



(12) 发明专利申请

(10) 申请公布号 CN 103677656 A

(43) 申请公布日 2014. 03. 26

(21) 申请号 201310052357. X

(22) 申请日 2013. 02. 06

(30) 优先权数据

13/600,644 2012. 08. 31 US

(71) 申请人 株式会社日立制作所

地址 日本东京都

(72) 发明人 出口彰

(74) 专利代理机构 北京市金杜律师事务所

11256

代理人 鄭迅

(51) Int. Cl.

G06F 3/06 (2006. 01)

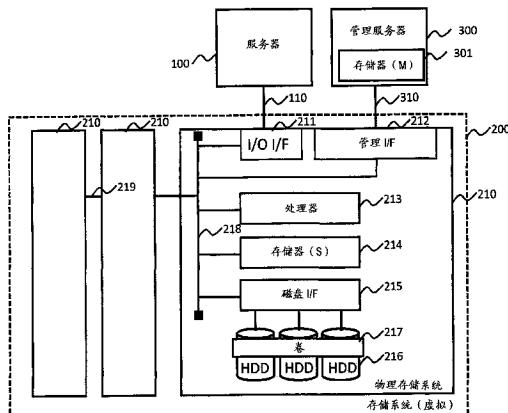
权利要求书4页 说明书14页 附图31页

(54) 发明名称

虚拟存储系统和远程复制系统的管理方法

(57) 摘要

示例性的实施方式提供了管理包括远程复制系统的存储系统并且通过使复杂的操作进行自动化来改进可管理性的技术。在一个实施方式中，计算机包括存储器和控制器。控制器可操作以：管理要向服务器提供的虚拟卷；管理从多个存储系统提供的多个逻辑卷；管理对虚拟卷所要求的条件，该条件与要发送到虚拟卷的数据的存储位置有关；管理多个逻辑卷中的每个逻辑卷的位置信息，逻辑卷的位置信息是基于逻辑卷的位置来限定的；以及基于虚拟卷的条件和逻辑卷的位置信息来进行控制，以将虚拟卷映射到多个逻辑卷中的逻辑卷。



1. 一种计算机,包括存储器和控制器,所述控制器可操作以:

管理要向服务器提供的虚拟卷;

管理从多个存储系统提供的多个逻辑卷;

管理对所述虚拟卷所要求的条件,所述条件与要发送到所述虚拟卷的数据的存储位置有关;

管理所述多个逻辑卷中的每个逻辑卷的位置信息,逻辑卷的所述位置信息是基于所述逻辑卷的位置来限定的;以及

基于所述虚拟卷的所述条件和所述逻辑卷的所述位置信息来进行控制,以将所述虚拟卷映射到所述多个逻辑卷中的任一逻辑卷。

2. 根据权利要求 1 所述的计算机,

其中,所述条件要求要发送到所述虚拟卷的数据的所述存储位置所在的站点与作为所述虚拟卷的提供对象的所述服务器所在的站点相同;并且

其中,要发送到所述虚拟卷的数据的所述存储位置是被映射到所述虚拟卷的所述逻辑卷的所述位置。

3. 根据权利要求 1 所述的计算机,

其中,所述虚拟卷被映射到主逻辑卷,所述主逻辑卷被映射到辅逻辑卷,所述主逻辑卷和所述辅逻辑卷是远程复制对并且处于不同的位置处。

4. 根据权利要求 3 所述的计算机,

其中,所述控制器可操作以基于所述逻辑卷的所述位置信息、对所述辅逻辑卷所要求的条件和所述多个逻辑卷所在的站点之间的距离的站点距离信息将所述主逻辑卷自动地映射到所述辅逻辑卷,其中对所述辅逻辑卷所要求的所述条件包括空闲容量和相对于所述主逻辑卷的连接性。

5. 根据权利要求 4 所述的计算机,

其中,所述控制器可操作以自动地创建用于存储所述主逻辑卷的日志的主日志卷并创建辅日志卷,将所述日志从所述主日志卷复制到所述辅日志卷以及将所述日志从所述辅日志卷复制到所述辅逻辑卷以实现异步远程复制。

6. 根据权利要求 3 所述的计算机,

其中,所述主逻辑卷被映射到所述辅逻辑卷和第三逻辑卷,所述主逻辑卷、所述辅逻辑卷和所述第三逻辑卷具有远程复制关系并且处于三个不同的位置处;并且

其中,所述控制器可操作以基于所述逻辑卷的所述位置信息、对所述辅逻辑卷所要求的条件、对所述第三逻辑卷所要求的条件、所述多个逻辑卷所在的站点之间的距离的站点距离信息将所述主逻辑卷自动地映射到所述辅逻辑卷和所述第三逻辑卷,其中,如果所述远程复制是级联的远程复制,则对所述辅逻辑卷所要求的所述条件包括空闲容量和相对于所述主逻辑卷的连接性,并且对所述第三逻辑卷所要求的所述条件包括空闲容量和相对于所述辅逻辑卷的连接性,如果所述远程复制不是级联的远程复制,则对所述辅逻辑卷所要求的所述条件包括相对于所述主逻辑卷的连接性。

7. 根据权利要求 3 所述的计算机,

其中,所述控制器可操作以在所述主逻辑卷出现 I/O(输入 / 输出)故障时将 I/O 从所述服务器与所述主逻辑卷之间的 I/O 路径自动地改变为所述服务器与所述辅逻辑卷之间

的另一 I/O 路径 ;并且

其中,所述辅逻辑卷的状态从 I/O 不可接收改变为 I/O 可接收。

8. 根据权利要求 3 所述的计算机,

其中,所述控制器可操作以在所述主逻辑卷出现 I/O( 输入 / 输出 ) 故障时将用于运行使用来自所述服务器的数据的应用软件的虚拟机自动地迁移到目的地服务器,并且将 I/O 从所述服务器与所述主逻辑卷之间的 I/O 路径自动地改变为所述目的地服务器与所述辅逻辑卷之间的另一 I/O 路径 ;并且

其中,所述辅逻辑卷的状态从 I/O 不可接收改变为 I/O 可接收。

9. 根据权利要求 3 所述的计算机,

其中,所述控制器可操作以在所述主逻辑卷出现 I/O( 输入 / 输出 ) 故障时基于所述逻辑卷的所述位置信息、对所述自动恢复辅逻辑卷所要求的条件和所述多个逻辑卷所在的站点之间的距离的站点距离信息将所述辅逻辑卷自动地映射到自动恢复辅逻辑卷作为自动恢复远程复制对,其中对所述自动恢复辅逻辑卷所要求的所述条件包括空闲容量和相对于所述辅逻辑卷的连接性。

10. 根据权利要求 1 所述的计算机,

其中,所述虚拟卷被映射到具有一个或多个主逻辑卷的主一致性组中的主逻辑卷,并且所述主一致性组被映射到具有一个或多个辅逻辑卷的辅一致性组,所述主一致性组和所述辅一致性组提供一个或多个远程复制对并且处于不同的位置处 ;并且

其中,所述控制器可操作以基于所述逻辑卷的所述位置信息、对所述辅逻辑卷所要求的条件、所述多个逻辑卷所在的站点之间的距离的站点距离信息将所述主一致性组中的所述主逻辑卷自动地映射到所述辅一致性组中的辅逻辑卷,其中对所述辅逻辑卷所要求的所述条件包括空闲容量和相对于所述主逻辑卷的连接性。

11. 一种系统,包括计算机和多个存储系统,所述计算机包括存储器和控制器,所述控制器可操作以 :

管理要向服务器提供的虚拟卷 ;

管理从所述多个存储系统提供的多个逻辑卷 ;

管理对所述虚拟卷所要求的条件,所述条件与要发送到所述虚拟卷的数据的存储位置有关 ;

管理所述多个逻辑卷中的每个逻辑卷的位置信息,逻辑卷的所述位置信息是基于所述逻辑卷的位置来限定的 ;以及

基于所述虚拟卷的所述条件和所述逻辑卷的所述位置信息来进行控制,以将所述虚拟卷映射到所述多个逻辑卷中的任一逻辑卷。

12. 根据权利要求 11 所述的系统,

其中,所述条件要求要发送到所述虚拟卷的数据的所述存储位置所在的站点与作为所述虚拟卷的提供对象的所述服务器所在的站点相同 ;并且

其中,要发送到所述虚拟卷的数据的所述存储位置是被映射到所述虚拟卷的所述逻辑卷的所述位置。

13. 根据权利要求 11 所述的系统,

其中,所述虚拟卷被映射到主逻辑卷,所述主逻辑卷被映射到辅逻辑卷,所述主逻辑卷

和所述辅逻辑卷是远程复制对并且处于不同的位置处；并且

其中，所述控制器可操作以基于所述逻辑卷的所述位置信息、对所述辅逻辑卷所要求的条件和所述多个逻辑卷所在的站点之间的距离的站点距离信息将所述主逻辑卷自动地映射到所述辅逻辑卷，其中对所述辅逻辑卷所要求的所述条件包括空闲容量和相对于所述主逻辑卷的连接性。

14. 根据权利要求 11 所述的系统，

其中，所述虚拟卷被映射到主逻辑卷，所述主逻辑卷被映射到辅逻辑卷，所述主逻辑卷和所述辅逻辑卷是远程复制对并且处于不同的位置处；

其中，所述控制器可操作以在所述主逻辑卷出现 I/O(输入 / 输出) 故障时将用于运行使用来自所述服务器的数据的应用软件的虚拟机自动地迁移到目的地服务器，并且将 I/O 从所述服务器与所述主逻辑卷之间的 I/O 路径自动地改变为所述目的地服务器与所述辅逻辑卷之间的另一 I/O 路径；并且

其中，提供所述辅逻辑卷的所述存储系统被配置为将所述辅逻辑卷的状态从 I/O 不可接收改变为 I/O 可接收。

15. 一种计算机可读存储介质，其存储用于控制数据处理器以管理数据存储的多个指令，所述多个指令包括：

使得所述数据处理器管理要向服务器提供的虚拟卷的指令；

使得所述数据处理器管理从多个存储系统提供的多个逻辑卷的指令；

使得所述数据处理器管理对所述虚拟卷所要求的条件的指令，所述条件与要发送到所述虚拟卷的数据的存储位置有关；

使得所述数据处理器管理所述多个逻辑卷中的每个逻辑卷的位置信息的指令，逻辑卷的所述位置信息是基于所述逻辑卷的位置来限定的；以及

使得所述数据处理器基于所述虚拟卷的所述条件和所述逻辑卷的所述位置信息来进行控制以将所述虚拟卷映射到所述多个逻辑卷中的任一逻辑卷的指令。

16. 根据权利要求 15 所述的计算机可读存储介质，

其中，所述条件要求要发送到所述虚拟卷的数据的所述存储位置所在的站点与作为所述虚拟卷的提供对象的所述服务器所在的站点相同；并且

其中，要发送到所述虚拟卷的数据的所述存储位置是被映射到所述虚拟卷的所述逻辑卷的所述位置。

17. 根据权利要求 15 所述的计算机可读存储介质，

其中，所述多个指令还包括：

使得所述数据处理器将所述虚拟卷映射到主逻辑卷的指令，所述主逻辑卷被映射到辅逻辑卷，所述主逻辑卷和所述辅逻辑卷是远程复制对并且处于不同的位置处；以及

使得所述数据处理器基于所述逻辑卷的所述位置信息、对所述辅逻辑卷所要求的条件和所述多个逻辑卷所在的站点之间的距离的站点距离信息将主逻辑卷自动地映射到所述辅逻辑卷的指令，其中对所述辅逻辑卷所要求的所述条件包括空闲容量和相对于所述主逻辑卷的连接性。

18. 根据权利要求 15 所述的计算机可读存储介质，其中，所述多个指令还包括：

使得所述数据处理器将所述虚拟卷映射到主逻辑卷的指令，所述主逻辑卷被映射到辅

逻辑卷和第三逻辑卷,所述主逻辑卷、所述辅逻辑卷和所述第三逻辑卷具有远程复制关系并且处于三个不同的位置处;以及

使得所述数据处理器基于所述逻辑卷的所述位置信息、对所述辅逻辑卷所要求的条件、对所述第三逻辑卷所要求的条件、所述多个逻辑卷所在的站点之间的距离的站点距离信息将所述主逻辑卷自动地映射到所述辅逻辑卷和所述第三逻辑卷的指令,其中,如果所述远程复制是级联的远程复制,则对所述辅逻辑卷所要求的所述条件包括空闲容量和相对于所述主逻辑卷的连接性,并且对所述第三逻辑卷所要求的所述条件包括空闲容量和相对于所述辅逻辑卷的连接性,如果所述远程复制不是级联的远程复制,则对所述辅逻辑卷所要求的所述条件包括相对于所述主逻辑卷的连接性。

19. 根据权利要求 15 所述的计算机可读存储介质,其中,所述多个指令还包括:

使得所述数据处理器将所述虚拟卷映射到主逻辑卷的指令,所述主逻辑卷被映射到辅逻辑卷,所述主逻辑卷和所述辅逻辑卷是远程复制对并且处于不同的位置处;以及

使得所述数据处理器在所述主逻辑卷出现 I/O(输入 / 输出)故障时将用于运行使用来自所述服务器的数据的应用软件的虚拟机自动地迁移到目的地服务器,并且将 I/O 从所述服务器与所述主逻辑卷之间的 I/O 路径自动地改变为所述目的地服务器与所述辅逻辑卷之间的另一 I/O 路径的指令。

20. 根据权利要求 15 所述的计算机可读存储介质,其中,所述多个指令还包括:

使得所述数据处理器将所述虚拟卷映射到主逻辑卷的指令,所述主逻辑卷被映射到辅逻辑卷,所述主逻辑卷和所述辅逻辑卷是远程复制对并且处于不同的位置处;以及

使得所述数据处理器在所述主逻辑卷出现 I/O(输入 / 输出)故障时基于所述逻辑卷的所述位置信息、对所述自动恢复辅逻辑卷所要求的条件、所述多个逻辑卷所在的站点之间的距离的站点距离信息将所述辅逻辑卷自动地映射到自动恢复辅逻辑卷以作为自动恢复远程复制对的指令,其中对所述自动恢复辅逻辑卷所要求的所述条件包括空闲容量和相对于所述辅逻辑卷的连接性。

## 虚拟存储系统和远程复制系统的管理方法

### 技术领域

[0001] 本发明一般涉及存储系统, 具体而言, 涉及远程复制系统和远程复制系统的技术的易用性。

### 背景技术

[0002] US2007/0079088 公开了一种用于向服务器提供具有相互连接的一个或多个物理存储系统的一个虚拟化存储设备的技术。具体而言, 该技术通过相互连接实现了其它物理存储系统的资源的使用。此外, 在虚拟化的存储系统(一个或多个物理存储系统)中提供了针对存储区域(卷)的唯一 ID。唯一的 ID 被称作全局卷 ID。因此, 即使在物理存储系统之间移动卷, 向服务器提供的卷 ID 也不会改变。

[0003] US2008/0034005 公开了虚拟存储系统。很多存储系统产品具有用于灾难恢复的远程复制功能。远程复制功能将主存储系统中存储的数据复制到辅存储系统。操作者不得不为辅物理存储系统做准备, 其包括: 在初始化远程复制之前创建复制目的地卷、远程复制路径配置等。

[0004] 第 7,680,919 号美国专利公开了与虚拟服务器有关的技术。这是用于在一个物理机器上创建一个或多个虚拟服务器的技术。此外, 该专利公开了用于在物理服务器之间移动虚拟服务器的技术。

[0005] 传统的虚拟化存储系统中的卷创建是按以下顺序操作的: 在物理存储系统中使用本地卷 ID 进行卷创建, 设置全局卷 ID 以及向服务器提供该卷。因此, 操作者必须知道物理存储系统以进行配置。用于自动地确定物理存储区域的技术不是基于指定的全局卷 ID 来提供的。因此, 操作成本很高。此外, 因为其未考虑超出数据中心(DC)的虚拟化存储系统, 因此不存在选择安装在最佳数据中心中的物理存储系统(例如, 安装在其中安装了服务器的相同的数据中心中的物理存储系统)的技术。如果传统的虚拟化存储系统应用于多个数据中心环境, 则将会出现超出 DC 的处理。这使得性能下降并且消耗 DC 之间的网络资源。

[0006] 传统的远程复制是使用两组存储系统来配置的。因此, 当仅向虚拟化的存储系统环境应用传统的远程复制时, 远程复制是使用两组虚拟化的存储系统来配置的。在该情况下, 对于两组虚拟化的存储系统而言, 需要设置远程复制。具体而言, 设置包括物理存储系统之间的有线连接、建立 I/F(接口)以进行数据传送、创建复制目的地卷、识别复制目的地的卷 ID 等。如果实现了超出数据中心的虚拟化存储系统环境, 则在该环境中, 远程复制功能可以用作本地复制功能。此外, 如果对于虚拟化的存储系统而言存在连接, 则用于远程复制的有线连接和 I/F 的建立将变得不必要。此外, 如果存储系统通过使用虚拟的卷技术自动地创建卷, 则创建复制目的地卷将变得不必要。此外, 如果使得复制源卷和复制目的地的 ID 是相同的, 则识别复制目的地卷 ID 将变得不必要。

### 发明内容

[0007] 本发明的示例性的实施方式提供了管理包括远程复制系统的存储系统以及通过

使复杂的操作自动化来改进可操作性的技术。

[0008] 第一实施方式涉及一种用于在虚拟存储系统中创建卷的方法。具体而言，该技术指派来自物理存储系统的物理存储容量，该物理存储系统安装在其中安装有服务器的相同的数据中心中。

[0009] 第二实施方式涉及一种自动远程复制环境的配置方法。一旦针对服务器使用的卷建立了远程复制属性，存储系统就自动地获取远程复制处理所需的资源，例如，复制目的地卷（辅卷、SVOL）或者存储系统之间的带宽。此外，使得复制源卷（主卷 PVOL）的卷 ID 和提供给服务器的 SVOL 的卷 ID 是相同的。当 PVOL 发生故障时，可以仅通过将 I/O 路径从 PVOL 改变为 SVOL 来重新启动 I/O 处理。因为不必如传统技术中一样将 I/O 处理改变为具有另一个 ID 的 SVOL，因此对应用没有影响。此外，可以使得一个存储配置中的服务器的故障恢复处理和多存储配置中的处理是相同的。

[0010] 第三实施方式涉及远程复制的自动恢复方法和将上述技术应用于三数据远程复制配置的方法。

[0011] 根据本发明的一个方面，计算机包括存储器和控制器。该控制器可操作以：管理要向服务器提供的虚拟卷；管理从多个存储系统提供的多个逻辑卷；管理对虚拟卷所要求的条件，该条件与要发送到虚拟卷的数据的存储位置有关；管理多个逻辑卷中的每个逻辑卷的位置信息，逻辑卷的位置信息是基于逻辑卷的位置来限定的；以及基于虚拟卷的条件和逻辑卷的位置信息来进行控制，以将虚拟卷映射到多个逻辑卷中的任一逻辑卷。

