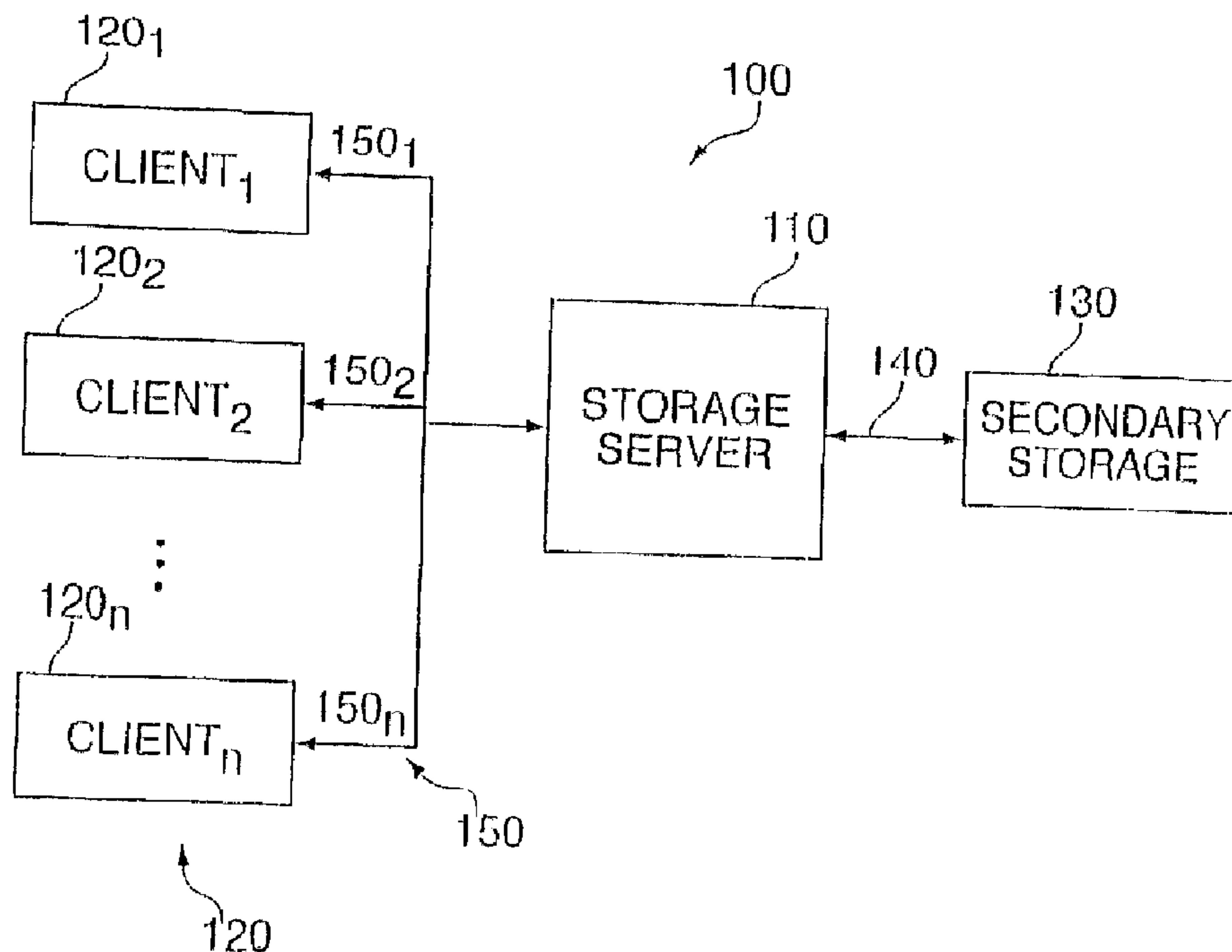




(22) Date de dépôt/Filing Date: 2002/05/09
(41) Mise à la disp. pub./Open to Public Insp.: 2002/11/21
(62) Demande originale/Original Application: 2 701 621
(30) Priorité/Priority: 2001/05/14 (US09/854,839)

(51) Cl.Int./Int.Cl. *H04N 21/218* (2011.01),
H04N 21/222 (2011.01), *H04N 21/235* (2011.01),
H04N 21/6377 (2011.01)
(71) Demandeur/Applicant:
COMCAST IP HOLDINGS I, LLC, A DELAWARE
LIMITED LIABILITY COMPANY, US
(72) Inventeurs/Inventors:
TAYLOR, CLEMENT G., US;
CHIN, DANNY, US;
LERMAN, JESSE S., US;
GOODE, CHRISTOPHER W. B., US
(74) Agent: RICHES, MCKENZIE & HERBERT LLP

(54) Titre : ARCHITECTURE DE SERVEUR DE STOCKAGE MODULAIRE A GESTION DE DONNEES DYNAMIQUE
(54) Title: MODULAR STORAGE SERVER ARCHITECTURE WITH DYNAMIC DATA MANAGEMENT



(57) **Abrégé/Abstract:**

The architecture of the present invention employs dynamic data management methods, which determine whether data should reside on disk or secondary storage (130), on which disk drives data should be stored, and how data should be replicated and/or migrated to new disk drives based on observed user (120) access patterns. These methods also migrate popular data to faster disk tracks to reduce average access time and thus improve performance. User access requests are assigned to modules based on the data stored at each module, and each module's current load (the number of requests waiting to be serviced). If the requested data is not on a disk drive, the data is retrieved from secondary storage (130), and may be stored on the disk drives for rapid subsequent access. When a requested data item on the disk drive is replicated, load balancing is performed by assigning the request to the module holding the data with the lowest load. In addition, user access requests waiting to retrieve replicated data may be dynamically and seamlessly migrated to another module based on changes in module loads.



ABSTRACT

The architecture of the present invention employs dynamic data management methods, which determine whether data should
5 reside on disk or secondary storage (130), on which disk drives data should be stored, and how data should be replicated and/or migrated to new disk drives based on observed user (120) access patterns. These methods also migrate popular data to faster disk tracks to reduce average access time and thus improve
10 performance. User access requests are assigned to modules based on the data stored at each module, and each module's current load (the number of request: waiting to be serviced). If the requested data is not on a disk driven, the data is retrieved from secondary storage (130), and may be stored on the disk
15 drives for rapid subsequent access. When a requested data item on the disk drive is replicated, load balancing is performed by assigning the request to the module holding the data with the lowest load. In addition, user access requests waiting to retrieve replicated data may be dynamically and seamlessly
20 migrated to another module based on changes in module loads.

**MODULAR STORAGE SERVER ARCHITECTURE WITH
DYNAMIC DATA MANAGEMENT**

This application is a divisional of Canadian patent
5 application serial number 2,701,621, which is a divisional of
Canadian patent application serial number 2,446,602 which is
the national phase of International application PCT/2002/14871
filed May 9, 2002 (09.05.2002) and published 21 November 2002
(21.11.2002) under International Publication No. WO 02/093298
10 A3.

FIELD OF THE INVENTION

The present invention generally relates to a modular
storage server architecture for retrieving data in response to
15 user access requests. In particular, the invention relates to a
server architecture in which data is dynamically distributed
over a plurality of disks, a plurality of processors are
assigned to particular disks, and data access requests are
assigned to particular processors in order to provide good data
20 access performance and server fault tolerance.

BACKGROUND OF THE DISCLOSURE

A storage server allows users to efficiently retrieve
information from large volumes of data stored on a plurality
25 of disks and secondary storage (e.g., magnetic tape or an
optical-disk jukebox). For example, a video server is a
storage server that accepts user requests to view a particular
movie from a video library, retrieves the requested program
from disk, and delivers the program to the appropriate
30 user(s). Such a video server is disclosed in U.S. patent
5,671,377, entitled "System For Supplying 30 Streams Of Data
To Multiple Users By Distributing A Data Stream To Multiple
Processors And Enabling Each User To Manipulate Supplied Data

-1a-

Stream" issued to Bleidt et al. on September 23, 1997.

The foregoing storage server employs one or more 35
processors that access data that is stored across an array
of disk drives using fault tolerant storage technique such
5 as RAID (Redundant Array of Inexpensive Disks). While such
architectures provide uniform non-blocking access to all of
the data stored on the disk drives, they do not facilitate a

-2-

modular architecture. Since data is striped across all of the disk drives in the array, adding or removing disk drives to/from the server requires that all of the data be re-striped across the new set of disk drives. Because the
5 servers are not modular, it is therefore inconvenient to increase or decrease storage capacity by adding or removing disk drives.

There is therefore a need in the art for a storage server architecture that is modular and can acceptably
10 resolve content blocking issues.

SUMMARY OF THE INVENTION

The disadvantages associated with the prior art are overcome by the present invention of a server comprising a
15 plurality of modules, each of which contains a single processor and a cluster of, for example, 16 disk drives, and a host controller that communicates with and assigns data requests to each of the modules. Data is written to the disk drives by striping the data across the 16-disk drive
20 cluster of a single module according to a RAID-5 protocol, with parity and spares distributed amongst the disk drives in the cluster.

The architecture of the present invention employs dynamic data management methods, which determine whether
25 data should reside on disk or secondary storage, on which disk drives data should be stored, and how data should be replicated and/or migrated to new disk drives based on observed user access patterns. These methods also migrate popular data to faster disk tracks to reduce average access
30 time and thus improve performance.

User access requests are assigned to modules based on the data stored at each module, and each module's current load (the number of requests waiting to be serviced). If the requested data is not on a disk drive, the data is
35 retrieved from secondary storage, and may be stored on the disk drives for rapid subsequent access. When a requested data item on the disk drive is replicated, load balancing is performed by assigning the request to the module holding the

data with the lowest load. In addition, user access requests waiting to retrieve replicated data may be dynamically and seamlessly migrated to another module based on changes in module loads.

5 Accordingly, in one aspect, the present invention provides a computer implemented method, comprising: receiving a request from at least one client for a title not resident in a storage server, said title comprising a play track, said play track comprising a plurality of chapters; initiating the
10 retrieval from a secondary storage device of play track portions proximate chapter delineation points; and in the case of a client request to begin presentation of said title at one of said chapters, initiating the streaming to said client of retrieved portions of said play track chapter, and initiating
15 the retrieval from said secondary storage device of at least unretrieved portions of said play track chapter and subsequent play track portions.

 In a further aspect, the present invention provides a computer implemented method, comprising: receiving a request
20 from at least one client for a title not resident in a storage server, said title comprising a play track, said play track comprising a plurality of chapters; initiating the retrieval from a secondary storage device of play track portions proximate chapter delineation points; determining bandwidth
25 capacity and quality-of-service (QoS) parameters associated with said secondary storage device; and in the case of a client request to begin presentation of said title at one of said chapters, initiating the streaming to said client of retrieved portions of said play track chapter, masking latency
30 associated with said secondary storage device, and initiating retrieval of at least unretrieved portions of said play track

-3a-

chapter and subsequent play track portions from said secondary storage device.

In a still further aspect, the present invention provides a computer implemented method comprising: for a first
5 publication of a product, storing in a primary storage device promotional information and meta data associated with the first published product; for a plurality of products, storing in the primary storage device product elements determined to be accessed more frequently based on user access patterns; and
10 for the plurality of products, removing at least a portion of a play track of one of the plurality of products from the primary storage device determined to be accessed less frequently based on user access patterns.

In a further aspect, the present invention provides a
15 computer implemented method comprising: in response to receiving a request for a video program from a client, streaming an initial portion of the video program to the client from primary storage of a video server; contemporaneously with the streaming of the initial portion,
20 retrieving a subsequent portion of the video program from secondary storage, wherein a duration of streaming the initial portion is configured to be longer than a latency of retrieving the subsequent portion of the video program from the secondary storage; and streaming the subsequent portion of
25 the video program to the client upon completion of the streaming of the initial portion.

In a still further aspect, the present invention provides a computer implemented method comprising: receiving requests for a plurality of titles from a plurality of clients; based
30 on the requests, determining which titles of a library stored in secondary storage are most likely to be requested by a client; determining a remaining capacity of a storage device

-3b-

of a streaming video server; removing, from the storage device of the streaming video server, one or more assets related to titles not determined most likely to be requested by a client; retrieving from the secondary storage one or more assets of the titles determined most likely to be requested by a client; and storing, to the storage device of the streaming video server, the retrieved one or more assets.

In a further aspect, the present invention provides in a storage server comprising a plurality of modules, each of said modules comprising a processor and a plurality of disks and a host controller, a method comprising the steps of: receiving data requests issued from a plurality of clients, sending vote solicitations from said host controller to each of said modules, each of said vote solicitations containing the name of the requested data, receiving votes from said modules that have said requested data on the module's respective disks, each of said votes containing a module ID, if one or more of said votes are received, selecting one of said votes according to a voting algorithm, forwarding said data request to the module that sent the selected vote, and after said module retrieves said requested data from the respective disks, receiving said data from said module and sending said data to said clients issuing said data request.

