

[19] 中华人民共和国国家知识产权局

[51] Int. Cl.

G06F 9/46 (2006.01)

G06F 13/00 (2006.01)



[12] 发明专利说明书

专利号 ZL 200410064190. X

[45] 授权公告日 2007 年 2 月 14 日

[11] 授权公告号 CN 1300688C

[22] 申请日 2004. 8. 24

[21] 申请号 200410064190. X

[30] 优先权

[32] 2003. 8. 29 [33] US [31] 10/652,021

[73] 专利权人 国际商业机器公司

地址 美国纽约

[72] 发明人 肯尼思·W·博伊德

柯尔比·G·达曼 肯尼思·F·戴

菲利普·M·道特玛斯

约翰·J·沃尔夫冈

[56] 参考文献

US6549992B1 2003. 4. 15

US6145034A 2000. 11. 7

JP2002244880A 2002. 8. 30

审查员 刘静_2

[74] 专利代理机构 中国国际贸易促进委员会专利
商标事务所

代理人 康建峰

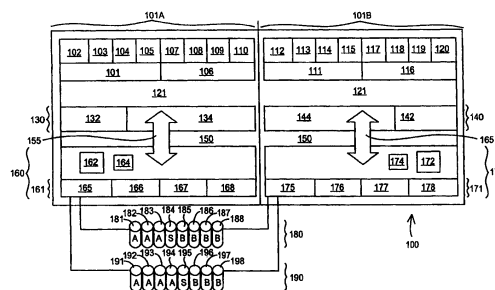
权利要求书 3 页 说明书 20 页 附图 7 页

[54] 发明名称

调整数据传输速率的装置和方法

[57] 摘要

一种用来调整(N)个初级备份设备之一的数据传输速率的方法。该方法由第一初级备份设备形成至少一个一致事务集。第一初级备份设备从其他(N-1)个初级备份设备中的每一个接收第(n)状态信号以及第(n+1)状态信号。该方法计算(N)个初级备份设备中的每一个的第(n)有效带宽、(N)个初级备份设备中的每一个的第(n)完成时间、以及所有(N)个初级设备的第(n)有效合计带宽。如果第一初级备份设备的第(n)完成时间大于其他(N-1)个初级备份设备中的每一个的第(n)完成时间,则该方法无延迟地将来自第一初级备份设备的至少一个一致事务集提供给第一次级备份设备。



1. 一种用来调整(N)个初级备份设备之一的数据传输速率的方法, 其中所述(N)个初级备份设备中的每一个能够与一个或多个第一数据存储和检索系统以及一个或多个次级备份设备通信, 其中所述一个或多个次级备份设备能够与一个或多个第二数据存储和检索系统通信, 所述方法包括以下步骤:

提供所述(N)个初级备份设备中的第一初级备份设备, 其中, 所述第一初级备份设备能够与所述一个或多个次级备份设备中的第一次级备份设备通信;

由所述第一初级备份设备形成至少一个包括从所述一个或多个第一数据存储和检索系统接收的信息的一致事务集;

从其他(N-1)个初级备份设备中的每一个接收第(n)状态信号;

从其他(N-1)个初级备份设备中的每一个接收第(n+1)状态信号;

计算所述(N)个初级备份设备中的每一个的第(n)有效带宽;

计算所述(N)个初级备份设备中的每一个的第(n)完成时间;

计算所有(N)个初级设备的第(n)有效合计带宽;

确定所述第一初级备份设备的第(n)完成时间是否大于其他(N-1)个初级备份设备中的每一个的第(n)完成时间;

如果所述第一初级备份设备的第(n)完成时间大于其他(N-1)个初级备份设备中的每一个的第(n)完成时间, 则操作以便无延迟地将来自所述第一初级备份设备的至少一个一致事务集的全部或部分提供给所述第一次级备份设备。

2. 如权利要求 1 所述的方法, 还包括以下步骤:

如果所述第一初级备份设备的第(n)完成时间不大于其他(N-1)个初级设备中的每一个的第(n)完成时间, 则操作以便确定第(n)合计带宽是否小于第(n-1)合计带宽;

如果第(n)合计带宽不小于第(n-1)合计带宽, 则操作以便确立所述第一初级备份设备的第(n)延迟;

使用所述第(n)延迟将来自所述第一初级备份设备的至少一个一致事务集的全部或部分提供给所述次级备份设备。

3. 如权利要求2所述的方法, 还包括以下步骤:

如果第(n)合计带宽小于第(n-1)合计带宽, 则操作以便确定所述第一初级备份设备的第(n-1)延迟是否大于所述第一初级备份设备的第(n-2)延迟;

如果所述第一初级备份设备的第(n-1)延迟大于所述第一初级备份设备的第(n-2)延迟, 则操作以便设置所述第一初级备份设备的第(n)延迟等于所述第一初级备份设备的第(n-2)延迟;

使用所述第(n)延迟将来自所述第一初级备份设备的至少一个一致事务集的全部或部分提供给所述次级备份设备。

4. 如权利要求3所述的方法, 还包括以下步骤:

如果所述第一初级备份设备的第(n-1)延迟不大于所述第一初级备份设备的第(n-2)延迟, 则操作以便确立所述第一初级备份设备的第(n)延迟;

使用所述第一初级备份设备的所述第(n)延迟将来自所述第一初级备份设备的至少一个一致事务集的全部或部分提供给所述次级备份设备。

5. 如权利要求2所述的方法, 其中所述确立步骤还包括以下步骤:

确定所有(N)个初级备份设备的第(n)平均完成时间;

确定所述第一初级备份设备的第(n)完成时间是否基本上等于所述第(n)平均完成时间;

如果所述第一初级备份设备的第(n)完成时间基本上等于所述第(n)平均完成时间, 则操作以便设置所述第一初级备份设备的第(n)延迟等于所述第一初级备份设备的第(n-1)延迟。

6. 如权利要求5所述的方法, 其中所述确立步骤还包括以下步骤:

提供标准延迟调整;

确定所述第一初级备份设备的第(n)完成时间是否大于所述第(n)平均完成时间;

如果所述第一初级备份设备的第(n)完成时间大于所述第(n)平均完成

时间，则操作以便设置所述第一初级备份设备的所述第(n)延迟等于所述第一初级备份设备的第(n-1)延迟减去所述标准延迟调整。

7. 如权利要求 6 所述的方法，还包括以下步骤：

如果所述第一初级备份设备的第(n)完成时间小于所述第(n)平均完成时间，则操作以便设置所述第一初级备份设备的所述第(n)延迟等于所述第一初级备份设备的第(n-1)延迟加上所述标准延迟调整。

8. 如权利要求 5 所述的方法，其中所述确立步骤还包括以下步骤：

提供延迟调整函数；

确定所述第(n)平均完成时间与所述第一初级备份设备的第(n)完成时间之间的差值；

根据所述差值和所述延迟调整函数确立所述第一初级备份设备的第(n)延迟。

9. 如权利要求 5 所述的方法，其中所述确立步骤还包括以下步骤：

提供标准延迟调整；

使用所述第(n)平均完成时间、(N-1)个其余初级备份设备中的每一个的第(n)完成时间和所述标准延迟调整，预测所述(N-1)个其余初级备份设备中的每一个的第(n)延迟；

根据所述(N-1)个其余初级备份设备中的每一个的预测延迟确立第一初级备份设备的第(n)延迟。

10. 如权利要求 2 所述的方法，其中所述确立步骤还包括以下步骤：

提供延迟调整函数；

使用所述延迟调整函数预测(N-1)个延迟值，其中所述预测延迟值中的每一个与其余(N-1)个初级备份设备不同之一相关联；

根据所述(N-1)个预测延迟值确立第一初级备份设备的第(n)延迟。

调整数据传输速率的装置和方法

技术领域

本发明涉及一种调整由多个备份设备中的每一个使用的数据传输速率的装置和方法。

背景技术

很多数据处理系统需要大量数据存储以用于高效存取、修改和再存储数据。数据存储典型地分成若干不同级别，每一个级别显现不同的数据存取时间或数据存储成本。第一或最高层数据存储涉及电子存储器，通常是动态或静态随机存取存储器(DRAM 或 SRAM)。电子存储器采取半导体集成电路的形式，其中数百万字节的数据可以存储在每个电路上，其中对这些数据字节的存取以纳秒测量。由于存取是完全电子式的，因此电子存储器提供最快的数据存取。

第二级数据存储通常涉及直接存取存储设备(DASD)。DASD 存储例如包括磁盘和/或光盘。数据比特作为盘表面上微米大小的磁性或光学改变的斑点来存储，从而表示组成数据比特二进制值的“一”和“零”。磁性 DASD 包括覆盖有残余磁性材料的一个或多个盘。这些盘旋转性地安装在受保护环境内。每个盘分成很多同心轨道或者紧密圆圈。数据沿着每个轨道逐比特地连续存储。

