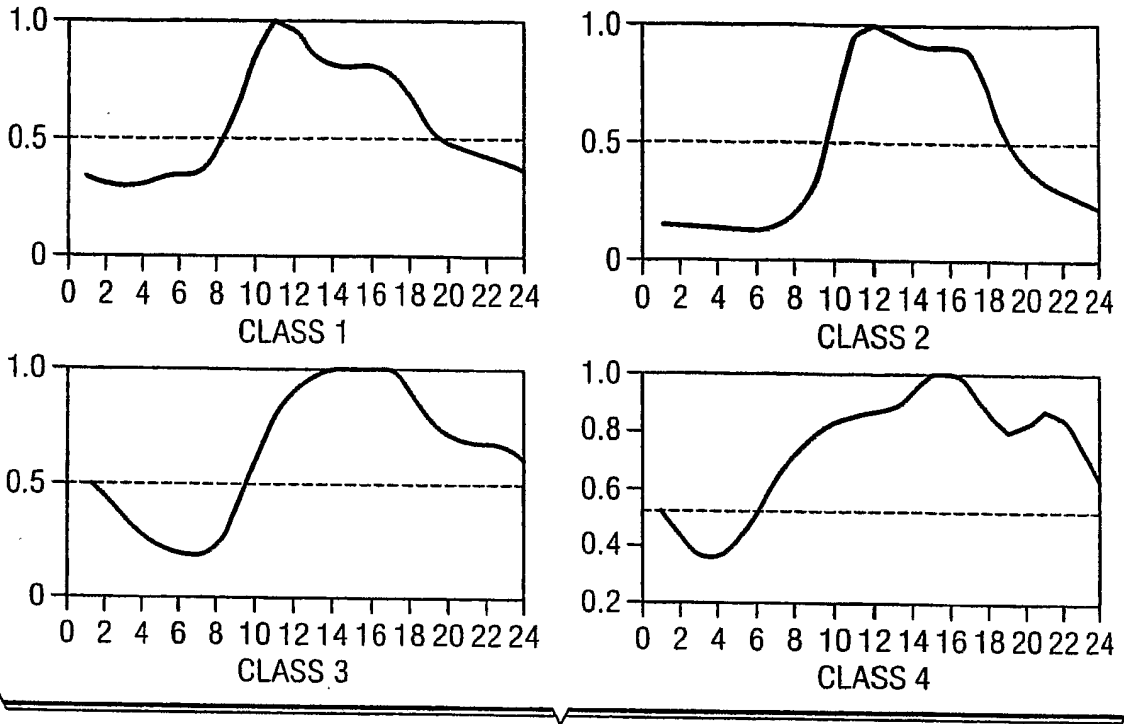


FIG. 9