In a still further aspect, the present invention provides a method comprising: in response to receiving a request for a video program from a client, streaming an initial portion of the video program to the client from primary storage; contemporaneously with the streaming of the initial portion, retrieving a subsequent portion of the video program from secondary storage; and streaming the subsequent portion of the video program to the client upon completion of the streaming of the initial portion.

-3c-

In a further aspect, the present invention provides a method comprising: receiving requests for a plurality of titles from a plurality of clients; based on the requests for the plurality of titles, determining which titles of a library stored in a secondary storage are most likely to be requested by a client; determining a remaining capacity of a primary storage device; and removing, from the primary storage device, one or more assets related to titles not determined most likely to be requested by a client; retrieving from the secondary storage one or more assets of the titles determined most likely to be requested by a client; and storing, to the primary storage device, the retrieved one or more assets.

In a still further aspect, the present invention provides an apparatus comprising: a storage device; and a controller device configured to: in response to receiving a request for a video program from a client, stream an initial portion of the video program to the client from the storage device; contemporaneously with the streaming of the initial portion, retrieve a subsequent portion of the video program from secondary storage; and stream the subsequent portion of the video program to the client upon completion of the streaming of the initial portion.

In a still further aspect, the present invention provides a method comprising: for a first publication of a product of a plurality of products, storing in a storage device promotional information associated with the product; storing in the storage device elements of the plurality of products determined to be accessed more frequently based on user access patterns; and removing an interval of a play track of one of the plurality of products from the storage device determined to be accessed less frequently based on the user access patterns.

-3d-

In a still further aspect, the present invention provides an apparatus comprising: a storage device; and a controller configured to: for a first publication of a product of a plurality of products, store in the storage device promotional information associated with the product; store in the storage device elements of the plurality of products determined to be accessed more frequently based on user access patterns; and remove an interval of a play track of one of the plurality of products from the storage device determined to be accessed less frequently based on the user access patterns.

In a further aspect, the present invention provides a method comprising: streaming, in response to receiving a request for a video program from a client, an initial portion of the video program to the client from primary storage; retrieving each remaining portion of the video program from secondary storage while bypassing the primary storage, including retrieving a first of the remaining portions contemporaneously with the streaming of the initial portion; and streaming the first of the remaining portions of the video program to the client upon completion of the streaming of the initial portion.

In a still further aspect, the present invention provides an apparatus comprising: a storage server comprising an array of storage devices; and a controller configured to: stream, in response to receiving a request for a video program from a client, an initial portion of the video program to the client from the storage server; retrieve, contemporaneously with the streaming of the initial portion, a subsequent portion of the video program from secondary storage while bypassing the storage server; and stream the subsequent portion of the video program to the client upon completion of the streaming of the initial portion.

-3e-

Further aspects of the invention will become apparent upon reading the following detailed description and drawings, which illustrate the invention and preferred embodiments of the invention.

5

BRIEF DESCRIPTION OF THE DRAWINGS

The teachings of the present invention can be readily understood by considering the following detailed description in conjunction with the accompanying drawings, in which:

10 FIG. 1 depicts a high-level block diagram of a data retrieval system that includes a storage server incorporating the present invention;

FIG. 2 depicts a detailed diagram of the architecture of the storage server;

15 FIG. 3 depicts a flowchart specification of the Data Initialization Protocol;

FIG. 4 depicts a flowchart specification of a more general version of the Data Initialization Protocol;

20 FIG. 5 depicts a flowchart specification of the Data Retrieval Protocol;

FIG. 6 depicts a flowchart specification of the Access Request Assignment Protocol; and

FIG. 7 depicts a flowchart specification of the Access Request Migration Protocol;

25 FIG. 8 depicts a high level block diagram of a data retrieval system including a plurality of storage servers;

FIG. 9 depicts a graphical representation of a plurality of assets forming a title;

30 FIG. 10 depicts a flow diagram of a method for managing the publication of new products;

FIG. 11 depicts a flow diagram of a method for managing server space.

-3f-

FIG. 12 depicts a flow diagram of an asset retrieval method suitable for use in an information distribution system utilizing primary and secondary storage.

To facilitate understanding, identical reference numerals 5 have been used, where possible, to designate identical elements that are common to the figures.

DETAILED DESCRIPTION

FIG. 1 depicts a client/server data retrieval system 100 that employs a storage server 110 to accept user data access requests from clients $120_1, 120_2, 120_3, \dots, 120_n$ (collectively referred to as clients 120) via bi-directional data paths $150_1, 150_2, \dots, 150_n$ (collectively referred to as paths 150). Server 110 retrieves the requested data from the disk drives within the server, and, when necessary, from secondary storage 130 via bi-directional path 140, and outputs the requested data for distribution to the appropriate client(s) along data path(s) 150. A detailed description of a video-on-demand system that finds the server of the present invention particularly useful is described in commonly assigned U.S. Patent No. 6,253,375, filed December 3, 1997.

FIG. 2 depicts a detailed diagram of the architecture of storage server 110. The storage server comprises a host controller 210 and a plurality of modules $220_1, 220_2, \dots, 220_n$ (collectively referred to as modules 220) that are coupled to the host controller by bi-directional data paths $230_1, 230_2, \dots, 230_n$ (collectively referred to as paths 230). In the preferred embodiment of the invention, data paths 230 are part of an ethernet network. Each of the modules 220 comprise a processor 240, a cluster of 16 disk drives 250, and a request queue 260 in which access requests wait to be serviced by the disk drives 250. The host controller 210 accepts incoming user access requests from data path 212, and assigns these requests to modules 220 by forwarding the requests along data paths 230.

A given data item (e.g., a video program) is stored on disk by striping the data across the 16 disk drives of one of the modules 220 according to the RAID 5 protocol. RAID 5 is well known in the art as a protocol that provides fault-tolerant error-correcting storage and retrieval on an array of disks.

A data item is therefore stored at a single module on a 16-disk drive cluster, in contrast to non-modular server

architectures in which data is striped across the entire set of disk drives. In order to manage data in this modular architecture, protocols are necessary for determining which module shall store the data at any given time. In addition, 5 it is possible in a modular architecture to store multiple copies of data within various modules, a technique known as data replication, and to dynamically migrate data from one module to another to accommodate changes in user access patterns. Finally, a substantial data library is stored in 10 secondary storage 130, e.g., content server, a magneto-optical storage array or other form of data storage, and selected data items are recalled from secondary storage 130 stored on disk 250 at one or more modules 220 to enable rapid access of the data. As such, special data management 15 routines are necessary to facilitate data movement from secondary storage 130 to the storage modules 220.

The Data Initialization Protocol determines which data items are retrieved from the secondary storage 130 and stored on the disk drives 250, and at which the module(s) 20 220 data are stored. The basic strategy of the Data Initialization Protocol is to repeatedly retrieve the most popular data items from the secondary storage 130 and store them on the disk drives 250, until there is insufficient disk space available to hold any additional data items. At 25 each iteration of the Data Initialization Protocol, the candidate modules 220 that have sufficient disk space to hold the data are identified, and the data is stored at the candidate module with the lowest load, where "load" refers to the total bandwidth requirement for the requests waiting 30 in the module's queue 260. Since at initialization, no user requests have yet been submitted, each data item is assigned an initial popularity value, and the load is estimated as the sum of the popularity values of the data items already stored on the module disk drive cluster 250. By selecting 35 the module 220 with the lowest load, the Data Initialization Protocol provides a significant load-balancing effect that keeps the popular data items evenly distributed over all the

-6-

modules 220, thus enhancing the performance of the server 110.

FIG. 3 depicts a flow diagram of the Data Initialization Protocol 300. At each iteration of the protocol, the most popular data item in secondary storage that can fit in at least one module's 16-disk drive cluster is selected at step 320. If no such data item exists, as determined by the query at step 330, the protocol terminates at step 370; otherwise, the set of candidate modules (M) with sufficient disk space is determined at step 340. At step 350, the module (M) having the lowest present load is selected from the set of modules M having sufficient space to store data d. At step 360, the data item is then stored on the disk drive cluster of the candidate module with the lowest load and the next iteration begins at step 320.

Note that the above specification of the Data Initialization Protocol does not have any provisions for multiple copies of data items stored on the disk drives (i.e., data replication.) In some applications such as a video server; however, it is desirable to replicate frequently requested data so that these data items don't act as "bottlenecks" when a substantial number of users request the same data item concurrently. Thus, the preferred embodiment of the invention employs a more general version of the Data Initialization Protocol in which data can be stored on the disk drive clusters of multiple modules simultaneously. In this more general version of the protocol, the popularity value of a data item is an integer denoting the desired degree of replication for that data item (i.e., how many copies of the data item should be stored.) Note that multiple copies of a data item should never be stored within the same module, as this offers no performance advantage while unnecessarily consuming additional disk space.

FIG. 4 depicts a flow diagram of the Data Initialization Protocol 400. At each iteration of the protocol 400, the data item with the highest replication count (denoted by copies(d)) is selected at step 420 and,

-7-

as in the previous version of the protocol, this data item is stored at the module with sufficient free disk space having the lowest load as determined in steps 440, 450, and 460. After the data item is stored, the replication count
5 copies(d) is decremented by 1, at step 470, and the next iteration begins at step 420.

When a user requests a data item that is not stored on the disk drives, the storage server retrieves the data from secondary storage, loads it onto a selected module disk
10 drive cluster, and outputs the data for delivery to the requesting user(s). In order to determine at which disk drive cluster the data should be stored, a Data Retrieval Protocol is used. If there is no available disk space on any of the clusters, and there is at least one data item
15 with more than a single copy on disk, the Data Retrieval Protocol will select one such data item and remove a copy from disk to make room for the new data item. The basic strategy of the Data Retrieval Protocol is to store the new data item on the least loaded module that has sufficient
20 free disk space to accommodate the new data item. If no modules have available disk space, however, then the Data Retrieval Protocol will make room for the new data item by replacing one of the data items currently on the disk drives. Alternatively, if no free space is available, the
25 system may send the data item to the user without storing the data item.

Selecting the data item on disk to be replaced is similar to the selection of virtual memory pages to be replaced in an operating system. The Data Retrieval
30 Protocol borrows from two page replacement algorithms, called Least Recently Used and Not Frequently Used, which are well known in the operating system arts. Specifically, the Data Retrieval Protocol maintains a view count for each data item which indicates the number of times the data item
35 has been requested, and a timestamp for each data item indicating the time at which the data was last requested. The protocol then identifies the set of data items which are both inactive and replicated, and from this set selects the

-8-

item that has been least frequently requested (i.e., having the smallest view count), breaking ties by selecting the item that was least recently used (i.e., having the earliest timestamp). A data item is active, and thus ineligible for
5 the aforementioned set of replacement candidates, in any one of the four following situations: (1) the data item is being read from the disk drives in response to a user access request, (2) the data item is being retrieved from secondary storage and stored on the disk drives, (3) the data item is
10 being migrated to another module disk drive cluster, and (4) the data item is in the process of being replicated.

FIG. 5 depicts a flow diagram of a Data Retrieval Protocol 500. This Protocol begins at step 510 and then checks (at step 520) whether the requested data item d is
15 stored on the disk drives, and if the data is stored on disk, the data item d is retrieved at step 585. If the data is not locally available, the protocol retrieves data item d from secondary storage at step 530, and then, at step 540, determines the set C of inactive data items having more than
20 one copy. If set C is empty as determined by the query at step 550, data item d is output for delivery to the requesting user(s) at step 590. Otherwise, at step 555, the protocol tries to select a subset C' of set C of data items in the same module, each having a low view count and/or
25 timestamp. If such a subset C' does not exist as determined at step 560, data item d is forwarded to the user without being stored on the disk drives at step 590. Otherwise, at step 570, the protocol removes set C' from the disk drives, stores data item d on disk in the vacated spot at step 580,
30 and forwards d to the host controller at step 590. The protocol ends at step 595.

In addition to the protocols disclosed, the present invention employs dynamic data migration to maintain server performance as user access patterns change over time. A
35 low-priority, non-real-time thread migrates data to new modules based on monitored access patterns, with the objective of evenly re-distributing popular data over all modules. A second low-priority non-real-time thread

migrates data to new tracks on the same disk drives, also based on observed access patterns, with the objective of locating popular data on the faster outer tracks of disk drives. The difference in seek and transfer rates between the outer and inner tracks of a disk are substantial; for example, the internal transfer rate of a Seagate Cheetah disk varies from 152 Mbps on inner tracks to 231 Mbps on outer tracks, and the seek time can vary from 2ms to 13ms depending on how far away the next segment of data on the disk drive is. The movement of popular data by the second thread to faster tracks can therefore significantly increase disk bandwidth, and consequently, overall server performance.