具有备份数据副本对于数据丢失将是灾难性的很多商业机构而言是强制性的。恢复丢失数据所需的时间也是重要的恢复考虑。通过磁带或库备份，初级数据周期性地通过在磁带或库存储装置上生成副本来备份。

另外，还需要保护以在整个系统或者甚至是场所被诸如地震、火灾、爆炸、飓风等的灾难破坏的情况下恢复数据。典型数据处理系统的灾难恢复保护要求在次级或远程位置上备份存储在初级 DASD 上的初级

数据。初级和次级位置之间的物理相隔距离可以根据对于用户可接受的风险级别来设置，并且可以从数公里变至数千公里。

次级场所不仅必须足够远离开于初级场所，而且必须能够实时备份初级数据。当初级数据被更新时，次级场所需要实时备份初级数据，其中只有某一极小的延迟。次级场所所需的困难任务在于次级数据必须是“次序一致”的，也就是，次级数据以需要大量系统考虑的与初级数据相同的顺序次序(顺序一致性)来拷贝。顺序一致性由于在数据处理系统中存在均控制多个 DASD 的多个存储控制器而复杂化。在没有顺序一致性的情况下，将产生与初级数据不一致的次级数据，从而破坏灾难恢复。

在某些数据处理应用中，将信息提供给一个或多个形成一个或多个一致事务集的初级备份设备。这些一个或多个初级备份设备通常位于初级存储场所或其附近。周期性地，一个或多个初级备份设备中的每一个通过公共通信链路将一致事务集的全部或部分提供给位于一个或多个远程存储场所的一个或多个次级备份设备。所需的是一种自主调整多个备份设备中的每一个的数据传输速率以最大化利用公共通信链路的可用数据传输带宽的方法。

发明内容

本申请人的发明包括一种用来调整(N)个初级备份设备之一的数据传输速率的装置和方法，其中这些(N)个初级备份设备中的每一个能够与一个或多个第一数据存储和检索系统以及第二备份设备通信，其中第二备份设备能够与一个或多个第二数据存储和检索系统通信。本申请人的方法提供所述(N)个初级备份设备中的第一初级备份设备，其中，该第一初级备份设备能够与多个次级备份设备中的第一次级备份设备通信。

该方法由第一初级备份设备形成至少一个包括从一个或多个第一数据存储和检索系统接收的信息的一致事务集。第一初级备份设备从其他(N-1)个初级备份设备中的每一个接收第(n)状态信号，然后从其他(N-1)个初级备份设备中的每一个接收第(n+1)状态信号。

该方法计算(N)个初级备份设备中的每一个的第(n)有效带宽，计算

(N)个初级备份设备中的每一个的第(n)完成时间, 并且计算所有(N)个初级设备的第(n)有效合计带宽。该方法然后确定第一初级备份设备的第(n)完成时间是否大于其他(N-1)个初级备份设备中的每一个的第(n)完成时间。如果第一初级备份设备的第(n)完成时间大于其他(N-1)个初级备份设备中的每一个的第(n)完成时间, 则该方法无延迟地将来自第一初级备份设备的至少一个一致事务集的全部或部分提供给第一次级备份设备。

附图说明

通过阅读下面结合附图的详细描述, 本发明将会得到更好的理解, 其中, 相同的附图标记用来指定相同的单元, 并且其中:

图 1 是示出本申请人的数据存储和检索系统的一个实施例的各组件的方框图;

图 2 是示出本申请人的数据存储和检索系统的第二实施例的各组件的方框图;

图 3 是示出本申请人的数据存储和检索系统的第三实施例的各组件的方框图;

图 4 是示出本申请人的远程拷贝数据存储和检索系统的各组件的方框图;

图 5 是概述本申请人的方法中的特定初始步骤的流程图;

图 6 是概述本申请人的方法中的特定附加步骤的流程图; 以及

图 7 是概述本申请人的方法中的特定附加步骤的流程图。

具体实施方式

参照附图在下面描述中以多个优选实施例描述本发明, 其中相同的标号表示相同或类似的单元。

图 4 示出本申请人的系统的各组件。现在参照图 4, 主机计算机 390 通过通信链路 402 与初级数据存储和检索系统 410、430 和 450 进行互连和通信。在某些实施例中, 通信链路 402 从包括串行互连如 RS-232 电缆或 RS-432 电缆、以太网互连、SCSI 互连、光纤通道互连、

ESCON 互连、FICON 互连、局域网(LAN)、私有广域网(WAN)、公用广域网、存储区域网(SAN)、传输控制协议/网际协议(TCP/IP)、因特网及其组合的组中选择。

初级数据存储和检索系统 410 将信息从初级信息存储介质 413 提供到次级数据存储和检索系统 425 以通过初级备份设备 415 和次级备份设备 420 拷贝到次级信息存储介质 428。信息存储和检索系统 410 还包括控制器 411 以及可选地数据高速缓冲存储器 412。信息存储和检索系统 425 还包括控制器 426 以及可选地数据高速缓冲存储器 427。

在某些实施例中，信息存储介质 413 包括 DASD。在某些实施例中，信息存储介质 413 包括一个或多个 RAID 阵列。在某些实施例中，信息存储介质 413 包括多个便携式信息存储介质，例如包括多个单独位于便携式容器例如磁带盒中的磁带。

在某些实施例中，信息存储介质 428 包括 DASD。在某些实施例中，信息存储介质 428 包括一个或多个 RAID 阵列。在某些实施例中，信息存储介质 428 包括多个便携式信息存储介质，例如包括多个单独位于便携式容器例如磁带盒中的磁带。

在某些实施例中，初级备份设备 415 与初级数据存储和检索系统 410 集成在一起。在图 4 的所示实施例中，初级备份设备 415 居于初级数据存储和检索系统 410 的外部，并且通过通信链路 414 与初级数据存储和检索系统 410 通信。在某些实施例中，通信链路 414 从包括串行互连如 RS-232 电缆或 RS-432 电缆、以太网互连、SCSI 互连、光纤通道互连、ESCON 互连、FICON 互连、局域网(LAN)、私有广域网(WAN)、公用广域网、存储区域网(SAN)、传输控制协议/网际协议(TCP/IP)、因特网及其组合的组中选择。

在某些实施例中，次级备份设备 420 与次级数据存储和检索系统 425 集成在一起。在图 4 的所示实施例中，次级备份设备 420 居于次级数据存储和检索系统 425 的外部，并且通过通信链路 429 与次级数据存储和检索系统 425 通信。在某些实施例中，通信链路 429 从包括串行互连如 RS-232 电缆或 RS-432 电缆、以太网互连、SCSI 互连、光纤通道

互连、ESCON 互连、FICON 互连、局域网(LAN)、私有广域网(WAN)、公用广域网、存储区域网(SAN)、传输控制协议/网际协议(TCP/IP)、因特网及其组合的组中选择。

初级数据存储和检索系统 430 将信息从初级信息存储介质 433 提供到次级数据存储和检索系统 445 以通过初级备份设备 435 和次级备份设备 440 拷贝到次级信息存储介质 448。信息存储和检索系统 430 还包括控制器 431 以及可选地数据高速缓冲存储器 432。信息存储和检索系统 445 还包括控制器 446 以及可选地数据高速缓冲存储器 447。

在某些实施例中，信息存储介质 433 包括 DASD。在某些实施例中，信息存储介质 433 包括一个或多个 RAID 阵列。在某些实施例中，信息存储介质 433 包括多个便携式信息存储介质，例如包括多个单独位于便携式容器例如磁带盒中的磁带。

在某些实施例中，信息存储介质 448 包括 DASD。在某些实施例中，信息存储介质 448 包括一个或多个 RAID 阵列。在某些实施例中，信息存储介质 448 包括多个便携式信息存储介质，例如包括多个单独位于便携式容器例如磁带盒中的磁带。

在某些实施例中，初级备份设备 435 与初级数据存储和检索系统 430 集成在一起。在图 4 的所示实施例中，初级备份设备 435 居于初级数据存储和检索系统 430 的外部，并且通过通信链路 434 与初级数据存储和检索系统 430 通信。在某些实施例中，通信链路 434 从包括串行互连如 RS-232 电缆或 RS-432 电缆、以太网互连、SCSI 互连、光纤通道互连、ESCON 互连、FICON 互连、局域网(LAN)、私有广域网(WAN)、公用广域网、存储区域网(SAN)、传输控制协议/网际协议(TCP/IP)、因特网及其组合的组中选择。

在某些实施例中，次级备份设备 440 与次级数据存储和检索系统 445 集成在一起。在图 4 的所示实施例中，次级备份设备 440 居于次级数据存储和检索系统 445 的外部，并且通过通信链路 449 与次级数据存储和检索系统 445 通信。在某些实施例中，通信链路 449 从包括串行互连如 RS-232 电缆或 RS-432 电缆、以太网互连、SCSI 互连、光纤通道

互连、ESCON 互连、FICON 互连、局域网(LAN)、私有广域网(WAN)、公用广域网、存储区域网(SAN)、传输控制协议/网际协议(TCP/IP)、因特网及其组合的组中选择。