[0012] 在一些实施方式中，该条件要求要发送到虚拟卷的数据的存储位置所在的站点与作为虚拟卷的提供对象的服务器所在的站点相同；以及要发送到虚拟卷的数据的存储位置是被映射到虚拟卷的逻辑卷的位置。虚拟卷被映射到主逻辑卷，主逻辑卷被映射到辅逻辑卷，主逻辑卷和辅逻辑卷是远程复制对并且处于不同的位置处。控制器可操作以基于逻辑卷的位置信息、对辅逻辑卷所要求的条件和多个逻辑卷所在的站点之间的距离的站点距离信息将主逻辑卷自动地映射到辅逻辑卷，其中对辅逻辑卷所要求的条件包括空闲容量和相对于主逻辑卷的连接性。控制器可操作以自动地创建用于存储主逻辑卷的日志的主日志卷并创建辅日志卷，将日志从主日志卷复制到辅日志卷以及将日志从辅日志卷复制到辅逻辑卷以实现异步远程复制。

[0013] 在具体实施方式中，主逻辑卷被映射到辅逻辑卷和第三逻辑卷，主逻辑卷、辅逻辑卷和第三逻辑卷具有远程复制关系并且处于三个不同的位置处。控制器可操作以基于逻辑卷的位置信息、对辅逻辑卷所要求的条件、对第三逻辑卷所要求的条件、多个逻辑卷所处的站点之间的距离的站点距离信息将主逻辑卷自动地映射到辅逻辑卷和第三逻辑卷，其中，如果远程复制是级联的远程复制，则对辅逻辑卷所要求的条件包括空闲容量和相对于主逻辑卷的连接性，并且对第三逻辑卷所要求的条件包括空闲容量和相对于辅逻辑卷的连接性，如果远程复制不是级联的远程复制，则对辅逻辑卷所要求的条件包括相对于主逻辑卷的连接性。

[0014] 在一些实施方式中，控制器可操作以在主逻辑卷出现 I/O（输入 / 输出）故障时将 I/O 从服务器与主逻辑卷之间的 I/O 路径自动地改变为服务器与辅逻辑卷之间的另一 I/O 路径；并且辅逻辑卷的状态从 I/O 不可接收改变为 I/O 可接收。

[0015] 在具体实施方式中，控制器可操作以在主逻辑卷出现 I/O（输入 / 输出）故障时将

用于运行使用来自服务器的数据的应用软件的虚拟机自动地迁移到目的地服务器，并且将 I/O 从服务器与主逻辑卷之间的 I/O 路径自动地改变为目的地服务器与辅逻辑卷之间的另一 I/O 路径；并且辅逻辑卷的状态从 I/O 不可接收改变为 I/O 可接收。

[0016] 在一些实施方式中，控制器可操作以在主逻辑卷出现 I/O(输入 / 输出)故障时基于逻辑卷的位置信息、对自动恢复辅逻辑卷所要求的条件和多个逻辑卷所处的站点之间的距离的站点距离信息将辅逻辑卷自动地映射到自动恢复辅逻辑卷作为自动恢复远程复制对，其中对自动恢复辅逻辑卷所要求的条件包括空闲容量和相对于辅逻辑卷的连接性。

[0017] 在具体实施方式中，虚拟卷被映射到具有一个或多个主逻辑卷的主一致性组中的主逻辑卷，并且主一致性组被映射到具有一个或多个辅逻辑卷的辅一致性组，主一致性组和辅一致性组提供一个或多个远程复制对并且处于不同的位置处。控制器可操作以基于逻辑卷的位置信息、对辅逻辑卷所要求的条件、多个逻辑卷所在的站点之间的距离的站点距离信息将主一致性组中的主逻辑卷自动地映射到辅一致性组中的辅逻辑卷，其中对辅逻辑卷所要求的条件包括空闲容量和相对于主逻辑卷的连接性。

[0018] 根据本发明的另一个方面，系统包括计算机和多个存储系统，计算机包括存储器和控制器。控制器可操作以：管理要向服务器提供的虚拟卷；管理从多个存储系统提供的多个逻辑卷；管理对虚拟卷所要求的条件，条件与要发送到虚拟卷的数据的存储位置有关；管理多个逻辑卷中的每个逻辑卷的位置信息，逻辑卷的位置信息是基于逻辑卷的位置来限定的；以及基于虚拟卷的条件和逻辑卷的位置信息来进行控制，以将虚拟卷映射到多个逻辑卷中的任一逻辑卷。

[0019] 在一些实施方式中，虚拟卷被映射到主逻辑卷，主逻辑卷被映射到辅逻辑卷，主逻辑卷和辅逻辑卷是远程复制对并且处于不同的位置处。控制器可操作以在主逻辑卷出现 I/O(输入 / 输出)故障时将用于运行使用来自服务器的数据的应用软件的虚拟机自动地迁移到目的地服务器，并且将 I/O 从服务器与主逻辑卷之间的 I/O 路径自动地改变为目的地服务器与辅逻辑卷之间的另一 I/O 路径。提供辅逻辑卷的存储系统被配置为将辅逻辑卷的状态从 I/O 不可接收改变为 I/O 可接收。

[0020] 本发明的另一个方面涉及存储用于控制数据处理器以管理数据存储设备的多个指令的计算机可读存储介质。多个指令包括：使得数据处理器管理要向服务器提供的虚拟卷的指令；使得数据处理器管理从多个存储系统提供的多个逻辑卷的指令；使得数据处理器管理对虚拟卷所要求的条件的指令，条件与要发送到虚拟卷的数据的存储位置有关；使得数据处理器管理多个逻辑卷中的每个逻辑卷的位置信息的指令，逻辑卷的位置信息是基于逻辑卷的位置来限定的；以及使得数据处理器基于虚拟卷的条件和逻辑卷的位置信息来进行控制以将虚拟卷映射到多个逻辑卷中的任一逻辑卷的指令。

[0021] 对于本领域普通技术人员来说，根据具体实施方式的以下详细描述，本发明的这些和其它特征和优点将变得显而易见。

## 附图说明

[0022] 图 1 示出了在其中可以应用本发明的方法和装置的系统的硬件配置的示例。

[0023] 图 2 是示出了根据第一实施方式的管理服务器中的存储器 (M) 的示例的详细框图。

- [0024] 图 3 是在管理服务器中管理的服务器站点管理表的示例。
- [0025] 图 4 是在管理服务器中管理的存储站点管理表的示例。
- [0026] 图 5 是示出了根据第一实施方式的物理存储系统中的存储器 (S) 的示例的详细框图。
- [0027] 图 6 是在物理存储系统中管理的全局 ID 管理表的示例。
- [0028] 图 7 是示出了卷创建处理流程的流程图的示例。
- [0029] 图 8 示出了用于同步远程复制处理的远程复制配置的示例。
- [0030] 图 9 示出了用于异步远程复制处理的远程复制配置的示例。
- [0031] 图 10 是示出了根据第二实施方式的物理存储系统中的存储器 (S) 的示例的详细框图。
- [0032] 图 11 是示出了根据第二实施方式的管理服务器中的存储器 (M) 的详细框图。
- [0033] 图 12 是物理存储系统的存储器 (S) 中的主存储对表的示例。
- [0034] 图 13 是物理存储系统的存储器 (S) 中的辅存储对表的示例。
- [0035] 图 14 是在物理存储系统中管理的流入表的示例。
- [0036] 图 15 是在管理服务器中管理的站点距离表的示例。
- [0037] 图 16 是用于数据复制的设置屏幕的示例。
- [0038] 图 17 是示出了对创建处理流程的流程图的示例。
- [0039] 图 18 是示出了从图 17 中的对创建程序调用的辅选择程序的处理流程的流程图的示例。
- [0040] 图 19 是示出了用于异步远程复制的对创建处理流程的流程图的另一个示例。
- [0041] 图 20 是示出了用于计算 JVOL 所需的容量的 JVOL 容量程序的处理流程的流程图的示例。
- [0042] 图 21 是示出了辅选择程序的处理流程的流程图的另一个示例。
- [0043] 图 22 是示出了在 I/O 故障之后的处理的概念图。
- [0044] 图 23 是示出了在 PVOL 出现故障之后向 SVOL 发起的写入处理的流程图的示例。
- [0045] 图 24 是示出了在 PVOL 出现故障以后服务器的变化和涉及 SVOL 的处理的后续变化的概念图。
- [0046] 图 25 是示出了在 PVOL 出现故障以后的虚拟机器的迁移处理和用于确定目的地的处理的流程图的示例。
- [0047] 图 26 是示出了远程复制对的一致性组的概念图。
- [0048] 图 27 是示出了考虑了一致性组的辅选择程序的处理流程的流程图的另一个示例。
- [0049] 图 28 是示出了远程复制配置的自动恢复方法的概念图。
- [0050] 图 29 是示出了配置恢复程序的处理流程的流程图的示例。
- [0051] 图 30 是示出了三数据中心远程复制的级联配置的概念图。
- [0052] 图 31 是示出了三数据中心远程复制的多目标配置的概念图。
- [0053] 图 32 示出了三数据中心远程复制的设置屏幕的另一个示例。
- [0054] 图 33 是示出了用于考虑三数据中心远程复制来确定复制目的地的三数据中心目的地选择程序的处理流程的流程图的示例。

## 具体实施方式

[0055] 在本发明的以下详细描述中参照附图，这些附图构成了本公开内容的一部分，并且在附图中通过举例说明而非限制性的方式示出了可以通过其实现本发明的示例性的实施方式。在附图中，相似的数字描述贯穿多个视图的基本上相似的组件。此外，应当注意的是，虽然详细描述提供了如下面所描述的并且在附图中示出的各个示例性的实施方式，但是本发明不限于本文所描述和示出的实施方式，而是可以扩展到本领域技术人员将知道或者即将知道的其它实施方式。在说明书中提及“一个实施方式”、“该实施方式”或者“这些实施方式”意味着结合实施方式所描述的特定的特征、结构或特性包括在本发明的至少一个实施方式中，并且这些短语在说明书中的各个位置的出现不一定均是指代相同的实施方式。此外，在下面的详细描述中，给出了大量具体细节以提供对本发明的彻底理解。然而，对于本领域普通技术人员而言显而易见的是，为了实现本发明，不是所有这些具体细节都是必需的。在其它环境中，公知的结构、材料、电路、过程和接口已经被详细描述，和 / 或可以以框图的形式被示出，以不必模糊本发明。

[0056] 此外，根据计算机中的操作的算法和符号表示给出了下面的详细描述的一些部分。这些算法描述和符号表示是数据处理领域中的技术人员用于最有效地向本领域其它技术人员传达其创新的实质的方式。算法是导致期望的最终状态或结果的一系列定义的步骤。在本发明中，所执行的步骤需要物理操纵有形的量以实现有形的结果。通常，虽然不是必需的，但这些量具有能够被存储、传送、组合、比较和以其它方式操纵的电信号或磁信号或者指令的形式。已经证明主要由于共同使用的原因，有时将这些信号称作比特、值、要素、符号、字符、术语、数字、指令等是便利的。然而，应当记住的是，所有这些和类似的术语将与适当的物理量相关联，并且仅仅是应用于这些量的便利的标签。除非另外专门声明，通过以下讨论显而易见的是，应当理解，贯穿说明书使用诸如“处理”、“计算”、“运算”、“确定”、“显示”等术语的讨论可以包括操纵计算机系统的寄存器和存储器中的表示为物理（电子）量的数据并将这些数据转换为类似地表示为计算机系统的存储器或寄存器或者其它信息存储、传输或显示设备中的其它数据的计算机系统或其它信息处理设备的动作和过程。

[0057] 本发明还涉及一种用于执行本文的操作的装置。该装置可以被专门构造用于所需的目的，或者其可以包括由一个或多个计算机程序选择性地激活或重新配置的一个或多个通用计算机。此类计算机程序可以存储在计算机可读存储介质中，所述计算机可读存储介质包括非瞬态介质，例如但不限于：光盘、磁盘、只读存储器、随机存取存储器、固态设备和驱动器或者适合于存储电子信息的任何其它类型的介质。本文呈现的算法和显示器未固有地涉及任何特定的计算机或其它装置。各个通用系统可以与根据本文的教导的程序和模块一起使用，或者可以证明构造更专用的装置来执行期望的方法步骤是便利的。此外，本发明不是参照任何特定的编程语言来描述的。将清楚的是，各种编程语言可以用于执行本文所描述的本发明的教导。编程语言的指令可以由诸如中央处理单元 (CPU)、处理器或控制器的一个或多个处理设备来执行。

[0058] 下面将更详细描述的本发明的示例性实施方式提供了用于管理包括远程复制系统的存储系统以及通过对复杂的操作进行自动化来改进可管理性的装置、方法和计算机程序。

[0059] 在下面的描述中,将描述处理并且在一些情况下将“程序”作为主题来操作。在程序由处理器执行的情况下,执行预定的处理。因此,处理的主题也可以是处理器。当程序作为主题被操作时所公开的处理还可以是由执行程序的处理器或者具有处理器(例如,控制设备、控制器和存储系统)的装置执行的处理。此外,当处理器执行程序时被执行的处理的一部分或全部也可以由作为处理器的替代物的硬件电路或者由除了处理器以外的硬件电路来执行。

#### [0060] 第一实施方式

[0061] 图1示出了在其中可以应用本发明的方法和装置的系统的硬件配置的示例。该系统包括服务器100、管理服务器300和形成虚拟存储系统200的一个或多个物理存储系统210(PDKC)。

[0062] 服务器具有CPU、存储器、I/O(输入/输出)I/F等。服务器通过执行诸如数据库管理系统的操作系统和软件来提供服务。数据库管理系统处理的数据被存储在存储系统200中。服务器100经由网络110被耦合到存储系统200。

[0063] 被配置为管理存储系统200的管理服务器300经由网络310耦合到存储系统200。管理服务器300被提供有CPU、存储器(M)301、输入/输出部分和管理I/F。存储器(M)301存储用于管理的程序。CPU执行管理程序以执行管理处理。输入/输出部分例如是由鼠标、键盘、显示器等配置的。输入/输出部分从执行管理的操作者接收各种指令输入,并且在显示器上显示各种信息。管理端口作为与存储系统200通信的中介。管理服务器300可以连接到服务器100并且还可以管理服务器100。

[0064] 在物理存储系统210中,I/O I/F211经由网络110耦合到服务器100,并且作为与服务器100通信的中介。管理I/F212经由网络310耦合到管理服务器300,并且作为与管理服务器300通信的中介。处理器213通过执行已经存储在存储器(S)214的程序单元2143中的各种程序来执行各种处理。此外,处理器213通过使用已经存储在存储器(M)214的控制信息单元2141中的各种信息来执行各种处理。例如,存储器单元(S)214是由例如至少一个存储器设备配置的,并且具有被配置为存储控制信息的控制信息单元2141、被配置为存储程序的程序单元2143和作为被配置为缓存数据的缓存存储器的示例的高速缓存单元2145。通常,与卷217的容量相比,高速缓存部分2145的容量更小。磁盘I/F215经由总线耦合到作为物理存储设备的示例的至少一个HDD216。例如,被配置为管理数据的卷217是由HDD216的至少一个存储区域配置的。物理存储设备不限于HDD并且还可以是SSD(固态驱动器)或DVD。可以在奇偶校验组的单元中收集至少一个HDD216,并且还可以使用诸如RAID(独立磁盘冗余阵列)的高可靠性技术。物理存储系统210可以具有上述资源中的一个或多个。这些资源经由内部总线218相互连接。

[0065] 在虚拟存储系统200中,物理存储系统210经由网络219相互连接。可以为了冗余来复用网络219。将相互连接的两个或更多个物理存储系统210作为一个存储系统(虚拟存储系统)200提供给服务器100。为了实现对物理存储系统210的虚拟化,需要唯一的卷ID分配方法和两个或更多个物理存储系统210之间的资源共享技术。这种技术是已知的。

[0066] 图2是示出了根据第一实施方式的管理服务器300中的存储器(M)301的示例的详细框图。控制信息单元302包含服务器站点管理表303和存储站点管理表304。这些表

存储稍后解释的由处理来使用的信息。在图 3 和图 4 中示出了这些表的细节。使用这些表来检测相同的数据中心中的服务器和存储。程序单元 305 包含用于向存储系统 200 指导卷创建的卷创建程序（管理）306。

[0067] 图 3 是在管理服务器 300 中管理的服务器站点管理表 303 的示例。服务器 ID 是用于在计算机系统中唯一地标识服务器的标识。站点 ID 是用于在计算机系统中唯一地标识站点（数据中心）的标识。

[0068] 图 4 是在管理服务器 300 中管理的存储站点管理表 304 的示例。PDKC ID 是用于在计算机系统中唯一地标识物理存储设备的标识。站点 ID 是用于在计算机系统中唯一地标识站点（数据中心）的标识。

[0069] 图 5 是示出了根据第一实施方式的物理存储系统 210 中的存储器（S）214 的示例的详细框图。控制信息单元 2141 包含全局 ID 管理表 2142。该表存储稍后解释的由处理使用的信息。图 6 中示出了该表的细节。程序单元 2143 包含用于创建新卷的卷创建程序（存储设备）2144。提供高速缓存单元 2145 以用于高速缓存。

[0070] 图 6 是在物理存储系统 210 中管理的全局 ID 管理表 2142 的示例。全局卷 ID 是用于在虚拟化的存储系统 200 中唯一地标识卷的卷标识。PDKC ID 是用于在计算机系统中唯一地标识物理存储系统 210 的标识。本地卷 ID 是用于在物理存储系统 210 中唯一地标识卷的卷标识。卷大小是由全局卷 ID 指定的卷的大小。

[0071] 图 7 是示出了卷创建处理流程的流程图的示例。该处理是由卷创建程序（管理）306 和卷创建程序（存储设备）2144 实现的。在步骤 S100 中，卷创建程序（管理）306 经由管理服务器 300 的输入 / 输出部分从操作者接收卷创建请求。此时，卷创建程序将服务器 ID、全局卷 ID 和卷大小作为参数进行接收。在步骤 S101 中，卷创建程序（管理）306 通过使用服务器站点管理表 303 和存储站点管理表 304 来检测相同站点中的 PDKC。在步骤 S102 中，卷创建程序（管理）306 寻找具有指定大小的空闲空间的一个或多个存储系统。在步骤 S103 中，卷创建程序（管理）306 向 PDKC 指导卷创建。如果在步骤 S102 中找到两个或更多个 PDKC，则可以选择具有最大空闲空间的 PDKC 并且可以选择具有最低负载的 PDKC。可替换地，该程序可以认为一旦找到了满足条件的第一 PDKC 则步骤 S102 就完成，并且前进至步骤 S103。

