



(12) 发明专利

(10) 授权公告号 CN 109313593 B

(45) 授权公告日 2022.03.01

(21) 申请号 201680086779.0

泽田翔

(22) 申请日 2016.09.16

(74) 专利代理机构 北京尚诚知识产权代理有限公司 11322

(65) 同一申请的已公布的文献号
申请公布号 CN 109313593 A

代理人 龙淳

(43) 申请公布日 2019.02.05

(51) Int.Cl.

G06F 11/10 (2006.01)

(85) PCT国际申请进入国家阶段日
2018.12.14

(56) 对比文件

(86) PCT国际申请的申请数据
PCT/JP2016/077533 2016.09.16

CN 104246707 A, 2014.12.24

CN 104205059 A, 2014.12.10

(87) PCT国际申请的公布数据
W02018/051505 JA 2018.03.22

WO 2014170984 A1, 2014.10.23

US 2015067439 A1, 2015.03.05

WO 2015145724 A1, 2015.10.01

(73) 专利权人 株式会社日立制作所
地址 日本东京都

审查员 薛聪帆

(72) 发明人 筱塚研太 武田贵彦 黑川勇

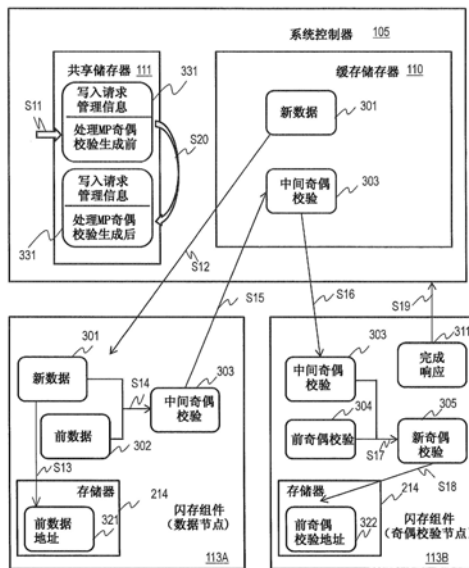
权利要求书3页 说明书11页 附图14页

(54) 发明名称

存储系统

(57) 摘要

控制器从主机接收写入请求,将与写入请求所示的指定地址对应的中间奇偶校验生成命令发送至多个存储设备中的第1存储设备,中间奇偶校验生成命令根据上述指定地址的新数据和被更新为上述新数据的前数据指示生成中间奇偶校验,中间奇偶校验生成命令包含存入有新数据的存储区域的第1地址,和用于存入中间奇偶校验的存储区域的第2地址,第1存储设备接收中间奇偶校验生成命令,从第1地址获取新数据,根据新数据和存储在第1存储设备的前数据生成中间奇偶校验,将中间奇偶校验存入第2地址。



1. 一种存储系统,其特征在于,包括:
具有存储区域的控制器;和
分别包括存储介质的多个存储设备,
所述控制器从主机接收写入请求,并且
将与所述写入请求所示的指定地址对应的中间奇偶校验生成命令发送至所述多个存储设备中的第1存储设备,

所述中间奇偶校验生成命令指示根据所述指定地址的新数据和要被更新为所述新数据的前数据生成中间奇偶校验,

所述中间奇偶校验生成命令包括:存入有所述新数据的所述存储区域的第1地址;和用于存入所述中间奇偶校验的所述存储区域的第2地址,

所述第1存储设备接收所述中间奇偶校验生成命令,并且

从所述第1地址获取所述新数据,

根据所述新数据和存储在所述第1存储设备的所述前数据生成所述中间奇偶校验,

将所述中间奇偶校验存入所述第2地址,

所述控制器包括多个处理器,

所述多个处理器中的第1处理器开始所述写入请求的处理,

所述第1处理器在由所述多个处理器共享的共享存储区域存入表示所述写入请求的处理阶段的写入请求管理信息,

所述多个处理器中的第2处理器响应于所述第1处理器的错误而接管所述写入请求的处理,

参照所述共享存储区域的所述写入请求管理信息,确定所述写入请求的处理的阶段。

2. 如权利要求1所述的存储系统,其特征在于:

所述控制器向所述多个存储设备中的第2存储设备发送新奇偶校验生成命令,

所述新奇偶校验生成命令指示根据所述中间奇偶校验和关于所述指定地址的前数据的前奇偶校验生成新奇偶校验,

所述新奇偶校验生成命令包括所述第2地址,

所述第2存储设备接收所述新奇偶校验生成命令,并且

从所述第2地址读取所述中间奇偶校验,

根据所述中间奇偶校验和存储在所述第2存储设备的所述前奇偶校验生成所述新奇偶校验。

3. 如权利要求2所述的存储系统,其特征在于:

在所述写入请求管理信息表示的是所述新奇偶校验的生成前时,所述第2处理器向所述第1存储设备和所述第2存储设备分别发送重置命令,该重置命令指示中止所述中间奇偶校验生成命令和所述新奇偶校验生成命令的处理。

4. 如权利要求3所述的存储系统,其特征在于:

所述第1存储设备管理:映射信息,其表示所述指定地址与存入有关于所述指定地址的数据的所述第1存储设备的存储介质的地址的对应关系;和存入有所述前数据的所述第1存储设备的存储介质的地址,

在所述第1存储设备的存储介质中,将所述新数据存入不同于所述前数据的地址,

在所述映射信息中,将与所述指定地址对应的地址更新为所述新数据的地址,响应所述重置命令,在所述映射信息中将与所述指定地址对应的地址写回到前数据的地址。

5. 如权利要求2所述的存储系统,其特征在于:

所述第1存储设备在所述第1存储设备的存储介质中将所述新数据存入不同于所述前数据的地址,

所述第2存储设备在所述第2存储设备的存储介质中将所述新奇偶校验存入不同于所述前奇偶校验的地址,

在所述写入请求管理信息表示的是所述新奇偶校验的生成前时,所述第2处理器从所述多个存储设备获取包含所述新数据的条带串的数据,根据所获取的所述条带串的数据生成奇偶校验来存入所述第2存储设备。

6. 一种存储系统中的数据存入的方法,所述存储系统包括具有存储区域的控制器和各自具有存储介质的多个存储设备,该方法的特征在于:

所述控制器从主机接收写入请求,

所述控制器将与所述写入请求所示的指定地址对应的中间奇偶校验生成命令发送至所述多个存储设备中的第1存储设备,

所述中间奇偶校验生成命令指示根据所述指定地址的新数据和要被更新为所述新数据的前数据生成中间奇偶校验,

所述中间奇偶校验生成命令包括:存入有所述新数据的所述存储区域的第1地址;和用于存入所述中间奇偶校验的所述存储区域的第2地址,

所述第1存储设备接收所述中间奇偶校验生成命令,

所述第1存储设备从所述第1地址获取所述新数据,

所述第1存储设备根据所述新数据和存储在所述第1存储设备的所述前数据生成所述中间奇偶校验,

所述第1存储设备将所述中间奇偶校验存入所述第2地址,其中,

所述控制器包括多个处理器,

所述多个处理器中的第1处理器开始所述写入请求的处理,

所述第1处理器将表示所述写入请求的处理阶段的写入请求管理信息存入在由所述多个处理器共享的共享存储区域,

所述多个处理器中的第2处理器响应于所述第1处理器的错误而接管所述写入请求的处理,

所述第2处理器参照所述共享存储区域的所述写入请求管理信息,确定所述写入请求的处理的阶段。

7. 如权利要求6所述的方法,其特征在于:

所述控制器向所述多个存储设备中的第2存储设备发送新奇偶校验生成命令,

所述新奇偶校验生成命令指示根据所述中间奇偶校验和关于所述指定地址的前数据的前奇偶校验生成新奇偶校验,

所述新奇偶校验生成命令包括所述第2地址,

所述第2存储设备接收所述新奇偶校验生成命令,

所述第2存储设备从所述第2地址读取所述中间奇偶校验，

所述第2存储设备根据所述中间奇偶校验和存储在所述第2存储设备的所述前奇偶校验生成所述新奇偶校验。

8. 如权利要求7所述的方法，其特征在于：

在所述写入请求管理信息表示的是所述新奇偶校验的生成前时，所述控制器向所述第1存储设备和所述第2存储设备分别发送重置命令，该重置命令指示中止所述中间奇偶校验生成命令和所述新奇偶校验生成命令的处理。

9. 如权利要求8所述的方法，其特征在于：

所述第1存储设备管理：映射信息，其表示所述指定地址与存入有关于所述指定地址的数据的所述第1存储设备的存储介质的地址的对应关系；和存入有所述前数据的所述第1存储设备的存储介质的地址，

在所述第1存储设备的存储介质中，所述第1存储设备将所述新数据存入不同于所述前数据的地址，

在所述映射信息中，所述第1存储设备将与所述指定地址对应的地址更新为所述新数据的地址，

所述第1存储设备响应所述重置命令，在所述映射信息中将与所述指定地址对应的地址写回到前数据的地址。

10. 如权利要求7所述的方法，其特征在于：

所述第1存储设备在所述第1存储设备的存储介质中将所述新数据存入不同于所述前数据的地址，

所述第2存储设备在所述第2存储设备的存储介质中将所述新奇偶校验存入不同于所述前奇偶校验的地址，

在所述写入请求管理信息表示的是所述新奇偶校验的生成前时，所述第2处理器从所述多个存储设备获取包含所述新数据的条带串的数据，根据所获取的所述条带串的数据生成奇偶校验来存入所述第2存储设备。

存储系统

技术领域

[0001] 本发明涉及存储系统。

背景技术

[0002] 作为本发明的背景技术,例如已知有国际公开第2015/145724号。国际公开第2015/145724号例如在摘要中公开有以下的结构。

[0003] 存储系统在对存储装置的条带内的多个不连续区域存入写入数据时,在将含有确定该多个不连续区域的信息的新数据发送命令和写入数据发送至存储装置后,从存储装置接收根据多个写入数据和该多个写入数据的更新前数据生成的中间奇偶校验,向存入奇偶校验的存储装置发送中间奇偶校验发送命令和中间奇偶校验。之后,当向多个存储装置发送含有确定多个不连续区域的信息的数据确定命令时,存入奇偶校验的存储装置根据所接收的中间奇偶校验和与该中间奇偶校验对应的更新前奇偶校验生成更新后奇偶校验,将更新后奇偶校验存入存储介质。