In addition to data management, in a modular server architecture it is necessary to assign user access requests to modules, a task that is particularly important when the requested data item is located at multiple modules. This task is accomplished in the present invention by an Access Request Assignment Protocol that is performed by the host controller. The object of the protocol is to assign requests to modules such that outstanding user requests are evenly distributed among the module processors, thus resulting in steady-state module loads. The Access Request Assignment Protocol 600 is depicted as a flow diagram in FIG. 6. The protocol 600 is essentially a conventional voting algorithm (voting algorithms are well-known in the distributed system arts.) The process begins at step 610 and proceeds to step 620 wherein each of the module processors determines, in parallel, whether the requested data item is stored on its disk drive cluster. At step 630, all processors that have the requested data item submit a vote to the host controller and, at step 640, the host controller forwards the access request to the module with the lightest load. The protocol 600 ends at step 650.

After an access request has been forwarded to a module by the Access Request Assignment Protocol 600, the request waits in the module queue for disk drive servicing. Since module loads may become uneven over time, the present

-10-

invention employs an Access Request Migration Protocol which attempts to dynamically re-balance module loads by moving requests waiting for service from one module's queue to another module's queue. Of course, the protocol can only
5 migrate a request to another module if the requested data item is replicated (i.e., another complete copy of the data is located on some other module disk cluster.)

The basic strategy of the Access Request Migration Protocol is to find the two modules in the server with the
10 maximum and minimum loads, and attempt to migrate a user request from the maximum-loaded module queue to the minimum-loaded module. The protocol repeats this process periodically, running as a background thread.

FIG. 7 depicts a flow diagram of the Access Request
15 Migration Protocol 700. The protocol begins at step 705 and proceeds to step 710 wherein the protocol first finds MAX, the module with the largest load, and MIN, the module with the lightest load. The queue of outstanding requests in MAX's queue is then examined at step 720 to see if there is
20 some request in the queue waiting to access a data item with another copy at MIN, such that the disk in MIN's cluster needed for the next service period has a free slot. If such a request is found, it is migrated at steps 730 and 740 from MAX to MIN). Optionally, in the case where multiple such
25 requests are found, it is advantageous for the protocol to migrate the request for the largest data item (i.e., the data that will take longest to read from disk). This process is then repeated indefinitely, starting at step 710.

The foregoing discloses a detailed description of a
30 modular architecture and data management methods for a general-purpose storage server. In video servers, there are a number of additional issues that arise due to the specific nature of video programs. In particular, it is advantageous to have the Data Initialization Protocol store only "leader
35 tracks" containing the first few minutes of a movie on the disk drives, and then, once a movie has been requested for viewing, retrieve the remainder of the movie from secondary storage to the disk drives while the user is watching the

-11-

leader track. This allows the video server to deliver a much larger volume of programs in real-time, rather than initially storing complete movies on the disk drives. As in the case of general-purpose data, leader tracks should be
5 evenly distributed over the server modules, and when retrieving the remainder of a movie from secondary storage, it is not necessary to store the data at the same module at which the leader is located. Finally, in the Access Request Migration Protocol, if there are multiple outstanding
10 request candidates in the maximum-loaded module, the optional selection procedure should choose the request for the movie with the longest remaining playing time and least recently manipulated, i.e., least recent use of fast forward, rewind and the like.

15 FIG. 8 depicts a high level block diagram of a data retrieval system including a plurality of storage servers. Specifically, FIG. 8 depicts a client/server data retrieval system 800 that employs, illustratively, a plurality of storage servers 810₁ through 810_N (collectively storage
20 servers 810) to accept user data access requests from respective client groups 815₁ through 815_N (collectively client groups 815). It is noted that each client group 815 comprises a plurality of clients, such as client 820₁₁ through 820_{1M}, which are shown in FIG. 8 as comprising client
25 group 815₁. Each server 810 communicates with its respective client group 815 via a respective bi-directional data path 150 (e.g., data paths 150₁₁ through 150_{1N}). Each server 810 retrieves the client requested data from disk drives within the server and, when necessary, from a secondary storage
30 device 830.

It is noted that secondary device 830 of FIG. 8 communicates with each of the servers 810₁ through 810_N via a respective bi-directional data path 140₁ through 140_N.

The data retrieval system 800 of FIG. 8 is especially
35 useful within the context of a networked secondary storage solution. The network implementing communications for such a secondary storage solution may comprise point-to-point links, an internet protocol (IP) network, a fiber channel

-12-

network or any other network suitable for carrying high bandwidth traffic. Preferably, the quality of service (QoS) provided to a customer requesting content via data paths 150 is not dependent upon a particular QoS level within the network 140. That is, as long as the network 140 delivers content from the secondary storage device 830 to the appropriate server 810 in a timely manner, irrespective of QoS deficiencies, the server 810 will be able to stream the requested content to a requesting subscriber via a respective network 150 in a manner providing a high quality presentation of the requested content to the subscriber.

However, where a server 810 has insufficiently cached requested content, the QoS of the requested content will also depend upon the QoS of network 140. That is, the controller must operate within the quality of service (QoS) associated with the network, such as latency, packet jitter and/or loss and other QoS factors. Thus, in one embodiment, a server 810 requesting content from the secondary storage device 830 via a network 140 performs various processing functions to insure that any degradation induced by the QoS of network 140 is either compensated for or masked prior to the serving of a requested title to a client 820.

Various processing techniques or algorithms may be utilized by a server 810 to address latency issues and other issues. For example, in one embodiment a server 810 is preloaded with one or more portions of a title such that subscriber access latency to the title is minimized. In this embodiment of the invention, a server 810 includes, for example, an initial portion of a title, which is streamed to a client immediately upon receiving a request for the title from the client. Additionally, other preloaded portions such as initial portions of title chapters or scenes may be preloaded such that user requests for title access at a particular chapter or scene may be rapidly satisfied. Contemporaneous to the start of streaming a preloaded portion, the server 810 accesses the secondary storage device 830 via the network 140 to retrieve the remaining portion of the title. In this manner, the server 810

operates as a caching server for at least the remaining portion of the title. Advantageously, the present invention operates to selectively provide push and/or pull provisioning techniques to control the type and/or amount of content stored within the servers 810.

In one embodiment, a title from the secondary storage unit 830 is loaded at a location other than the start of the title or other asset, such that latency to a particular portion of the title or other asset is minimized. That is, a title or other asset is loaded in a manner providing one or more entry points to the title or asset other than an initial entry point. Each of the one or more entry points comprises a portion of the title that may be accessed rapidly. For example, in the case of a title divided into a plurality of chapters such as used in digital versatile disks (DVD) storage media, a title loaded onto the secondary server may be loaded as a plurality of chapters, where each chapter may be rapidly entered (i.e., retrieved) based on client requests. This random entry of predefined portions of a title helps reduce latency. By storing several portions of an asset within a server 810, where the stored portions are proximate logical play track or content entry points (i.e., chapter points or scene change points), the server 810 is more likely to be able to immediately satisfy a customer request.

In one embodiment, a video reservation method is implemented. In this embodiment, a client request for a particular title or other asset is queued by the corresponding server 810. The server 810 retrieves the requested title or other asset from the secondary storage device 830 via the network 140. When the requested title or other asset is retrieved, a notification is sent to the client that the requested title or asset is available for delivery and presentation to the client.

In one embodiment, a preview, promotion or advertisement is used to mask latency associated with the transfer of a title or other asset from the secondary storage 830 to the server 810. In this embodiment, various

-14-

"eye candy" is presented to a client while a client request is processed by the server. It is noted that further refinement of this embodiment is allowing a client to skip the preview, promotion or advertisement in the case of the title or requested asset having at least an initial portion preloaded on the server, as discussed above.

The secondary storage interconnect network 140 has associated with it various capacity and QoS parameters, while title or asset availability depends upon the provisioning of secondary storage 830 and the server 810 serving a particular client. These parameters heavily influence the time required to acquire a title, product or other asset stored in secondary storage 830. As such, the close controller within server 810 (or other controllers within the data retrieval system operating to satisfy a client request) must evaluate provisioning, QoS and/or bandwidth availability information to responsively determine an appropriate latency masking and/or client message strategy. In each instance, the goal is to enhance the client-side interaction experience. This experience is enhanced by timely providing requested content, by masking latency associated with satisfying such requests, by directing appropriate promotional material to clients and/or by providing pre-loaded content or content portions on the server 130 or random access to appropriate entry points within content titles.

In one embodiment, a server 810 requiring content to satisfy a user request retrieves the required content from the secondary storage device 830 or, optionally, from another storage server 810. Content may be transferred between storage servers 810 via the secondary storage device 830 using the network 140. Optionally, a second or auxiliary network 840 connecting one or more of the storage servers 810 is provided for this purpose. The auxiliary network 840 optionally includes connectivity with the secondary storage device 830. The auxiliary network 840 is especially useful where two storage servers 810 are used to store respective portions of available content, such that

-15-

each storage server 810 operates as a secondary storage device with respect to a portion of available content. Optionally, each server 810₁ through 810_N has associated with it a respective network connection 845₁ through 845_N used to
5 connect the storage server to other equipment (not shown). Such other equipment may comprise local area network (LAN), wide area network (WAN), additional secondary and/or primary storage devices and the like. Essentially, each storage server 810 may be used to provide additional services and/or
10 coordinate with other equipment.

The above-described embodiments are associated with a managed server storage model, in which the service operator and server demand determine the provisioning of content, either wholly or in part, within secondary and/or primary
15 server modules.

The inventors have determined that for some applications a demand storage model is appropriate. A demand storage model of managing content may comprise, for example, a push model and a pull model. A pull model is
20 where provisioning is adapted in response to client demand. A push model is where provisioning is adapted to storage and/or bandwidth availability, along with decisions regarding specific content to be promoted. As an example of a push model, a push from the server 830 to the caches 810
25 may be selected using service decisions made by programming personnel, as well as service decisions based upon aggregate purchasing behavior of network customers as evaluated against the characteristics of the asset(s) to be pushed. The total bandwidth usage of the network 140 may be
30 minimized by using a broadcast (i.e., a multi-cast) model for the push implementation. In this case, illustratively all servers 810 receive the push content but each server implements a level of filtering of received content according to a server-specific evaluation of the importance
35 of the pushed content made according to the needs and/or preferences of the subtended customers associated with the respective server 810. Thus, content pushed from secondary storage 830 to the servers 810 is blended with pull criteria

-16-

based upon client groups 815 served by respective servers 810.

Within the context of the storage server 800 of FIG. 8, or any server network comprising secondary storage communicating with multiple streaming servers, optimal provisioning and other data management functions may be realized by balancing push and pull methodologies within a demand storage model.

10

STORAGE MANAGEMENT

Within the concept of data retrieval systems utilizing primary and secondary storage (especially where high bandwidth information such as video information is stored), it is important to efficiently store and delete content from at least the primary storage devices. For purposes of this discussion, all content titles are primarily stored in secondary storage and a demand pull algorithm is used to decide when to bring a product, title and/or asset to a primary storage device, such as a streaming server. The following definitions are applied: A "product" is a subscriber purchasable item, such as a movie, package of movies, subscription and the like. Products are defined by meta data which includes description, price, discount rules and the like. Products contain assets such as the play track of a movie, fast play and fast reverse tracks, navigation related assets such as navigation screens, promotional videos and the like.

Information associated with a product comprises at least product rights (or information) enabling a host controller to determine what a customer may do with a particular product, such as viewing time, use time, price, purchase rules, encryption methodologies and the like. A "title" is a viewable entity, such as a movie, a television episode and the like. An "asset" is an elementary component of a title or product, such as an MPEG play track, a preview track, a JPEG still image associated with a title, fast forward and/or rewind tracks and the like.

-17-

Products are typically stored on secondary storage by either the network operator or third party program providers. The availability of such stored products is "published" to the application servers that present a user interface to the customers. The products become elements of the overall user interface, either as on screen graphical or video imagery or web page information. Methods and apparatus for providing such on-screen graphical or video imagery within the context of a user interface are disclosed in U.S. Patent No. 6,208,335 (Attorney Docket No. 006), granted on March 27, 2001.

When a customer or a client requests a product that is not currently resident on the respective storage server or primary storage device, the host controller must retrieve the product definition of the requested product, which will provide, inter alia, a list of titles or products comprising the requested product. The product definition allows the application server to adapt the presented user interface of a customer to achieve a next level of selection, where a list of titles is presented to the customer for final title selection.