[0072] 接收指导的卷创建程序（存储设备）2144 创建卷，并且将该卷映射到连接到由服务器 ID 指定的服务器的 I/O I/F（步骤 S200）。该创建的卷可以是虚拟卷。当从服务器 100 写入数据时，虚拟卷技术分配物理存储区域。在很多存储产品中实现该技术。在卷创建之后，卷创建程序（存储设备）2144 向调用方发送完成消息（步骤 S201）。在步骤 S104，卷创建程序（管理）306 从 PDKC 接收完成消息并且终止处理。

[0073] 当服务器仅通过前面提到的处理来指定全局卷 ID 时，确定最佳的物理存储区域。此外，当虚拟化的存储系统在不同的数据中心中包括 PDKC 时，不选择引起 I/O 延迟下降的 PDKC。因此，降低了用于进行配置的设计成本和操作成本。

#### [0074] 第二实施方式

[0075] 第二实施方式涉及一种从构成虚拟化的存储系统的物理存储系统中自动地选择用于远程复制的资源的方法。该方法还执行用于远程复制的自动配置（例如，远程复制路径配置）并且创建传送数据存储区域。通过使用图 8 和图 9 描述了远程复制功能。将远程

复制分为两类：同步远程复制和异步远程复制。

[0076] 图 8 示出了用于同步远程复制处理的远程复制配置的示例。主存储系统 600 是远程复制的源存储系统，并且包括存储由服务器使用的数据的主卷 (PVOL) 601。辅存储系统 700 是远程复制的目的地存储系统，并且包括存储 PVOL601 的复制数据的辅卷 (SVOL) 701。PVOL601 与 SVOL701 之间的关系被称作远程复制对。主存储系统 600 和辅存储系统 700 分别具有用于数据传送的一个或多个 I/O I/F602、702。对于同步远程复制处理而言，主存储系统 600 从服务器 100 接收写入请求。然后，主存储系统 600 向 PVOL601 写入写入数据，并且将写入数据传送到辅存储系统 700。主存储系统 600 不向服务器发送完成消息，直到接收到来自辅存储系统 700 的完成消息为止。接收到复制请求的辅存储系统 700 将经传送的数据写入 SVOL701，并且向主存储系统 600 发送完成消息。最后，主存储系统 600 向服务器 100 发送完成消息。在这种同步远程复制中，在 PVOL601 与 SVOL701 之间实现了数据的高一致性。因此，即使主存储系统 600 出现故障，也可以使由此产生的数据丢失最小化。然而，随着主存储系统与辅存储系统之间的距离增加，存在着延长写入响应时间的缺点。

[0077] 图 9 示出了用于异步远程复制处理的远程复制配置的示例。主存储系统 600 从服务器 100 接收写入请求。然后，主存储系统 600 在 PVOL601 和 JVOL608 中写入数据作为日志 (JNL)。JNL 是传送数据，并且 JVOL608 是用于暂时地存储 JNL 的存储区域。JVOL 不必是卷。JNL 可以存储在高速缓存单元 2145 中。此后，主存储系统 600 向服务器发送完成消息。在辅存储系统 700 中向 JVOL703 异步传送 JNL。最后，向 SVOL701 复制 JVOL703 中存储的 JNL。因此，与同步远程复制相比，写入响应时间更短。然而，还存在以下缺点：与同步远程复制相比，如果主存储系统 600 出现故障，则存在更容易数据丢失和更大量的数据丢失的缺点。在现有的远程复制系统中，操作必须配置辅存储系统、SVOL、JVOL、用于远程复制的 I/O I/F、服务器（辅）、远程复制对等。这些操作非常复杂。

[0078] 图 10 是示出了根据第二实施方式的物理存储系统 210 中的存储器 (S) 214 的示例的详细框图。与图 5 中所示的第一实施方式相比，添加了主存储对表 2146、辅存储对表 2147、p- 对创建程序 2149、s- 对创建程序 2150、写入程序 (存储设备) 2151 和流入表 2148。主存储对表 2146 和辅存储对表 2147 是用于管理 PVOL601 与 SVOL701 之间的关系的表。在图 12 和图 13 中示出了这些表的细节。p- 对创建程序 2149 和 s- 对创建程序 2150 是用于在主存储系统和辅存储系统中创建远程复制对的程序。在图 17 中示出了这些程序的处理流程。写入程序 (存储设备) 2151 是用于向目标卷存储写入数据的程序。在该实施方式中，在图 22 和图 24 中描述了在 PVOL 故障以后接收 SVOL701 写入的两种情况。流入表 2148 是用于管理写入到 PVOL601 中的写入数据的量的表。提供该信息以考虑在主存储系统与辅存储系统之间所需的带宽。

[0079] 图 11 是示出了根据第二实施方式的管理服务器 300 中的存储器 (M) 301 的详细框图。与图 2 中所示的第一实施方式相比，添加了站点距离表 307、对创建程序 308、辅选择程序 309、JVOL 容量程序 311 和迁移站点程序 312。站点距离表 307 是用于管理数据中心之间的距离的表。在图 15 中示出了表的细节。对创建程序 308 和辅选择程序 309 是用于选择复制目的地并且向存储系统指导远程复制对的创建的程序。在图 17 至图 19 中示出了这些程序的处理流程。JVOL 容量程序 311 是用于估计所需的 JVOL 容量的程序。在图 20 中示出了程序的处理流程。迁移站点程序 312 是用于确定迁移的目的地以进行服务器处理的程

序。在图 24 和图 25 中示出了程序的概念和处理。

[0080] 图 12 是物理存储系统 210 中的存储器 (S) 214 中的主存储对表 2146 的示例。图 13 是物理存储系统 210 的存储器 (S) 214 中的辅存储对表 2147 的示例。当物理存储系统 210 被配置为远端复制的主存储系统时, 使用主存储对表 2146。当物理存储系统 210 被配置为远端复制的辅存储系统时, 使用辅存储对表 2147。主卷 ID 是作为复制源卷的卷的 ID。辅 PDKC ID 是具有 SVOL701 的辅物理 PDKC700 的 ID。辅卷 ID 是作为 PVOL601 的复制目的地卷的本地卷 ID。主 PDKC ID 是具有 PVOL601 的主物理 PDKC600 的 ID。

[0081] 图 14 是在物理存储系统 210 中管理的流入表 2148 的示例。该表管理去往每一个卷的写入数据的流入。本地卷 ID 存储每一个卷的 ID, 并且流入存储针对每个卷的每秒写入数据量。

[0082] 图 15 是在管理服务器 300 中管理的站点距离表 307 的示例。站点 ID 的意义与图 3 中的站点 ID 的意义相同。距离存储由站点 ID 指定的数据中心之间的距离。

[0083] 通过使用图 16 至图 21 来描述自动远程复制配置处理。

[0084] 图 16 是用于数据复制的设置屏幕 800 的示例。HA(高可靠性)意味着本地复制。RC(同步)意味着同步类型的远程复制。RC(异步)意味着异步类型的远程复制。设置屏幕 800 具有复选框。通过复选该框来应用诸如远程复制的功能。在所示的示例中, 向具有卷 ID1 的卷应用异步远程复制。

[0085] 图 17 是示出了对创建处理流程的流程图的示例。在步骤 S300 中, 对创建程序 308 调用辅选择程序 309 以选择辅 PDKC。在图 18 中示出了辅选择程序 309 的处理流程。在步骤 S301 中, 已经从辅选择程序 309 接收到辅信息的对创建程序 308 向所选择的存储系统指导卷创建。所选择的存储系统中的卷创建程序(存储设备)2144 的处理与图 7 的步骤 S200 和 S201 相同。在步骤 S302 中, 对创建程序 308 向主 PDKC 指导对创建, 并且等待完成消息。此时, 将操作者指定的全局卷 ID 和在步骤 S301 中选择的本地卷 ID 以及在步骤 S300 中选择的辅 PDKC700 的 ID 作为参数进行发送。

[0086] 接收对创建请求的主 PDKC600 执行 p- 对创建程序 2149。p- 对创建程序 2149 确认是否可以创建对(步骤 S400)。例如, 其包括确定主存储对表 2146 是否具有空闲条目, 主 PDKC600 是否连接到辅 PDKC700 等等。在确认以后, p- 对创建程序 2149 向在步骤 S300 中选择的辅 PDKC 指导对创建。将全局卷 ID、PVOL 和 SVOL 的本地卷 ID 和主 PDKC 的 ID 作为参数进行发送。

[0087] 接收对创建请求的辅 PDKC700 执行 s- 对创建程序 2150。s- 对创建程序 2150 确认该对是否可以被创建并且更新辅存储对表 2147(步骤 S500 和 S501)。在步骤 S502 中, s- 对创建程序 2150 更新全局 ID 管理表 2142 并且向主 PDKC600 发送完成消息。所接收的全局卷 ID 被存储在全局卷字段中, 并且所接收的本地卷 ID 被存储在本地卷 ID 字段中。其自己的 ID 被存储在 PDKC ID 字段中。由本地 ID 指定的卷的大小被存储在卷大小字段中。

[0088] 在步骤 S401 中, p- 对创建程序 2149 在从辅 PDKC700 中接收到完成消息以后, 更新主 PDKC600 中的主存储对表 2146。在步骤 S402 中, p- 对创建程序 2149 向调用方(对创建程序 308)发送完成消息。在步骤 S303 中, 对创建程序 308 终止该处理。

[0089] 图 18 是示出了从图 17 中的对创建程序 308 调用的辅选择程序 309 的处理流程的流程图的示例。该程序用于选择最佳的复制目的地存储设备。在步骤 S600 中, 辅选择程序

309 确认该请求是否是远程复制请求。如果该请求不是远程复制请求，则程序前进至步骤 S601 并且选择其自己的 PDKC 作为复制目的地 PDKC（步骤 S601），并且然后在步骤 S607 中向调用方报告 PDKC ID。这是因为该请求针对本地复制。如果请求是远程复制请求，则辅选择程序 309 标识直接连接到具有 PVOL 的主 PDKC 的存储系统（步骤 S602）。然后，辅选择程序 309 确认该请求是否是同步远程复制请求（步骤 S603）。如果该请求是同步远程复制请求，则程序前进至步骤 S604 并且寻找与指定的主 PDKC 相距约 100km 的存储系统（步骤 S604）。100km 的意义是可以大体上提供可允许的 I/O 响应时间以进行同步远程复制的距离。可替换地，可以应用可以提供可允许的 I/O 响应时间的一些其它距离。如果该请求不是同步远程复制请求，则辅选择程序 309 寻找与指定的主 PDKC 相距约 1000km 或者更远的存储系统（步骤 S605）。1000km 的意义是可以避免主数据中心灾难的后果的距离。可替代地，可以应用可以避免灾难的后果的某一其它距离。在步骤 S604 或者 S605 之后，在步骤 S606 中，辅选择程序 309 选择满足条件（例如，针对 SVOL 的空闲容量、连接性等等）的 PDKC。最后，辅选择程序 309 向调用方报告决定的 PDKC ID（步骤 S607）。

[0090] 通过使用图 19 和图 20 来描述异步远程复制的情况。

[0091] 图 19 是示出了用于异步远程复制的对创建处理流程的流程图的另一个示例。在步骤 S700 中，对创建程序 308 调用 JVOL 容量程序 311 以计算 JVOL608、703 所需的容量。对创建程序 308 调用辅选择程序 309 并且指导卷创建（步骤 S701 和 S702）。步骤 S701 和 S702 与图 17 中的步骤 S300 和 S301 相同。卷创建程序（存储设备）2144 的处理与图 7 的步骤 S200 和 S201 相同。然后，对创建程序 308 向主 PDKC 和辅 PDKC 指导 JVOL 创建（步骤 S703）。将在步骤 S701 中计算出的 JVOL 大小作为参数进行发送。接收 JVOL 创建请求的每一个 PDKC 创建具有指定的大小的卷，并且该创建的卷可以是虚拟的卷（步骤 S800）。PDKC 中的每一个向调用方发送完成消息（步骤 S801）。最后，在步骤 S704 中，对创建程序 308 向主 PDKC 指导对创建。步骤 S704 之后的后续步骤与图 17 中的步骤 S302 之后的后续步骤相同。

[0092] 图 20 是示出了用于计算 JVOL608、703 所需的容量的 JVOL 容量程序 311 的处理流程的流程图的示例。下面将描述 JVOL608、703 的容量。在例如主 PDKC600 与辅 PDKC700 之间发生路径故障的情况下，对于 JVOL608 而言，JNL 开始维持。在此后解决路径故障的情况下，JVOL608 中存储的 JNL 被发送到 JVOL703。因此，在 JVPL608 具有足以存储已经生成的 JNL 的容量的情况下，可以在无需暂停和重新同步远程复制对的情况下继续远程复制。根据用户希望抵抗路径故障的所需时间来设计 JVOL608 的容量。因此，可以通过使用所需的时间和从服务器 100 写入的数据的量来计算 JVOL 容量。JVOL 容量程序 311 从主 PDKC600 获得针对 PVOL601 的每秒写入数据量（步骤 S900）。然后，当发生路径故障时，JVOL 容量程序 311 获得所需的持续时间（步骤 S901）。通过使用输入 / 输出部分向管理服务器 300 输入所需的持续时间。最后，JVOL 容量程序 311 计算 JVOL 容量（步骤 S902）。通过写入数据的量乘以持续时间来计算 JVOL 容量。

[0093] 图 21 是示出了辅选择程序 309 的处理流程的流程图的另一个示例。该示例考虑主 PDKC 与辅 PDKC 之间的带宽。仅解释与图 18 不同的步骤。在步骤 S602 之后执行额外的步骤 S1000 和 S1001，并且在步骤 S605 或者步骤 S604 之后执行步骤 S1002 而不是步骤 S606。在步骤 S1000 中，辅选择程序 309 从主 PDKC600 获得针对 PVOL 的每秒写入数据量 (MB/s)。

在步骤 S1001 中, 辅选择程序 309 计算主 PDKC 与辅 PDKC 之间所需的带宽。在辅 PDKC 选择步骤 (S1002) 中, 辅选择程序 309 考虑计算出的带宽。

[0094] 图 22 是示出了在 I/O 故障之后的处理的概念图。服务器 100 使用其卷 ID 为 0 的卷。服务器 100 具有交替路径软件。在图 22 的示例中, 交替路径软件具有去往卷 0 的两个路径。虚拟化的存储系统 200 被提供给服务器 100。在存储系统层中, 向物理存储系统 600 发起来自服务器 100 的 I/O, 并且向 (物理存储系统 700 的) 卷 820 复制在 (物理存储系统 600 的) 卷 810 中存储的数据。当主 PDKC 出现故障时, 交替路径软件改变 I/O 路径并且再次发起 I/O。在存储系统层中, 向物理存储系统 700 发起来自服务器 100 的 I/O。在物理存储系统 700 的卷 820 中存储写入数据。因为向服务器 100 提供的卷 ID 是相同的, 因此服务器 100 没有意识到物理存储系统的改变。通常, 为了一致性, 远程复制的 SVOL 拒绝服务器 I/O。在 PVOL 故障以后, SVOL 的状态必须变为可以接收和处理服务器 I/O 的状态。在图 23 中示出了该处理。

[0095] 图 23 是示出了在 PVOL 出现故障以后向 SVOL 发起的写入处理的流程图的示例。该写入程序 (存储设备) 2151 接收去往 SVOL 的 I/O (步骤 S1100) 并且确认 PVOL 故障是否已经发生 (步骤 S1101)。如果 PVOL 故障未发生, 则写入程序 (存储设备) 2151 向服务器 100 报告错误 (步骤 S1102)。如果 PVOL 故障已经发生, 则写入程序 (存储设备) 2151 将 SVOL 状态改变为 I/O 可接收的状态 (步骤 S1103)。然后, 写入程序 (存储设备) 执行 I/O 并且终止处理 (步骤 S1104 和 S1105)。

[0096] 在图 22 中示出了服务器未经改变的情况。在本发明的配置的情况下, 因为安装在不同的数据中心中的物理存储系统 210 可以构成虚拟化的存储系统 200, 因此 I/O 目的地物理存储系统可以改变为与其中安装了服务器的数据中心不同的数据中心。因此, 需要超出数据中心的 I/O 并且性能下降。为了避免这种情况, 还可以改变服务器。

[0097] 图 24 是示出了在 PVOL 出现故障以后的服务器变化和涉及 SVOL 的处理中的后续变化的概念图。仅解释与图 22 的区别。服务器 400 作为迁移的源被安装, 并且服务器 500 作为迁移的目的地被安装。服务器 400 连接到服务器 500。在物理服务器 400 和 500 中的每一个中创建虚拟机 410。通过虚拟机来为应用提供服务。可以从一个物理服务器向另一个物理服务器迁移虚拟机。可以通过传统的技术来执行迁移。在卷 810 故障以后或者在 PDKC600 和 I/O 目的地改变为 PDKC700 中的卷 820 以后, 虚拟机 410 从服务器 400 迁移到处于与 PDKC700 相同的数据中心中的服务器 500。因此, 可以避免超出数据中心的 I/O。

[0098] 图 25 是在 PVOL 出现故障以后的虚拟机 410 的迁移处理以及用于确定目的地的处理的流程图的示例。在步骤 S1200 中, 写入程序 (服务器) 检测 I/O 故障。与预定的值相比更长的响应时间也可以是检测的条件。写入程序 (服务器) 向管理服务器 300 查询迁移的目的地。接收查询的管理服务器 300 调用迁移站点程序 312。迁移站点程序 312 向具有 PVOL 的 PDKC 查询具有 SVOL 的 PDKC 的信息 (步骤 S1300)。PDKC 可以使用主存储对表 2146 来对其进行确认。然后, 迁移站点程序 312 向服务器报告 PDKC 的所获得的信息 (步骤 S1301)。最后, 写入程序 (服务器) 向所报告的 PDKC 移动虚拟机 410。

[0099] 通常, 远程复制功能具有由一个或多个远程复制对构成的一致性组 609。图 26 是示出了远程复制对的一致性组 609 的概念图。由一个应用使用的一个或多个远程复制对应当属于相同的一致性组 609。操作者可以针对一致性组 609 执行远程复制操作。因此,

操作者不需要操作每一个远程复制对。此外,使用远程复制的异步类型,在恢复处理中向 SVOL701 确认去往属于相同的一致性组 609 的 PVOL601 的写入顺序。在远程复制对创建处理时,应当从相同的数据中心中选择属于相同的一致性组 609 的 PVOL 的复制目的地。否则,故障以后的 I/O 将超出数据中心。

[0100] 图 27 是示出了考虑一致性组 609 的辅选择程序 309 的处理流程的流程图的另一个示例。仅解释与图 21 的区别。在步骤 S1001 之后执行额外的步骤 S1400、S1401 和 S1402。首先,当使用复选框应用远程复制时,由操作者经由设置屏幕 800 输入一致性组 ID。在步骤 S1400 中,辅选择程序 309 从主 PDKC 获得指定的一致性组信息。在步骤 S1401 中,辅选择程序 309 确定指导的远程复制对是否是一致性组的第一对。如果是,则辅选择程序 309 执行步骤 S603 和 S607 并且确定辅 PDKC。否则,如果指导的远程复制对不是一致性组的第一对,则辅选择程序 309 前进至步骤 S1402,并且选择具有已经属于指定的一致性组的 SVOL 的辅 PDKC(步骤 S1402)。然后,辅选择程序 309 前进至步骤 S607,并且向调用方报告 PDKC ID。