[0004] 现有技术文献

[0005] 专利文献

[0006] 专利文献1:国际公开第2015/145724号

发明内容

[0007] 发明所要解决的问题

[0008] 在上述存储系统中,储存控制器与多个存储装置例如按SCSI标准连接,在奇偶校验更新时执行以下的步骤。(1)从控制器向存储装置1发出新数据发送命令(新数据从控制器发送至存储装置1)。(2)从存储装置1发出针对新数据发送命令的完成响应(完成命令)。(3)从存储装置1向控制器发出中间奇偶校验发送命令(中间奇偶校验从存储装置1发送至控制器)。

[0009] (4)从控制器发出针对中间奇偶校验发送命令的完成响应(完成命令)。(5)从控制器向存储装置2发出中间奇偶校验发送命令。(6)从存储装置2发出针对中间奇偶校验发送命令的完成响应(完成命令)。(7)从存储装置2发出对控制器的新奇偶校验生成完成响应(完成命令)。

[0010] 上述存储系统在每次伴随来自主机的数据写入的奇偶校验更新时,发出步骤(1)~(7)的命令。命令发出数量多时,成为处理器的负荷。因而,为了抑制奇偶校验更新时的处理器的负荷,提高存储系统的处理性能,期望发出命令数量的削减。

[0011] 用于解决问题的方式

[0012] 本发明的代表性的一个例子为一种存储系统,其包括:具有存储区域的控制器;和分别包括存储介质的多个存储设备,所述控制器从主机接收写入请求,并且将与所述写入请求所示的指定地址对应的中间奇偶校验生成命令发送至所述多个存储设备中的第1存储设备,所述中间奇偶校验生成命令指示根据所述指定地址的新数据和要被更新为所述新数

据的前数据生成中间奇偶校验,所述中间奇偶校验生成命令包括:存入有所述新数据的所述缓存区域的第1地址;和用于存入所述中间奇偶校验的所述缓存区域的第2地址,所述第1存储设备接收所述中间奇偶校验生成命令,并且从所述第1地址获取所述新数据,根据所述新数据和存储在所述第1存储设备的所述前数据生成所述中间奇偶校验,将所述中间奇偶校验存入所述第2地址。

[0013] 发明的效果

[0014] 根据本发明,能够提高奇偶校验生成的效率而抑制处理器的负荷,实现存储系统的性能的提高。

附图说明

[0015] 图1表示计算机系统的结构例。

[0016] 图2表示闪存组件的结构例。

[0017] 图3表示基于来自自主计算机的写入请求的数据更新的正常流程例。

[0018] 图4表示基于来自自主计算机的写入请求的数据更新的正常流程例。

[0019] 图5表示写入请求管理信息的结构例。

[0020] 图6表示中间奇偶校验生成命令的结构例。

[0021] 图7表示新奇偶校验生成命令的结构例。

[0022] 图8表示响应了来自自主计算机的写入请求的数据更新中产生障碍的情况下的流程例。

[0023] 图9A表示重置命令的结构例。

[0024] 图9B表示重置命令的结构例。

[0025] 图10表示响应了来自自主计算机的写入请求数据更新中产生了障碍的情况下的另一流程例。

[0026] 图11表示响应了来自自主计算机的写入请求的数据更新的流程例。

[0027] 图12表示响应了来自自主计算机的写入请求的数据更新中产生了障碍的情况下的另一流程例。

[0028] 图13表示响应了来自自主计算机的写入请求的数据更新的另一正常流程例。

[0029] 图14表示响应了来自自主计算机的写入请求的数据更新的另一正常流程例。

[0030] 图15表示图13和图14所示的正常流程中产生了障碍的情况下的流程例。

具体实施方式

[0031] 实施例1

[0032] 以下,参照附图说明实施例。但是,本实施例仅是用于实现发明的一个例子,并不限定发明的技术范围。此外,对各图中共同的结构标注相同的参照附图标记。

[0033] 另外,在以后的说明中以“表”这样的表达说明本发明的信息,但是这些信息并不是一定以表形成的数据结构来表现,也可以以“列表”、“DB(数据库)”,“队列”等数据结构及其以外的方式表现。因此,为了表示不依赖于数据结构,对于“表”、“列表”、“DB”、“队列”等简称为“信息”。此外,在对各信息的内容进行说明时,能够使用“识别信息”、“识别符”、“名”、“名称”、“ID”这样的表现,并能够对它们彼此替换。

[0034] 在以后的说明中,存在以“程序”为主语进行说明的情况,但是由于程序是一边利用存储器和通信端口(通信控制装置)一边进行由处理器执行而确定的处理,所以也可以采取以处理器为主语的说明,还可以采取以控制器为主语的说明。此外,以程序为主语所公开的处理也可以为管理装置或信息处理装置等的计算机进行的处理。既可以利用专用硬件实现程序的一部分或全部,此外,也可以模块化。各种程序也可以由程序分送伺服器及存储介质安装于各计算机。

[0035] 此外,近年来,在计算机及存储系统中,为了进行大量数据的高速解析和高速I/O处理,需要有大容量的存储区域。例如在计算机中需要存储器DB那样的应用。但是,可装载于装置的DRAM容量受限于成本方面的原因和电气安装的限制。因此,作为缓和措施,出现了利用比DRAM慢但比HDD高速的NAND闪存存储器等的半导体存储介质的动向。

[0036] 这些半导体存储介质以SSD(Solide State Drive:固态硬盘)的名称来称呼,通过SATA(Seiral ATA:串行ATA)、SAS(Serial Attached SCSI:串行连接SCSI)等磁盘I/O接口连接及其协议与计算机和储存控制器连接,来加以利用。

[0037] 但是,对计算机的性能提高而言,经由这些磁盘I/O接口和协议的访问的开销大而且延迟较长。因此,近些年来出现了能够装载在能够与处理器直接连接的作为通用总线的PCI-Express(PCIe)(接口)上,且为了充分利用其高速性而使用新制定的NVMe(Non Volatile Memory Express:非易失性储存器接口规范)协议、能够以低延迟进行访问的PCIe连接型SSD(PCIe-SSD或PCIe-Flash)。

[0038] 本实施方式例如使用在奇偶校验生成中能够从存储设备直接访问控制器上的存储器的NVMe命令。在NVMe中,用于数据发送和接收的支持的I/O命令非常简洁,支持必须命令仅有Write(写)、Read(读)、Flush(闪存)这3个。此外,在SAS等现有磁盘I/O协议中,主机为主体,将命令和数据发送至装置侧,与此相对,在NVMe中,仅从主机将已制作了命令的意思通知给装置,命令本身的获取和数据的转送以装置侧为主体来实施。

[0039] 即变更为从装置侧的行动而实施。例如,在装置所获取的命令的内容为Write(写入)的情况下,在现有方式中主机向装置发送Write数据(写入数据),而在NVMe中,通过装置Read(读取)主机的数据的动作来实现。相反,在命令的内容为Read的情况下,Read命令的处理通过装置向主机的存储器Write数据的动作来实现。

[0040] 即,在NVMe中,装置侧将动作契机与命令的接收、数据的读取·写入转送一起掌握,因此不需要用于随时都能够受理来自主机的请求的、额外的资源确保。

[0041] 由此,在本实施例中,通过在发出命令中含有控制器的地址信息、储存器件的地址信息来提高奇偶校验生成的效率,抑制处理器的负荷,由此提高存储系统的处理性能。

[0042] 图1表示实施方式的计算机系统100的结构例。计算机系统100包括主计算机101、管理装置102、存储系统104。主计算机101、管理装置102、存储系统104通过网络103相互连接。网络103例如为Storage Area Network(SAN(存储区域网络))。管理装置102也可以通过不同于网络103的管理网络与其它装置连接。

[0043] 主计算机101是执行应用程序的计算机,经由网络103访问存储系统104的逻辑存储区域。主计算机101例如包括输入器件、输出器件、CPU(Central Processing Unit:中央处理器)、存储器、磁盘适配器、网络适配器和二次存储设备(未图示)。

[0044] 主计算机101执行用户使用的应用程序、进行与存储系统104的接口控制的存储系

统控制程序。主计算机101使用存储系统104提供的的数据卷。主计算机101通过对所提供的的数据卷发出读取命令和写入请求,来访问存入在数据卷中的数据。

[0045] 管理装置102管理存储系统104,例如进行存储系统104的存储区域的构成。管理装置102执行用于进行存储系统104的的管理的管理程序。管理装置102与通用的计算机同样,具有键盘和显示器等输入输出器件、CPU、存储器、网络适配器和二次存储设备。

[0046] 存储系统104包括系统控制器105和多个闪存组件113。存储系统104将数据存入闪存组件113的存储区域。存储系统104对主计算机101提供1个以上的数据卷。

[0047] 系统控制器105包括主机接口(I/F)106、维护I/F、驱动器I/F108、MP(Micro Processor:微处理器)109、存储器单元119、多个闪存组件113。这些构成要素通过总线112相互连接。

[0048] 主机I/F106是用于存储系统104与主计算机101的启动器(initiator)进行通信的接口器件。主计算机101为了访问数据卷而发出的请求(读取请求、写入请求等)送到主机I/F106。存储系统104从主机I/F106向主计算机101返送信息(例如读取数据)。

[0049] 维护I/F107是用于存储系统104与管理装置102进行通信的接口器件。来自管理装置102的命令送到维护I/F107。存储系统104从维护I/F107向管理装置102反馈信息。

[0050] 图1表示主机I/F106和维护I/F107均与网络103连接的结构,但是主机I/F106连接的网络和维护I/F107连接的网络也可以为不同的网络。

[0051] MP109在系统控制器105内设置有1个或多个,执行用于进行存储系统104的各种控制的程序。存储器单元119由缓存存储器110、共享存储器111、本地存储器118构成。缓存存储器110、共享存储器111、本地存储器118既可以物理上在一个存储器内将区域分割使用,也可以分别由物理上独立的存储器构成。