FIG. 9 depicts a graphical representation of a plurality of assets forming a title. Specifically, the exemplary title 900 of FIG. 9 comprises a play track 901, a fast forward track 902, a rewind track 903, a promotional track 904, a still image 905 and meta data (i.e., data related to the title or product) 906. It is noted that the relative sizes of the various tracks and other asset elements are not to scale. However, it is noted that the amount of disk space required to store the play track 901 is significantly greater than the amount of space required to store the fast forward track 902 or rewind track 903. Moreover, the amount of disk space required to store the fast forward/rewind tracks is significantly greater (typically) than the amount of time required to store the promotional track 904. The still imagery 905 and meta data 906 require even less memory to store. As such, multi-level

-13-

5 caching preferably favors the removal of high storage requirement assets, such as play tracks and FF/REW tracks, prior to the removal of other assets. As an example, a 100 minute movie encoded at 4 megabits per second results in, approximately, 3 gigabytes of data. Depending upon the speed of the trick play tracks, the FF/REW tracks comprise approximately 20% of the total file size. The preview track comprises, typically, 60 megabytes of data. The majority of the three gigabyte file is occupied by the play track.

10 To provide information sufficient to enable a customer to select an individual title, a subset of the assets associated with the various titles may be provided. For example, in one embodiment, a promotional track is streamed in response to user interaction with an interactive program
15 guide. Optionally, still imagery associated with a potentially selected title may be provided. JPEG files and meta data that describe individual titles may be used to create electronic program guide screens which are presented to a subscriber via the subscriber's presentation or display
20 device. These screens may be presented as video information or as web pages, such as HTML pages. A customer may choose to view a promotional or preview track of a title or product, in which case the server will retrieve the appropriate promotional track from a secondary storage
25 device or from the primary storage server devices.

It is important to note that the storage on a primary storage server is generally much less than the storage within a secondary storage device. As such, the host controller must actively manage the server storage to insure
30 a high quality interactive experience for each of the subscribers within the system. Storage management functionality within the context of the present invention includes a number of elements, such as:

Complete titles as well as the individual assets
35 forming complete titles are individually managed. Thus, when a product is first published the promotional information and meta data (e.g., usage rules) are retrieved and stored in primary storage. Some of the information

-19-

related to the individual titles, such as JPEG files or other still imagery, may be stored in primary storage such that latency is reduced to a user interacting via an interactive program guide.

5 Where elements of the product are expected to be frequently viewed, such elements may also be retrieved and stored in a primary storage device. As previously discussed, a push method may be used to distribute the various assets associated with a product. For example, when
10 a product is first published and heavily promoted, preview information associated with individual titles may be retrieved from secondary storage and stored within primary storage prior to customer requests. Additionally, information published with the product, such as Motion
15 Picture Association of America (MPAA) ratings or other ratings may also be placed in primary storage as meta data. Such information is useful in enabling subscribers to select appropriate content.

Storage management algorithms according to the
20 invention adapt to the usage pattern of customers, the storage requirements and/or capacity of the primary and secondary storage devices, and the expected demand for content.

In this method, the controller waits for new products
25 to be published. Product meta data may be retrieved (not shown) and, from that meta data, decisions are made by the controller regarding the assets to be transferred. For example, product promotional assets are always transferred, since these assets are generally required for a customer to
30 access the product. Thus, if no space is available in a storage server, space will be made for at least the promotional assets. A decision is then made to transfer titles and assets that may be frequently requested, though these assets are only transferred if space is available.
35 Otherwise, this requirement is queued to be rechecked at a future time. In transferring frequently used titles and assets, the actual transfers will depend upon the available space within the storage server. A space management method

will be described in more detail below with respect to FIG. 11.

FIG. 10 depicts a flow diagram of a method for managing the publication of new products. The method 1000 of FIG. 10 is suitable for use in, for example, a controller (such as the host controller 210 of FIG. 2) operating with a server (such as the servers 110 or 810 of FIGS. 1 and 8).

The method 1000 is entered at step 1005 where the controller waits for new product. Upon receiving new product, a query is made at step 1010 as to whether any promotional assets exist in the new product. If no promotional assets exist, the method 1000 returns to step 1005. If promotional assets exist, then a determination is made at steps 1015 and 1020 as to whether space on the storage server 110/810 is available. If space is not available, then at step 1025 sufficient space within the storage server is freed or released.

At step 1030, the promotional assets associated with the new product are transferred to the existing (or newly released) space on the storage server. At step 1035, a query is made as to whether any high usage titles are included with the new product. If no high usage titles are included, then the method 1000 returns to step 1005 to await further new products. If high usage titles are included, then at step 1040 a determination is made as to whether sufficient space is available to store the high usage titles. At step 1045, a query is made as to whether the determination at step 1040 indicates that sufficient space is available. If available, at step 1050 the assets that will fit are transferred to the storage server. Otherwise, a request for the high usage title assets is placed in a request queue. The method 1000 then returns to step 1005.

FIG. 11 depicts a flow diagram of a method for managing server space. The space management method 1100 of FIG. 11 is suitable for use in, for example, a controller (such as the host controller 210 of FIG. 2) operating with a server (such as the storage servers 110 of FIG. 1 or 810 of FIG. 8). The space management algorithm 1100 takes into account

-21-

the need to minimize the time to restore a title and provide full VOD capability. The flow diagram depicts some of the considerations, such as not deleting bookmarked titles, deleting the play track first and the like. Other
5 modifications to the method 1100 of FIG. 11 include deleting the end of a play track before deleting an entire play track and other modifications.

At step 1105, the available space within the storage server to be managed is checked. At step 1110, a query is
10 made as to whether sufficient space to store the assets associated with a product or title has been found. If sufficient space has been found, the method 1100 exits at step 1115.

If sufficient space has not been found at step 1110,
15 then at least one infrequently used title is found at step 1120. At step 1125, a query is made as to whether the title found at step 1120 has been bookmarked. If the infrequently used title of step 1120 has been bookmarked, then the method 1100 proceeds to step 1120 where the next infrequently used
20 title is found. The method iterates through this loop until an infrequently used title that has not been bookmarked is found, at which point the method 1100 proceeds to step 1135. If no unbookmarked title is found (i.e., bookmarked titles are not deleted) then the method 1100 proceeds to step 1130,
25 where the title least used or requested by subscribers is identified.

At step 1135, at least a portion of the play track of the non-bookmarked infrequently used title (or the least used title) is deleted. As previously noted, the play track
30 of a title comprises the video and associated audio information associated with the title, and is the largest single asset associated with the title. In one embodiment, only those portions of the play track not proximate chapter entry points and/or scene change points are deleted. In
35 this manner, customer requests for the entry into a play track at typical chapter entry points may be rapidly satisfied. Such requests will be discussed in more detail

-22-

below with respect to step 1235 in the method 1200 of FIG. 12. The method 1100 then proceeds to step 1140.

At step 1140, a query is made as to whether the remaining free space is sufficient to store a desired title. 5 If the available space is sufficient, then the method 1100 exits at step 1145. Otherwise, other assets associated with the infrequently used title (or least used title) are deleted at step 1150. At step 1155, a determination is made as to whether the available space is now sufficient to store 10 the desired title. If sufficient space is available, then the method 1100 exits at step 1145. If the space is insufficient, then the method 1100 performs the above steps beginning with step 1120 to free up additional space within the server. It is noted that where a bookmarked title is 15 found that the play track of the bookmarked title may be partially deleted to free space.

SERVICE MANAGEMENT/ASSET RETRIEVAL

Within the context of the present invention, the 20 storage management methodologies employed to effect efficient video asset caching are adapted in response to service considerations. That is, the retention or deletion of assets associated with products and/or titles depends upon the need to provide access to new titles or products as 25 well as the need to avoid the expense of storing infrequently used or otherwise less necessary titles and products on a primary storage device. To achieve appropriate service levels for customers or subscribers within the system, it is necessary to transfer video assets 30 from secondary storage to primary storage prior to streaming the video assets to a customer. Within the context of video assets including play tracks, fast forward tracks, fast reverse tracks, bookmark or chapter entry points and the like, this transfer must be intelligently handled. In one 35 embodiment, the transfer of video assets from primary to secondary storage occurs prior to any streaming of the video assets to a requesting subscriber. Unfortunately, this results in a relatively high degree of latency. To avoid

such latency, a smart retrieval method will now be described with respect to FIG. 1.

FIG. 12 depicts a flow diagram of an asset retrieval method suitable for use in an information distribution system utilizing primary and secondary storage, such as for the distribution of video information. Specifically, the method 1200 of FIG. 12 is especially useful within the context of titles or products comprising a plurality of video assets, such as described above with respect to FIG. 9.

The method 1200 of FIG. 12 is entered at step 1205 when a controller executing the method receives a title request. At step 1210, a query is made as to whether the requested title is resident within a primary storage device associated with a requesting user or subscriber. If the requested title is resident, then at step 1215 the primary storage device begins streaming the play track of the requested title to the requesting subscriber. Optionally, a promotional track or other asset associated with the requested title is provided to the subscriber. The method 1200 then exits at step 1220. It is noted that further user interactions requiring the streaming of fast forward/rewind or other trick play tracks are handled in the manner previously described above with respect to FIGS. 1-8.

If at step 1210 it is determined that the requested title is not resident with the primary storage, then at step 1225 the controller begins retrieving portions of the play track proximate chapter delineation points. That is, at step 1225, where a play track is divided into chapters or segments, such as chapters commonly found on digital versatile disk (DVD) media, several blocks or portions of the play track near each chapter point are retrieved prior to the retrieval of the remainder of the play track. In this manner, if a customer chooses to "jump" to a chapter point, at least some of the play track proximate that chapter point will be available for presentation, thereby reducing the latency experienced by the user.

-24-

At step 1230, the play track retrieval continues until all the portions proximate chapter delineation points are retrieved or until a user interaction indicates a request to jump to the presentation of retrieved content at a chapter point. If a jump request is received, then at step 1235 the retrieved play track is streamed to the user beginning at the requested chapter. Additionally, the controller causes the play track to be retrieved from the secondary storage device beginning with the play track portions necessary to provide the requested chapter and subsequent chapters in the title or product. The method 1200 then proceeds to step 1240. When the play track portions proximate chapter delineation points have been retrieved (per step 1230) or when at least the requested chapter play track has been retrieved (per step 1235), at step 1240 the controller begins retrieving play track portions supportive of trick play operation. Specifically, portions of the play track are associated with respective portions of trick play tracks, such as fast forward or fast rewind tracks. Transitions between play and trick play tracks are supported using indexing information, whereby streaming and presentation of a play track results in the advancement of an index associated with the play track, said index having a trick play index counterpart indicating the entry point into a trick play track should a user request fast forward or rewind of a currently streamed play track. Thus, given a play track being streamed to a user, it is desirable to retrieve from secondary storage at least portions of trick play tracks associated with presently streamed play track portions.

At step 1245, the retrieval of play track portions supportive of trick play operation is continued until play track portions capable of supporting at least limited trick play transitions have been retrieved or the user enters a trick play mode.

At step 1245, if the user requests a trick play mode, then the method proceeds to step 1250 where the controller begins retrieving the appropriate trick play track (i.e., FF

-25-

or REW) from the secondary storage device and streaming the
trick play track to the user beginning at the requested
entry point. At step 1255, when the user exits the trick
play mode, the controller begins retrieving play track data
5 and streaming the retrieved play track data at the play
track entry point associated with the trick play track exit
point. The method 1200 then proceeds to step 1260.

At step 1260, the controller begins retrieving the
remaining portions of the trick play tracks from the
10 secondary storage device. At step 1265, the trick play
track retrieval continues until all the trick play tracks
have been retrieved or until user input is received. Upon
retrieval of all the trick play tracks from the secondary
storage device, the controller begins retrieving the
15 remainder of the play track and any other assets not yet
retrieved at step 1270.

At step 1275, the controller waits for (or responds to)
user input, such as a jump request, a trick play mode entry
or exit request and the like. The method 1200 exits at step
20 1280.

Within the context of the present invention, it is
contemplated that fast play and fast reverse play tracks are
utilized in which compressed video streams (e.g., MPEG or
other compression schemes) including intracoded frames (I-
25 frames), forward predictively coded frames (i.e., P-frames)
and/or bidirectionally predicted frames (i.e., B-frames) are
utilized. In this manner, a predefined presentation speed
ratio may be established between the play track and the
trick play tracks, illustratively a 1x9 ratio. It is noted
30 that in the case of trick play tracks providing a 9x speed
increase with respect to standard play tracks that the
amount of time required to retrieve a quick play track from
a secondary storage device is approximately one fifth of the
time required to retrieve the play track from the secondary
35 storage device. Using the method 1200 of FIG. 12, where
selected blocks of the play track (e.g. less than 5% of the
total play track) are retrieved, the latency experienced by
a subscriber is greatly reduced. It is noted that the

-26-

blocks retrieved at step 1225 represent approximately 5% of the total play track blocks. These blocks are preferably equally spaced throughout the play track, each block containing at least one complete group of pictures (GOP).