[0101] 通过使用上面提到的技术,用户的操作将是指定全局卷 ID 以及应用远程复制功能。因此,准备 SVOL、配置远程复制路径等变得没有必要。

[0102] 图 28 是示出了远程复制配置的自动恢复方法的概念图。仅解释与图 24 的区别。通过使用图 24 所示的技术,在移动虚拟机 410 并且将 I/O 目的地改变到 PDKC700 中的卷 820 以后,计算机系统处于不应用远程复制的状态中。因此,可靠性较低。自动恢复方法基于在初始设置时输入的内容自动地确定复制目的地,并且重新启动远程复制处理。添加对于配置恢复所需的第三物理存储设备 900,如图 28 所示。第三物理存储设备 900 中的卷 830 是 PDKC700 中的卷 820 的复制目的地,并且也具有相同的全局卷 ID0。

[0103] 通过图 29 中所示的管理服务器 300 中的配置恢复程序来实现上面的处理。可以在图 25 中的迁移站点程序 312 的步骤 S1301 之后调用配置恢复程序。首先,配置恢复程序接收 PVOL 或主 PDKC600 故障的通知(步骤 S1500)。该配置恢复程序获得针对目标卷 820 的初始设置的内容(步骤 S1501)。在图 28 的示例中,目标卷是具有全局卷 ID0 的卷。例如,内容是“RC(同步)”、“RC(异步)”等。然后,配置恢复程序使用获得的内容调用对创建程序 308(步骤 S1502)并且终止处理(步骤 S1503)。

[0104] 在服务器 400 发生故障而不是 PDKC600 发生故障的情况下,远程复制的复制目的地可以是 PDKC600。在该情况下,卷 820 成为复制源卷,并且卷 810 成为复制目的地卷。在一般的远程复制功能中,只将不同的数据从卷 820 复制到卷 810。

#### [0105] 第三实施方式

[0106] 第三实施方式涉及远程复制配置的自动恢复方法和将上面的技术应用于三数据中心远程复制配置的方法。远程复制功能支持同步类型的远程复制和异步类型的远程复制的组合。在图 30 和图 31 中示出了两个示例。

[0107] 图 30 是示出了三数据中心远程复制的级联配置的概念图。级联的远程复制将 PVOL 中存储的数据同步地复制到 SVOL,并且将 SVOL 中存储的数据异步地复制到 TVOL(第三卷)。

[0108] 图 31 是示出了三数据中心远程复制的多目标配置的概念图。多目标远程复制将 PVOL 中存储的数据同步地复制到 SVOL,并且将 PVOL 中存储的数据异步地复制到 TVOL。

[0109] 级联和多目标的远程复制被简单地称作三数据中心远程复制(3DC 远程复制)。这

些功能既实现了作为同步类型的优点的无数据损失,又实现了作为异步类型的优点的更短的响应时间。

[0110] 图 32 是用于三数据中心远程复制的设置屏幕 800 的另一个示例。添加了用于级联配置和多目标配置的复选框。

[0111] 图 33 是示出了用于考虑到三数据中心远程复制来确定复制目的地的 3DC 目的地选择程序的处理流程的流程图的示例。3DC 目的地选择程序是由于到图 32 的设置屏幕 800 的输入而从对创建程序 308 调用的。

[0112] 在步骤 S1600 中,3DC 目的地选择程序识别直接连接到主 PDKC 的存储系统。该步骤与图 18 中的步骤 S602 是相同的。然后,3DC 目的地选择程序确定同步类型的远程复制的复制目的地。3DC 目的地选择程序查找与指定的主 PDKC 仅相距约 100km 的存储系统(步骤 S1601),并且选择满足作为同步远程复制的条件的辅 PDKC(步骤 S1602)。3DC 目的地选择程序确定指定的类型是否是级联的(步骤 S1603)。如果指定的类型不是级联的,则 3DC 目的地选择程序从在步骤 S1600 中获得的信息中查找与主 PDKC 和辅 PDKC 相距约 1000km 或者更远的存储系统,并且然后前进至步骤 S1607。如果指定的类型是级联的,则 3DC 目的地选择程序识别直接连接到所选择的辅 PDKC 的存储系统(步骤 S1605),并且从在前一步骤 S1605 中获得的信息中查找与主 PDKC 和辅 PDKC 相距约 1000km 或者更远的存储系统(步骤 S1606)。在(步骤 S1604 或步骤 1606 之后的)步骤 S1607 中,3DC 目的地选择程序选择满足诸如异步远程复制的条件的 PDKC。最后,3DC 目的地选择程序向调用方报告辅 PDKC ID 和第三 PDKC ID(步骤 S1608)。

[0113] 在级联的情况下,不能检测到最佳的 PDKC。因此,3DC 目的地选择程序返回步骤 S1602,并且在选择下一个最佳的辅 PDKC 之后再次执行步骤 S1605 至步骤 S1607。

[0114] 图 33 中的步骤可以作为辅选择程序 309 的一部分被执行(图 21 或图 27)。例如,这些步骤可以在图 27 中的步骤 S1001 之后被执行。在该情况下,针对三数据中心配置考虑 PDKC 之间的带宽等。

[0115] 当然,图 1、8、9、24、26、28、30 和 31 中所示的系统配置仅仅是在其中实现本发明的信息系统的示例,并且本发明并不限于特定的硬件配置。实现本发明的计算机和存储系统还可以具有已知的 I/O 设备(例如,CD 和 DVD 驱动器、软盘驱动器、硬盘驱动器等),该已知的 I/O 设备可以存储和读取用于执行上面所描述的发明的模块、程序和数据结构。可以在这些计算机可读介质上编码这些模块、程序和数据结构。例如,本发明的数据结构可以独立于其上存储有本发明中使用的程序的一个或多个计算机可读介质而存储在计算机可读介质上。可以通过数字数据通信的任何形式或介质(例如,通信网络)来对系统的组件进行互连。通信网络的示例包括局域网、诸如因特网的广域网、无线网络、存储区域网络等。

[0116] 在说明书中,给出了大量细节以用于解释的目的,从而提供对本发明的彻底理解。然而,对于本领域技术人员而言显而易见的是,为了实现本发明,不是所有这些具体细节都是必需的。还应当注意的是,可以将本发明描述为通常描绘为流程图、流程示意图、结构示意图或框图的过程。虽然流程图可以将操作描述为连续过程,但是可以并行地或并发地执行这些操作中的很多操作。此外,可以重新排列操作的顺序。

[0117] 如本领域中公知的,上文所描述的操作可以由硬件、软件或者软件和硬件的某一组合来执行。本发明的实施方式的各个方面可以使用电路和逻辑设备(硬件)来执行,而

其它方面可以使用存储在机器可读介质（软件）上的指令来执行，所述指令如果由处理器执行，则将使得处理器执行用于实现本发明的实施方式的方法。此外，本发明的一些实施方式可以单独地用硬件来执行，而其它实施方式可以单独地用软件来执行。此外，所描述的各个功能可以在单个单元中执行，或者可以以任意数量的方式分布在大量组件中。当由软件执行时，这些方法可以由诸如通用计算机的处理器基于存储在计算机可读介质上的指令来执行。如果期望的话，指令可以以压缩和 / 或加密的格式存储在介质上。

[0118] 根据上述内容将清楚的是，本发明提供了用于管理包括远程复制系统的存储系统并且通过使复杂的操作自动化来改进可管理性的方法、装置和存储在计算机可读介质上的程序。此外，虽然已经在本说明书中示出和描述了具体的实施方式，但是本领域普通技术人员将清楚的是，可以用被计算为实现相同的目的的任何安排来替代所公开的具体实施方式。本公开内容旨在涵盖本发明的任何或所有调节或变化，并且应当理解的是，下面的权利要求中使用的术语不应当理解为将本发明限制于说明书中公开的具体实施方式。更确切地说，本发明的范围应当完全由下面的权利要求进行确定，其应当根据所建立的权利要求解释的教导以及这些权利要求有权享有的等同形式的整个范围来进行理解。

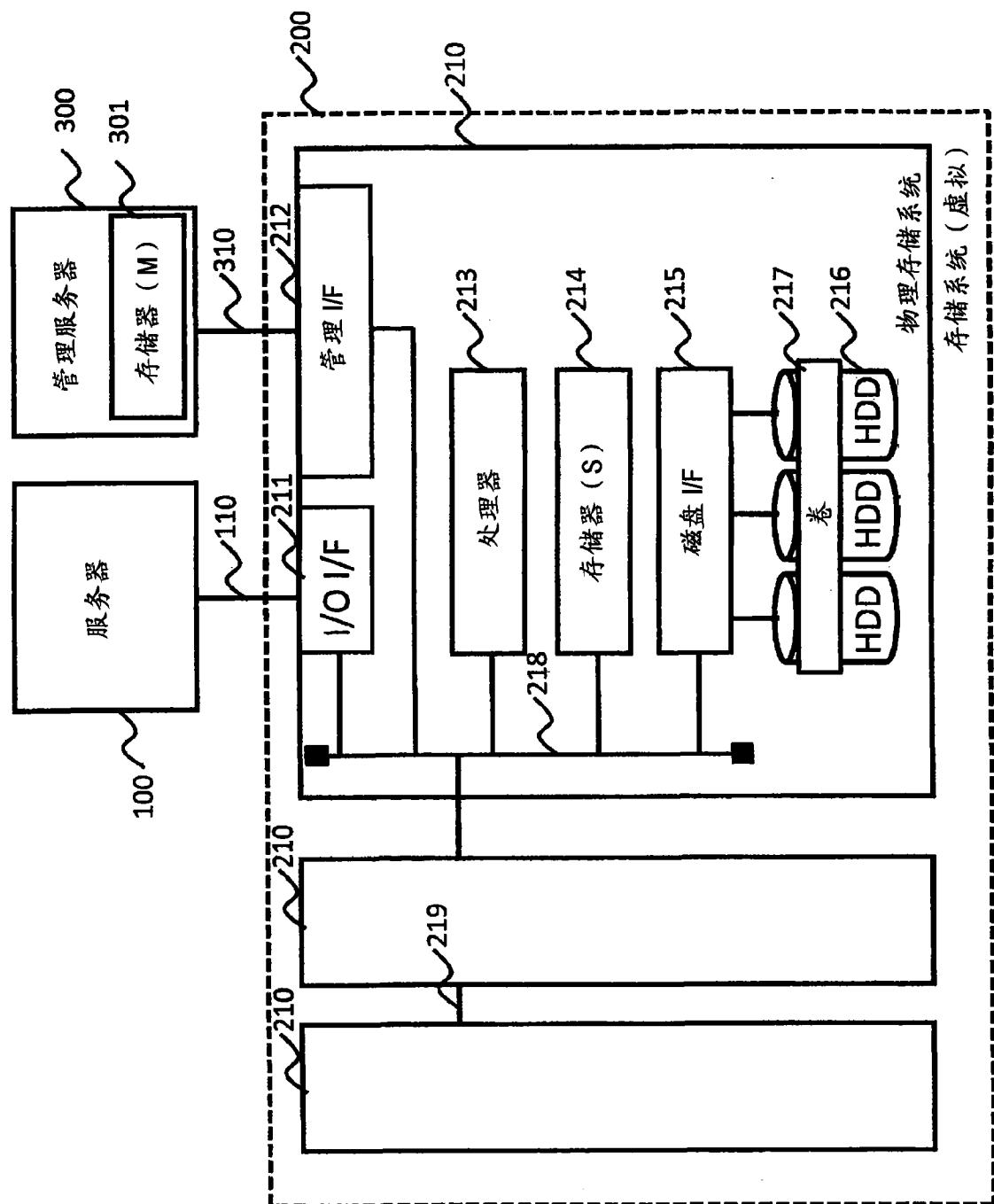


图 1

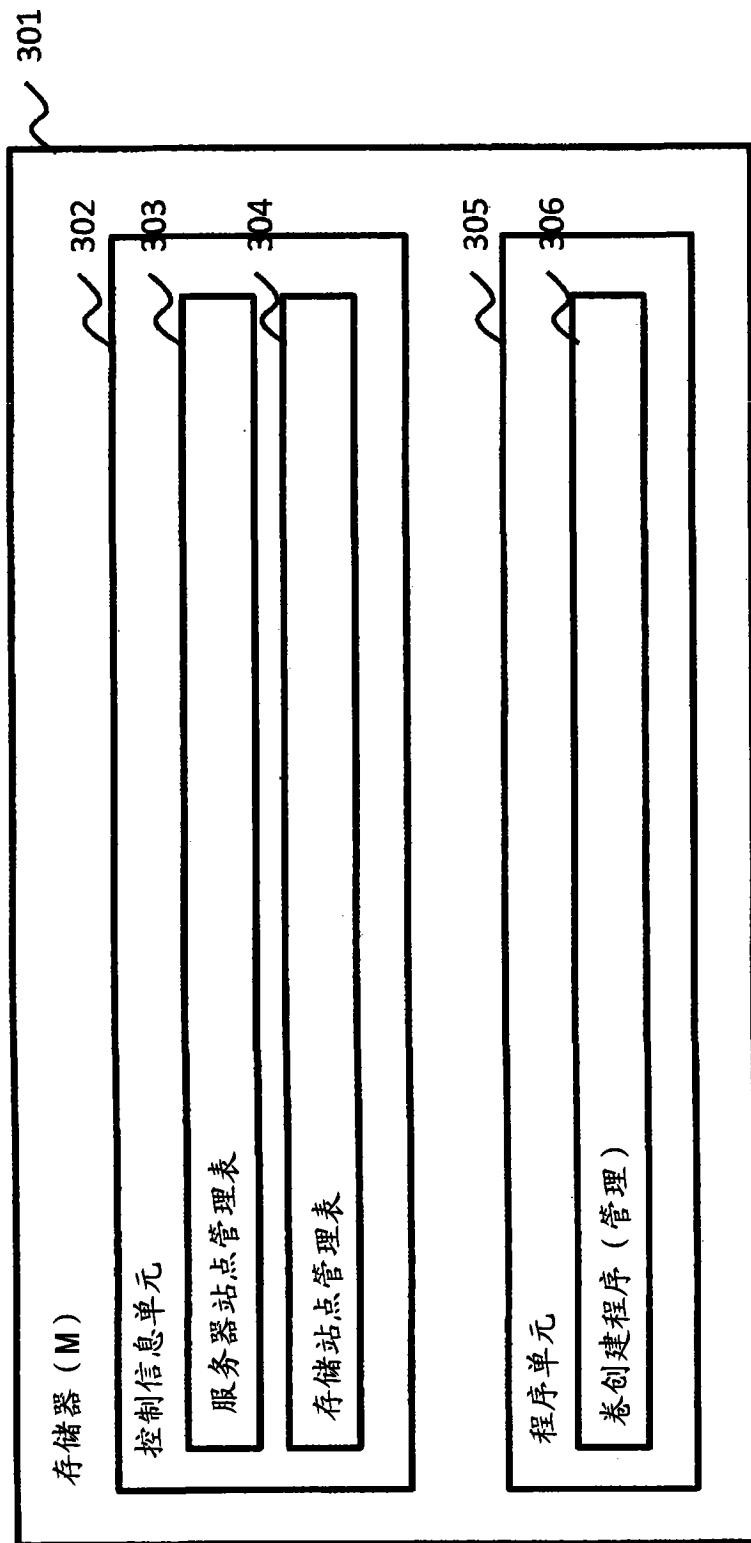


图 2

服务器 ID	站点 ID
0	0
1	0
2	0
3	0
...	...

图 3

PDKC ID	站点 ID
0	0
1	0
2	1
3	1
...	...