[0052] 缓存存储器110提供缓存区域。缓存存储器110例如由RAM(Random Access Memory:随机存取存储器)构成,临时地保存在闪存组件113被读写的数据。共享存储器111由硬盘和闪存存储器、RAM等构成,存入在储存控制器进行工作的程序和构成信息等。本地存储器118用于临时存入MP109执行的程序和MP109使用的信息。

[0053] 缓存存储器110作为临时存储区域使用,其临时存储对数据卷(存储设备)的写入数据或从数据卷(存储设备)读取的数据(读取数据)。

[0054] 具体而言,例如,在存储在缓存存储器110的读取数据被反馈主计算机101之后,也可以不消除该读取数据而存储在缓存存储器110不动。

[0055] 此外,例如在读取数据被返送给主计算机101之后,也可以从缓存存储器110删除该读取数据。在缓存存储器110使用DRAM(Dynamic Random Access Memory:动态随机存取存储器)、SRAM(Static Random Access Memory:静态随机存取存储器)等易失性存储介质,但也可以在缓存存储器110使用非易失性存储器。

[0056] 共享存储器111提供用于存入MP109使用的管理信息的共享存储区域。与缓存存储器110同样地在共享存储器111使用DRAM、SRAM等易失性存储介质,但也可以使用非易失性存储介质。

[0057] 闪存组件113是具有用于存储来自主计算机101的写入数据的非易失性存储介质的存储设备。闪存组件113作为存储介质使用闪存存储器等非易失性半导体存储器或其它存储介质。闪存组件113的一个例子为SSD。

[0058] 为了实现高可靠性,多个闪存组件113构成RAID (Redundant Array of Independent Disks:独立磁盘冗余阵列)组115。MP109具有即使一台闪存组件113发生故障也能够将该闪存组件113的数据恢复的RAID功能。

[0059] 在RAID组115中制作一个以上逻辑卷。一个逻辑卷与构成RAID组115的闪存组件113具有的物理的存储区域相关联。

[0060] RAID功能将从主计算机101接收的写入数据和冗余数据的集合分散地存入于RAID组115的闪存组件113。已知多个RAID级别。例如,RAID1将写入数据及其复制数据存入在不同的闪存组件113中。

[0061] 除此以外,RAID5将由写入数据和一个奇偶校验构成的冗余数据集合分散地存入于不同的闪存组件113,RAID5将由写入数据和二个奇偶校验构成的冗余数据分散地存入于不同的闪存组件113。

[0062] 另外,在本实施例中例示了主计算机101通过网络与系统控制器12连接的例子,但存储系统42的硬件结构也可以是与服务器同样的结构。例如也可以代替本实施例的存储系统104,在个人计算机等通用的计算机(以下,将其简称为“计算机”)装载(或连接)多个闪存组件113,在计算机上执行各种程序。在这种情况下,计算机从服务器受理I/O请求,进行数据向闪存组件113的存入或进行从闪存组件113的数据读取。

[0063] 此外,在计算机上执行各种程序的结构的情况下,也可以以在存储系统上执行的各种程序和由服务器执行的程序均在同一计算机上执行的方式构成。在这种情况下,例如通过在计算机上执行形成虚拟机的管理程序(Hypervisor、系统管理程序),在计算机上至少形成执行由服务器执行的程序的服务器虚拟机、和在存储系统上执行各种程序的储存虚拟机即可,读取请求或写入请求从服务器虚拟机发出。由此,以下的主计算机101也可以与“服务器虚拟机”相替换。

[0064] 图2表示闪存组件113的结构例。闪存组件113具有器件控制器210和闪存存储器280,该闪存存储器280是用于存储来自主计算机101的写入数据的存储介质。器件控制器210具有驱动器I/F211、器件处理器213、存储器214、闪存I/F215、逻辑电路216,它们通过内部网络212相互连接。逻辑电路216例如进行奇偶校验运算、密码处理和压缩处理。

[0065] 驱动器I/F211是用于与系统控制器105进行通信的接口器件。闪存I/F215是器件控制器210用于与闪存存储器280进行通信的接口器件。

[0066] 器件处理器213执行用于进行闪存组件113的控制的程序。在本例中,以下说明的闪存组件113进行的处理通过器件处理器213执行程序来进行。

[0067] 存储器214是用于保存器件处理器213执行的程序和器件处理器213使用的控制信息等的存储器。还包括临时保存在闪存存储器280被读写的数据的临时存入区域。

[0068] 器件处理器213从MP109受理请求,执行遵从所受理的命令的处理。例如,器件处理器213从MP109接收写入命令,将该写入命令的数据写入闪存存储器280。器件处理器213在将数据写入存储器214内的临时存入区域或闪存存储器280后,将写入命令的完成响应(响应命令)返送至MP109。

[0069] 器件处理器213从MP109接收读取命令,将该读取命令的数据从存储器214内的临时存入区域或闪存存储器280读取,返送至MP组件109。从MP109发送至闪存组件113的命令和闪存组件113对命令的处理的详细情况后述。

[0070] 存储系统104例如如日志构造文件系统那样在作为存储介质的闪存组件113内的闪存存储器280追加新数据。存储系统104当由来自主计算机101的请求所示的逻辑地址确定的数据有了更新请求时,不将更新数据覆盖写入前数据,而在闪存存储器280中、在与前数据不同的存储区域进行追加。闪存存储器280保持来自主计算机101的请求所示的逻辑地址的过去数据和当前数据。

[0071] 此外,存储系统104例如在闪存组件113的存储器214具有映射信息。映射信息是指将逻辑地址与物理地址对应地进行管理的信息。逻辑地址是指用于主计算机101进行访问的地址空间的地址,是对于数据卷的I/O地址(数据卷标识符与数据卷内地址)。物理地址是指在闪存存储器280内实际存入数据的地址。

[0072] 在从主计算机101收到数据的更新请求的情况下,存储系统104通过将关联从存入有与更新请求所示的逻辑地址关联的前数据的物理地址变更至存入有更新数据的物理地址,来进行映射信息的更新。映射信息既可以仅保存更新数据的物理地址,也可以保存更新数据的物理地址和前数据的物理地址的两者。

[0073] 在映射信息仅保存更新数据的物理地址的情况下,存储系统104也可以在将更新数据存入闪存存储器280时解除前数据的物理地址与逻辑地址的关联,即,删除前数据的物理地址而使更新数据的物理地址与逻辑地址关联,由此进行映射信息的更新。

[0074] 在映射信息保存更新数据和前数据的物理地址两者的情况下,存储系统104也可以在将更新数据存入闪存存储器280之后的规定的期间,使更新数据的物理地址和前数据的物理地址两者与逻辑地址关联。并且,在经过规定的期间后,存储系统104也可以将逻辑地址与前数据的物理地址的关联解除,即,将前数据的物理地址删除。关联解除的契机也可以并不限于时间的经过,而是在接收规定的命令或规定的完成响应之后。

[0075] 通过数据的追加,提高写入速度,特别是提高使用NAND闪存的闪存组件113的写入速度。闪存组件113或MP109进行用于追加的地址管理。

[0076] 图3和图4表示与来自主计算机101的写入请求相应的数据更新的正常流程例。系统控制器(图中为CTL)105从主计算机101接收写入请求和新数据301。在本实施例例中,预先确定与写入请求的地址相应地处理该写入请求的MP109。在本例中,MP__A109A为预先确定的处理MP。

[0077] MP__A109A在写入请求管理信息331登记关于所接收到的写入请求的信息(S11)。图5表示写入请求管理信息331的结构例。写入请求管理信息331被存入共享存储器111,表示写入请求、处理该写入请求的MP和该写入请求的处理的阶段。具体而言,写入请求管理信息331具有请求ID、处理MP、阶段和逻辑地址的栏。

[0078] 请求ID栏表示写入请求的标识符。对新的写入请求赋予唯一的值。处理MP栏表示处理该写入请求的MP的标识符。阶段栏表示该写入请求的处理的阶段。具体而言,表示“奇偶校验生成前”或“奇偶校验生成后”的任一阶段。在接收了新的写入请求的情况下,按“奇偶校验生成前”登记。逻辑地址栏表示作为写入请求所示的写入目的地的逻辑地址。

[0079] 来自主计算机101的写入请求的数据(新数据)存入在缓存存储器110。MP__A109A根据从主计算机101接收的写入请求生成中间奇偶校验生成命令。MP__A109A将中间奇偶校验生成命令发送至存入前数据302的闪存组件(数据节点)113A。

[0080] 中间奇偶校验生成命令指示存入新数据301,并且指示根据新数据301和前数据

302生成中间奇偶校验303。图6表示中间奇偶校验生成命令341的结构例。中间奇偶校验生成命令341具有OPECODE字段、地址信息1字段、地址信息2字段、命令详情字段和地址信息3字段。

[0081] OPECODE字段表示操作的种类,在本例中表示来自自主计算机101的写入请求的处理。地址信息1字段表示缓存存储器110中存入有新数据301的地址。地址信息2字段表示缓存存储器110中用于存入有新生成的中间奇偶校验303的地址。

[0082] 命令详情字段表示本命令是中间奇偶校验生成命令。地址信息3字段表示新数据301的逻辑地址。该逻辑地址也是前数据302的逻辑地址。

[0083] MP_A109A参照映射信息,从写入请求所示的逻辑地址鉴别存入有写入请求的前数据302的数据节点113A中的物理地址。缓存存储器110中用于存入新生成的中间奇偶校验的地址表示缓存存储器110的空白区域。

[0084] 数据节点113A参照所接收的中间奇偶校验生成命令341,从缓存存储器110读取存入在由地址信息1指定的地址的新数据301,存入在存储器214的临时存入区域(S12)。数据节点113A将新数据301存入闪存存储器280的空白区域,将存入在存储器214的映射信息更新,即,使逻辑地址与存入了新数据301的物理地址相关联。