初级数据存储和检索系统 450 将信息从初级信息存储介质 453 提供到次级数据存储和检索系统 465 以通过初级备份设备 455 和次级备份设备 460 拷贝到次级信息存储介质 468。信息存储和检索系统 450 还包括控制器 451 以及可选地数据高速缓冲存储器 452。信息存储和检索系统 465 还包括控制器 466 以及可选地数据高速缓冲存储器 467。

在某些实施例中，信息存储介质 453 包括 DASD。在某些实施例中，信息存储介质 453 包括一个或多个 RAID 阵列。在某些实施例中，信息存储介质 453 包括多个便携式信息存储介质，例如包括多个单独位于便携式容器例如磁带盒中的磁带。

在某些实施例中，信息存储介质 468 包括 DASD。在某些实施例中，信息存储介质 468 包括一个或多个 RAID 阵列。在某些实施例中，信息存储介质 468 包括多个便携式信息存储介质，例如包括多个单独位于便携式容器例如磁带盒中的磁带。

在某些实施例中，初级备份设备 455 与初级数据存储和检索系统 450 集成在一起。在图 4 的所示实施例中，初级备份设备 455 居于初级数据存储和检索系统 450 的外部，并且通过通信链路 454 与初级数据存储和检索系统 450 通信。在某些实施例中，通信链路 454 从包括串行互连如 RS-232 电缆或 RS-452 电缆、以太网互连、SCSI 互连、光纤通道互连、ESCON 互连、FICON 互连、局域网(LAN)、私有广域网(WAN)、公用广域网、存储区域网(SAN)、传输控制协议/网际协议(TCP/IP)、因特网及其组合的组中选择。

在某些实施例中，次级备份设备 460 与次级数据存储和检索系统 465 集成在一起。在图 4 的所示实施例中，次级备份设备 460 居于次级数据存储和检索系统 465 的外部，并且通过通信链路 469 与次级数据存储和检索系统 465 通信。在某些实施例中，通信链路 469 从包括串行互连如 RS-232 电缆或 RS-452 电缆、以太网互连、SCSI 互连、光纤通道

互连、ESCON 互连、FICON 互连、局域网(LAN)、私有广域网(WAN)、公用广域网、存储区域网(SAN)、传输控制协议/网际协议(TCP/IP)、因特网及其组合的组中选择。

初级备份设备 415、435 和 455 分别从初级数据存储和检索系统 410、430 和 450 接收信息。可替换地，任何初级备份设备可以从任何初级数据存储和检索系统接收信息。周期性地，每个初级备份设备形成一致事务集。在此所用的“一致事务集”是指这样一个事务集，即当在次级数据存储和检索系统控制器上应用该事务集中的所有事务时，次级存储看上去将相同于创建该事务集的时间点上的初级存储。

在某些实施例中，数据存储和检索系统 410、425、430、445、450 和/或 465 中的一个或多个包括数据存储和检索系统 100(图 1)。现在参照图 1，本申请人的信息存储和检索系统 100 包括第一群集 101A 和第二群集 101B。每个群集包括处理器部分 130/140 和输入/输出部分 160/170。每个群集的内部 PCI 总线分别通过远程 I/O 桥 155/165 连接在处理器部分 130/140 与 I/O 部分 160/170 之间。

信息存储和检索系统 100 还包括位于四个主机舱(host bay)101、106、111 和 116 内的多个主机适配器 102-105、107-110、112-115 和 117-120。每个主机适配器可以包括一个光纤通道端口、一个 FICON 端口、两个 ESCON 端口或者两个 SCSI 端口。每个主机适配器通过一个或多个公共平台互连总线 121 和 150 连接到两个群集，使得每个群集可以处理来自任何主机适配器的 I/O。

处理器部分 130 包括处理器 132 和高速缓冲存储器 134。在某些实施例中，处理器 132 包括基于 64 比特 RISC 的对称多处理器。在某些实施例中，处理器 132 包括内置故障和错误纠正功能。高速缓冲存储器 134 用来存储读取和写入数据以改善所附主机系统的性能。在某些实施例中，高速缓冲存储器 134 包括大约 4 吉字节。在某些实施例中，高速缓冲存储器 134 包括大约 8 吉字节。在某些实施例中，高速缓冲存储器 134 包括大约 12 吉字节。在某些实施例中，高速缓冲存储器 134 包括大约 16 吉字节。在某些实施例中，高速缓冲存储器 134 包括大约 32 吉字

节。

处理器部分 140 包括处理器 142 和高速缓冲存储器 144。在某些实施例中，处理器 142 包括基于 64 比特 RISC 的对称多处理器。在某些实施例中，处理器 142 包括内置故障和错误纠正功能。高速缓冲存储器 144 用来存储读取和写入数据以改善所附主机系统的性能。在某些实施例中，高速缓冲存储器 144 包括大约 4 吉字节。在某些实施例中，高速缓冲存储器 144 包括大约 8 吉字节。在某些实施例中，高速缓冲存储器 144 包括大约 12 吉字节。在某些实施例中，高速缓冲存储器 144 包括大约 16 吉字节。在某些实施例中，高速缓冲存储器 144 包括大约 32 吉字节。

I/O 部分 160 包括非易失性存储装置(“NVS”)162 和 NVS 电池 164。NVS 162 用来存储写入数据的第二副本以在发生群集故障的电源故障和该数据的高速缓冲存储器副本丢失的情况下确保数据完整性。NVS 162 存储提供给群集 101B 的写入数据。在某些实施例中，NVS 162 包括大约 1 吉字节的存储装置。在某些实施例中，NVS 162 包括四个独立存储卡。在某些实施例中，每对 NVS 卡具有由电池供电的充电系统，即使整个系统掉电，其也在高达 72 小时内保护数据。

I/O 部分 170 包括 NVS 172 和 NVS 电池 174。NVS 172 存储提供给群集 101A 的写入数据。在某些实施例中，NVS 172 包括大约 1 吉字节的存储装置。在某些实施例中，NVS 172 包括四个独立存储卡。在某些实施例中，每对 NVS 卡具有由电池供电的充电系统，即使整个系统断电，其也在高达 72 小时内保护数据。

在群集 101B 发生故障的情况下，故障群集的写入数据将驻留在位于正常群集 101A 内的 NVS 162 中。然后，该写入数据以高优先级降级 (destage) 到 RAID 等级(rank)。同时，正常群集 101A 将开始对于其自己的写入数据使用 NVS 162，从而确保仍然保持写入数据的两个副本。

I/O 部分 160 还包括多个设备适配器如设备适配器 165、166、167 和 168 和组织成两个 RAID 等级即 RAID 等级“A”和 RAID 等级“B”的十六个盘驱动器。在某些实施例中，RAID 等级“A”和“B”利用 RAID 5

协议。在某些实施例中，RAID 等级“A”和“B”利用 RAID 10 协议。

在某些实施例中，数据存储和检索系统 410、425、430、445、450 和/或 465 中的一个或多个包括数据存储和检索系统 200(图 2)。图 2 示出系统 200 的一个实施例。

系统 200 被安排用于响应来自一个或多个主机系统如主机计算机 390(图 4)的命令而存取便携式数据存储介质。系统 200 包括前壁 270 和后壁 290 上的多个存储架 260，用于存储容纳数据存储介质的便携式数据存储盒。系统 200 还包括：至少一个数据存储驱动器 250，用于对数据存储介质进行数据读取和/或写入；以及至少一个存取器(accessor)210，用于在多个存储架 260 与数据存储驱动器 250 之间运输数据存储介质。系统 200 可以可选地包括操作员面板 230 或其他用户接口如基于万维网(web)的接口，其允许用户与存储库进行交互。系统 200 可以可选地包括上方导入/导出台 240 和/或下方导入/导出台 245，其允许将数据存储介质插入到存储库中并且/或者从存储库中移走数据存储介质而不打断存储库操作。

存取器 210 包括升降伺服部件 212，其能够沿着 Z 轴进行双向移动。存取器 210 还包括至少一个机械爪组件(gripper assembly)216，其用于抓握(gripping)一个或多个数据存储介质。在图 2 的所示实施例中，存取器 210 还包括条形码扫描器 214 或其他阅读系统如智能卡阅读器或类似系统，以“阅读”有关数据存储介质的标识信息。在图 2 的所示实施例中，存取器 210 还包括位于升降伺服部件 212 上的第二机械爪机构 218。

在某些实施例中，系统 200 包括一个或多个存储框架(storage frame)，其中每一个都具有可由存取器 210 存取的存储架 260。存取器 210 在导轨(rail)205 上沿着 X 轴双向移动。在包括多个框架的存储库 100 的实施例中，这些单独框架的每一个中的导轨 205 被对齐成使得存取器 210 可以沿着邻接导轨系统从存储库的一端行驶到另一端。