图 4

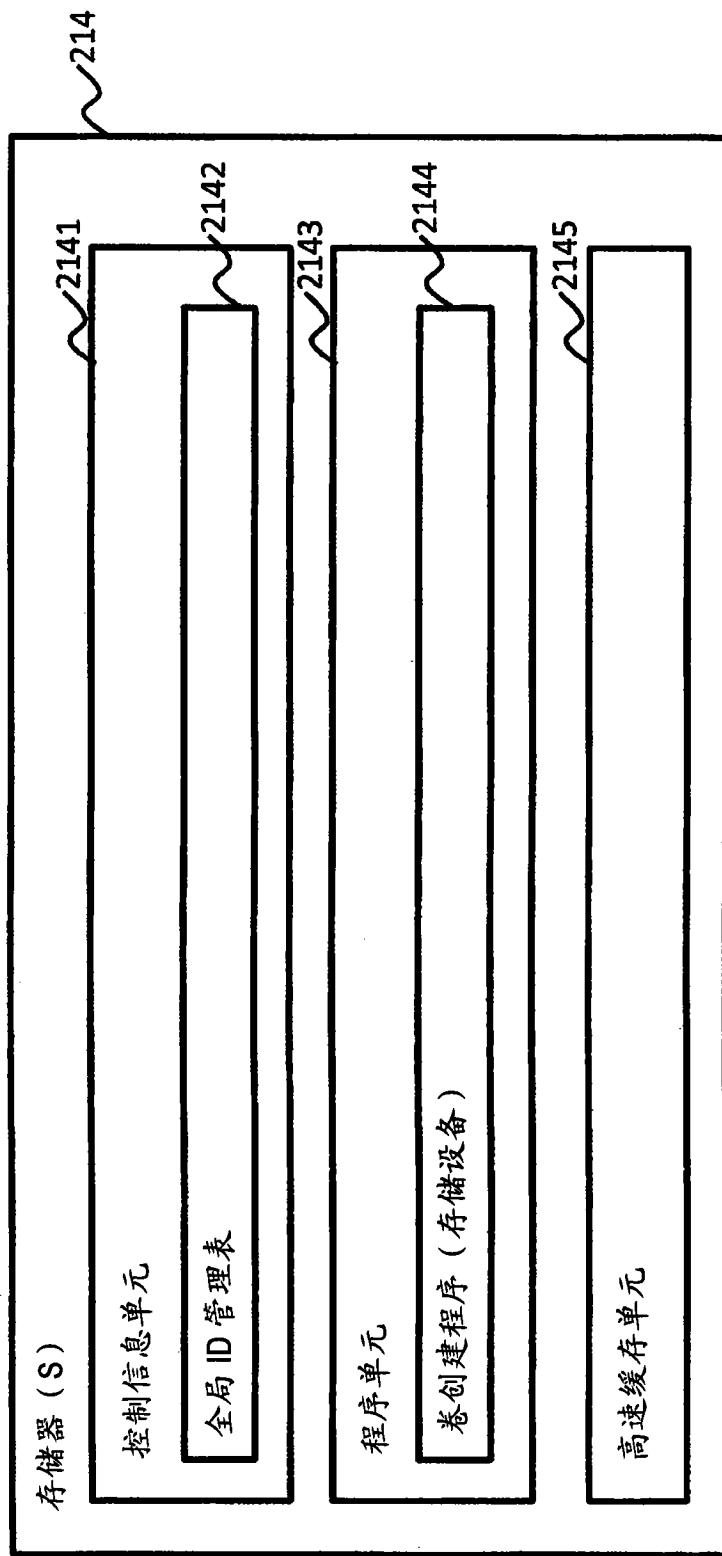
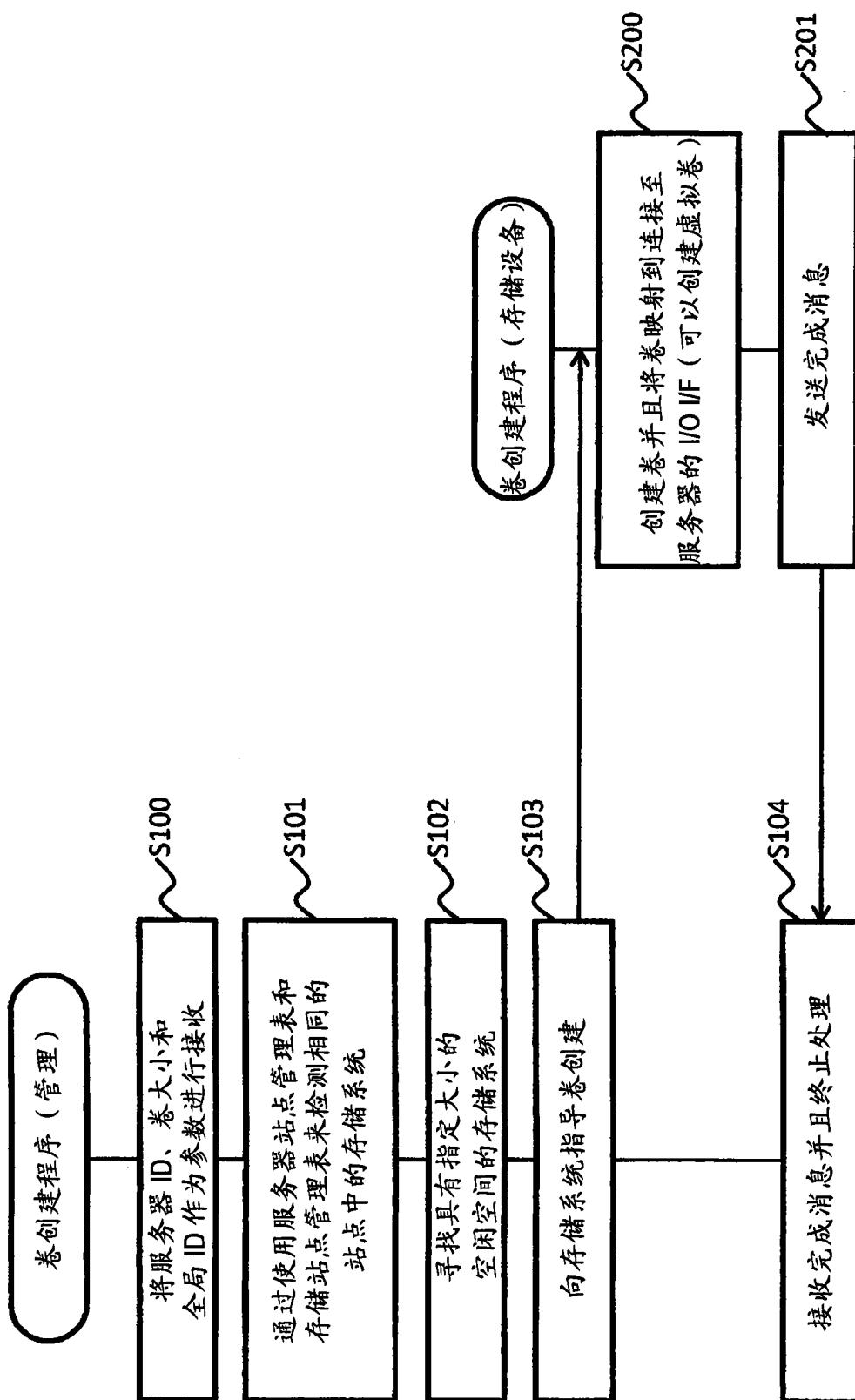


图 5

2142

全局卷 ID	PDKC ID	本地卷 ID	卷大小
0	0	0	500GB
1	0	1	800GB
2	2	0	300GB
3	3	5	500GB
...	...	...	...

图 6



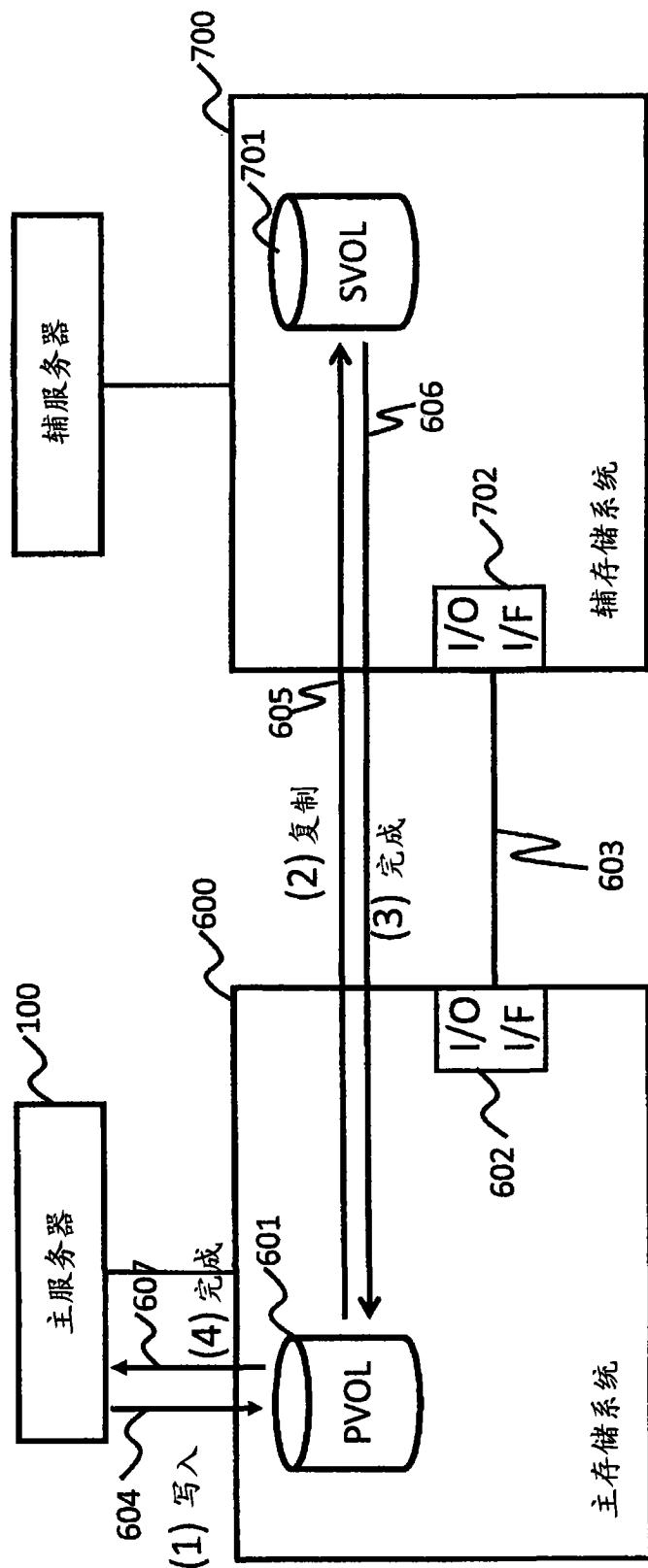


图 8

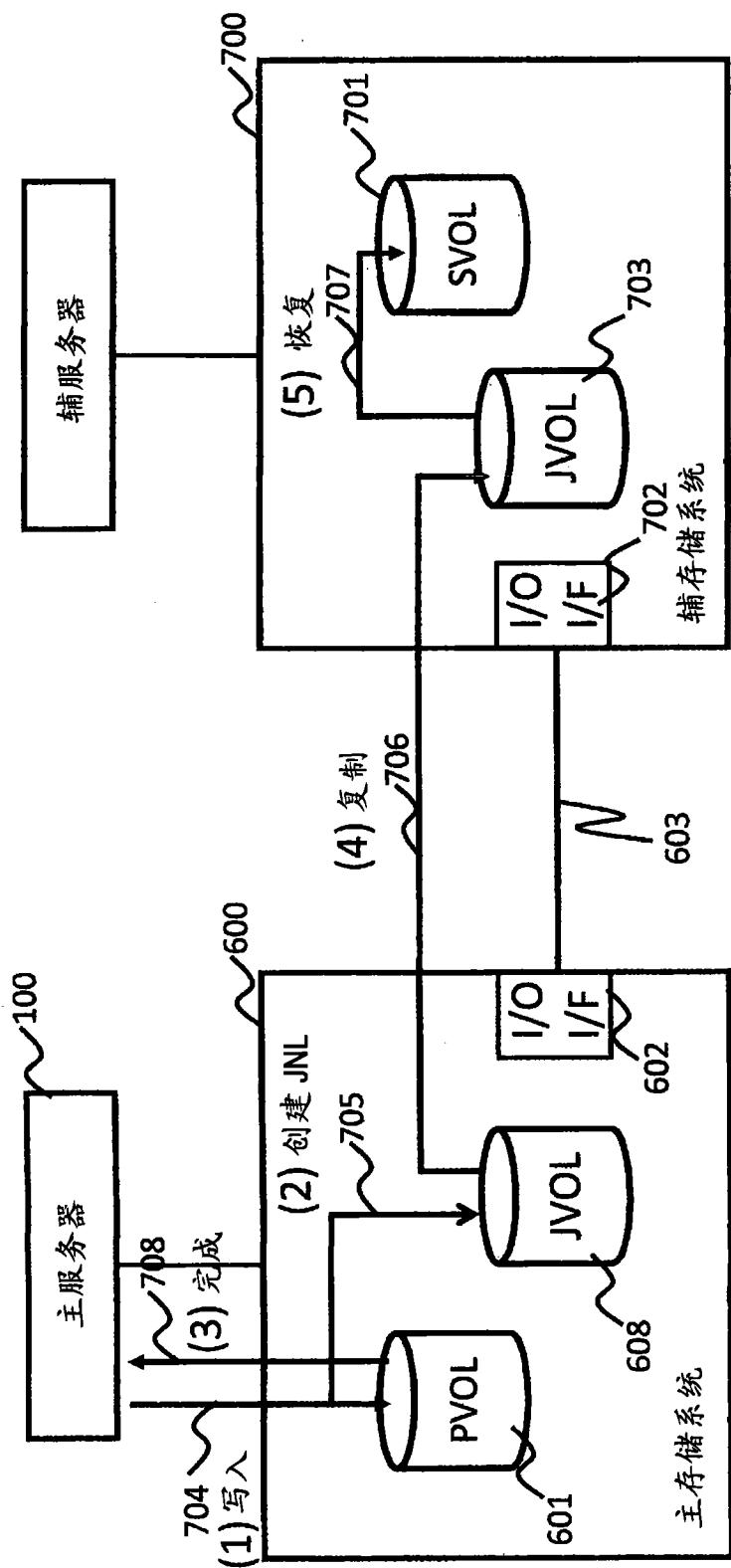


图 9

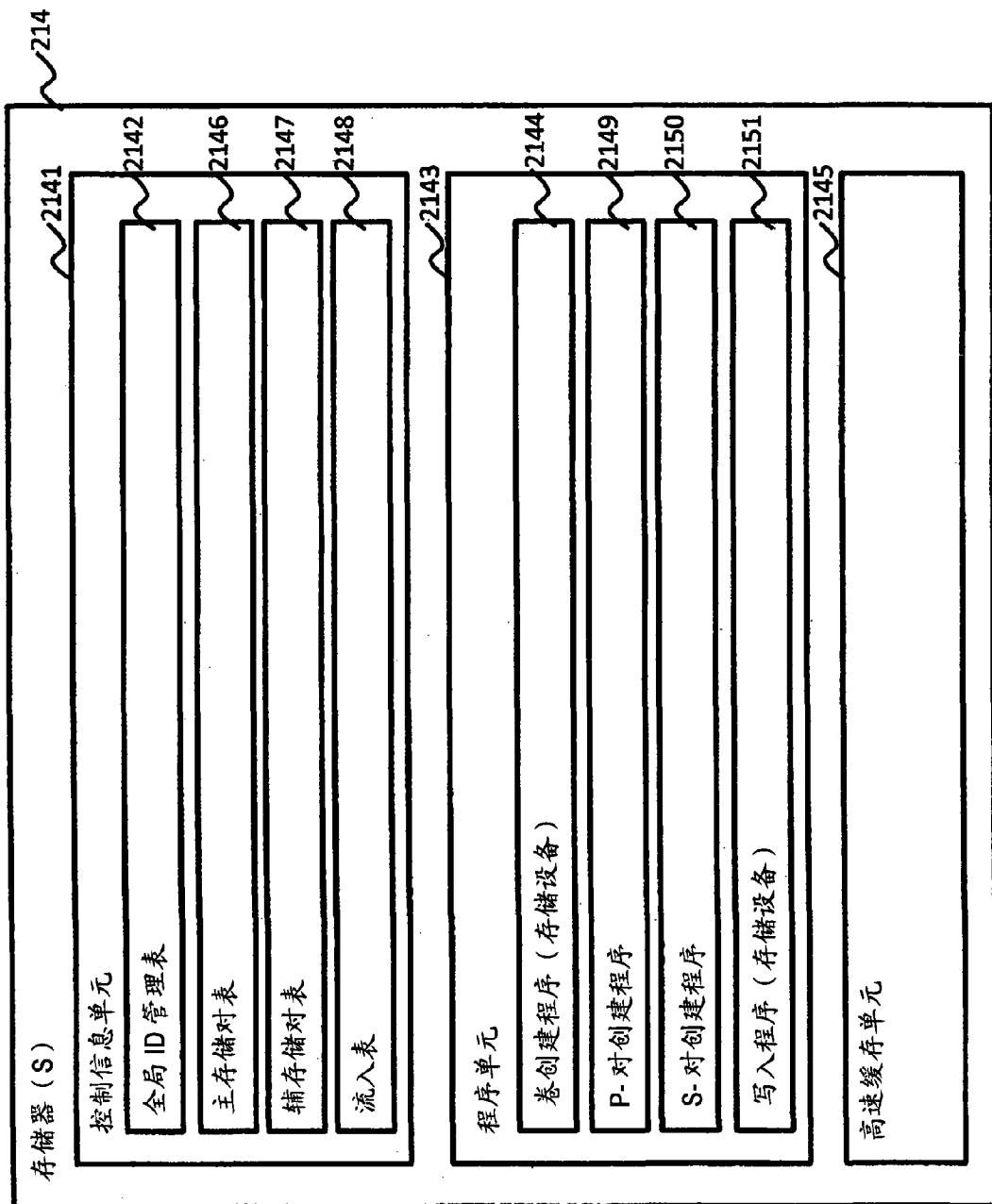


图 10

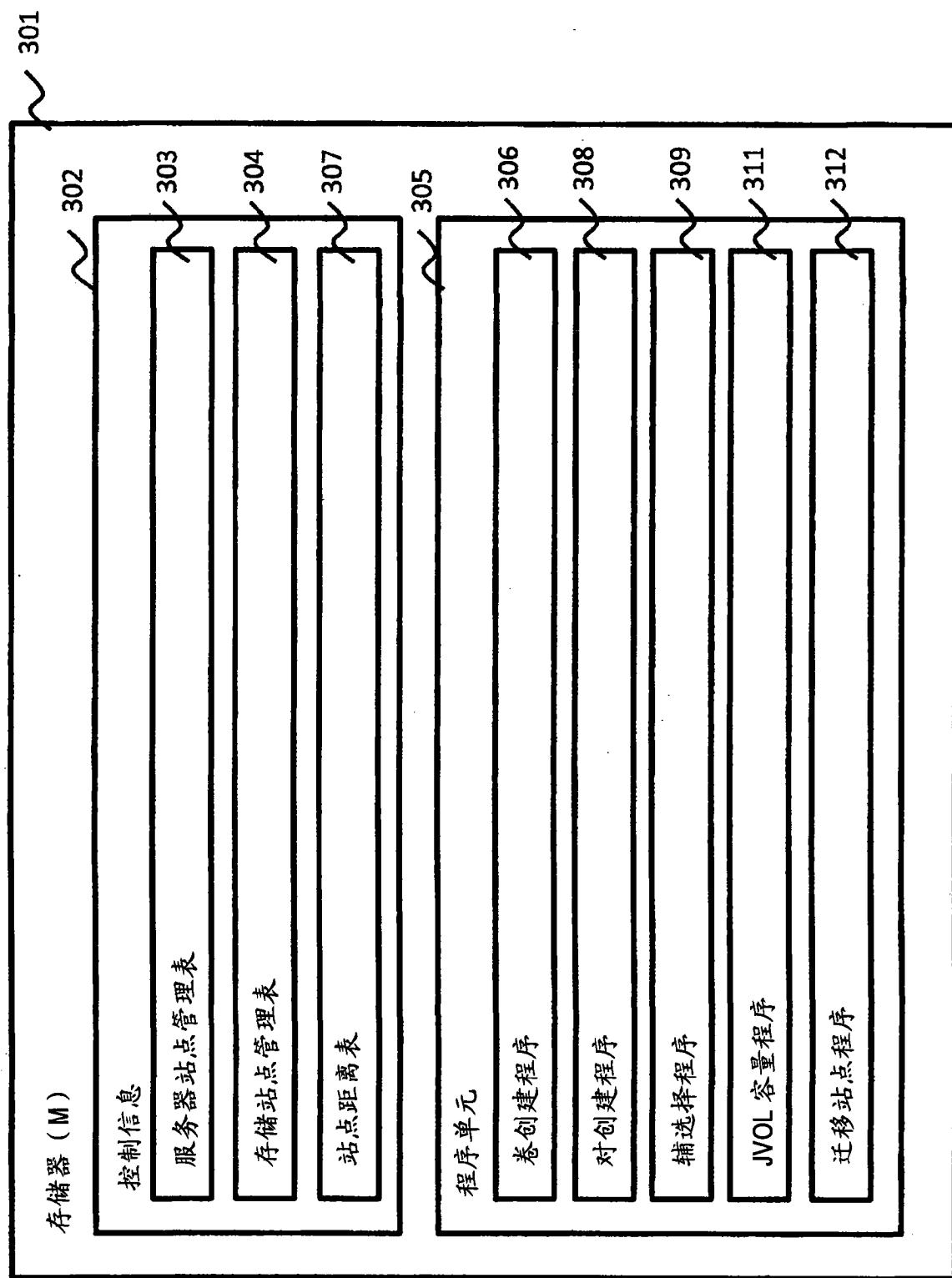


图 11

2146

主卷 ID	辅 PDKC ID	辅卷 ID
0	2	0
1	2	10
2	3	3
3	3	6
...	...	...