[0085] 另外,将新数据301存入在闪存存储器280的处理和映射信息的更新也可以不在该S12的时间进行。例如,也可以不与中间奇偶校验生成命令341的接收同步,而按规定的周期地进行,还可以在对器件处理器213的处理性能而言有余地的时间进行。

[0086] 此外,本实施例不以中间奇偶校验生成命令的接收为契机、即与中间奇偶校验生成命令的接收不同步地开始将新数据301存入闪存存储器280的处理和生成中间奇偶校验的处理。器件处理器213接收中间奇偶校验生成命令,在存储器214的临时存入区域写入数据后,将写入命令的完成响应(响应命令)反馈到MP109。

[0087] 另外,也可以以中间奇偶校验生成命令的接收为契机,即,与中间奇偶校验生成命令的接收同步地开始将新数据301存入闪存存储器280的处理和生成中间奇偶校验的处理。在这种情况下,器件处理器213在完成了新奇偶校验生成后向MP109反馈写入命令的完成响应(响应命令)。

[0088] 数据节点113A伴随映射信息的更新将与写入请求所示的新数据的逻辑地址对应的前数据302的物理地址、即前数据地址321存入存储器214(S13)。前数据地址321以可知与写入请求所示的逻辑地址的对应的方式存储即可。例如,前数据地址321既可以包含于映射信息,也可以以与映射信息不同的形式被保存。另外,存入存储器214的前数据地址321既可以与关于该写入请求的处理非同步地被消除,也可以从被存入起经过规定的时间后被消除。前数据地址321用于在产生障碍时参照。详细情况后述。

[0089] 如上所述,本例的闪存组件113不在前数据上覆盖写入新数据,而在不同的物理地址追加新数据。

[0090] 数据节点113A参照映射信息,从地址信息3指定的逻辑地址鉴别存入前数据302的数据节点113A的闪存存储器280的物理地址。数据节点113A读取前数据302,存入在存储器214。数据节点113A根据存入在存储器214的新数据301和前数据302生成中间奇偶校验303(S14)。

[0091] 数据节点113A参照中间奇偶校验生成命令341,在由地址信息2指定的缓存存储器

110的地址存入中间奇偶校验303 (S15)。

[0092] 如上所述,通过如本实施例那样,通过在缓存存储器110中存入新数据301的地址、在缓存存储器110中用于存入新生成的中间奇偶校验303的地址也包含于中间奇偶校验生成命令341,由此,由1个中间奇偶校验生成命令的发出实现新数据301从系统控制器105至闪存组件113的读取,以及中间奇偶校验303的生成和从闪存组件113至系统控制器105的写入。由此,能够通过少的命令和步骤有效地生成中间奇偶校验303。

[0093] 接着,MP__A109A根据从存储主计算机101接收到的写入请求生成新奇偶校验生成命令。MP__A109A将新奇偶校验生成命令发送至存入前奇偶校验304的闪存组件(奇偶校验节点)113B。本例展示一个奇偶校验节点,但是在存在多个奇偶校验节点的情况下,处理MP109将新奇偶校验生成命令分别发送至存入前奇偶校验的奇偶校验节点。

[0094] 中间奇偶校验生成命令根据中间奇偶校验303和前奇偶校验304指示生成新奇偶校验305。图7表示新奇偶校验生成命令342的结构例。新奇偶校验生成命令342具有OPECODE字段、地址信息1字段、命令详情字段和地址信息3字段。

[0095] OPECODE字段表示操作的种类,在本例中,表示来自主计算机101的写入请求的处理。地址信息1字段表示在缓存存储器110中存入中间奇偶校验303的地址。命令详情字段表示本命令为新奇偶校验生成命令。

[0096] 地址信息3字段表示关于有了写入请求的数据的新奇偶校验的逻辑地址。该逻辑地址也是前奇偶校验304的逻辑地址。

[0097] MP__A109A参照映射信息,从写入请求所示的逻辑地址鉴别存入有写入请求的前奇偶校验304的奇偶校验节点113B的物理地址。

[0098] 奇偶校验节点113B参照所接收的新奇偶校验生成命令342,从缓存存储器110读取存入在由地址信息1所指定的地址的中间奇偶校验303,并存入存储器214 (S16)。

[0099] 奇偶校验节点113B参照映射信息,从地址信息3指定的逻辑地址鉴别存入前奇偶校验304的奇偶校验节点113B的闪存存储器280中的物理地址。奇偶校验节点113B读取前奇偶校验304,并存入存储器214。

[0100] 奇偶校验节点113B根据存入存储器214的中间奇偶校验303和前奇偶校验304生成新奇偶校验305 (S17)。奇偶校验节点113B将新奇偶校验305存入闪存存储器280的空白区域,更新存入于存储器214的映射信息。

[0101] 奇偶校验节点113B伴随映射信息的更新将与该逻辑地址对应的前奇偶校验地址322存入存储器214 (S18)。前奇偶校验地址322在障碍发生时参照。详细情况后述。

[0102] 最后,奇偶校验节点113B将完成响应311反馈至系统控制器105 (S19)。

[0103] MP__A109A在接收到完成响应311时更新写入请求管理信息331 (S20)。具体而言,MP__A109A在已经处理了的写入请求的记录将阶段栏的“奇偶校验生成前”变更为“奇偶校验生成后”。

[0104] 如上所述,通过新奇偶校验生成命令342能够实现从系统控制器105至闪存组件113的中间奇偶校验303的读取以及中间新奇偶校验305的生成。由此,能够以少的命令和步骤高效地生成新奇偶校验305。

[0105] 上述例中,闪存组件113管理追加的数据的地址。在其它例中,也可以由系统控制器105 (MP109)管理追加的数据的地址。系统控制器105管理前数据地址和前奇偶校验地址。

[0106] 由系统控制器105发出的中间奇偶校验生成命令包含在缓存存储器110中存入新数据的地址和存入中间奇偶校验的地址。并且,包含闪存组件113中的新数据的存入地址和前数据的存入地址。闪存组件113从闪存存储器280的所指定的地址读取前数据,并且将新数据存入指定的地址。

[0107] 由系统控制器105发出的新奇偶校验生成命令包含在缓存存储器110中存入中间奇偶校验的地址、以及闪存组件113中的新奇偶校验的存入地址和前奇偶校验的存入地址。闪存组件113从闪存存储器280的所指定的地址读取前奇偶校验,并且将新奇偶校验存入指定的地址。

[0108] 图8表示与来自自主计算机101的写入请求相应的数据更新中、产生障碍的情况下的流程例。在本例中,在MP__A109A中产生障碍。以下主要说明与图4所示的流程例的差异。

[0109] 在图8中,在紧接着新数据301从系统控制器105的缓存存储器110被转送至数据节点113A(S12)之后,便在预先设定的处理MP__A109A中产生了障碍。处理MP从MP__A109A移交给MP__B109B。

[0110] MP__B109B在中断图4所示的流程后,从新数据301的转送(S12)起再执行。具体而言,MP__B109B参照写入请求管理信息331,获取从MP__A109A移交来的处理的信息(S31)。具体而言,处理MP获取表示MP__A109A的记录。在本例中,写入请求管理信息331表示从MP__A109A移交来的处理的阶段为“奇偶校验生成前”。

[0111] MP__B109B将本处理的中断通知给数据节点113A和奇偶校验节点113B。具体而言,MP__B109B向数据节点113A和奇偶校验节点113B分别发送重置命令(S31)。

[0112] 图9A和图9B表示重置命令的结构例。图9A表示发送至数据节点的重置命令的结构例351,图9A表示发送至奇偶校验节点的重置命令的结构例352。重置命令表示地址信息字段和命令详情字段。

[0113] 发送至数据节点的重置命令351中,地址信息字段与中间奇偶校验生成命令341中的地址信息3同样,表示在数据节点中存入新数据的逻辑地址。命令详情表示中间奇偶校验生成的重置。

[0114] 在发送至奇偶校验节点的重置命令352中,地址信息字段与新奇偶校验生成命令342中的地址信息3同样,表示存入新奇偶校验的逻辑地址。命令详情表示新奇偶校验生成的重置。

[0115] 接收到重置命令的数据节点113A和奇偶校验节点113B将中间奇偶校验生成处理和新奇偶校验生成处理中断。在图8的例子中,数据节点113A不开始中间奇偶校验生成处理,奇偶校验节点113B不开始新奇偶校验生成处理。数据节点113A和奇偶校验节点113B分别将对重置命令的完成响应反馈至MP__B109B(S32和S33)。

[0116] 当从数据节点113A和奇偶校验节点113B接收到完成响应时,MP__B109B更新写入请求管理信息331(S34)。具体而言,MP__B109B在写入请求管理信息331的该记录中将处理MP从MP__A变更为MP__B。流程的以后的步骤除了处理MP从MP__A109A变为MP__B109B这一点以外与图4相同。

[0117] 图10表示与来自自主计算机101的写入请求相应的数据更新中,产生了障碍的情况下的其它流程例。在本例中,在MP__A109A中产生障碍。以下主要说明与图8所示的流程例的差异。

[0118] 在图10中,在数据节点113A生成中间奇偶校验(S14)后,转送至系统控制器105之前,在预先设定的处理MP__A109A中产生了障碍。处理MP从MP__A109A移交给MP__B109B。

[0119] 数据节点113A和奇偶校验节点113B从MP__B109B接收重置命令(S31)。数据节点113A关于映射信息将与重置命令351的逻辑地址对应的物理地址反馈至在S13中存入在存储器214的前数据地址321,然后反馈完成响应(S41)。由此,能够在之后的步骤中生成正确的中间奇偶校验。

[0120] 图11表示与来自自主计算机101的写入请求相应的数据更新的流程例。在数据节点113A将中间奇偶校验发送至系统控制器105(S15)后,系统控制器105将中间奇偶校验转送至奇偶校验节点113B之前,在预先设定的处理MP__A109A中产生了障碍。处理MP从MP__A109A移交给MP__B109B。其它方面与图10的流程一样。