在某些实施例中，数据存储和检索系统 410、425、430、445、450 和/或 465 中的一个或多个包括数据存储和检索系统 300(图 3)。现在参

照图 3, 虚拟磁带服务器 300(“VTS”)300 通过后台程序(daemon)370、372 和 374 与一个或多个主机以及一个或多个虚拟磁带服务器通信。在图 3 的所示实施例中, 后台程序 370 通过通信链路 380 与第一主机通信。在图 3 的所示实施例中, 后台程序 372 通过通信链路 382 与第二主机通信。后台程序 374 通过通信链路 384 与例如初级备份设备如设备 415 通信。

VTS 300 还与直接存取存储设备(DASD)310、多个数据存储设备 330 和 340 通信。在某些实施例中, 数据存储设备 330 和 340 位于一个或多个数据存储和检索系统内。在某些实施例中, DASD 310 与主机 110 集成在一起(图 1)。在某些实施例中, DASD 310 与 VTS 300 集成在一起。在某些实施例中, DASD 310 与数据存储和检索系统集成在一起。在某些实施例中, DASD 310 外部于主机 110、VTS 300 以及与 VTS 300 通信的一个或多个数据存储和检索系统。

VTS 300 还包括存储管理器 320 如 IBM Adstar®分布式存储管理器。存储管理器 320 控制从 DASD 310 到安装在数据存储设备 330 和 340 中的信息存储介质的数据移动。在某些实施例中, 存储管理器 320 包括 ADSM 服务器 322 和 ADSM 分级式存储管理器客户端 324。可替换地, 服务器 322 和客户端 324 均可包括 ADSM 系统。来自 DASD 310 的信息通过 ADSM 服务器 322 和 SCSI 适配器 385 提供到数据存储设备 330 和 340。

VTS 300 还包括自主控制器 350。自主控制器 350 通过分级式存储管理器(HSM)客户端 324 控制 DASD 310 的操作, 以及 DASD 310 与数据存储设备 130 和 140 之间的数据传输。

回到图 4, 每个初级备份设备以不同于其他初级备份设备的速率从各个不同初级存储控制器接收数据。在这种情况下, 由初级备份设备形成的一致事务集的大小可能变化很大。序列号为 10/339,957、名称为“Method, System and Article of Manufacture for Creating a Consistent Copy”且转让给其共同受让人的未决专利申请描述了一种形成一致事务集的方法, 并且在此将其全文引作参考。

初级备份设备如设备 415、425 和 435 均通过公共通信链路如通信链路 470 分别提供一致事务集给其对应的次级备份设备如设备 420、430 和 440。在某些实施例中，通信链路 470 从包括串行互连如 RS-232 电缆或 RS-432 电缆、以太网互连、SCSI 互连、光纤通道互连、ESCON 互连、FICON 互连、局域网(LAN)、私有广域网(WAN)、公用广域网、存储区域网(SAN)、传输控制协议/网际协议(TCP/IP)、因特网及其组合的组中选择。

为了优化从多个初级备份设备如初级设备 415、435 和 455 到多个次级备份设备如设备 420、440 和 460 的数据传输，通信链路 470 的带宽应当保持完全利用。在某些实施例中，没有一个初级备份设备可以完全利用公共通信链路即通信链路 470 的带宽。在这些实施例中，最好是多个初级备份设备在提供一致事务集给一个或多个次级备份设备时利用通信链路。

另外，最好是每一个初级备份设备以近似相同的时间完成一致事务集的传输，因为这些传输是逐事务集地发生的。因此，让一个初级备份设备在其余初级设备完成其一致事务集传输之前完成一致事务集传输是无益的。

如果由于第一初级备份设备未被分配通信链路 470 的足够带宽而该第一设备的一致事务集传输完成时间("TTC")超过其余初级备份设备的 TTC，则该不同 TTC 将会不利地影响次级备份设备提供信息给次级数据存储和检索系统的速率。使用本申请人的方法，初级备份设备中的每一个使用互连这些初级备份设备中的每一个与多个次级备份设备的公共通信链路自主调整其 TTC。

本申请人的发明包括一种在本申请人的系统中自主调整每个初级备份设备的数据传输速率的方法。图 5 概述了自主调整(N)个初级备份设备中的每一个的数据传输速率的本申请人方法的各步骤。为了描述起见，图 5、6 和 7 的步骤在下面被描述为由一个初级备份设备即第一初级设备执行。在实现中，图 5、6 和/或 7 的步骤由(N)个初级备份设备中的每一个独立地即自主地执行。

现在参照图 5, 在步骤 510, 第一初级备份设备从其余(N-1)个初级备份设备中的每一个接收状态信号。每个初级备份设备包括第一初级备份设备周期性地有时称作“心跳”信号的状态信号发送到其他初级备份设备中的每一个。

这些状态信号中的每一个均包括通过公共通信链路如互连(N)个初级备份设备与次级备份设备中的每一个的通信链路 470(图 4)传输到一个或多个次级备份设备的第(n)信息量。

由每个初级备份设备在其第(n)状态信号中报告的第(n)信息量包括固定数据量。在某些实施例中, 第(n)信息量包括至少一个一致事务集。在某些实施例中, 第(n)信息量包括至少一个一致事务集的一部分。

本申请人的方法从步骤 510 转至步骤 515, 其中第一初级备份设备从其余(N-1)个初级备份设备中的每一个接收下一个状态信号即第(n+1)状态信号。这些第(n+1)状态信号中的每一个包括通过公共通信链路传输到一个或多个次级备份设备的第(n+1)信息量。

步骤 515 的第(n+1)信息量包括固定数据量。步骤 515 的第(n+1)信息量典型地小于步骤 510 的第(n)信息量。如果第一初级备份设备在第(n)状态信号和第(n+1)状态信号之间的间隔内不提供任何数据给一个或多个次级备份设备, 则第(n+1)信息量等于第(n)信息量。

本申请人的方法从步骤 515 转至步骤 520, 其中第一初级备份设备计算(N)个设备中的每一个的第(n)有效带宽。第一初级备份设备可以通过将发送到一个或多个次级设备的信息量除以发送该信息的时间间隔(time interval)来确定其第(n)有效带宽。第一初级备份设备通过将其余初级设备中的每一个的第(n)信息量与第(n+1)信息量之间的各自差值除以状态信号间隔时间来确定其余(N-1)个设备中的每一个的第(n)有效带宽。

本申请人的方法从步骤 520 转至步骤 525, 其中第一初级备份设备计算(N)个初级备份设备中的每一个的完成时间(“TTC”)。步骤 525 的 TTC 值包括(N)个初级备份设备中的每一个使用步骤 520 的第(n)有效带宽将其余信息量发送到一个或多个次级备份设备所需的时间。步骤 525

包括对于(N)个设备中的每一个将第(n+1)信息量除以步骤 520 的第(n)有效带宽。

本申请人的方法从步骤 525 转至步骤 530, 其中第一初级备份设备计算所有(N)个初级备份设备的第(n)合计带宽。本领域的技术人员应当理解, 步骤 530 包括平均步骤 520 的各个第(n)单独带宽的(N)个值。

本申请人的方法从步骤 530 转至步骤 535, 其中该方法确定第一初级备份设备是否具有(N)个初级备份设备中的最大第(n)TTC 时间。如果本申请人的方法在步骤 535 确定第一初级备份设备具有(N)个初级备份设备的最大第(n)TTC 时间, 则本申请人的方法从步骤 535 转至步骤 550, 其中该方法将第一初级备份设备的第(n)延迟设为 0。本申请人的方法从步骤 550 转至步骤 560, 其中该方法由第一初级备份设备使用第(n)延迟通过公共通信链路将包括其一致事务集的全部或部分的数据提供给一个或多个次级备份设备。

在此所用的“使用第(n)延迟”提供数据是指每个设备发送包括该设备的一致事务集的全部或部分的固定数据量, 在发送那个数据之后, 每个设备然后“睡眠”其第(n)延迟值。重复该过程直到第(n)延迟值发生改变或者该一致事务集的数据全被发送为止。

如果本申请人的方法在步骤 535 确定第一初级备份设备不具有(N)个初级备份设备的最大 TTC 时间, 则本申请人的方法从步骤 535 转至步骤 540, 其中该方法确定第(n)合计带宽是否小于第(n-1)带宽。如果本申请人的方法在步骤 540 确定第(n)合计带宽不小于第(n-1)带宽, 则本申请人的方法从步骤 540 转至步骤 555, 其中该方法确立第一初级设备的第(n)延迟。在本发明方法的初始迭代中, 即在(n)为 1 的情况下, 没有第(n-1)合计带宽。因此, 在(n)为 1 的情况下, 步骤 540 的确定结果一定是“否”。

如果本申请人的方法在步骤 540 确定第(n)合计带宽小于第(n-1)带宽, 则本申请人的方法从步骤 540 转至步骤 545, 其中该方法确定第一初级备份设备的第(n-1)延迟是否大于第一初级备份设备的第(n-2)延迟。如果本申请人的方法在步骤 545 确定第一初级备份设备的第(n-1)延迟大