图 12

2147

辅卷 ID	主 PDKC ID	主卷 ID
0	0	0
1	0	9
2	0	12
3	0	5
...	...	...

图 13

本地卷 ID	流入
0	5MB/s
1	3MB/s
2	10MB/s
...	...

图 14

站点 ID	站点 ID	距离
0	1	100km
0	2	1000km
0	3	300km
1	2	900km
1	3	200km
...	...	...

图 15

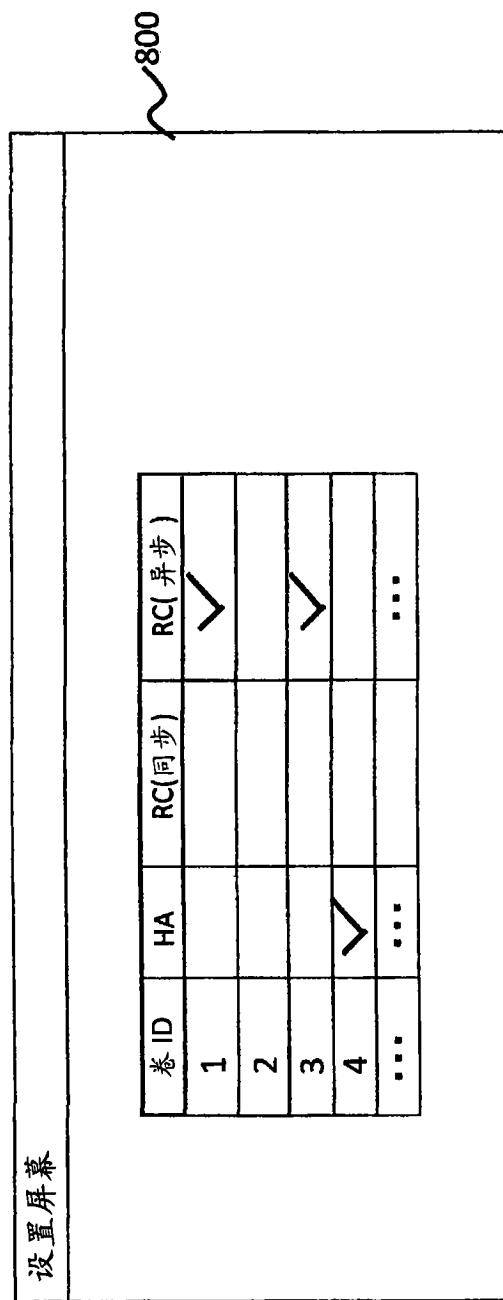


图 16

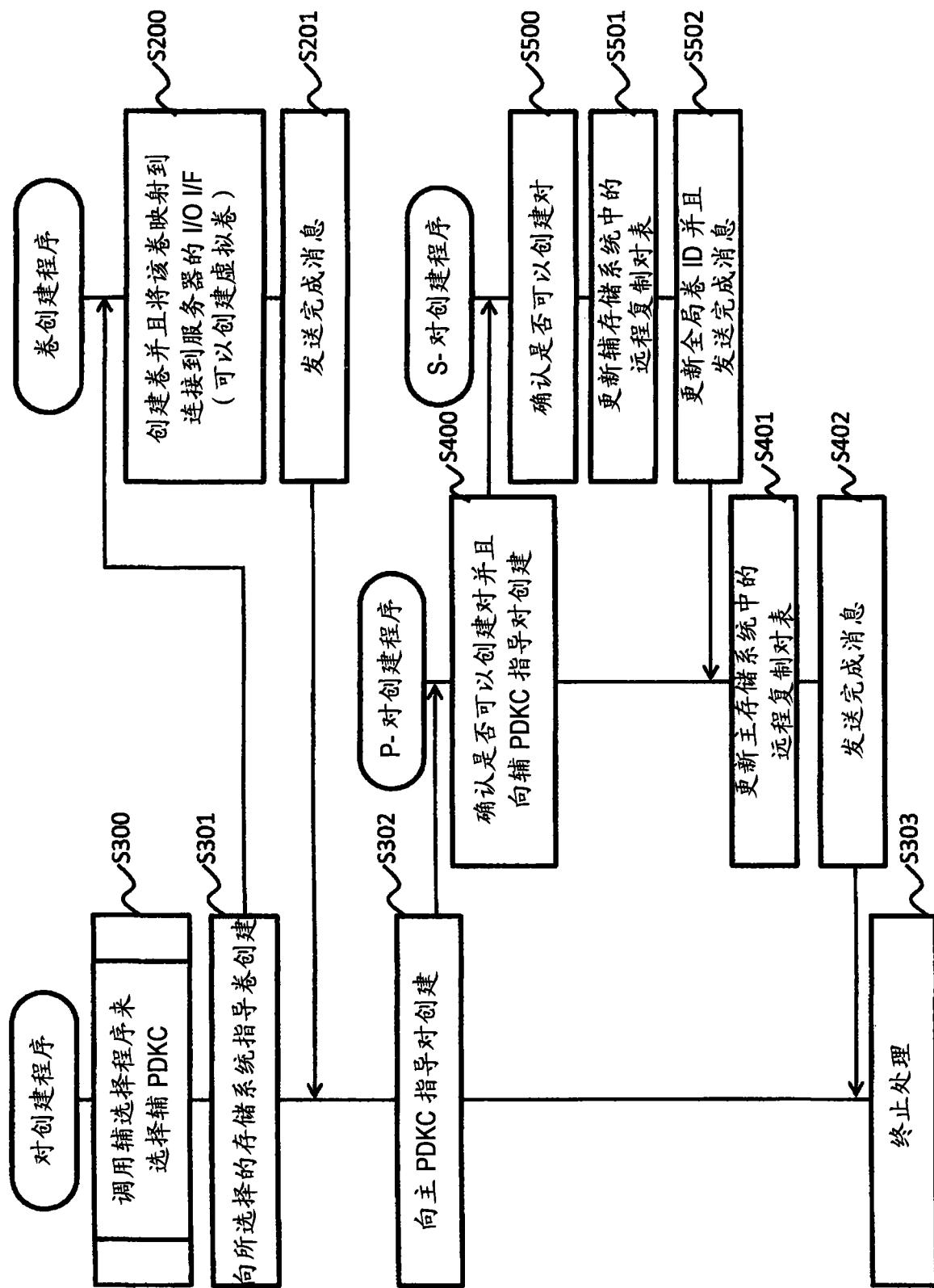


图 17

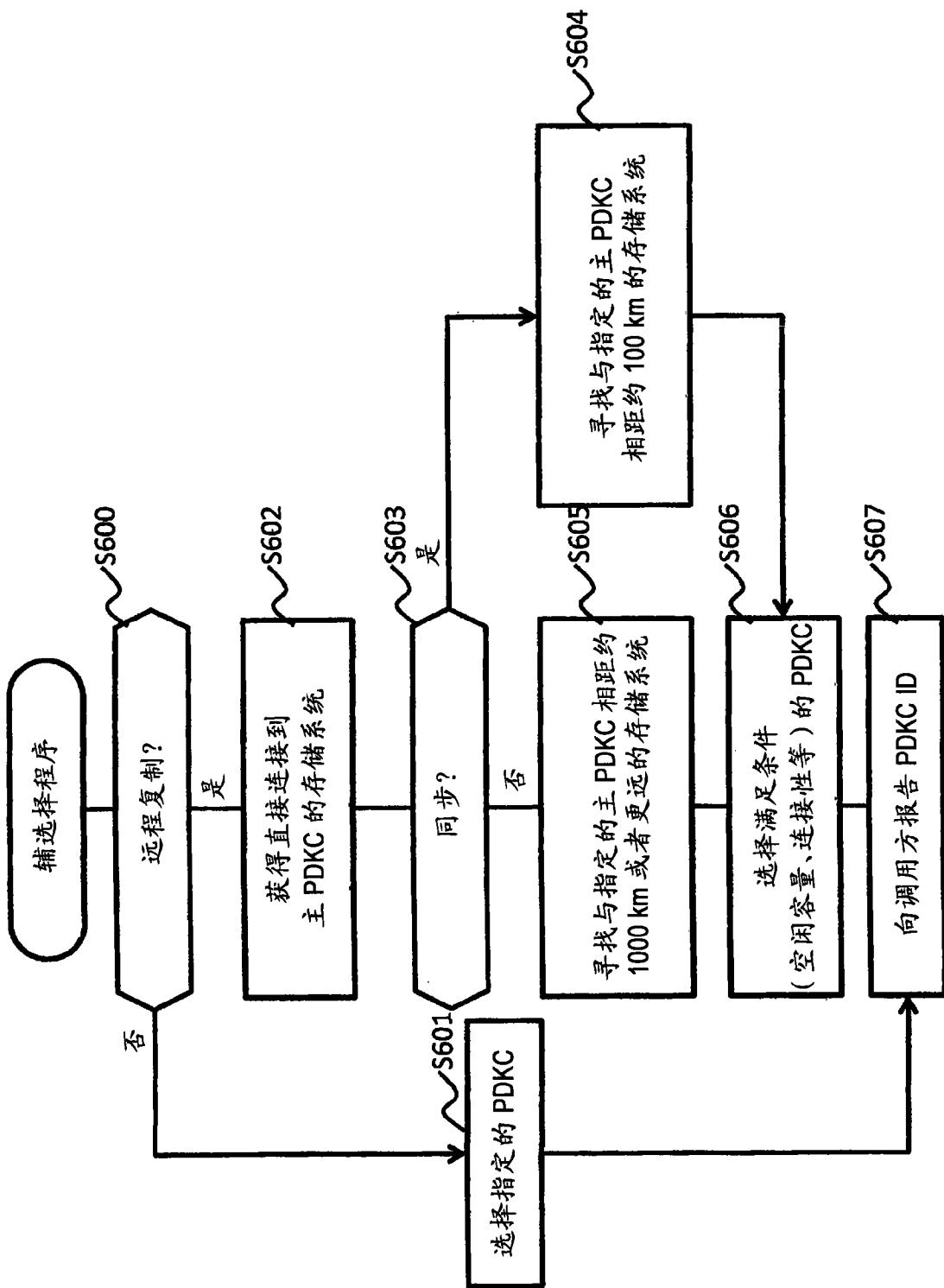


图 18

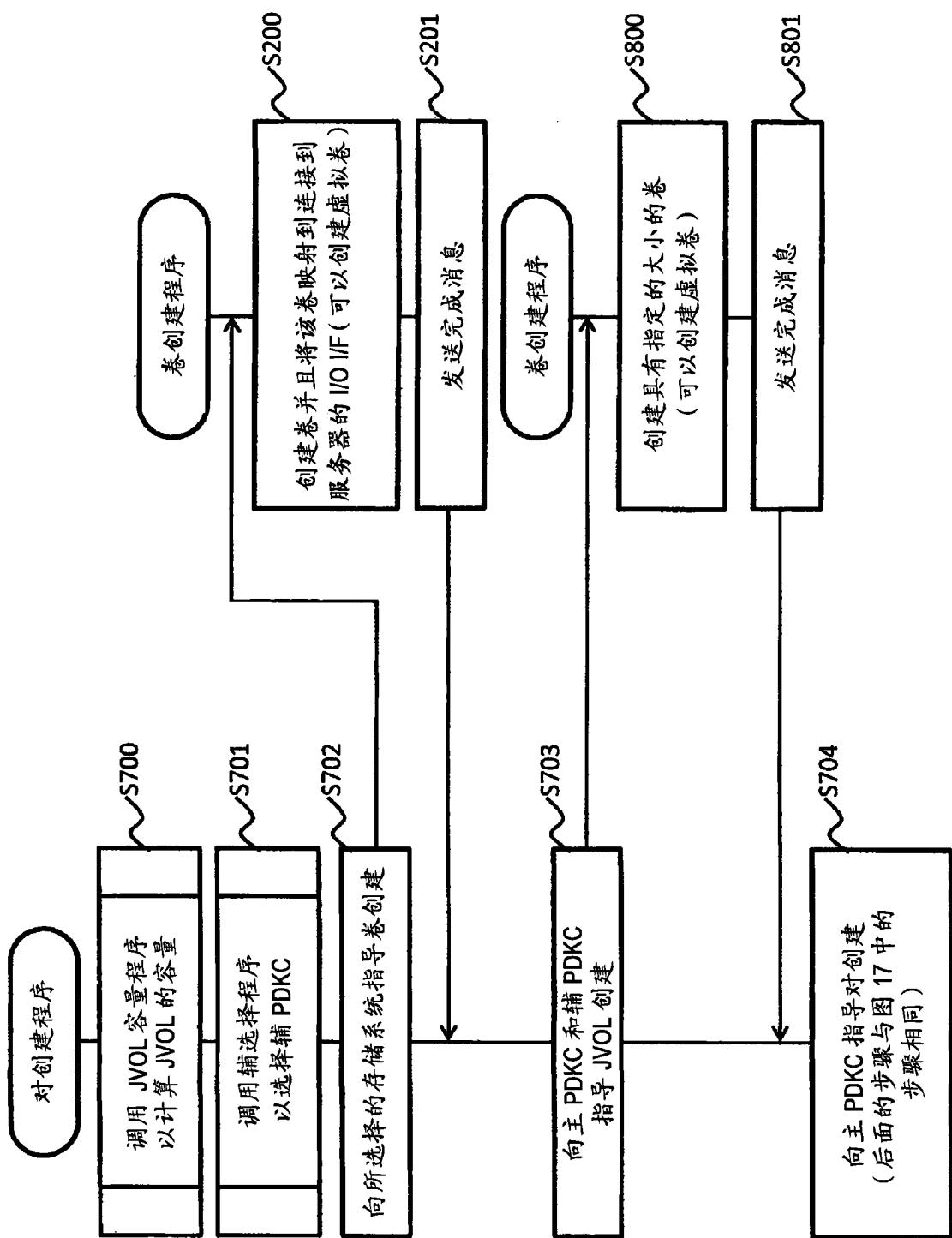