[0121] 图12表示与来自自主计算机101的写入请求相应的数据更新中产生了障碍的情况下的其它流程例。在奇偶校验节点113B生成新奇偶校验(S17)而将前奇偶校验的物理地址存入在存储器214(S18)后,将完成响应发送至系统控制器105(S19)前,在预先设定的处理MP__A109A中产生了障碍。处理MP从MP__A109A移交给MP__B109B。

[0122] 数据节点113A和奇偶校验节点113B从MP__B109B接收重置命令(S31)。

[0123] 数据节点113A将与映射信息的逻辑地址对应的物理地址反馈至在S13存入在存储器214中的前数据地址321,然后反馈完成响应(S41)。同样,奇偶校验节点113B关于映射信息将与重置命令352的逻辑地址对应的物理地址反馈至在S18中存入在存储器214的前奇偶校验地址322,然后反馈完成响应(S42)。由此,在之后的步骤中能够生成正确的中间奇偶校验和新奇偶校验。其它方面与图11的流程相同。

[0124] 与上述任一情况均不同,继续进行(继承)处理的MP__B109B所参照的写入请求管理信息331表示“奇偶校验生成后”的情况下,MP__B109B判断为完成了对写入请求的奇偶校验更新。这样,通过写入请求管理信息331,能够恰当地控制障碍发生时的处理。

[0125] 此外,如上所述,在写入请求的处理途中的任一步骤的MP障碍中,继承处理的MP均从新数据的转送开始,由此能够高效地构成系统控制器105。

[0126] 在上述例中,写入请求管理信息331仅管理“奇偶校验生成前”、“奇偶校验生成后”,因此能够使系统结构简洁。写入请求管理信息331也可以表示其它阶段。例如,写入请求管理信息331也可以表示“中间奇偶校验生成前”、“中间奇偶校验生成后”的阶段。继承处理的MP__B109B进行与阶段相应的处理。

[0127] 实施例2

[0128] 使用图13~15说明与实施例1不同的实施例。在本实施例中,对将前数据地址321和前奇偶校验地址322没有存入存储器214中的情况下的处理进行说明。

[0129] 图13和图14表示本实施例中的与来自自主计算机101的写入请求相应的数据更新的另一正常流程例。以下主要说明与图3和图4所示的流程例的差异。图3和图4所示的流程例的前数据地址321和前奇偶校验地址322以及它们的存入步骤S13、S18在图13和图14的流程中省略。其它方面图3和图4所示的流程与图13和图14所示的流程相同。

[0130] 接着,对图13和图14所示的正常流程中产生了障碍的情况下的流程例进行说明。在图15的流程例中,在数据节点113A将中间奇偶校验发送至系统控制器105(S15)后,系统控制器105将中间奇偶校验转送至奇偶校验节点113B(S16)前,在MP__A109A在产生了障碍。

[0131] 以下主要对与实施例1的图11和图15所示的流程例的差异进行说明。在预先设定的处理MP__A109A中产生了障碍,处理MP从MP__A109A移交给MP__B109B。

[0132] MP__B109B参照写入请求管理信息331,获取从MP__A109A移交来的处理的信息(S51)。在本例中,写入请求管理信息331表示从MP__A109A移交来的处理的阶段为“奇偶校验生成前”。

[0133] 在图15的例子中,本流程的与图11的差异在于不使用前数据地址321和前奇偶校验地址322,因此不存在其存入和回写的步骤。此外,没有写入请求管理信息331的重置指示,写入请求管理信息331保持不更新的状态、即没有图11的S31~S34地直接前进至奇偶校验更新处理(S12~S19)。该奇偶校验更新处理(S12~S19)由于写入请求管理信息331的阶段为“奇偶校验生成前”因而进行。

[0134] 当从奇偶校验节点113B接收到新奇偶校验生成的完成响应时(S19),MP__B109B要更新写入请求管理信息331时进行参照,但是当判定为存储在写入请求管理信息331的处理MP与MP__B109B不同时,判定为发生奇偶校验不匹配。因此,将写入请求管理信息331的处理MP更新为MP__B109B(S52),再次生成奇偶校验(S53)。

[0135] 另外,关于S53的奇偶校验再生成,在本实施例中,由于没有将前数据地址321和前奇偶校验地址322存入存储器214,因此不采用之前说明的那样的使用新数据、前数据和前奇偶校验生成新奇偶校验的方法,而采用所谓的顺序写入的奇偶校验生成方法。具体而言,MP__B109B在本流程中,从包含数据节点113A在内的多个数据节点接收存入新数据并生成新奇偶校验的条带串的数据块。MP__B109B根据所接收到的条带串的数据块生成奇偶校验。

[0136] 本例不参照上次数据地址和上次奇偶校验地址,此外,变更了处理MP,因此存在奇偶校验节点113B生成不正确的新奇偶校验的可能性。由此,通过在继承了处理的MP__B109B,根据新数据再生成奇偶校验,能够可靠地保存正确的奇偶校验。

[0137] 在上述结构中,能够将一部分省略。例如,存储系统104也可以不使用新奇偶校验生成命令而生成新奇偶校验。例如,也可以是系统控制器105从奇偶校验节点读取前奇偶校验,生成新奇偶校验来存入奇偶校验节点。或者,也可以从系统控制器105向奇偶校验节点发送中间奇偶校验发送命令,接收其完成响应和新奇偶校验生成完成响应。

[0138] 另外,本发明并不限定于上述的实施例,能够包括各种各样的变形例。例如,上述的实施例为了将本发明说明得容易明白而进行了详细的说明,但是并不一定限定于包括所说明的所有结构。此外,能够将一个实施例的结构的一部分替换为其它实施例的结构,此外,还能够在一个实施例的结构中加入其它实施例的结构。此外,能够对各实施例的结构的一部分进行其它结构的追加·削除·替换。

[0139] 此外,上述的各结构、功能、处理部等例如也可以通过利用集成电路等进行设计、用硬件实现其一部分或全部。此外,上述各结构、功能等也可以通过处理器解释、执行实现各个功能的程序而以软件实现。实现各功能的程序、图表、文件夹等信息能够存储于存储器、硬盘、SSD等记录装置或IC卡、SD卡、DVD等记录介质。