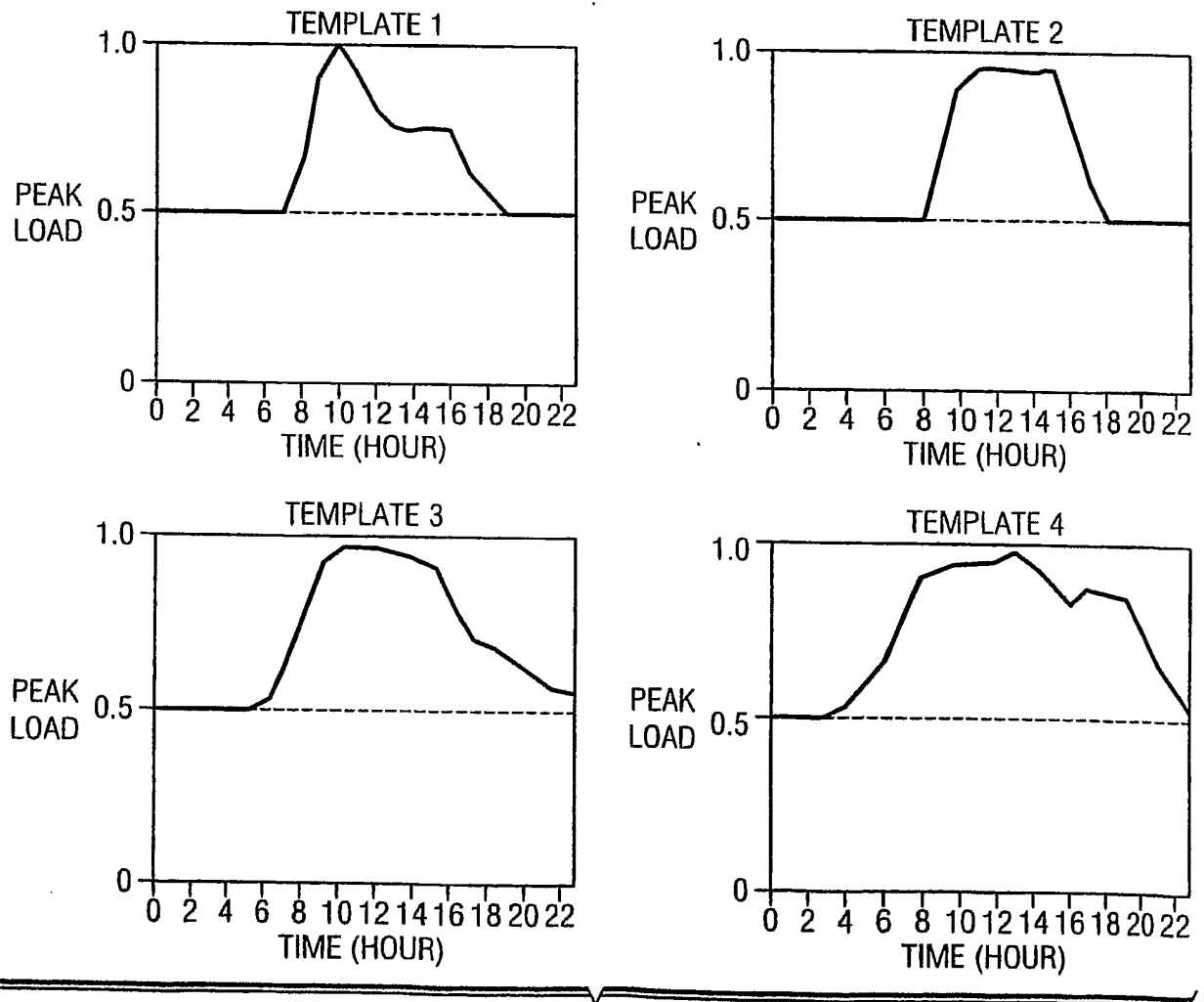
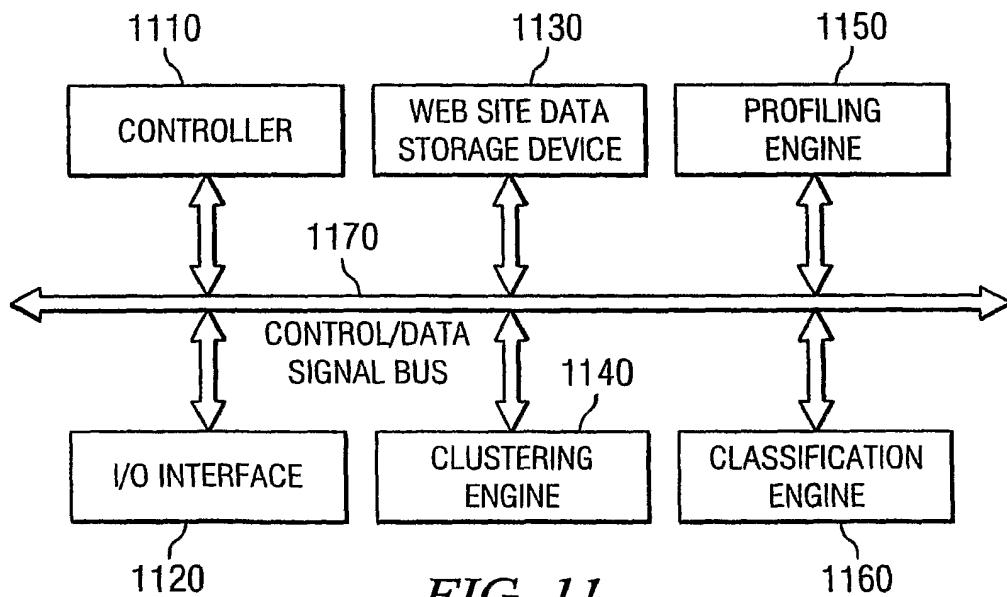