于第一初级备份设备的第 $(n-2)$ 延迟，则本申请人的方法从步骤 545 转至步骤 565，其中该方法将第 (n) 延迟设为第 $(n-2)$ 延迟。或者，如果本申请人的方法在步骤 545 确定第一初级备份设备的第 $(n-1)$ 延迟不大于第一初级备份设备的第 $(n-2)$ 延迟，则本申请人的方法从步骤 545 转至步骤 555。

作为第一例子，在本申请人方法的第二次迭代中，即在 (n) 为 2 的情况下，不可能存在第 $(n-2)$ 延迟，因此，在 (n) 为 2 的情况下，步骤 545 的确定结果一定是“否”，并且该方法从步骤 545 转至步骤 555。作为第二例子，在本申请人方法的第一次迭代中，即 (n) 为 1，为第一初级备份设备设置第一延迟，并且使用该第一延迟，本申请人的系统以第一合计带宽提供数据。然后在第二次迭代中，即 n 为 2，增加第一初级备份设备的延迟，并且使用第一初级备份设备的该第二延迟，本申请人的系统以第二合计带宽提供数据，其中，第二合计带宽小于第一合计带宽。在第三次迭代中，即 (n) 为 3，本申请人的方法在步骤 545 返回“是”的确定结果，并且转至步骤 565，其中，该方法将第一初级备份设备的第三延迟设为第一延迟值。

步骤 550、555 和 565 转至步骤 560，其中该方法使用第 (n) 延迟将数据从第一初级备份设备提供一个或多个次级备份设备。本申请人的方法从步骤 560 转至步骤 570，其中该方法增加 (n) 。该方法从步骤 570 转至步骤 515 并且继续。

图 6 概述了在步骤 555(图 5)确立第 (n) 延迟的本申请人方法的两个实施例的各步骤。现在参照图 6，在步骤 610，本申请人的方法计算所有 (N) 个初级备份设备的第 (n) 平均 TTC。在某些实施例中，步骤 610 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 610 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。本领域的技术人员应当理解，步骤 610 包括确定在步骤 520(图 5)确定的 (N) 个第 (n) 带宽中的每一个的平均值。

本申请人的方法从步骤 610 转至步骤 615，其中该方法确定第一初级备份设备的第 (n) TTC 是否基本上等于所有设备的第 (n) 平均 TTC。

“基本上相等”在此是指相差小于大约正负百分之十(10%)。在某些实施例中，步骤 615 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 615 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

如果本申请人的方法在步骤 615 确定第一初级备份设备的第(n)TTC基本上等于所有设备的第(n)平均 TTC，则该方法从步骤 615 转至步骤 625，其中该方法设置第一初级备份设备的第(n)延迟等于该设备的第(n-1)延迟。在某些实施例中，步骤 625 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 625 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

如果本申请人的方法在步骤 615 确定第一初级备份设备的第(n)TTC不基本上等于所有设备的第(n)平均 TTC，则该方法从步骤 615 转至步骤 620，其中该方法确定是否使用其余(N-1)个设备的预测延迟值设置第一初级备份设备的第(n)延迟。在某些实施例中，步骤 620 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 620 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。如果本申请人的方法选择使用其余(N-1)个设备的预测延迟值设置第一初级备份设备的第(n)延迟，则该方法从步骤 620 转至步骤 710(图 7)。

步骤 620 的决定基于先前作出且在固件中实现的策略决定。在某些实施例中，该策略决定由设备制造商作出，并且制造时在初级备份设备内的固件中实现。在某些实施例中，该策略决定由系统用户作出，并且使用例如操作员输入台在初级备份设备内的固件中实现。在某些实施例中，该操作员输入台与备份设备集成在一起。在某些实施例中，该操作员输入台与互连到备份设备的数据存储和检索系统集成在一起。在某些实施例中，该操作员输入台外部于备份设备以及与该备份设备互连的一个或多个数据存储和检索系统。

或者，如果本申请人的方法选择不使用其余(N-1)个设备的预测延迟值设置第一初级备份设备的第(n)延迟，则该方法从步骤 620 转至步骤 630，其中该方法确定是否使用标准延迟调整。步骤 630 的决定基于先

前作出且在固件中实现的策略决定。在某些实施例中，该策略决定由设备制造商作出，并且制造时在初级备份设备内的固件中实现。在某些实施例中，该策略决定由系统用户作出，并且使用例如操作员输入台在初级备份设备内的固件中实现。在某些实施例中，该操作员输入台与备份设备集成在一起。在某些实施例中，该操作员输入台与互连到备份设备的数据存储和检索系统集成在一起。在某些实施例中，该操作员输入台外部于备份设备以及与该备份设备互连的一个或多个数据存储和检索系统。

如果本申请人的方法在步骤 630 选择不使用标准延迟调整，则该方法从步骤 630 转至步骤 680。如果本申请人的方法在步骤 630 选择使用标准延迟调整，则该方法从步骤 630 转至步骤 640，其中该方法提供标准延迟调整。

在某些实施例中，步骤 640 的标准延迟调整设置在位于第一初级备份设备如初级备份设备 415(图 4)内的固件如固件 416 中。在某些实施例中，步骤 640 的标准延迟调整设置在位于第一初级备份设备如初级备份设备 415(图 4)的控制器如控制器 417 内的固件中。在某些实施例中，步骤 640 的标准延迟调整设置在位于初级数据存储和检索系统如数据存储和检索系统 410(图 4)的控制器如控制器 411(图 4)内的固件中。在某些实施例中，步骤 640 的标准延迟调整由主机计算机如主机计算机 390(图 3)提供。

本申请人的方法从步骤 640 转至步骤 650，其中该方法确定第一初级备份设备的第(n)TTC 是否大于步骤 610 的第(n)平均 TTC。在某些实施例中，步骤 650 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 650 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

如果本申请人的方法在步骤 650 确定第一初级备份设备的第(n)TTC 大于步骤 610 的第(n)平均 TTC，则该方法从步骤 650 转至步骤 660，其中该方法将第一初级备份设备的第(n)延迟设为第一初级备份设备的第(n-1)延迟减去步骤 630 的标准延迟调整。在某些实施例中，步骤 660 由

第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 660 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

如果本申请人的方法在步骤 650 确定第一初级备份设备的第(n)TTC 不大于步骤 610 的第(n)平均 TTC，则该方法从步骤 650 转至步骤 670，其中该方法将第一初级备份设备的第(n)延迟设为第一初级备份设备的第(n-1)延迟加上步骤 630 的标准延迟调整。在某些实施例中，步骤 670 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 670 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

在不存在第(n-1)延迟的情况下，即当(n)为 1 的情况下，本申请人的方法将第一初级备份设备的第(n)延迟设为标准延迟调整。

如果本申请人的方法在步骤 630 选择不使用标准延迟调整，则该方法从步骤 630 转至步骤 680，其中该方法提供延迟调整函数。在某些实施例中，步骤 680 的延迟调整函数由主机计算机如主机计算机 390(图 3)提供。在某些实施例中，步骤 680 的延迟调整函数设置在位于第一初级备份设备如初级备份设备 415(图 4)内的固件如固件 416(图 4)中。在某些实施例中，步骤 680 的延迟调整函数设置在位于第一初级备份设备的控制器如控制器 417(图 4)内的固件中。

在某些实施例中，步骤 680 的延迟调整函数包括一个查询表，其包括在步骤 610 确定的第(n)平均 TTC 与第一初级备份设备的第(n)TTC 之间的各个差值即 $TTC_{agg}-TTC_{(n)}$ 所对应的特定延迟值。在某些实施例中，步骤 680 的延迟调整函数包括方程(1)：

$$\text{延迟}=a(TTC_{agg}-TTC_{(n)})+b \quad (1)$$

本申请人的方法从步骤 680 转至步骤 690，其中该方法使用步骤 680 的延迟调整函数设置第(n)延迟。在某些实施例中，步骤 680 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 680 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

如果本申请人的方法在步骤 620 选择使用其余(N-1)个设备的预测延迟值设置第一初级备份设备的第(n)延迟，则该方法从步骤 620 转至步骤

710(图 7)。现在参照图 7, 在步骤 710, 本申请人的方法计算步骤 610 的第(n)平均 TTC 与其他(N-1)个初级备份设备中的每一个的第(n)TTC 之差。在某些实施例中, 步骤 710 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中, 步骤 710 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

在步骤 720, 本申请人的方法确定是否使用标准延迟调整来预测其余(N-1)个初级备份设备的延迟值。在某些实施例中, 步骤 720 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中, 步骤 720 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

步骤 720 的决定基于先前作出且在固件中实现的策略决定。在某些实施例中, 该策略决定由设备制造商作出, 并且制造时在初级备份设备内的固件中实现。在某些实施例中, 该策略决定由系统用户作出, 并且使用例如操作员输入台在初级备份设备内的固件中实现。在某些实施例中, 该操作员输入台与备份设备集成在一起。在某些实施例中, 该操作员输入台与互连到备份设备的数据存储和检索系统集成在一起。在某些实施例中, 该操作员输入台外部于备份设备以及与该备份设备互连的一个或多个数据存储和检索系统。