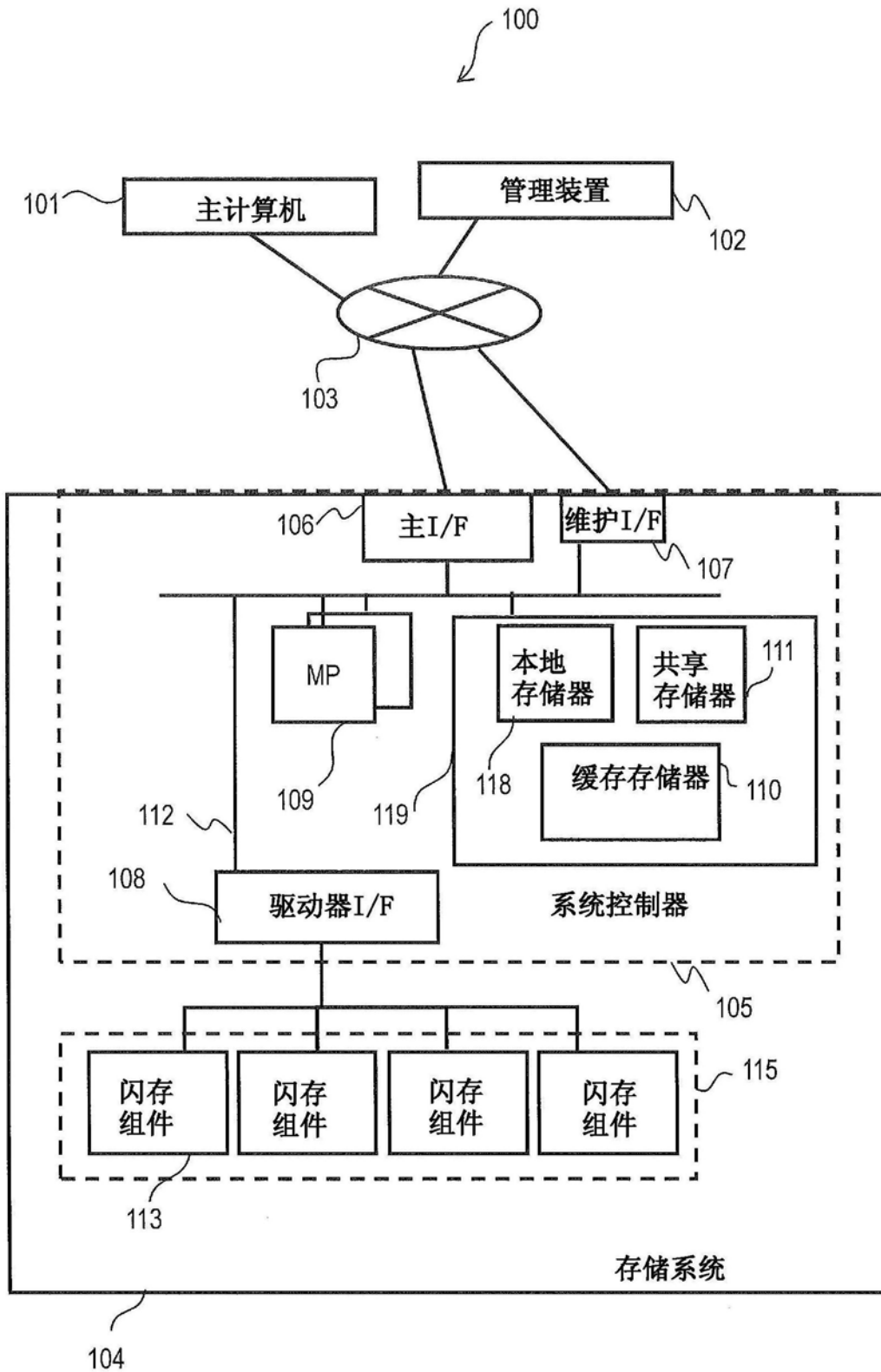


图1

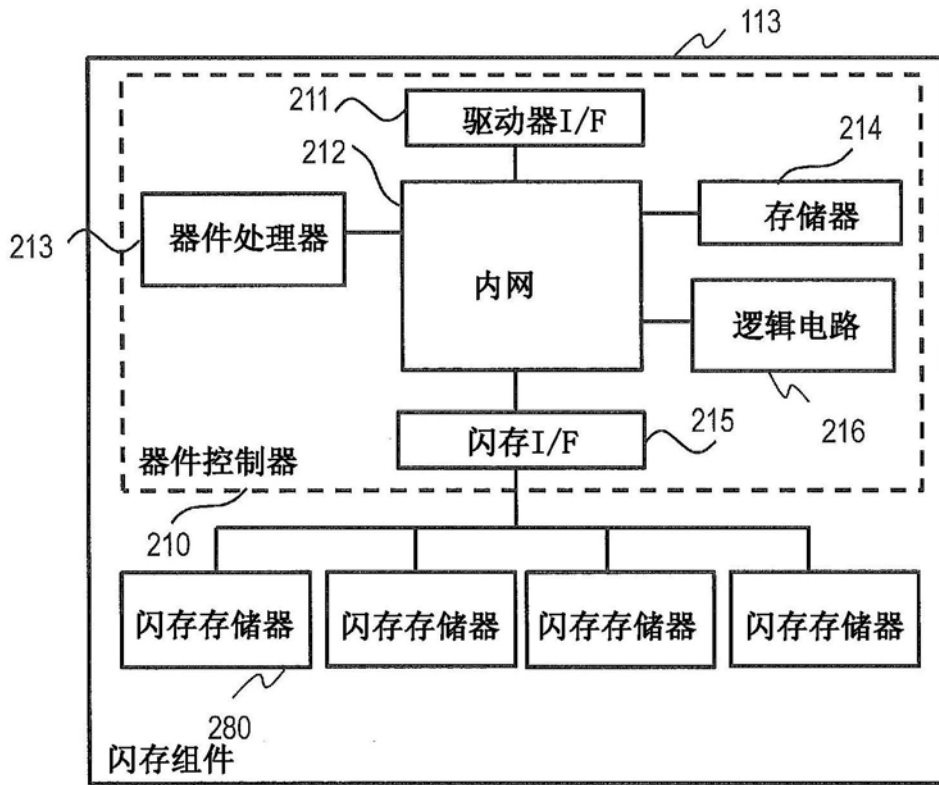


图2

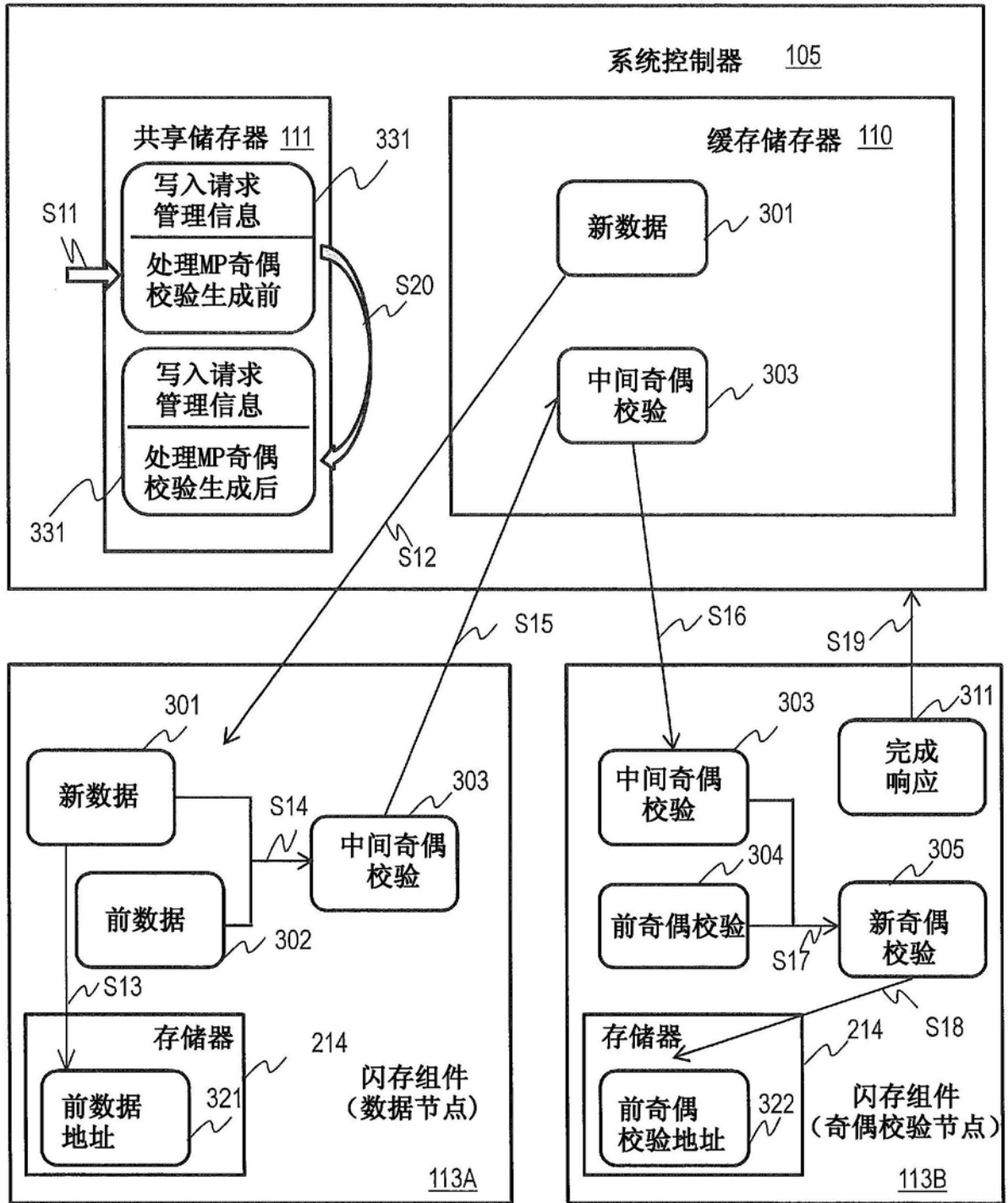


图3

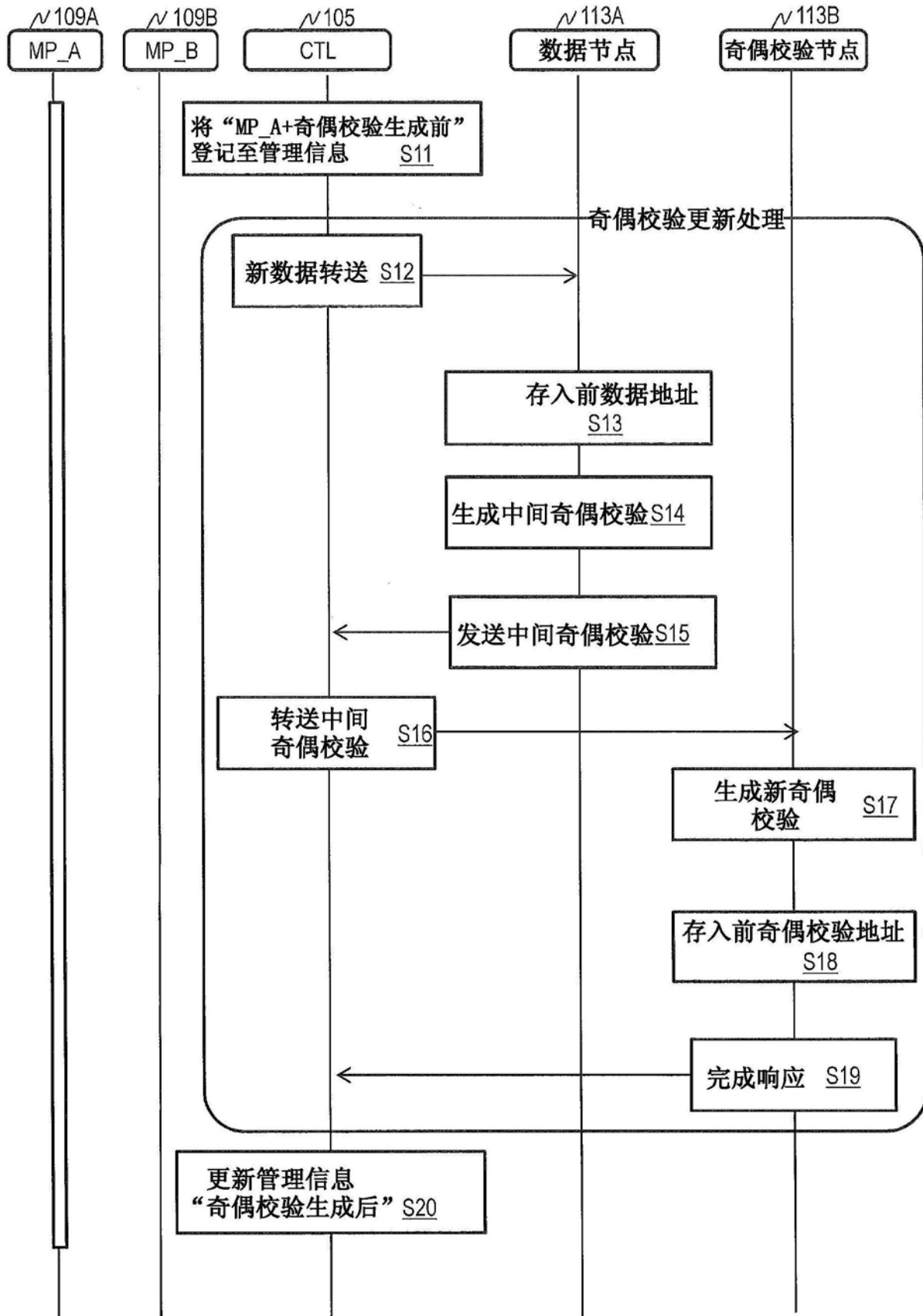


图4

写入请求管理信息 331			
请求ID	处理MP	阶段	逻辑地址
...
xxxx	MP_A	奇偶校验 生成后	aaaa
yyyy	MP_B	奇偶校验 生成前	bbbb