图 19

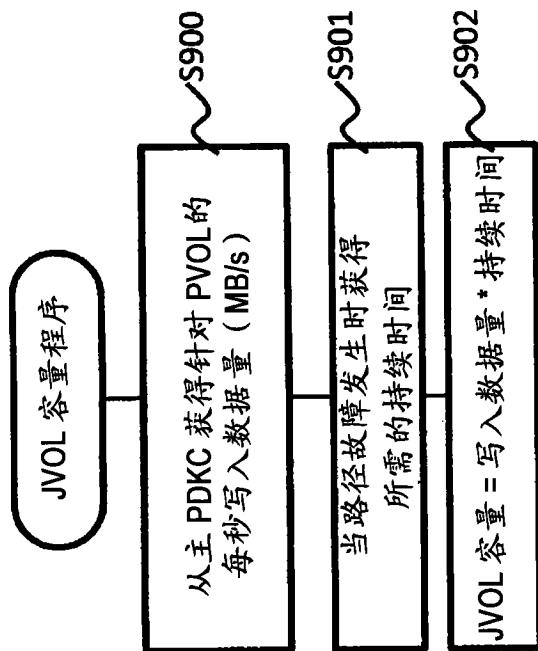


图 20

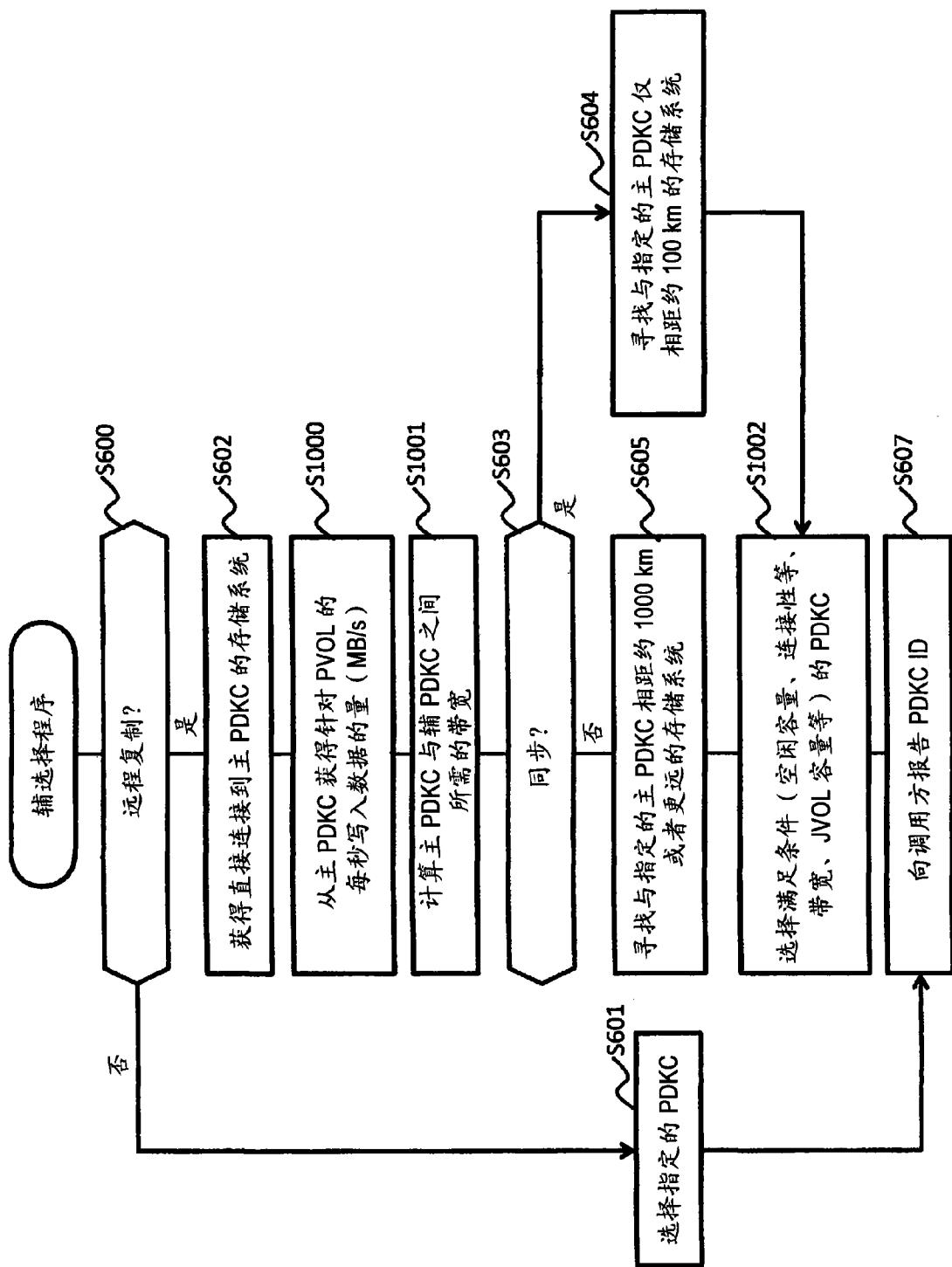


图 21

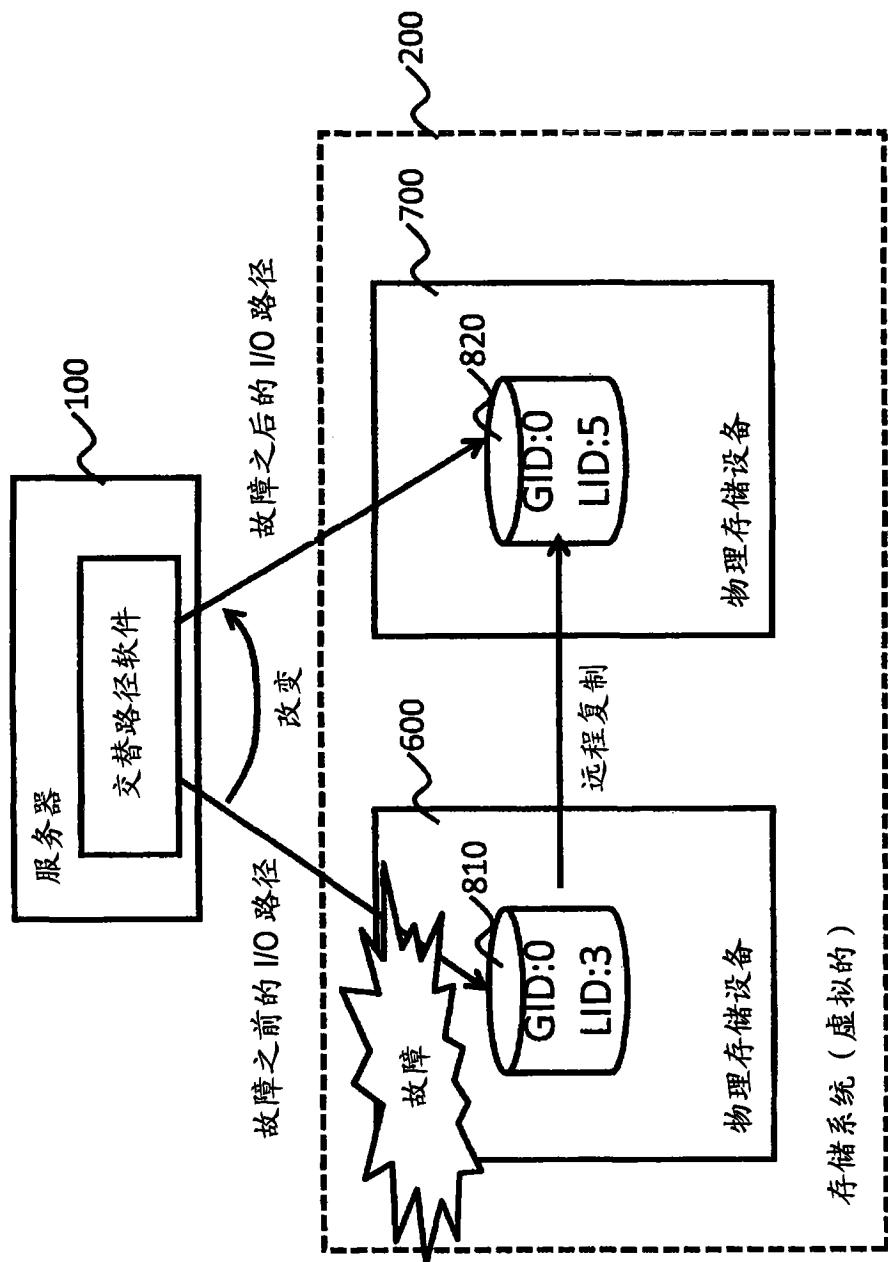


图 22

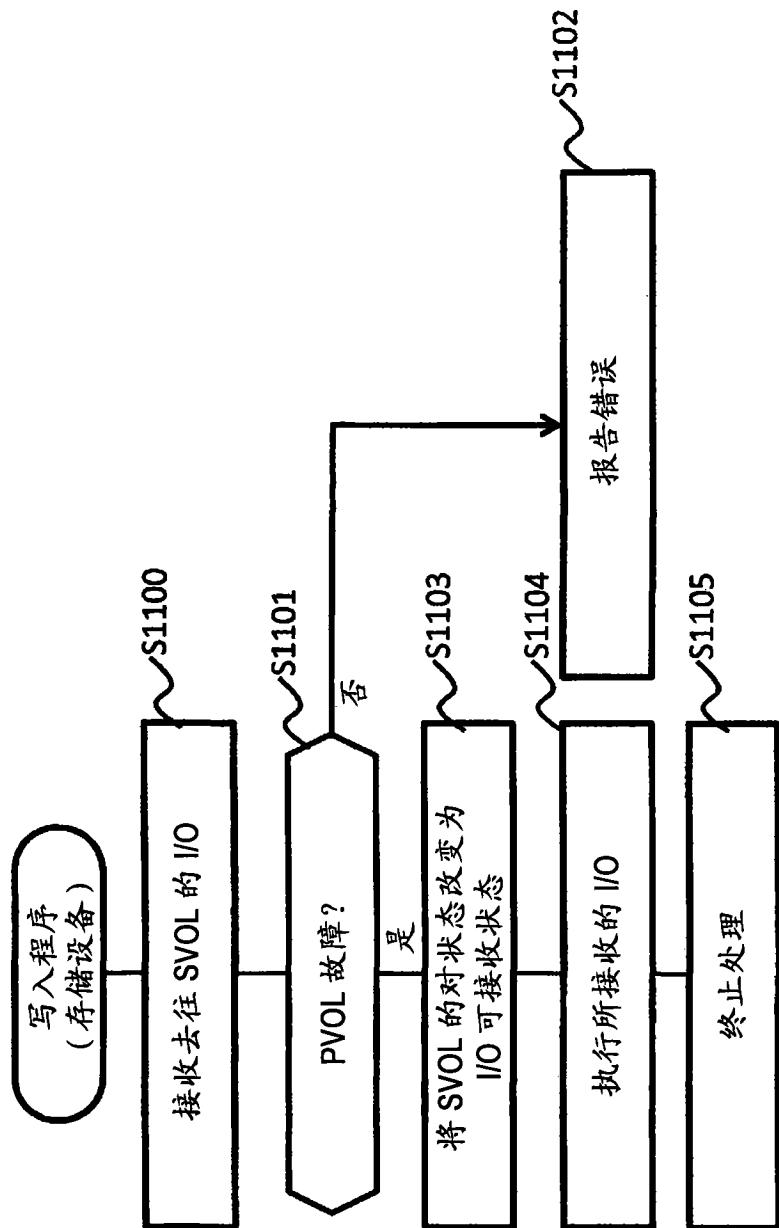


图 23

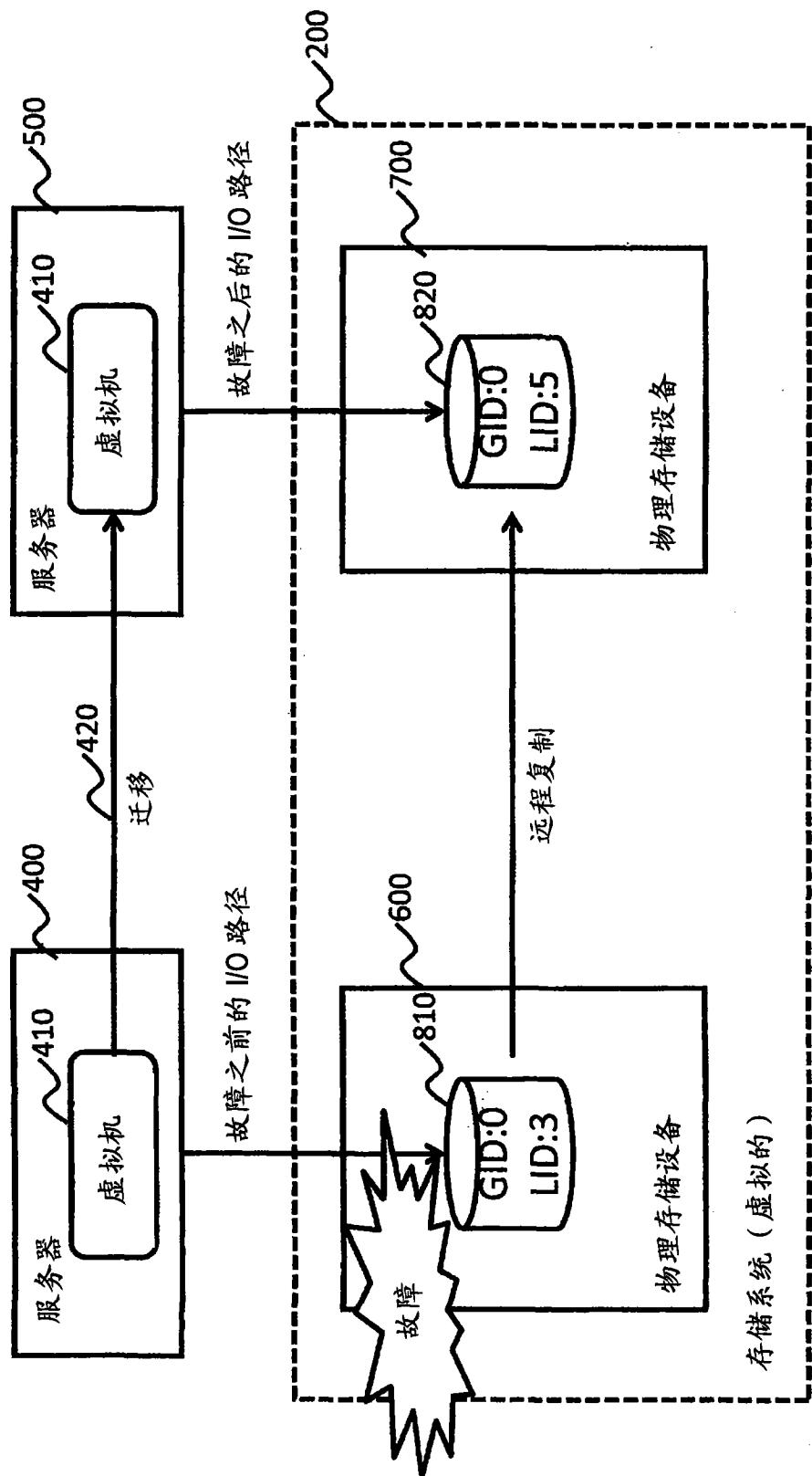


图 24

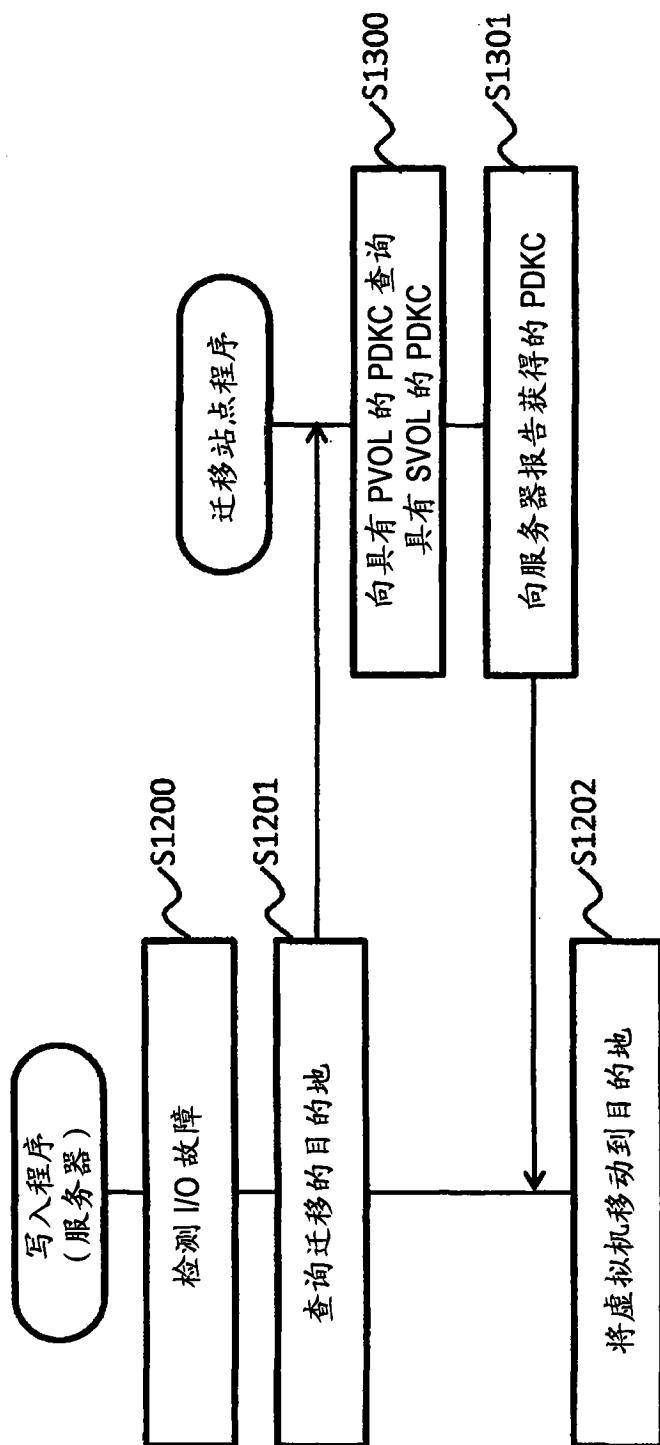


图 25

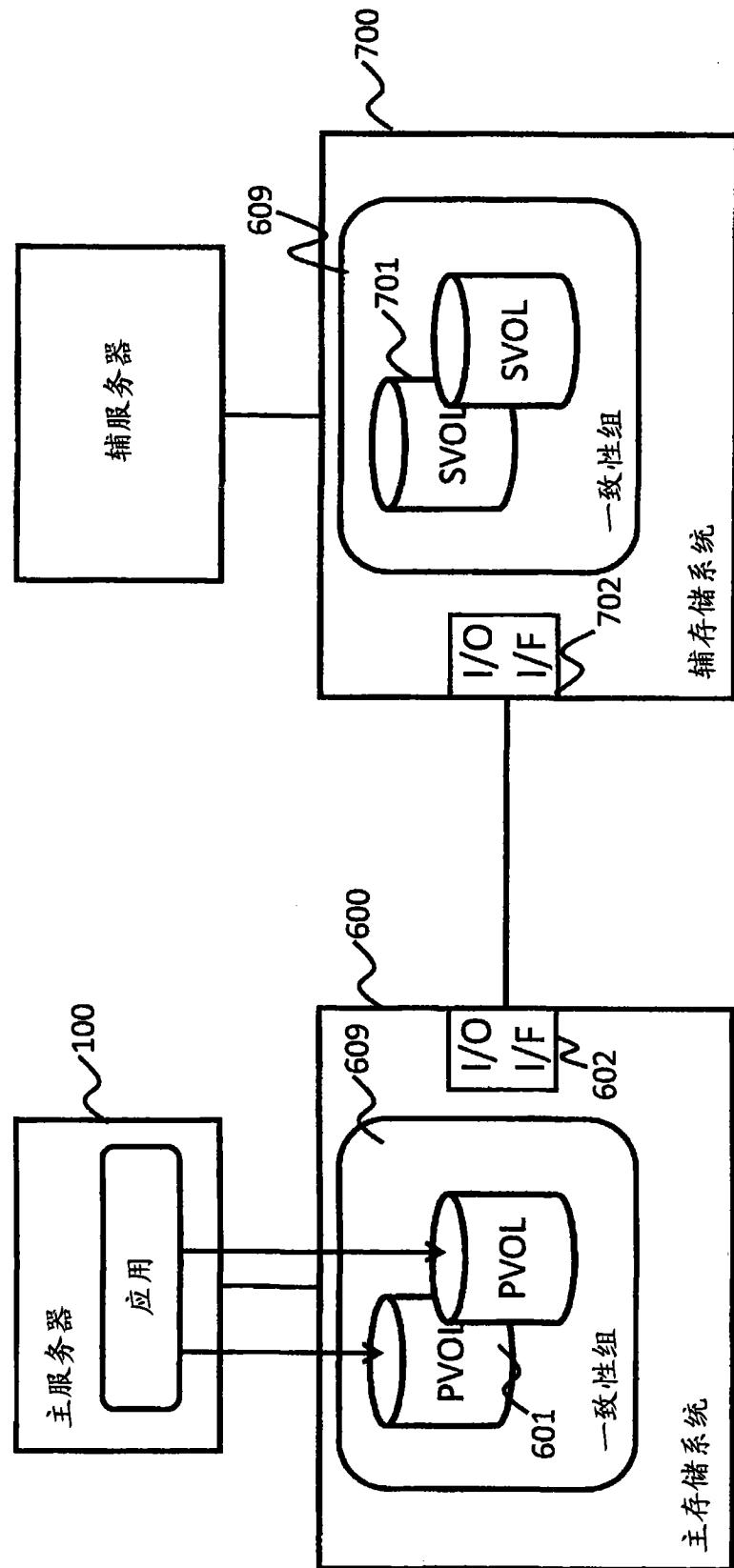


图 26

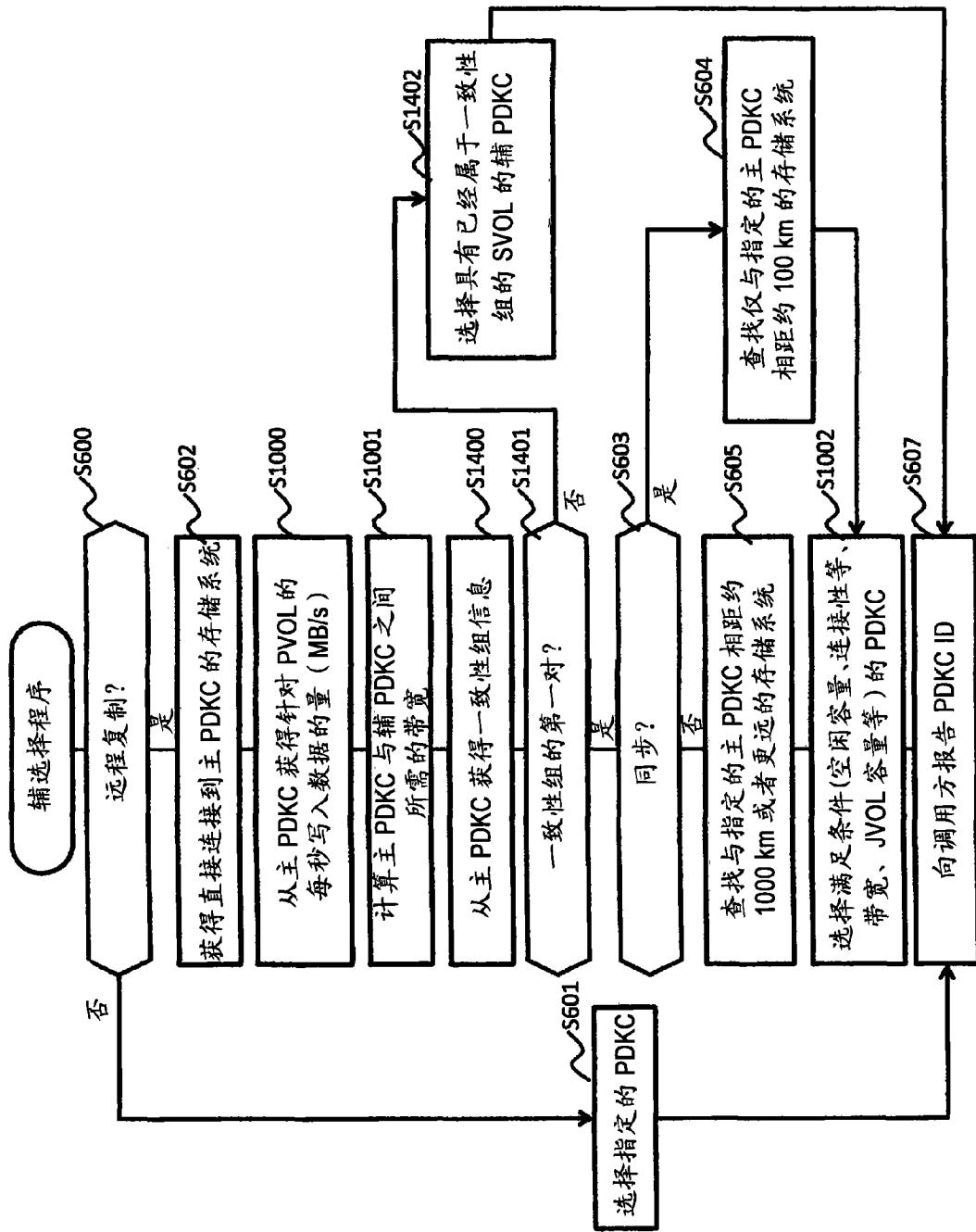