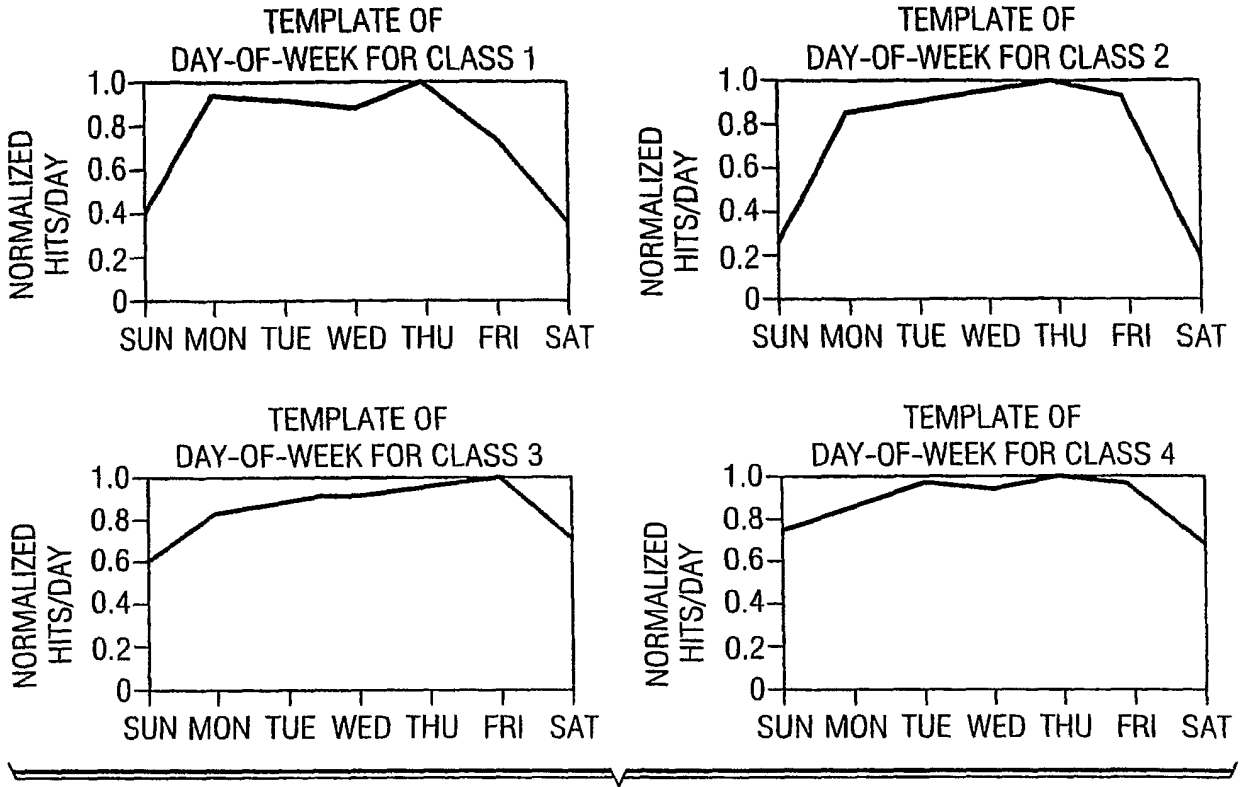
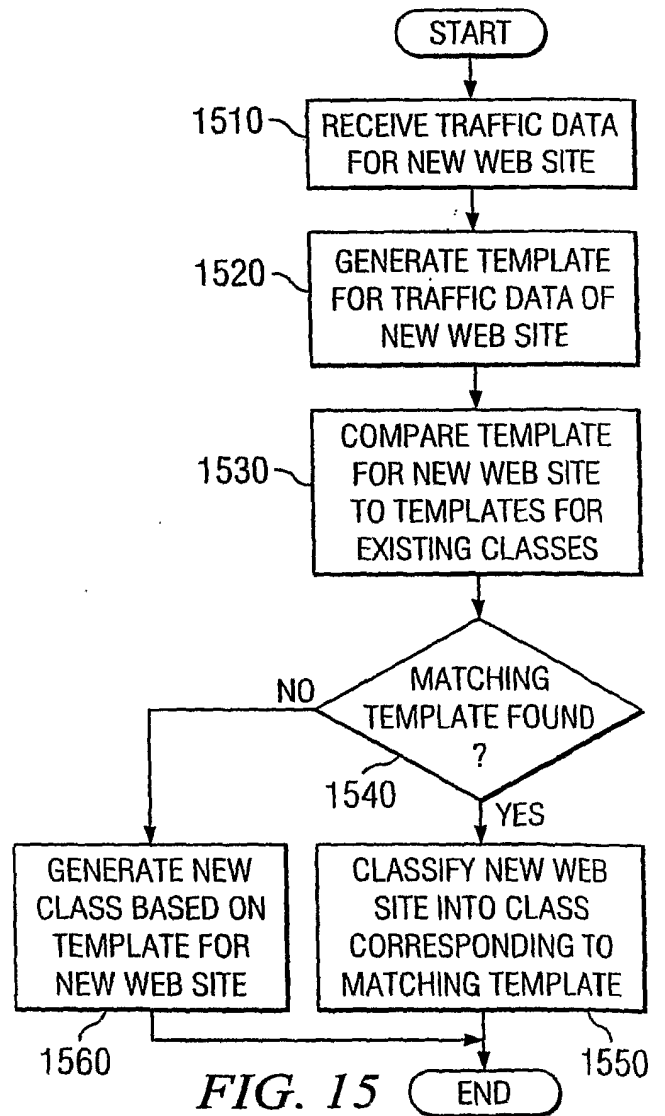