图5

中间奇偶校验生成命令 341	
字段类别	内容
OPECODE	WRITE
地址信息1	新数据的缓存存储器内地址
地址信息2	中间奇偶校验存入目的地的缓存存储器内地址
命令详情	中间奇偶校验生成
地址信息3	新数据的逻辑地址

图6

新奇偶校验生成命令 342	
字段类别	内容
OPECODE	WRITE
地址信息1	中间奇偶校验的缓存存储器内地址
命令详情	新奇偶校验生成
地址信息3	新奇偶校验的逻辑地址

图7

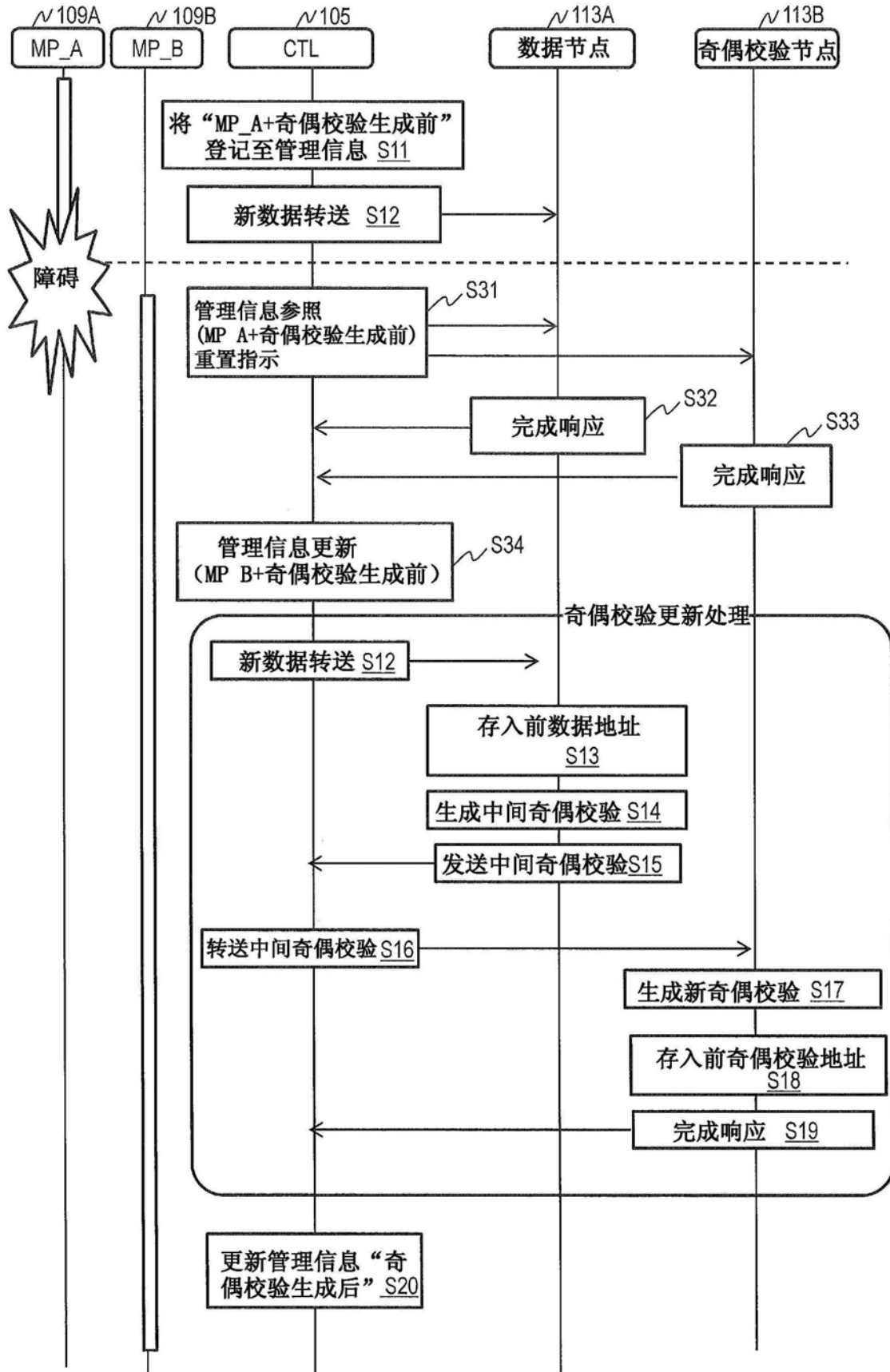


图8

重置命令 <u>351</u>	
字段类别	内容