图 27

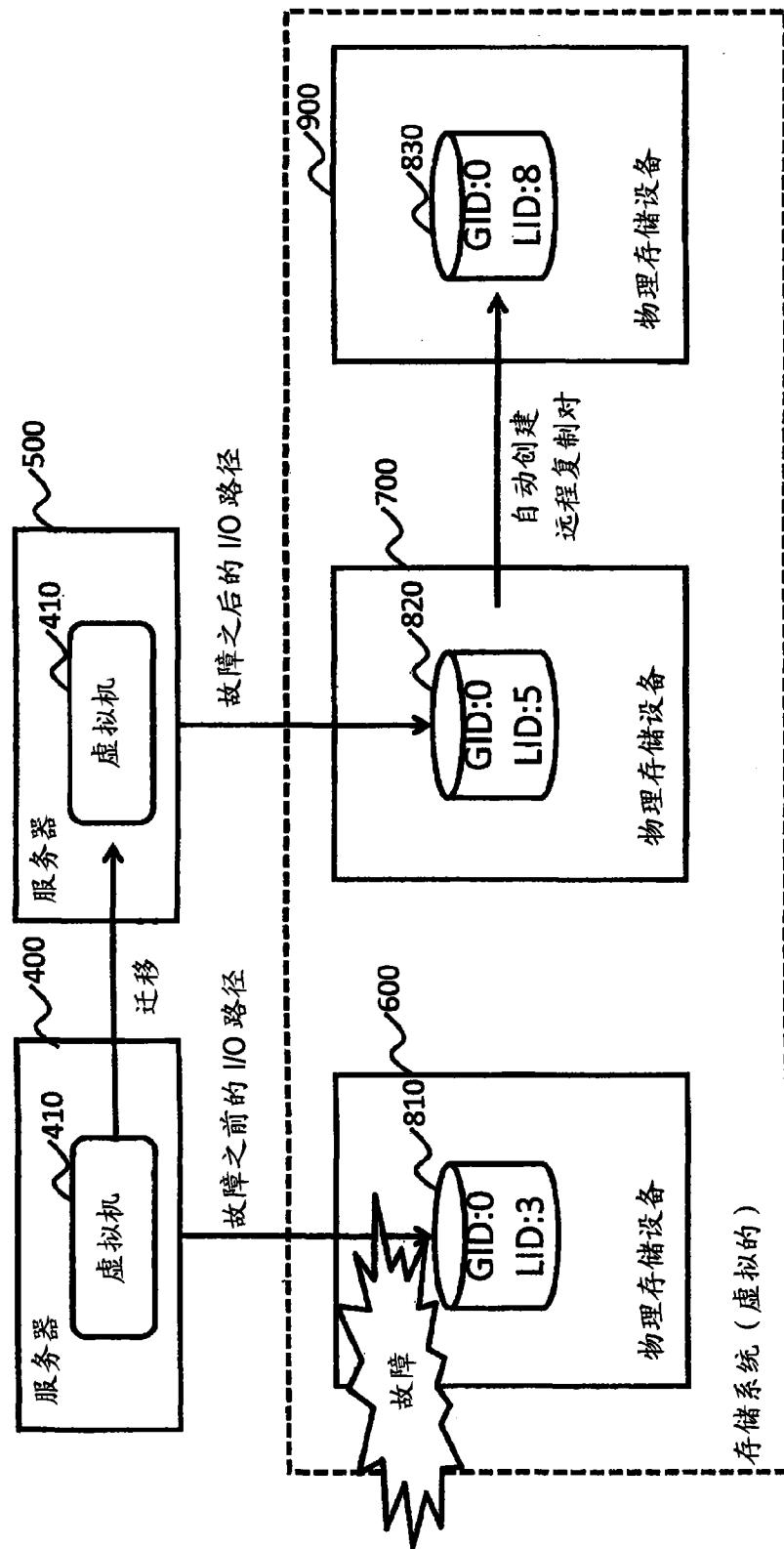


图 28

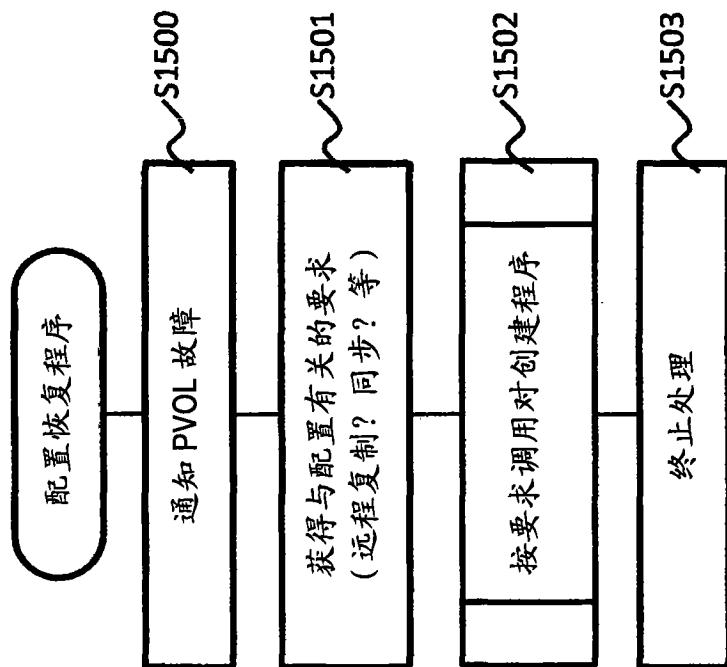


图 29

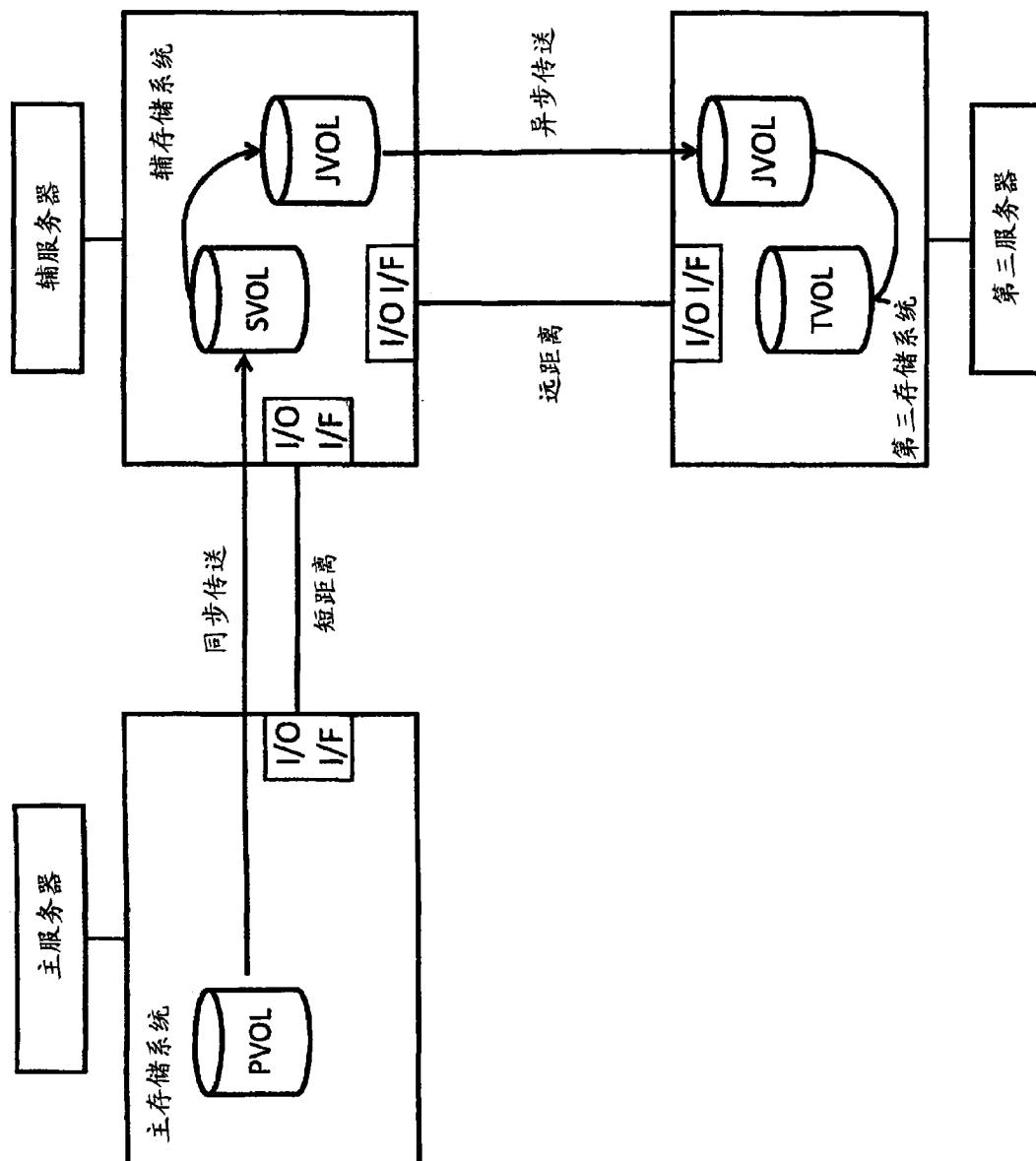


图 30

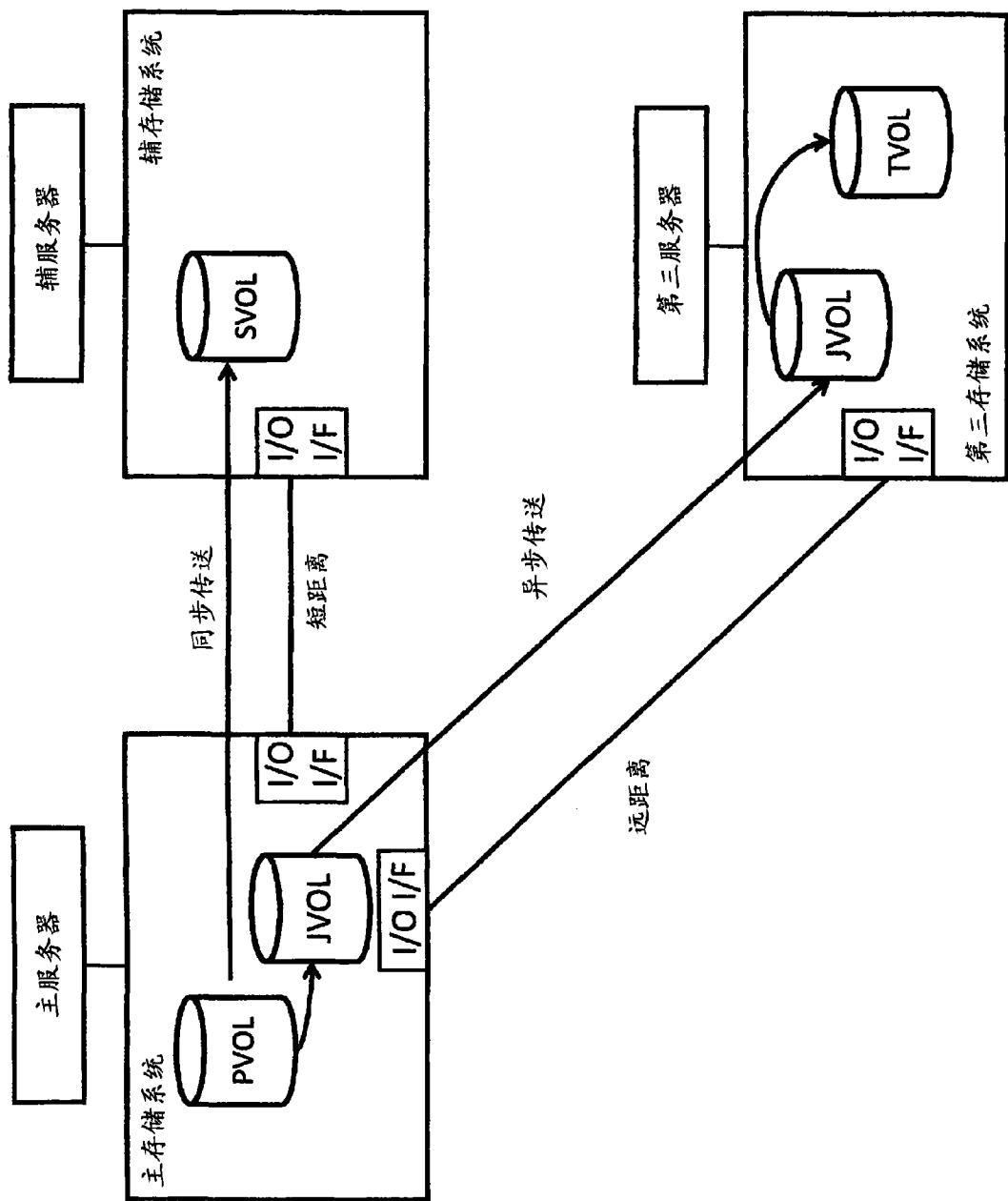


图 31

设置屏幕	ID	HA	RC(同步)	RC(异步)	级联	多目标
1			✓			
2				✓		
3					✓	
4	...	...	✓			
				...	...	...
					...	...
						...

图 32

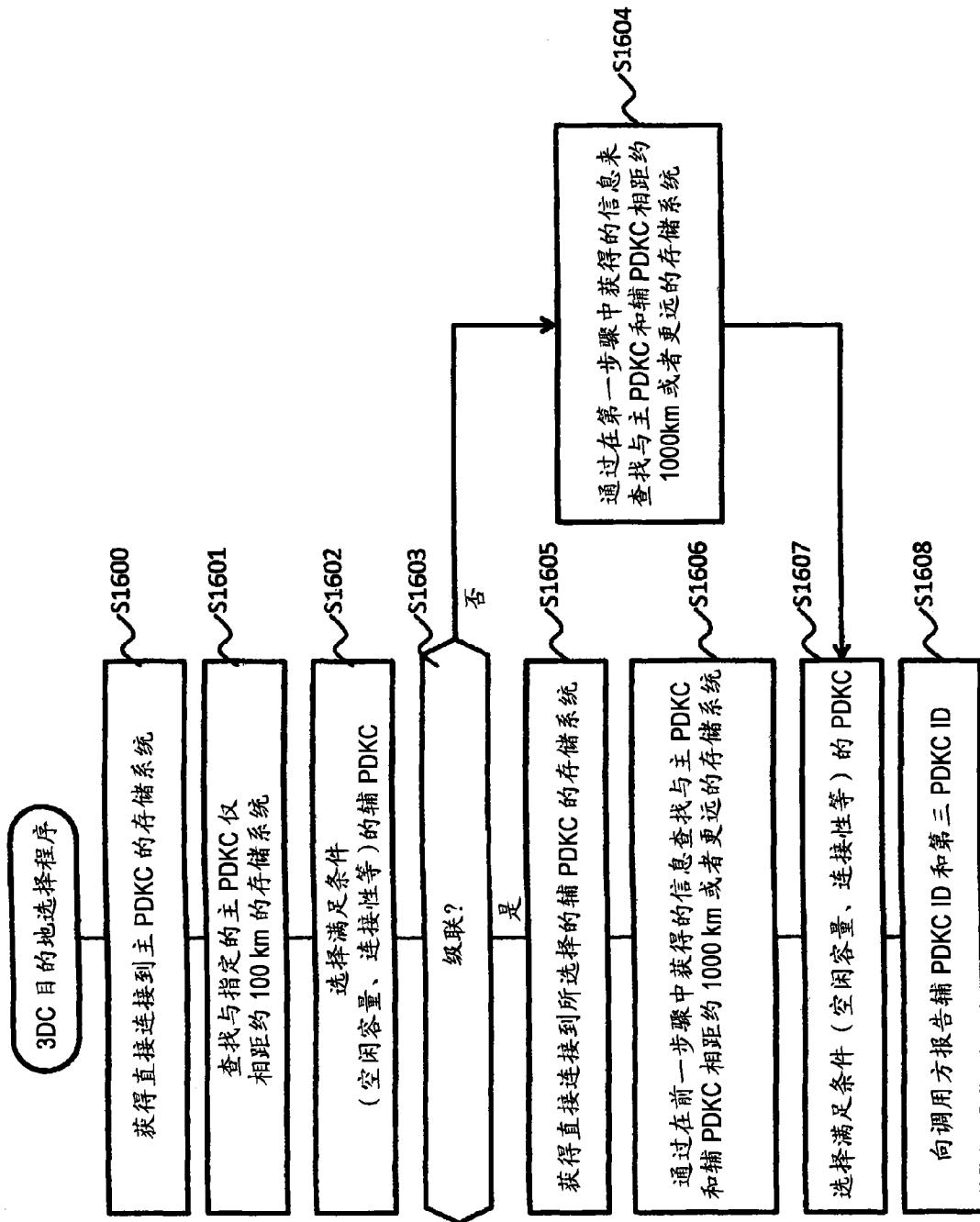


图 33