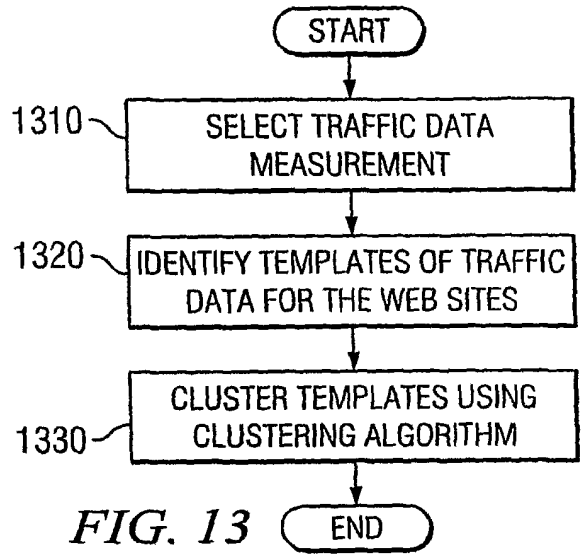
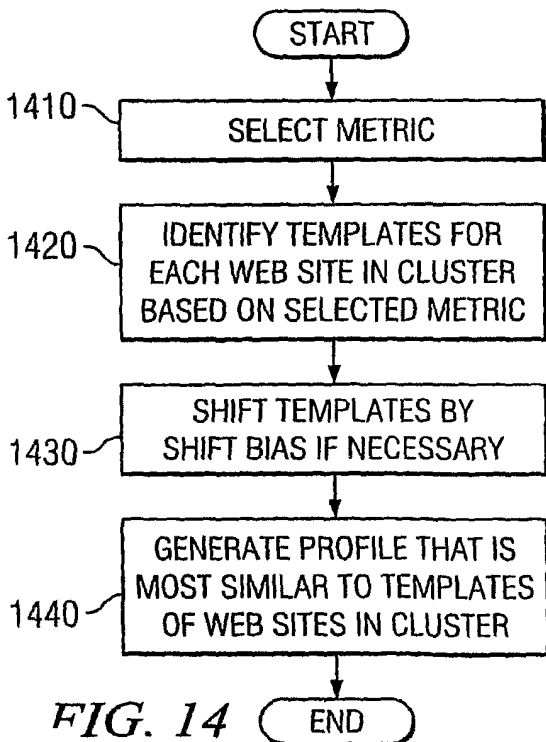
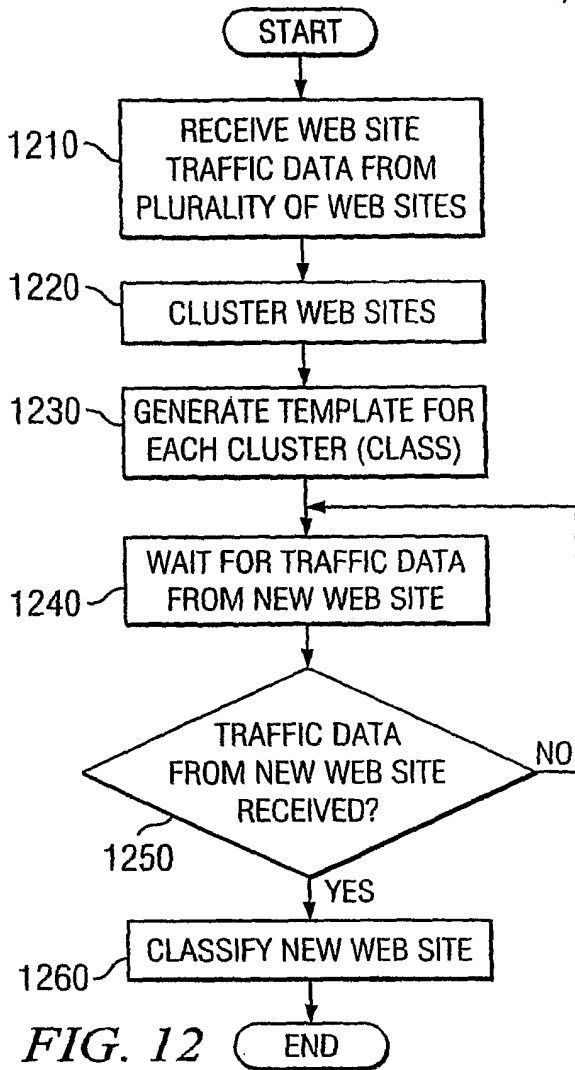