5 In one embodiment of the invention, imagery displayed upon the user's display device comprises a "progress bar," which serves to graphically represent the relative presentation position of a presently viewed play or quick play track. The progress bar represents the beginning of a
10 title on one end and the end of the title on another end, with a marker therebetween indicating a relative temporal presentation distance between the beginning and end portions of the title or product. In one embodiment of the invention, control information is received from the
15 subscriber indicative of subscriber selection of a portion of the control bar, said portion of the control bar being translated by the server module controller into a segment identification or chapter point, at which time the controller causes the streaming of content from that point
20 in either a play mode or trick play mode.

The method 1200 of FIG. 12 advantageously minimizes the latency experienced by a user within an interactive information distribution system by selectively retrieving portions of a title or product from a secondary storage
25 device in a manner calculated to meet user expectations. In one embodiment of the invention, one or more of the play and trick play tracks are stored in the secondary storage device multiple times at respective multiple quality levels. It is known that the amount of data required to represent video at
30 a relatively coarse quality level is less than the amount of data required to represent the same video at a relatively high quality level, where quality is defined in terms of spatial and/or temporal resolution. It is also noted that the above-described teachings with respect to video tracks
35 may also be applied to audio tracks associated with the video tracks. However, since the amount of information typically required to represent audio is far less than the amount of information required to represent corresponding

- 27 -

video, such processing of audio tracks is typically unnecessary, though such processing is contemplated by the inventors within the context of the present invention.

5 In one embodiment of the invention, the secondary storage device 830 of FIG. 8 may operate as a primary storage device. Thus, where the secondary storage device 830 operates as a primary storage device (such as a storage device 810), the primary storage device operations
10 described above with respect to storage device 810 and FIGS. 8-12 are applicable to at least a subset of the secondary storage device 830 operations.

 While this invention has been particularly shown and described with references to a preferred embodiment
15 thereof, the scope of the claims should not be limited by the preferred embodiments set forth in the examples, but should be given the broadest interpretation consistent with the description as a whole.

The embodiments of the invention in which an exclusive property or privilege is claimed are defined as follows:

1. A method comprising:

streaming, in response to receiving a request for a video program from a client, an initial portion of the video program to the client from primary storage;

retrieving each remaining portion of the video program from secondary storage while bypassing the primary storage, including retrieving a first of the remaining portions contemporaneously with the streaming of the initial portion; and

streaming the first of the remaining portions of the video program to the client upon completion of the streaming of the initial portion.

2. The method of claim 1, wherein a duration of the streaming of the initial portion is configured to be longer than a latency of the retrieving of the first of the remaining portions of the video program from the secondary storage.

3. The method of claim 1, wherein the remaining portions of the video program comprises a plurality of sequential segments, the method further comprising:

receiving, from the client, a second request for one segment of the plurality of sequential segments;

streaming, in response to receiving the second request for the one segment of the plurality of sequential segments, an initial portion of the one segment to the client;

retrieving, contemporaneously with the streaming of the initial portion of the one segment, a subsequent portion of the one segment from the secondary storage; and

streaming the subsequent portion of the one segment to the client upon completion of the streaming of the initial portion of the one segment.

4. The method of claim 3, wherein the plurality of sequential segments comprises video program chapters.

5. The method of claim 1, further comprising:

storing meta data related to the video program into the primary storage; and

sending, in response to the receiving of the request for the video program, the meta data to the client from the primary storage, wherein the meta data comprises one or more of a product description, product rights, a product price, product discount rules, still images, moving images, preview tracks, HTML links, title descriptions, MPAA ratings, and product information, and wherein the meta data is configured for interactive presentation within the client.

6. The method of claim 1, further comprising:

receiving requests for a plurality of video programs from a plurality of clients;

determining a pattern of the requests for the plurality of video programs;

retrieving, based on the pattern, initial portions of a subset of the plurality of video programs from the secondary storage; and

storing the initial portions of the subset into the primary storage.

7. The method of claim 1, further comprising:

transmitting the initial portion from the secondary storage to the primary storage.

8. The method of claim 1, wherein the primary storage includes a storage array storing the initial portion.

9. An apparatus comprising:

a storage server comprising an array of storage devices;
and

a controller configured to:

stream, in response to receiving a request for a video program from a client, an initial portion of the video program to the client from the storage server;

retrieve, contemporaneously with the streaming of the initial portion, a subsequent portion of the video program from secondary storage while bypassing the storage server; and

stream the subsequent portion of the video program to the client upon completion of the streaming of the initial portion.

10. The apparatus of claim 9, wherein a duration of the streaming of the initial portion is configured to be longer than a latency of the retrieving of the subsequent portion of the video program from the secondary storage.

11. The apparatus of claim 9, wherein the video program comprises a plurality of sequential segments, and wherein the controller is further configured to:

receive, from the client, a second request for one segment of the plurality of sequential segments;

stream, in response to receiving the second request for the one segment of the plurality of sequential segments, an initial portion of the one segment to the client from the storage server;

retrieve, contemporaneously with the streaming of the initial portion of the one segment, a subsequent portion of the one segment from the secondary storage; and

stream the subsequent portion of the one segment to the client upon completion of the streaming of the initial portion of the one segment.

12. The apparatus of claim 11, wherein the plurality of sequential segments comprises video program chapters.

13. The apparatus of claim 9, wherein the controller is further configured to:

store meta data related to the video program into the storage server; and

send, in response to the receiving of the request for the video program, the meta data to the client from the storage server.

14. The apparatus of claim 9, wherein the controller is further configured to:

receive requests for a plurality of video programs from a plurality of clients;

determine a pattern of the requests for the plurality of video programs;

retrieve, based on the pattern, initial portions of a subset of the plurality of video programs from the secondary storage; and

store the initial portions of the subset into the storage server.

15. The apparatus of claim 9, further comprising:

the secondary storage, wherein the secondary storage is configured to transmit the initial portion from the secondary storage to the storage server.