如果本申请人的方法在步骤 720 选择使用标准延迟调整来预测其余(N-1)个初级备份设备的延迟值, 则该方法从步骤 720 转至步骤 730, 其中该方法提供标准延迟调整。在某些实施例中, 步骤 730 的标准延迟调整设置在位于第一初级备份设备如初级备份设备 415(图 4)内的固件如固件 416(图 4)中。在某些实施例中, 步骤 730 的标准延迟调整设置在位于第一初级备份设备如初级备份设备 415(图 4)的控制器如控制器 417 内的固件中。在某些实施例中, 步骤 730 的标准延迟调整设置在位于初级数据存储和检索系统如数据存储和检索系统 410(图 4)的控制器如控制器 411(图 4)内的固件中。在某些实施例中, 步骤 730 的标准延迟调整由主机计算机如主机计算机 390(图 3)提供。

本申请人的方法从步骤 730 转至步骤 740, 其中该方法使用步骤 730 的标准延迟调整预测其余(N-1)个初级备份设备的延迟值。在某些实

施例中，步骤 740 包括对于其他(N-1)个初级备份设备中的每一个使用步骤 650(图 6)、660(图 6)和 670(图 6)。在某些实施例中，步骤 740 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 740 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

本申请人的方法从步骤 740 转至步骤 770，其中该方法使用步骤 740 的预测延迟值设置第一初级备份设备的第(n)延迟。在某些实施例中，步骤 770 由第一初级备份设备如初级备份设备 415(图 4)执行。在某些实施例中，步骤 770 由位于第一初级备份设备内的控制器如控制器 417(图 4)执行。

如果本申请人的方法在步骤 720 选择不使用标准延迟调整来预测其余(N-1)个初级备份设备的延迟值，则该方法从步骤 720 转至步骤 750，其中该方法提供延迟调整函数。在某些实施例中，步骤 750 的延迟调整函数包括一个查询表，其包括在步骤 610(图 6)确定的第(n)平均 TTC 与第一初级备份设备的第(n)TTC 之间的各个差值即 $TTC_{agg}-TTC_{(n)}$ 所对应的特定延迟值。在某些实施例中，步骤 680 的延迟调整函数包括二阶方程如方程(2)：

$$\text{延迟}=a(TTC_{agg} - TTC_{(n)})+b \quad (2)$$

本申请人的方法从步骤 750 转至步骤 760，其中该方法使用步骤 750 的延迟调整函数预测其他(N-1)个初级备份设备中的每一个的延迟值。本申请人的方法从步骤 760 转至步骤 770。

在某些实施例中，图 5、6 和/或 7 所述的各个步骤可以经过组合、删除或重新排序。

本申请人的发明还包括一种制造品，其包括一个计算机可用介质例如计算机可用介质 418、423、438、443、458 和/或 463，其中包含有用来使用图 5、6 和/或 7 所述的步骤调整数据传输速率的计算机可读程序代码。

本申请人的发明还包括一种可与可编程计算机处理器一起使用的计算机程序产品，例如计算机程序产品 419、424、439、444、459 和/或 464，其具有使用图 5、6 和/或 7 所述的步骤调整数据传输速率的计算

机可读程序代码。在某些实施例中，该计算机程序产品位于数据存储和检索系统中。在某些实施例中，该计算机程序产品位于备份设备中。在某些实施例中，该计算机程序代码实现图 5、6 和/或 7 所述的步骤。