FIG. 10A





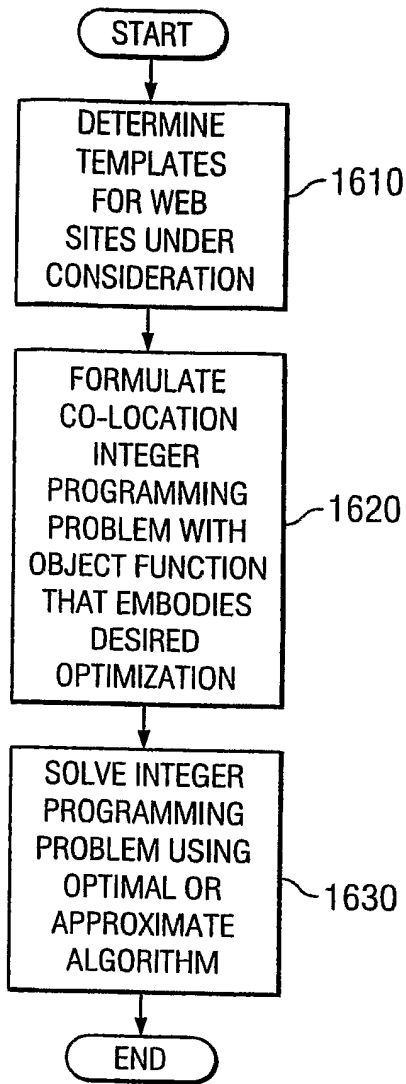


FIG. 16

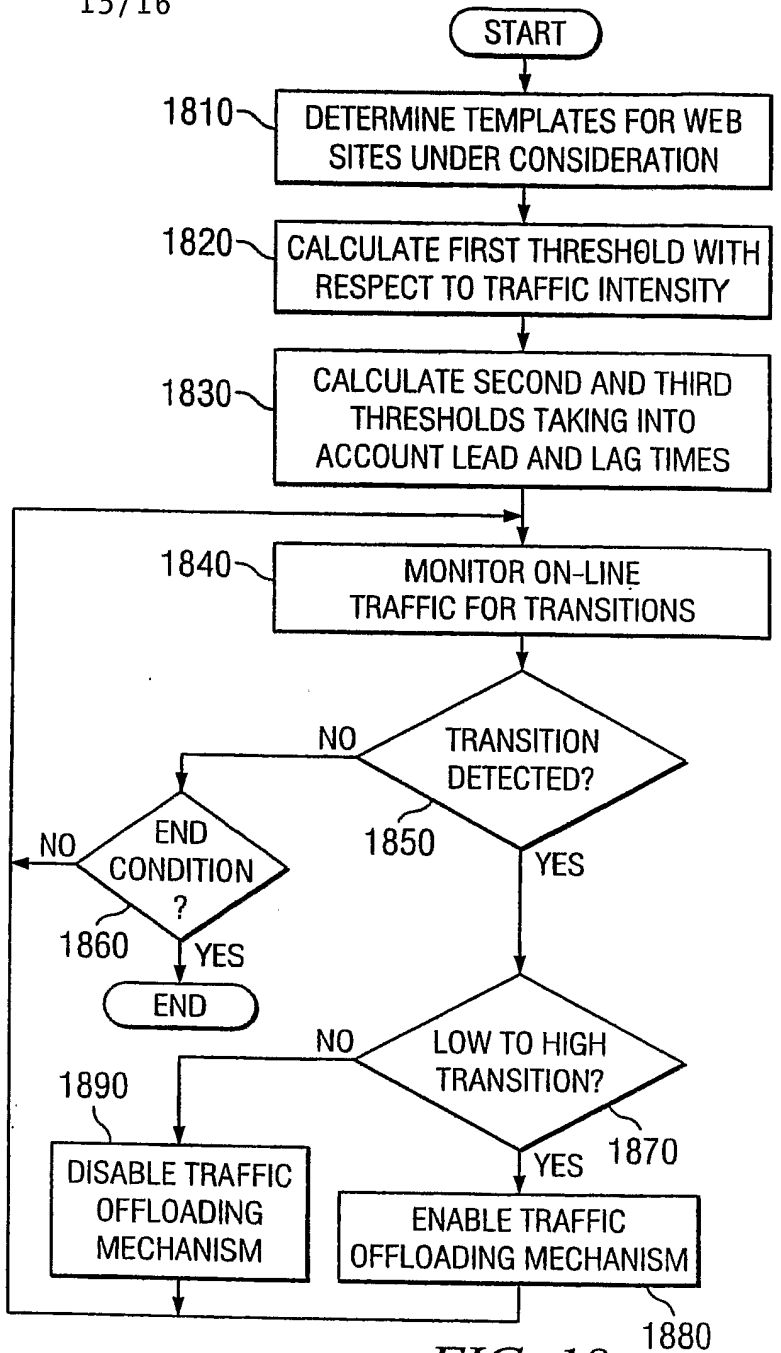


FIG. 18

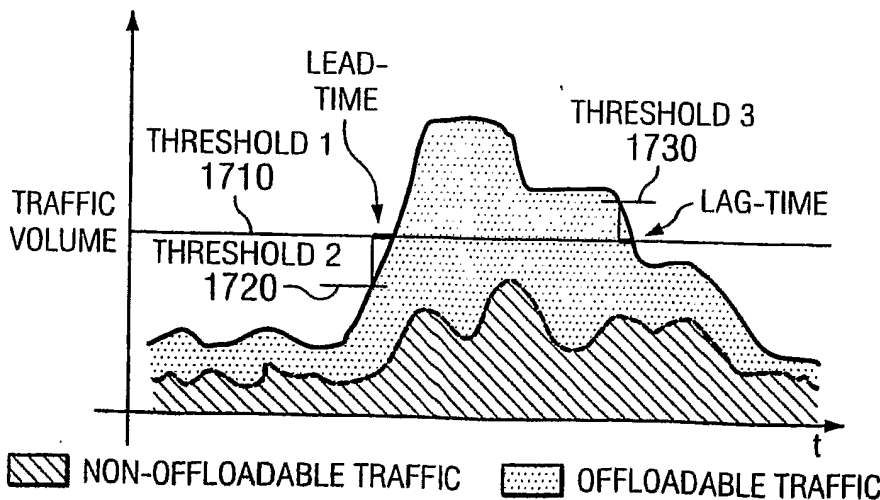


FIG. 17

16/16

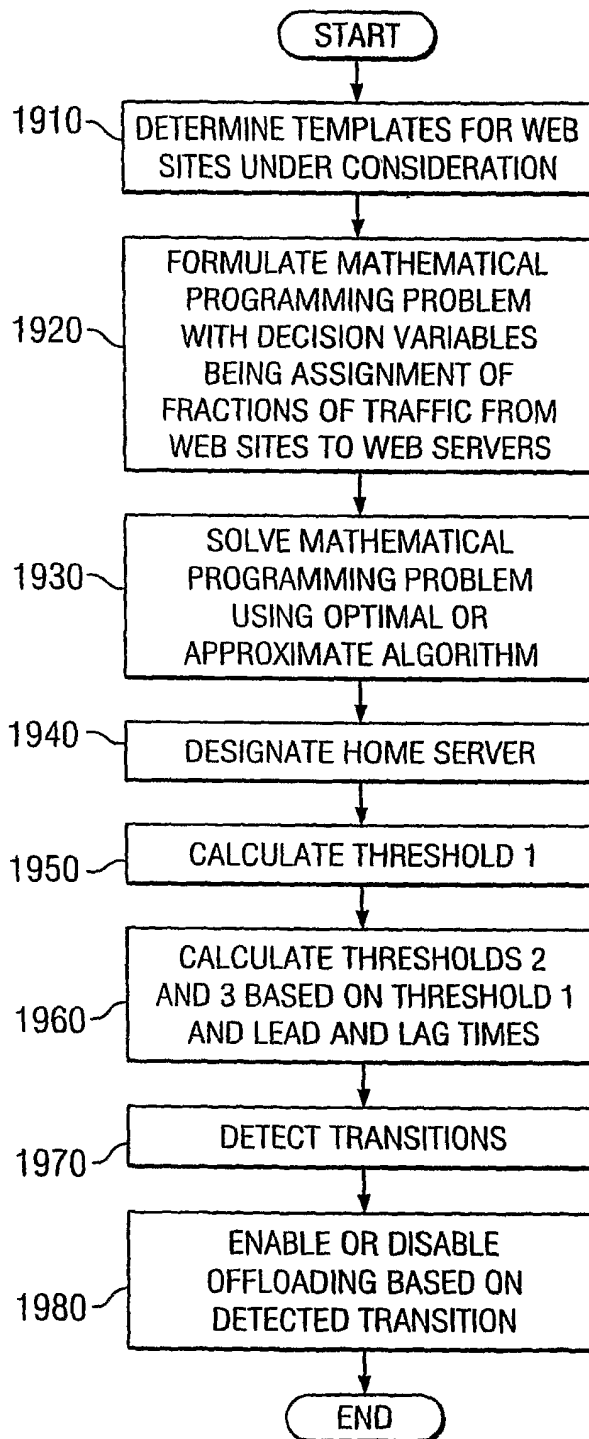


FIG. 19

INTERNATIONAL SEARCH REPORT

International Application No

PCT/EP 03/15016

A. CLASSIFICATION OF SUBJECT MATTER
 IPC 7 H04L29/06 H04L12/24

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
 IPC 7 H04L G06F

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

EPO-Internal, INSPEC

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category °	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 2002/143945 A1 (BALAJI KANNAN ET AL) 3 October 2002 (2002-10-03)	1,2, 7-14, 19-26, 31-37
Y	abstract paragraphs '0016!-'0034!, '0050!-'0055!, '0099!, '0100!; claims 1,8-10,12; figure 3 ----- -/--	3-6, 15-18, 27-30

Further documents are listed in the continuation of box C.

Patent family members are listed in annex.

° Special categories of cited documents :