16. The apparatus of claim 9, wherein the subsequent portion includes an each remaining portion of the video program.

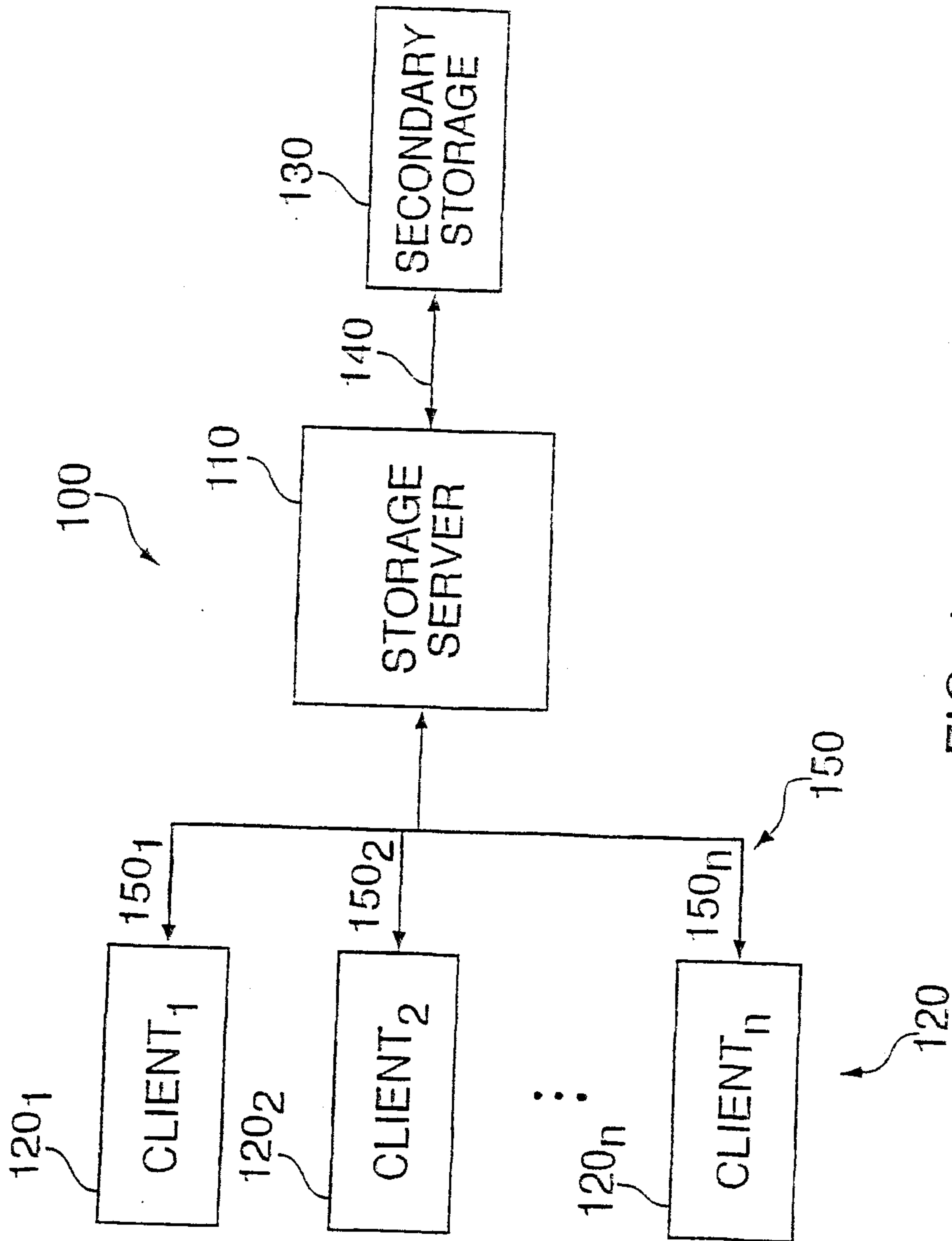


FIG. 1

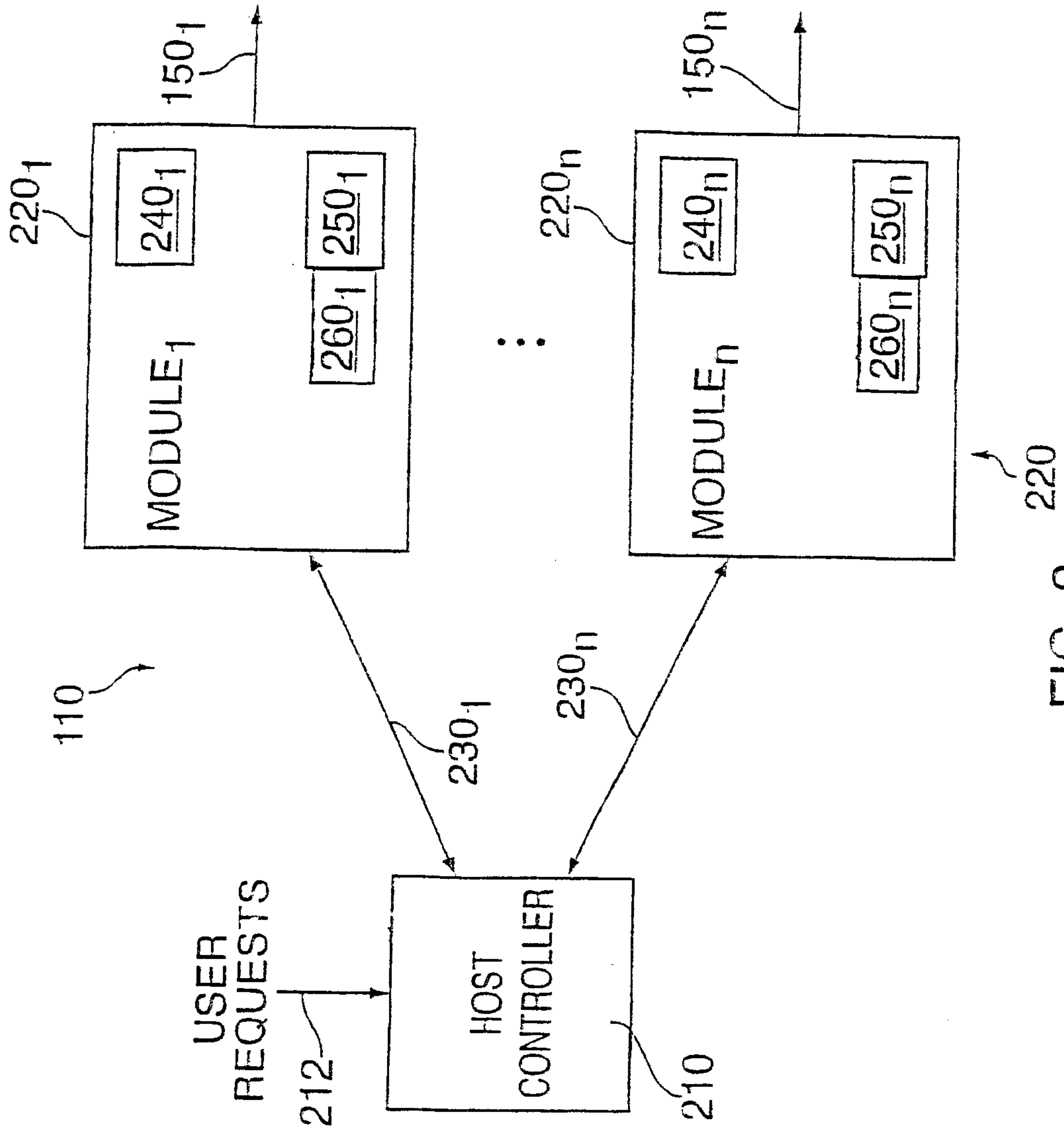


FIG. 2

3/12

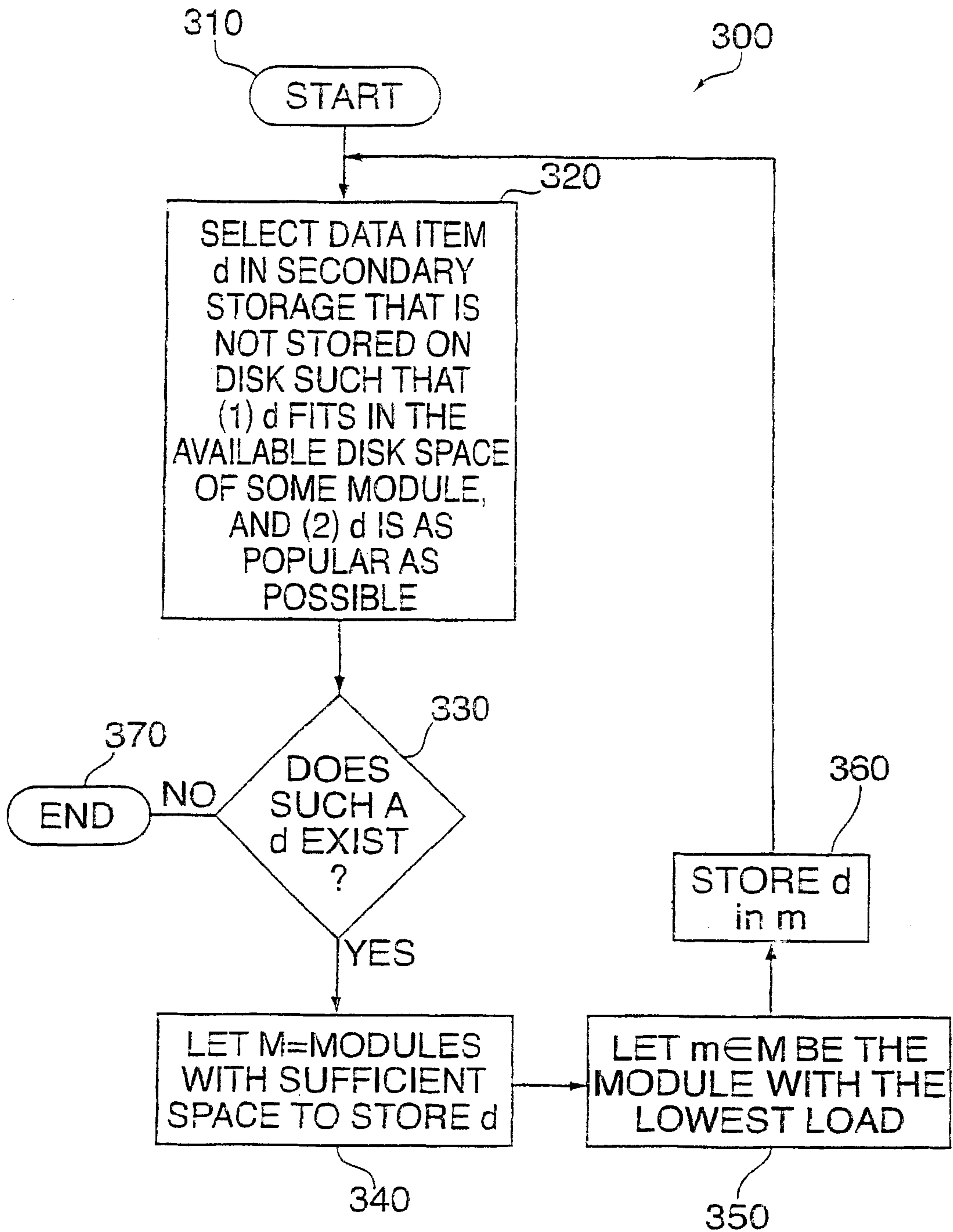