虽然详细地举例说明了本发明的优选实施例，但是本领域的技术人员应当清楚，在不脱离由所附权利要求限定的本发明的范围的情况下，可以对这些实施例进行各种修改和变动。

图1

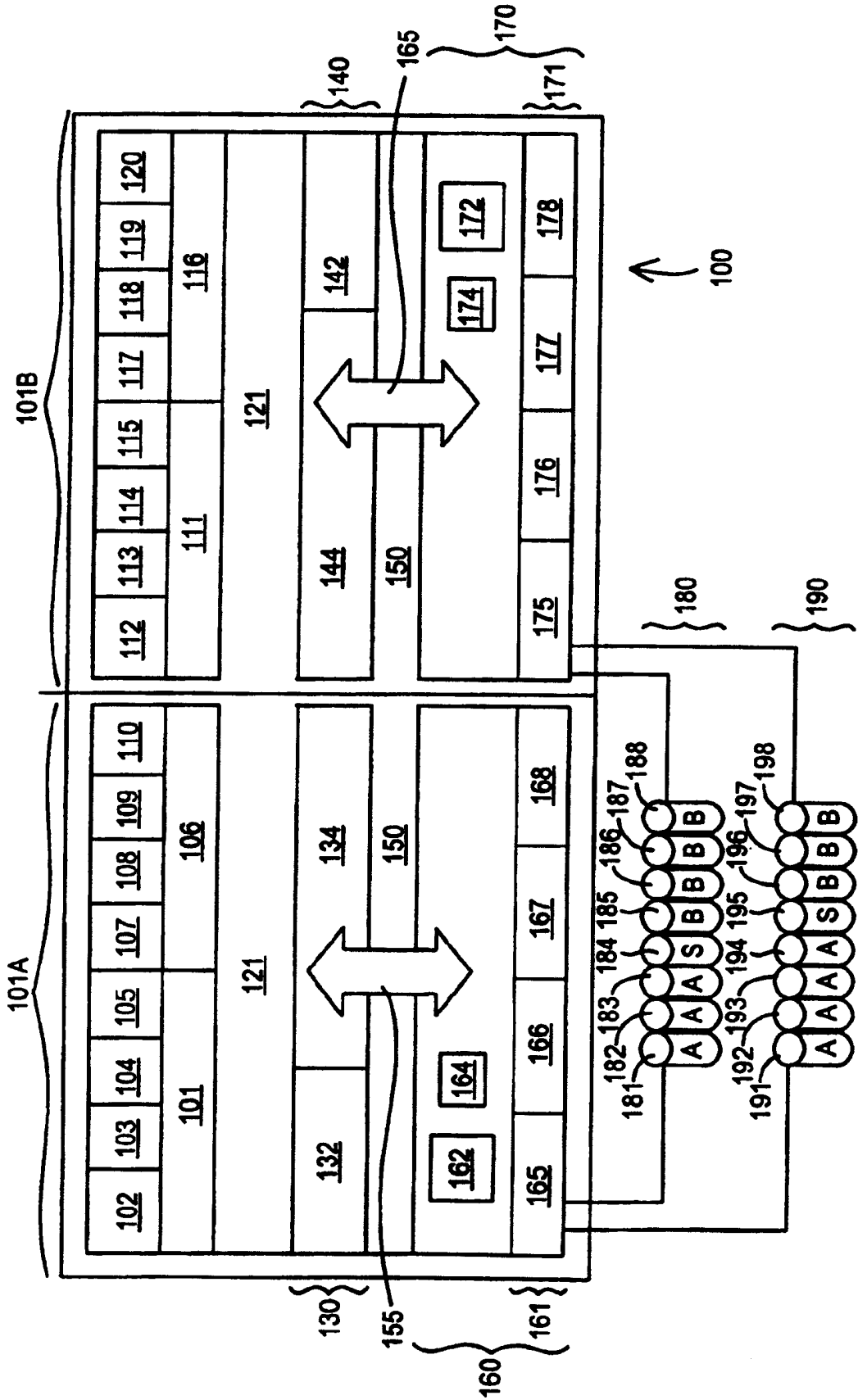


图2

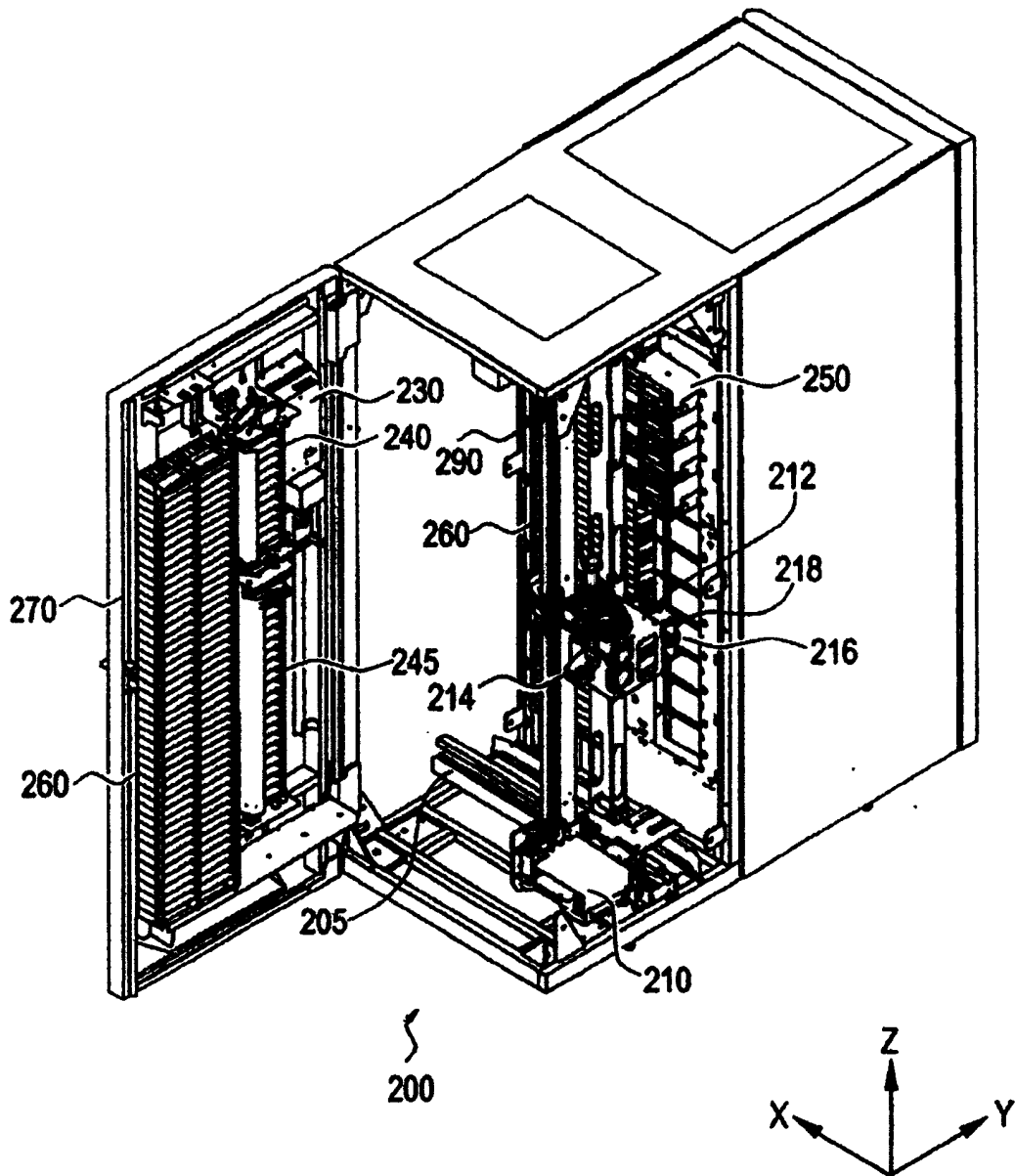


图3