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.

"&" document member of the same patent family

Date of the actual completion of the international search

23 March 2004

Date of mailing of the international search report

30/03/2004

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
 NL - 2280 HV Rijswijk
 Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
 Fax: (+31-70) 340-3016

Authorized officer

Goller, W

INTERNATIONAL SEARCH REPORT

International Application No
PCT/EP 03/15016

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT		
Category °	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
Y	<p>BAKER S M ET AL: "Distributed cooperative Web servers" COMPUTER NETWORKS, ELSEVIER SCIENCE PUBLISHERS B.V., AMSTERDAM, NL, vol. 31, no. 11-16, 17 May 1999 (1999-05-17), pages 1215-1229, XP004304550 ISSN: 1389-1286 abstract paragraphs '0001!, '0002!, '03.2!, '03.3!, '0004!, '04.1!, '04.2!</p>	<p>3-6, 15-18, 27-30</p>
A	<p>US 2002/112036 A1 (BOHANNON THOMAS A ET AL) 15 August 2002 (2002-08-15) abstract paragraphs '0034!-'0039!, '0119!, '0138!</p>	<p>3-6, 15-18, 27-30</p>
A	<p>US 6 374 297 B1 (WOLF JOEL L ET AL) 16 April 2002 (2002-04-16) abstract column 1, line 45 -column 2, line 60 column 3, line 22 -column 8, line 36; claim 3</p>	<p>3-6, 15-18, 27-30</p>
A	<p>US 6 351 775 B1 (YU PHILIP SHI-LUNG) 26 February 2002 (2002-02-26) abstract column 9, line 50 -column 10, line 42</p>	<p>3-6, 15-18, 27-30</p>

INTERNATIONAL SEARCH REPORT

Information on patent family members

International Application No

PCT/EP 03/15016

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 2002143945 A1	03-10-2002	JP 2002318791 A	31-10-2002
US 2002112036 A1	15-08-2002	AU 3649702 A CA 2430416 A1 EP 1344118 A2 WO 0244848 A2	11-06-2002 06-06-2002 17-09-2003 06-06-2002
US 6374297 B1	16-04-2002	NONE	
US 6351775 B1	26-02-2002	CN 1202772 A ,B JP 2970760 B2 JP 11004261 A TW 451117 B	23-12-1998 02-11-1999 06-01-1999 21-08-2001