FIG. 3

4/12

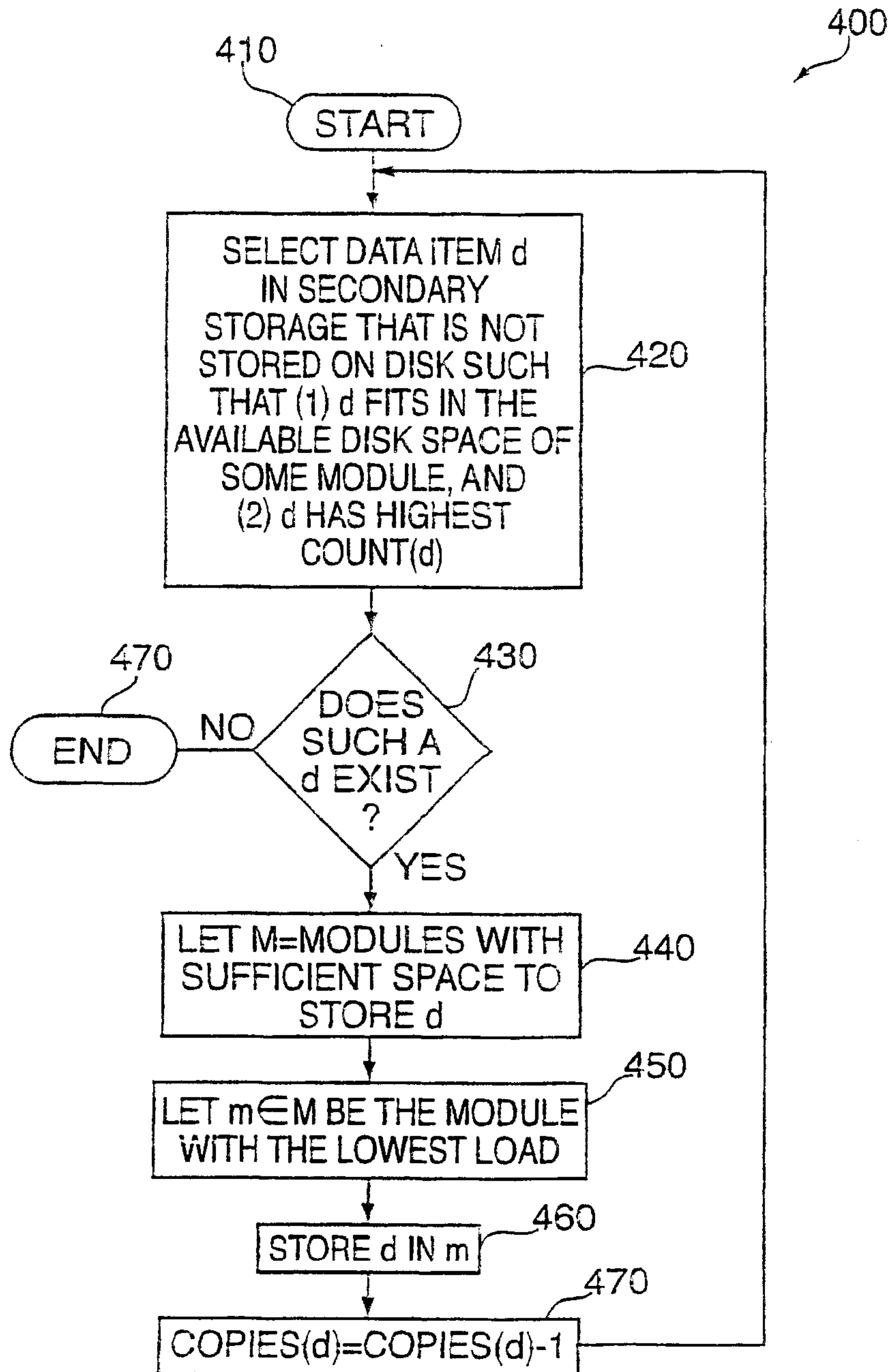


FIG. 4

5/12

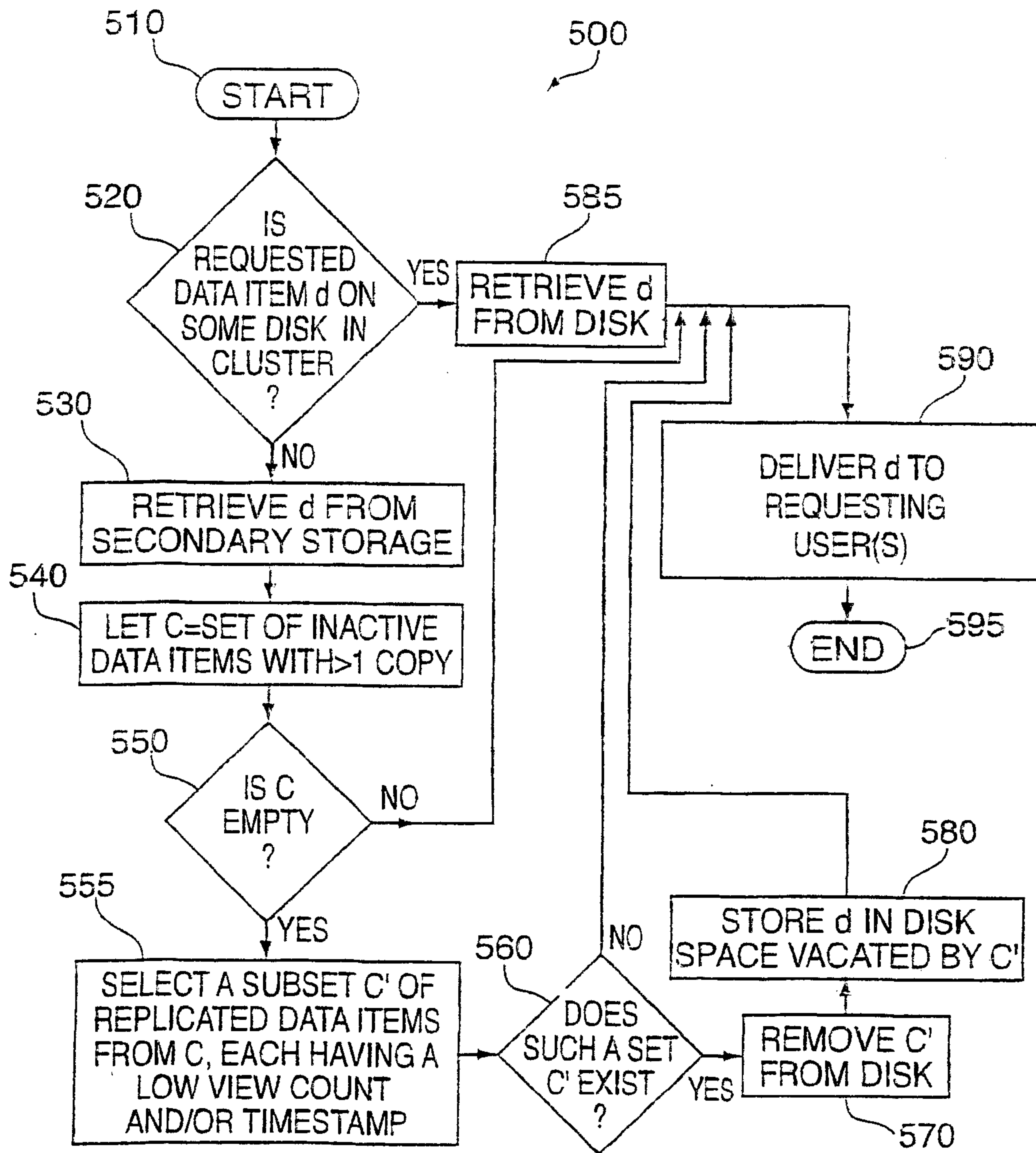


FIG. 5

6/12

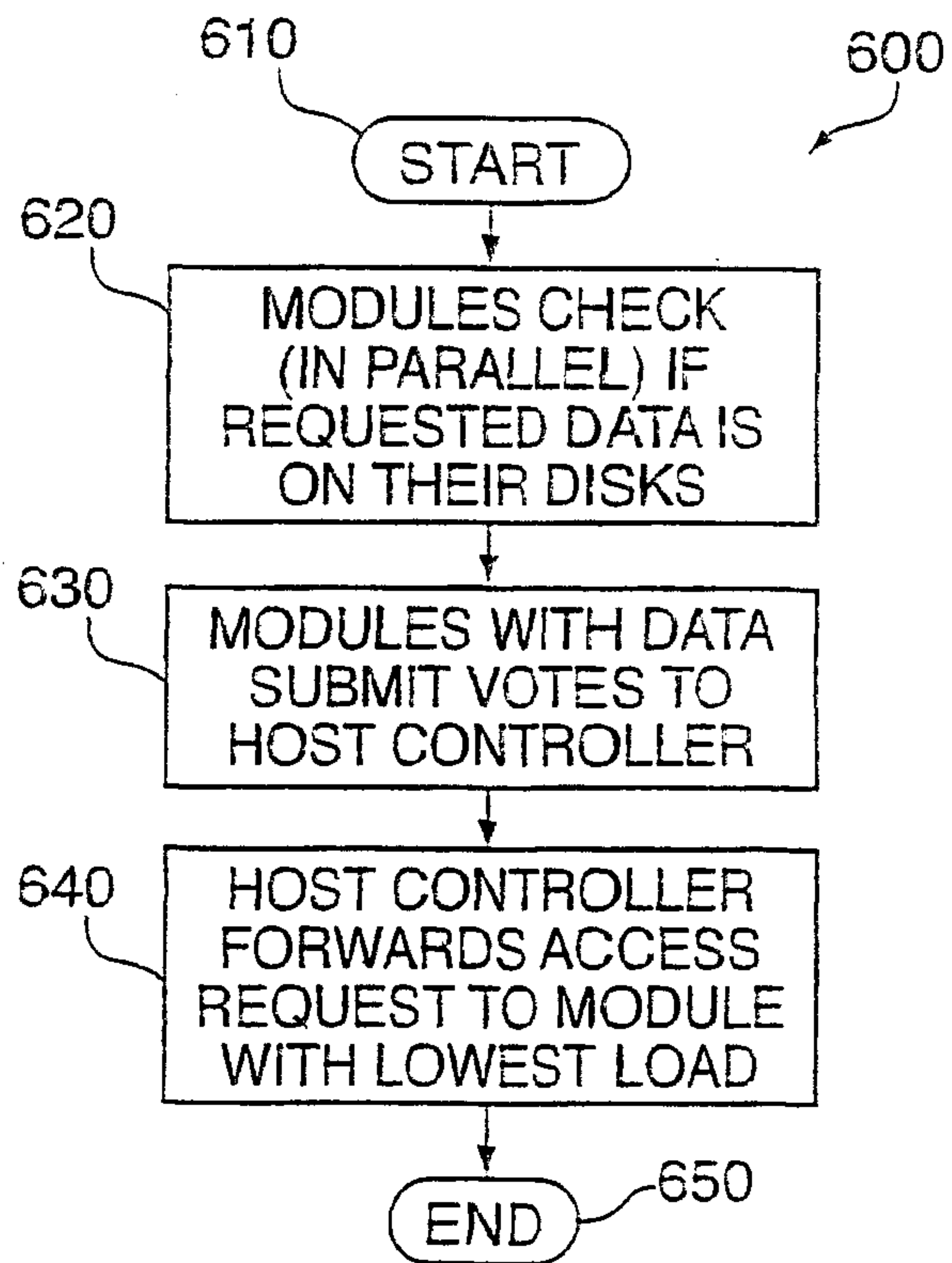


FIG. 6

7/12

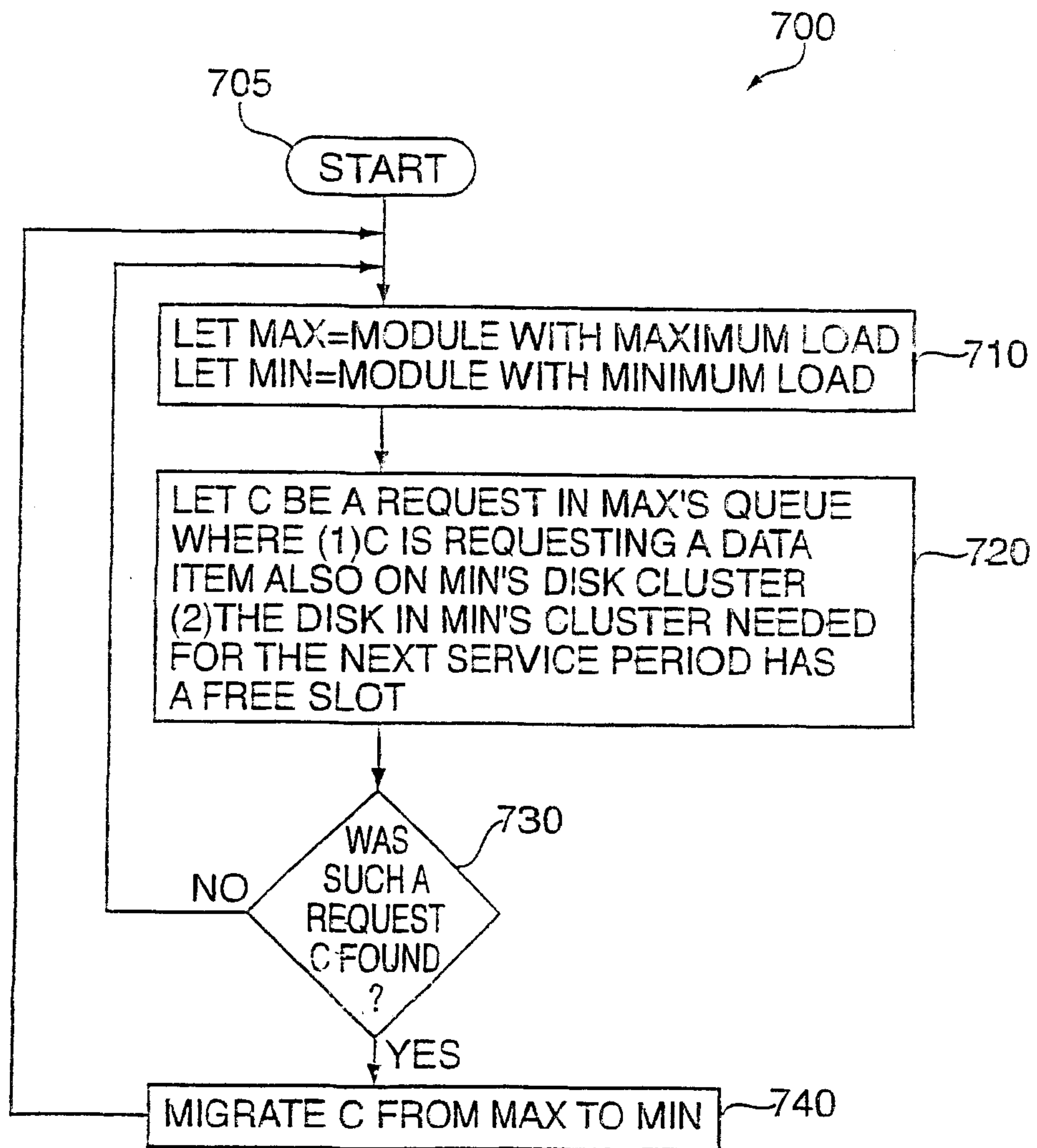


FIG. 7