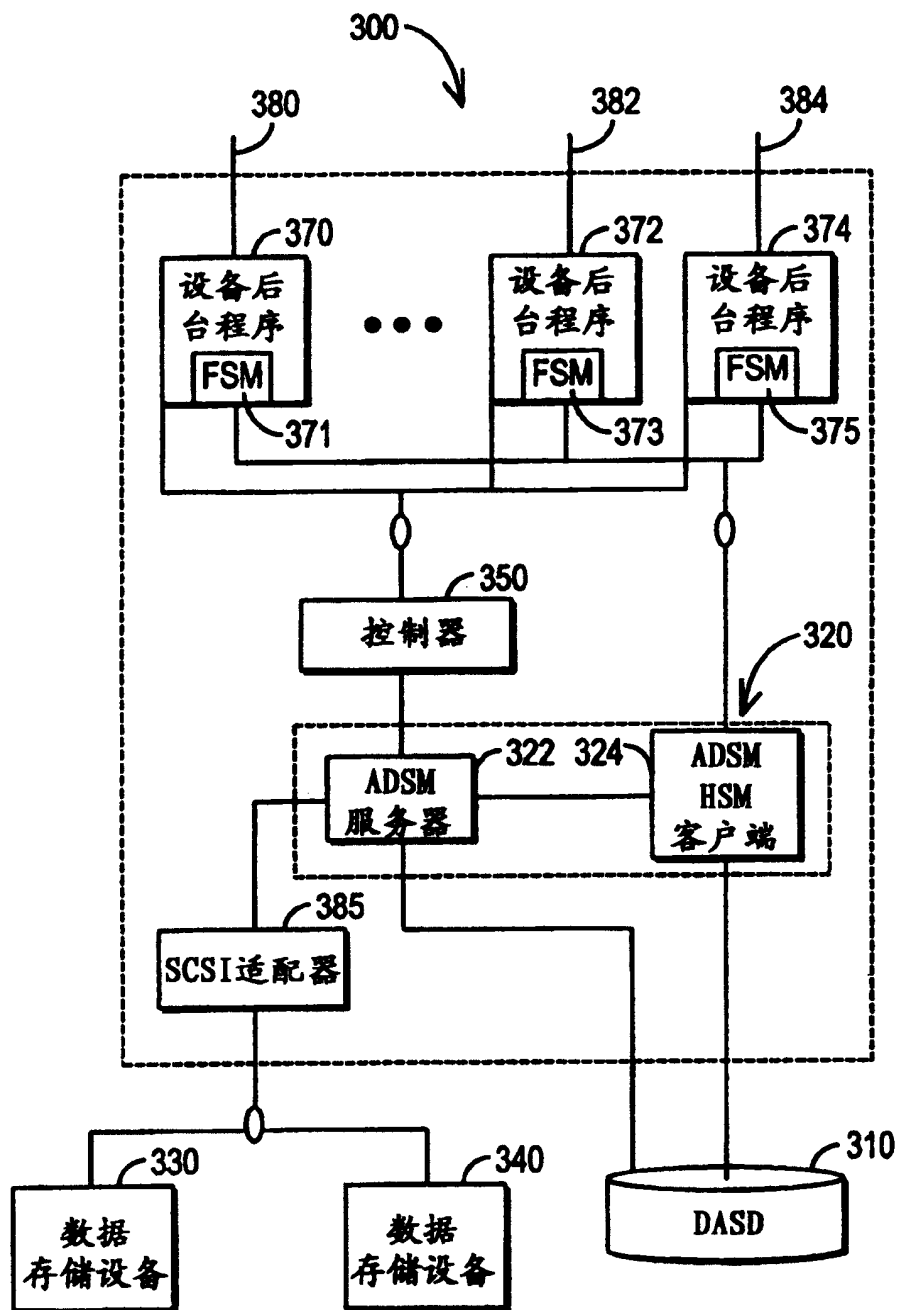


图 4

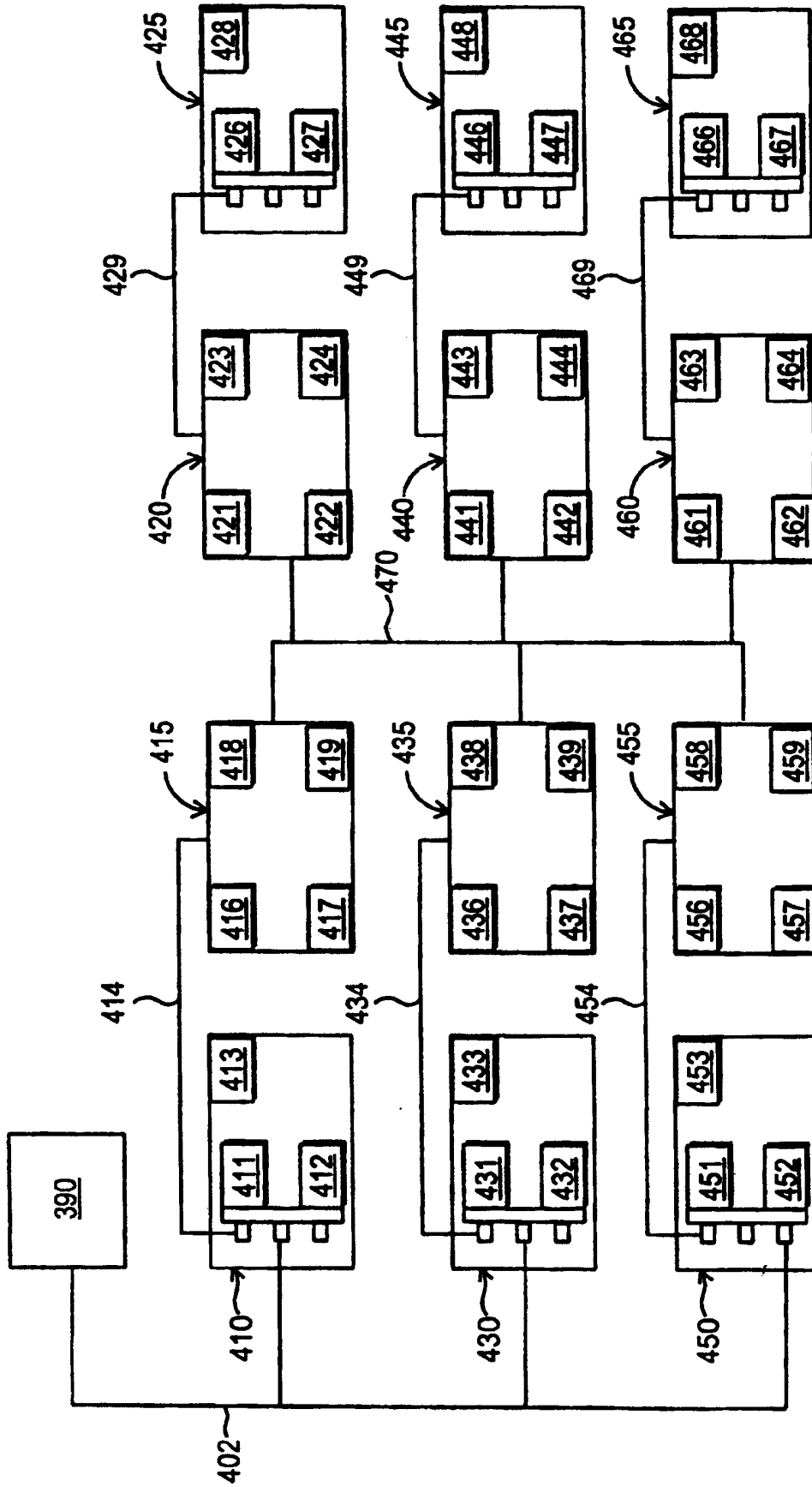


图5

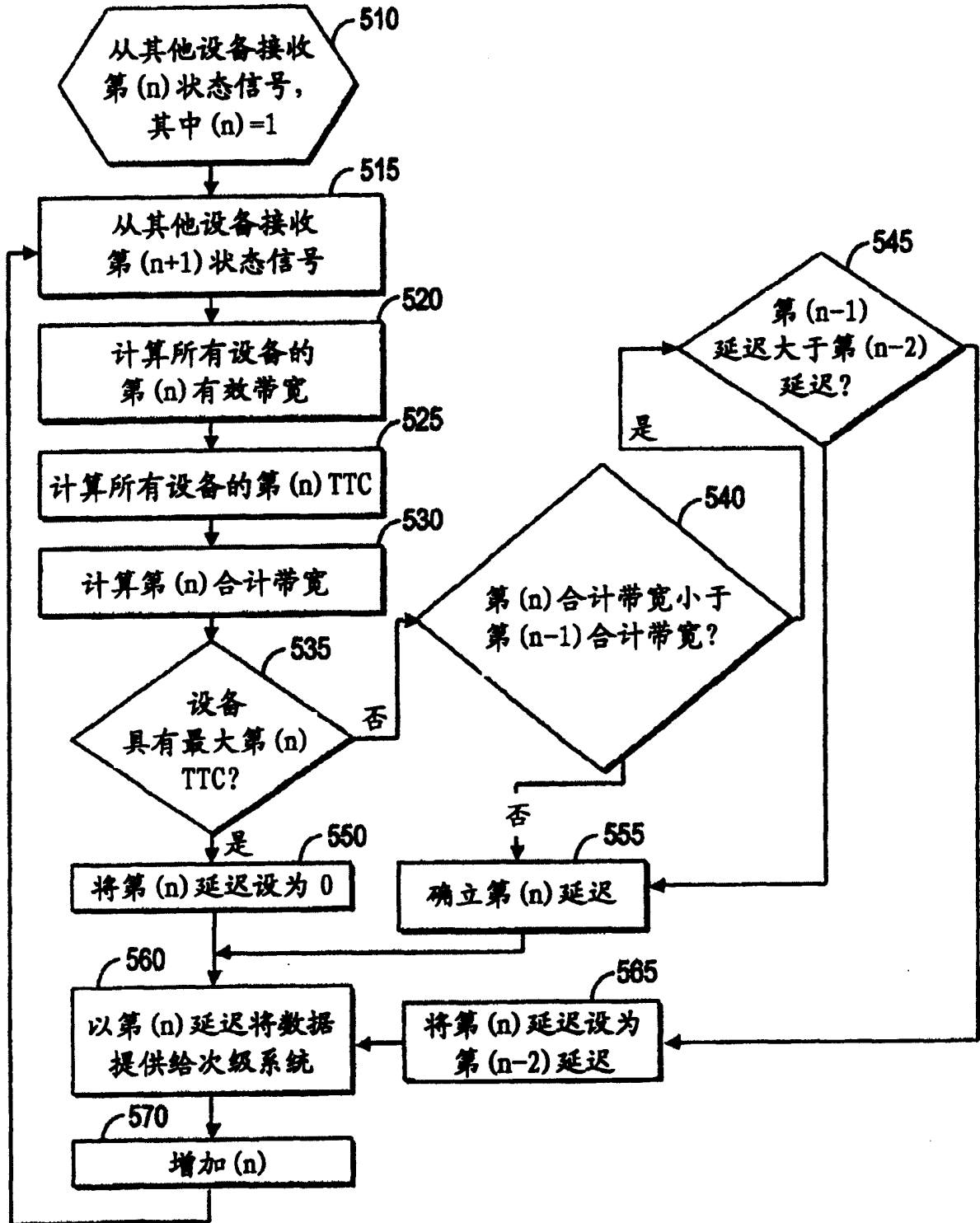


图6

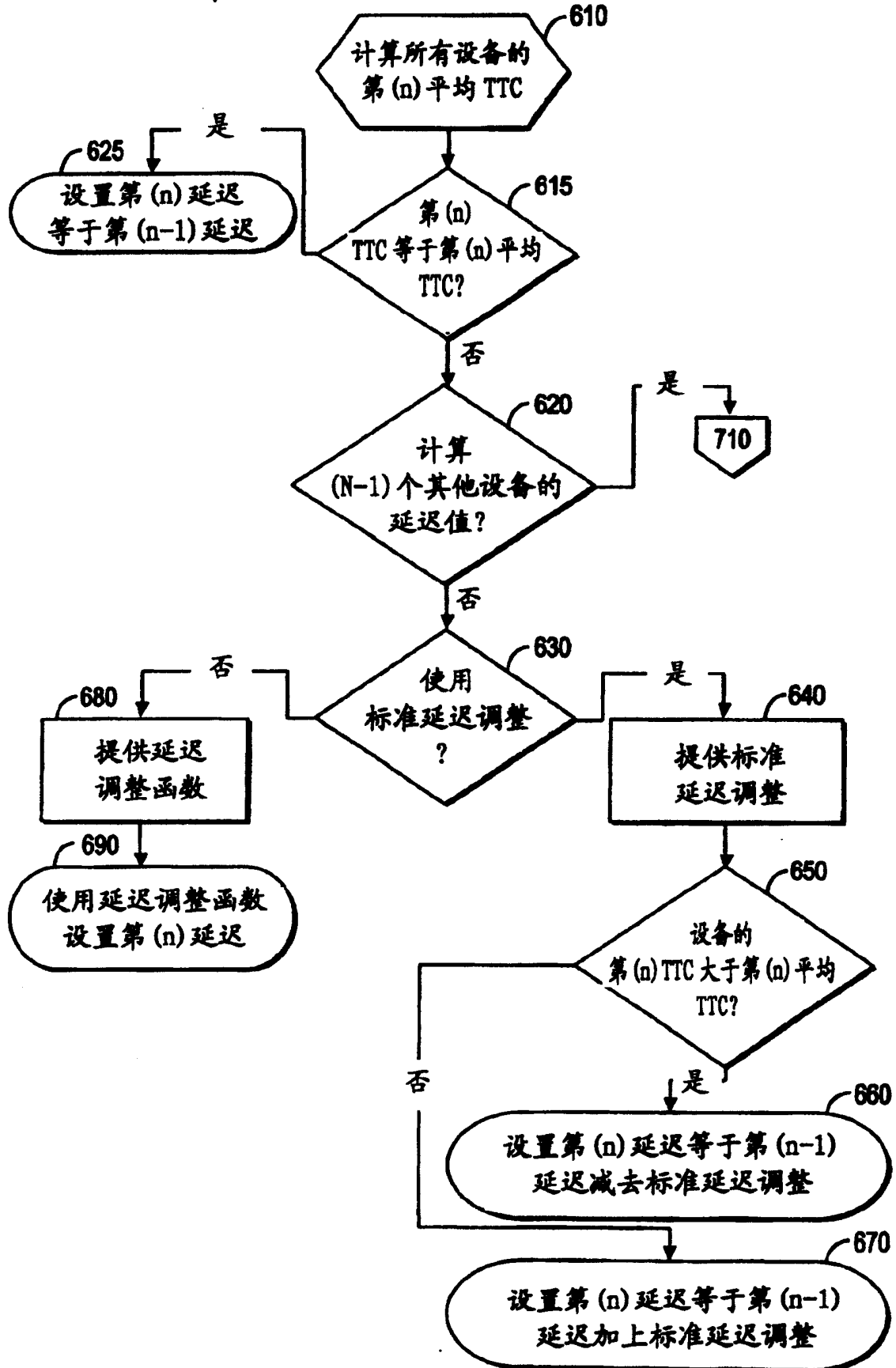


图 7