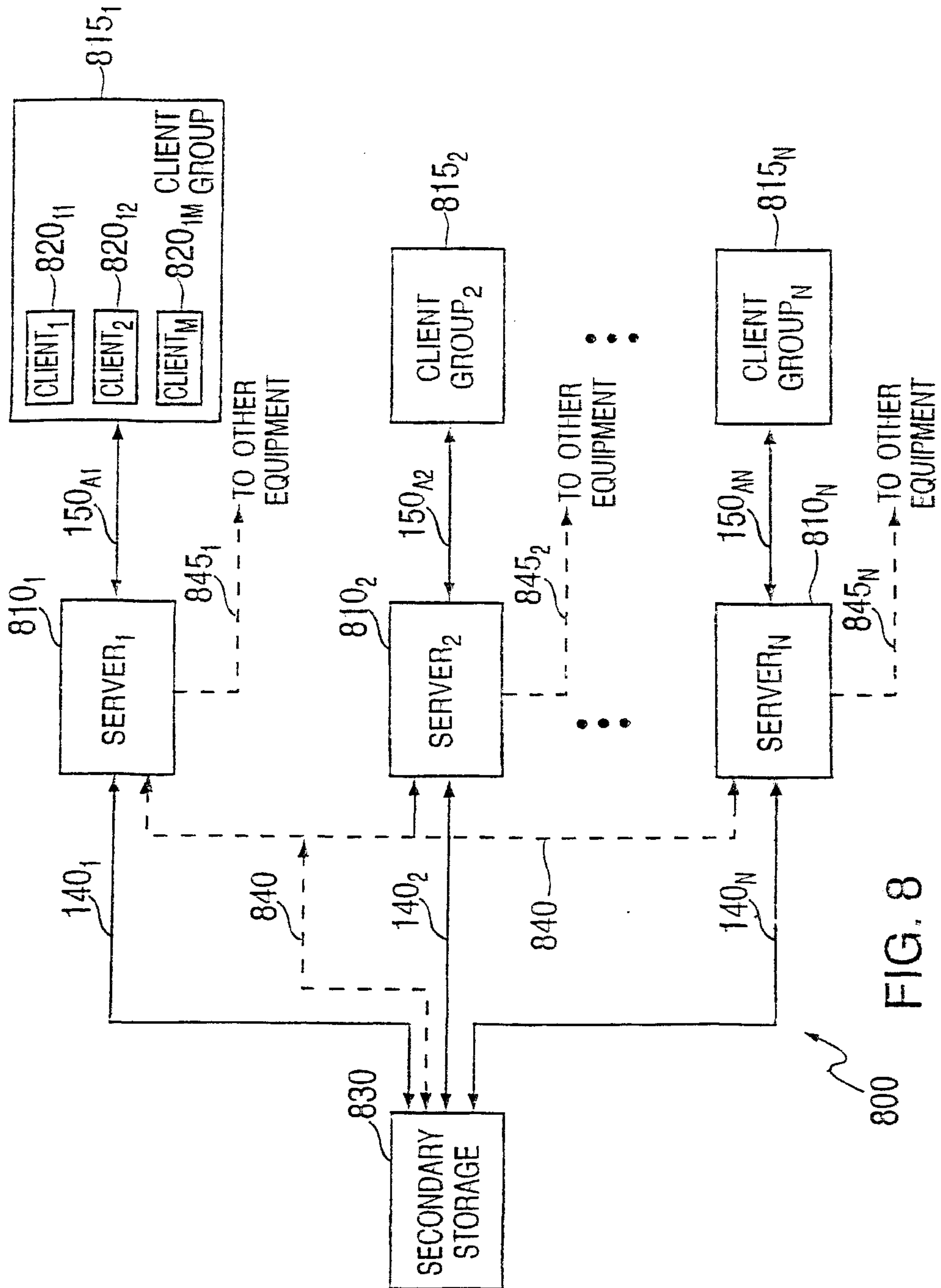


FIG. 8

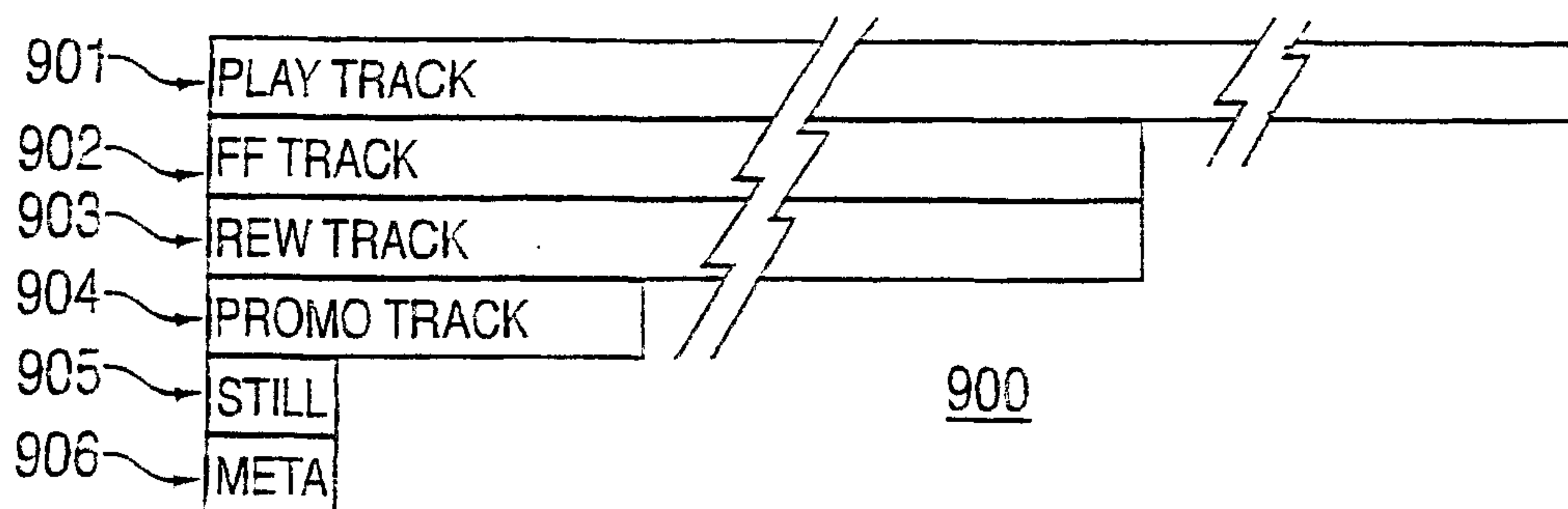


FIG. 9

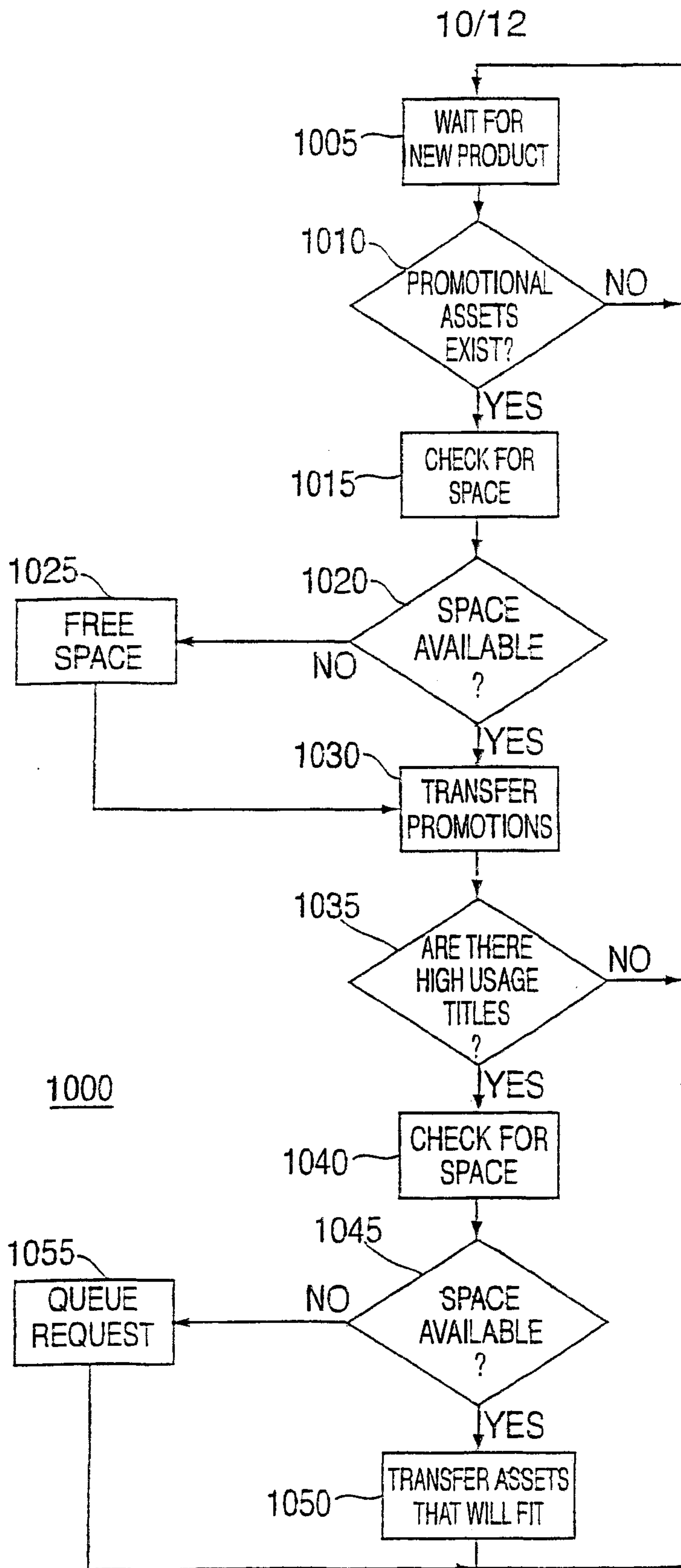


FIG. 10

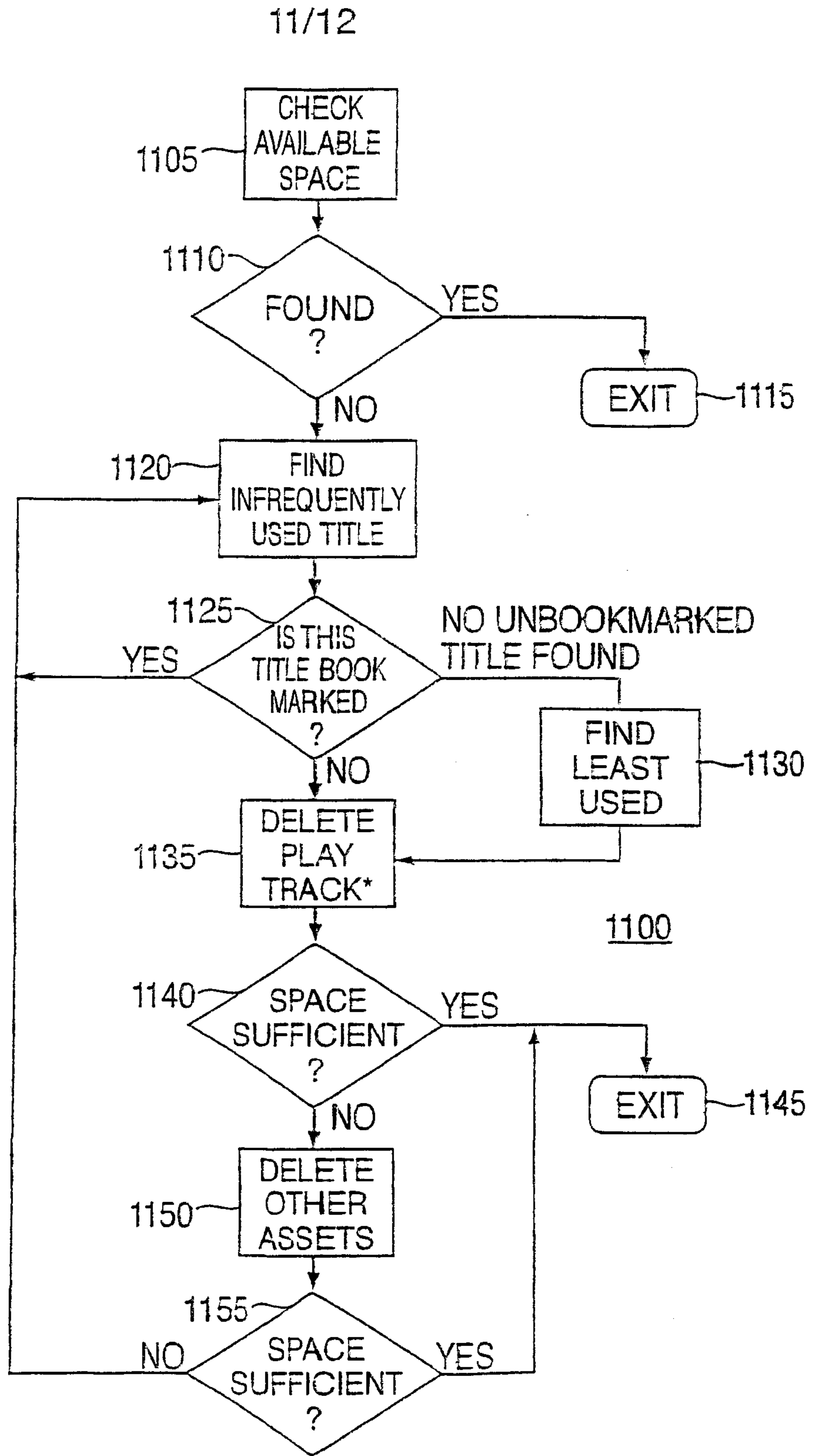


FIG. 11

12/12

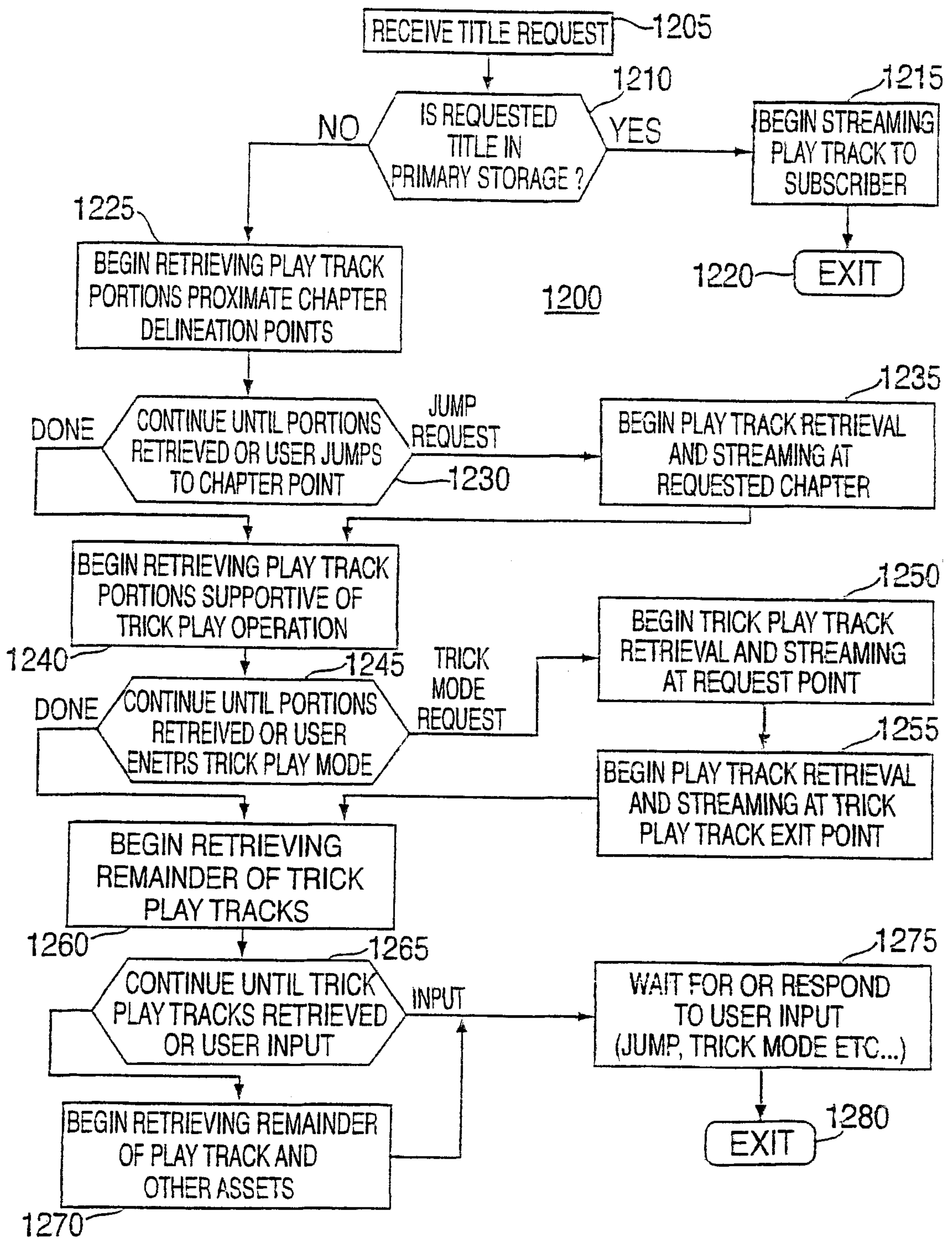


FIG. 12