地址信息	数据节点的逻辑地址
命令详情	中间奇偶校验生成的重置

图9A

重置命令 <u>352</u>	
字段类别	内容
地址详情	奇偶校验节点的逻辑地址
命令详情	新奇偶校验生成处理的重置

图9B

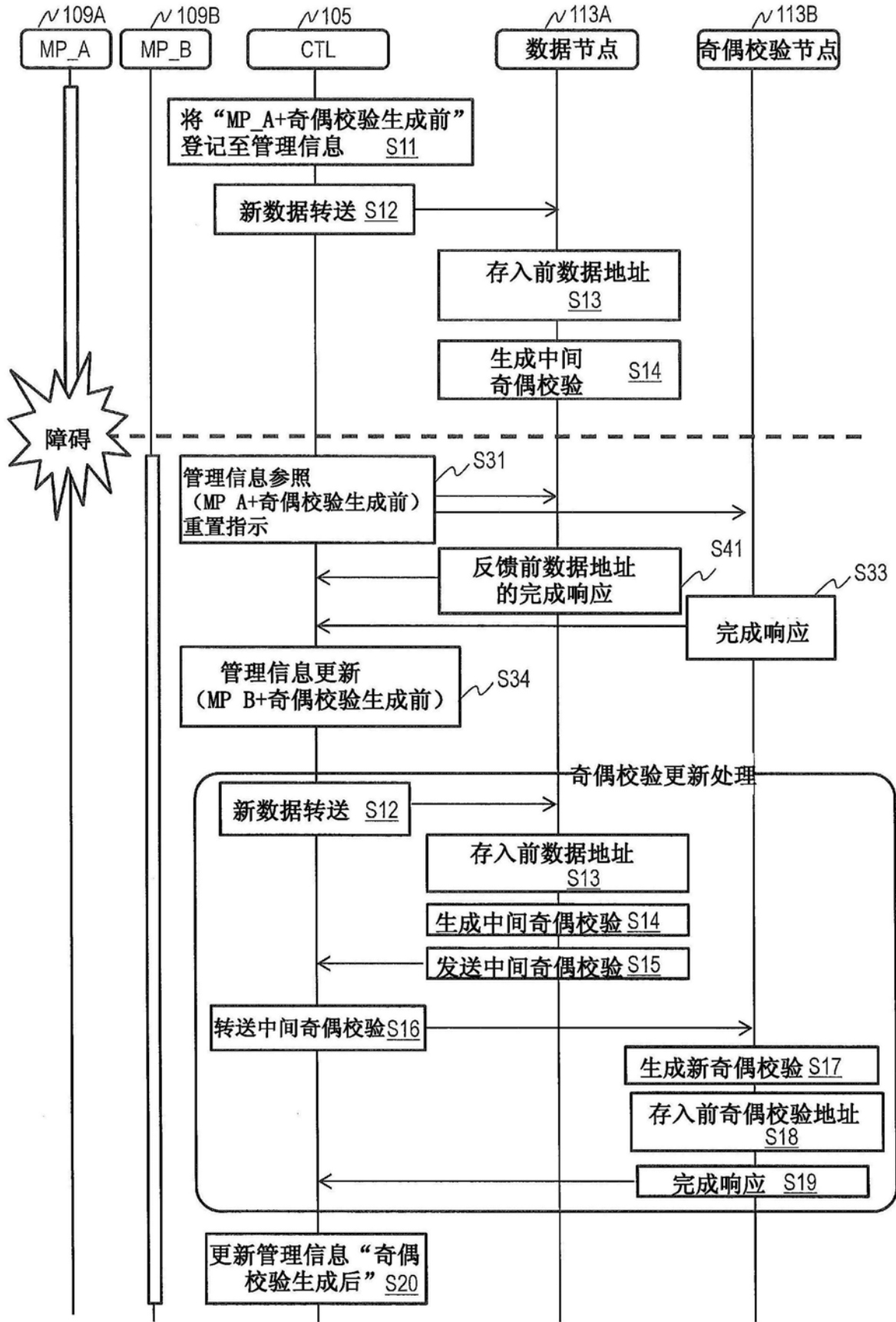


图10

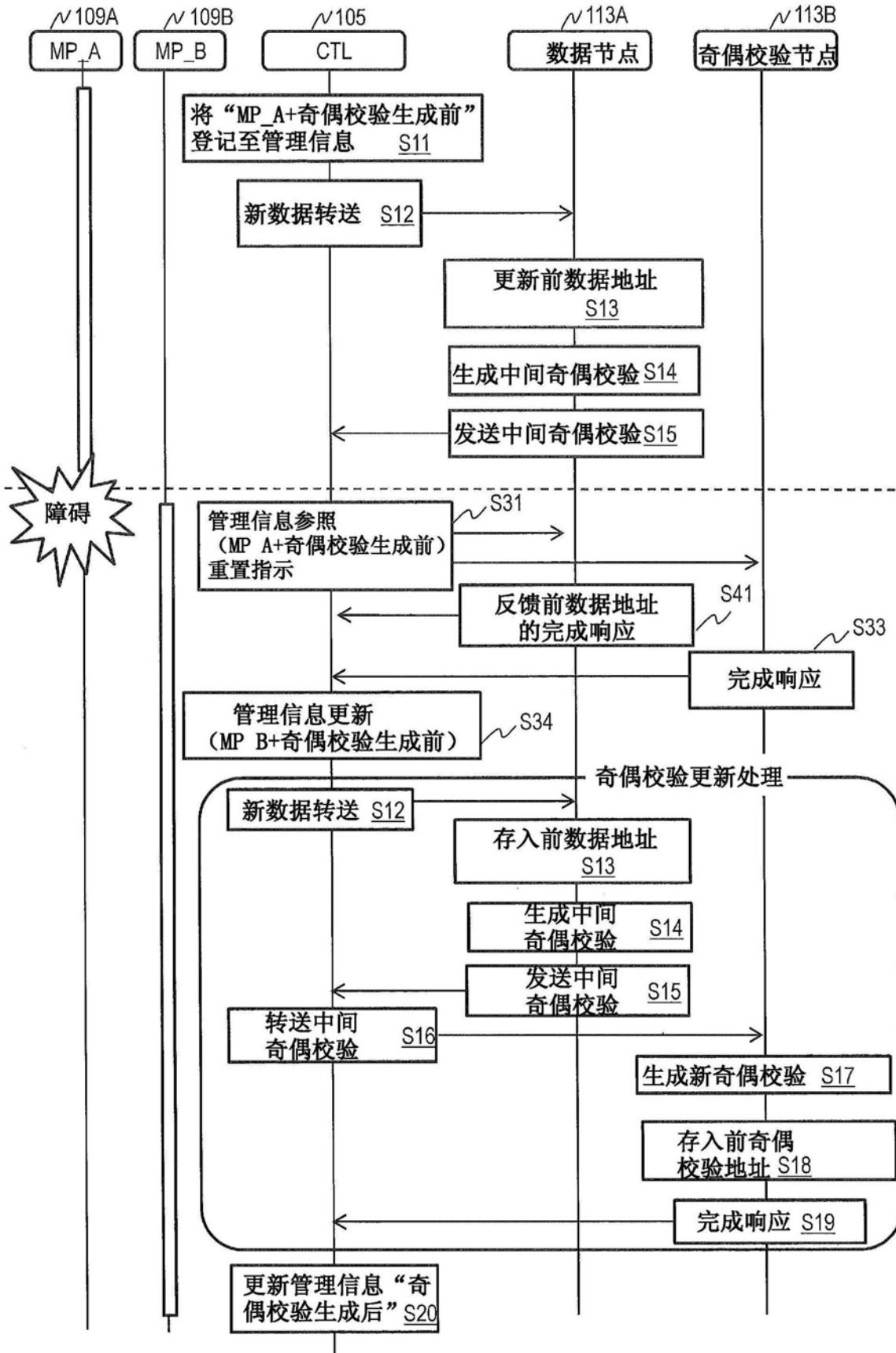


图11

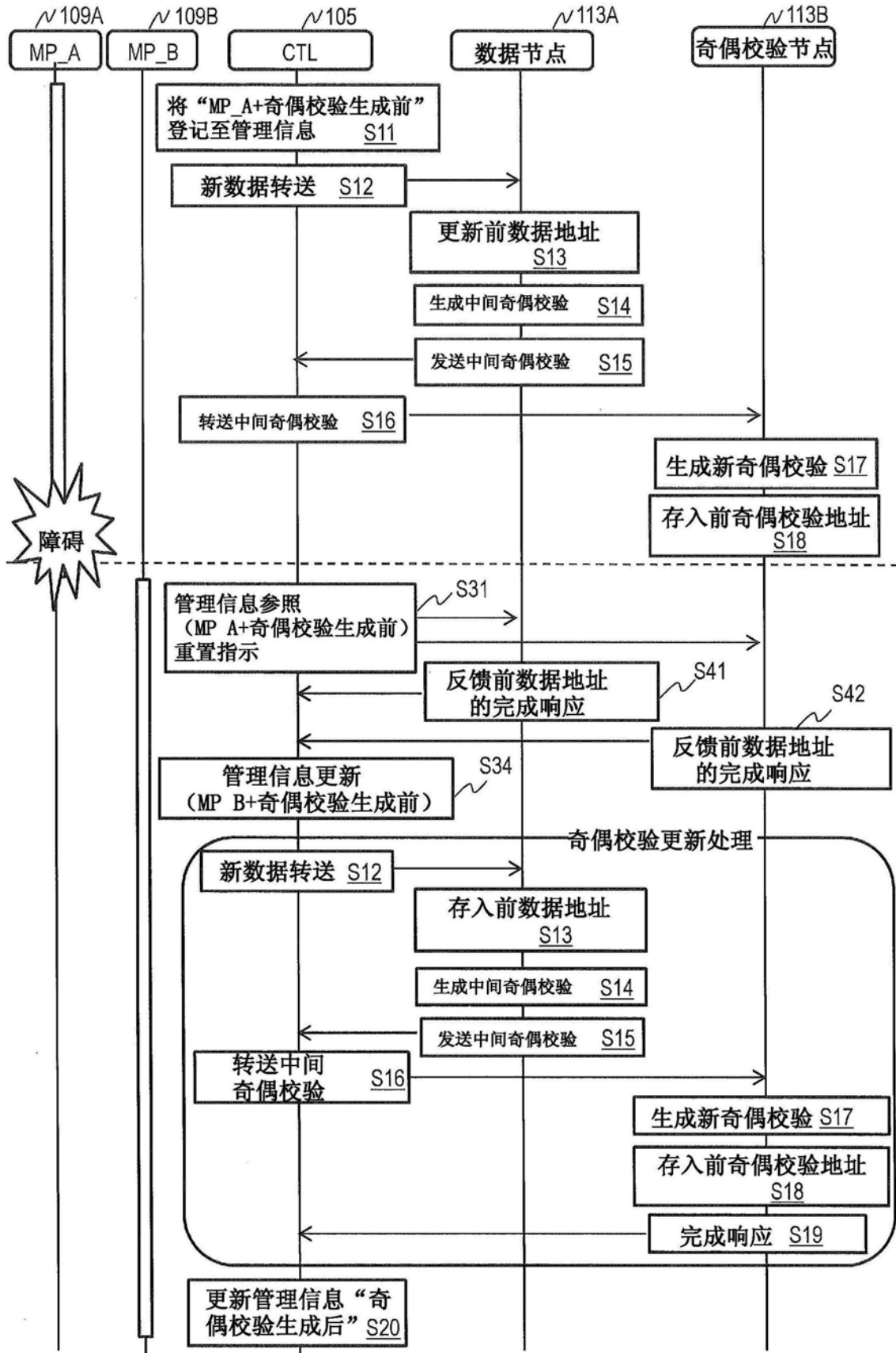


图12

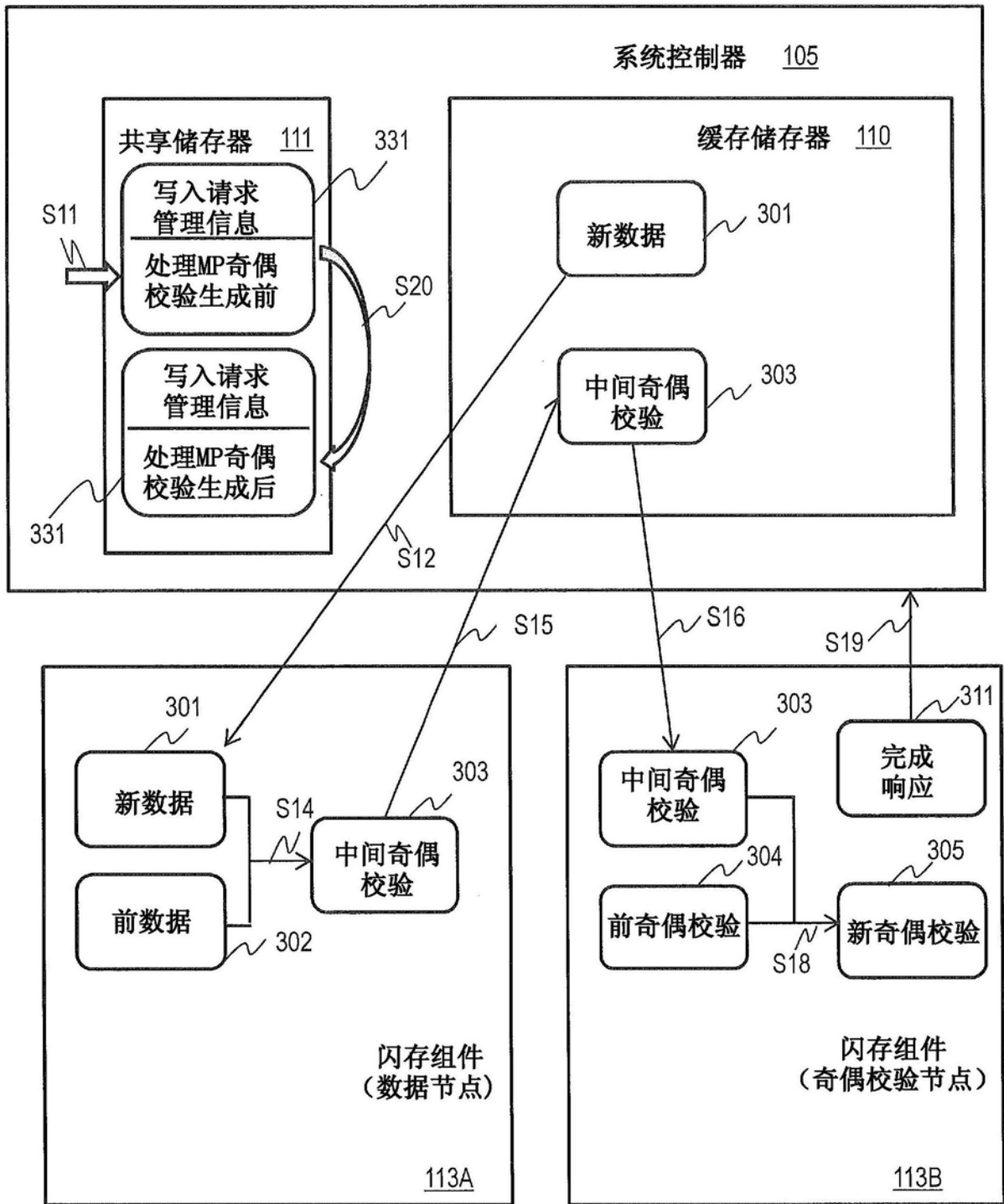


图13

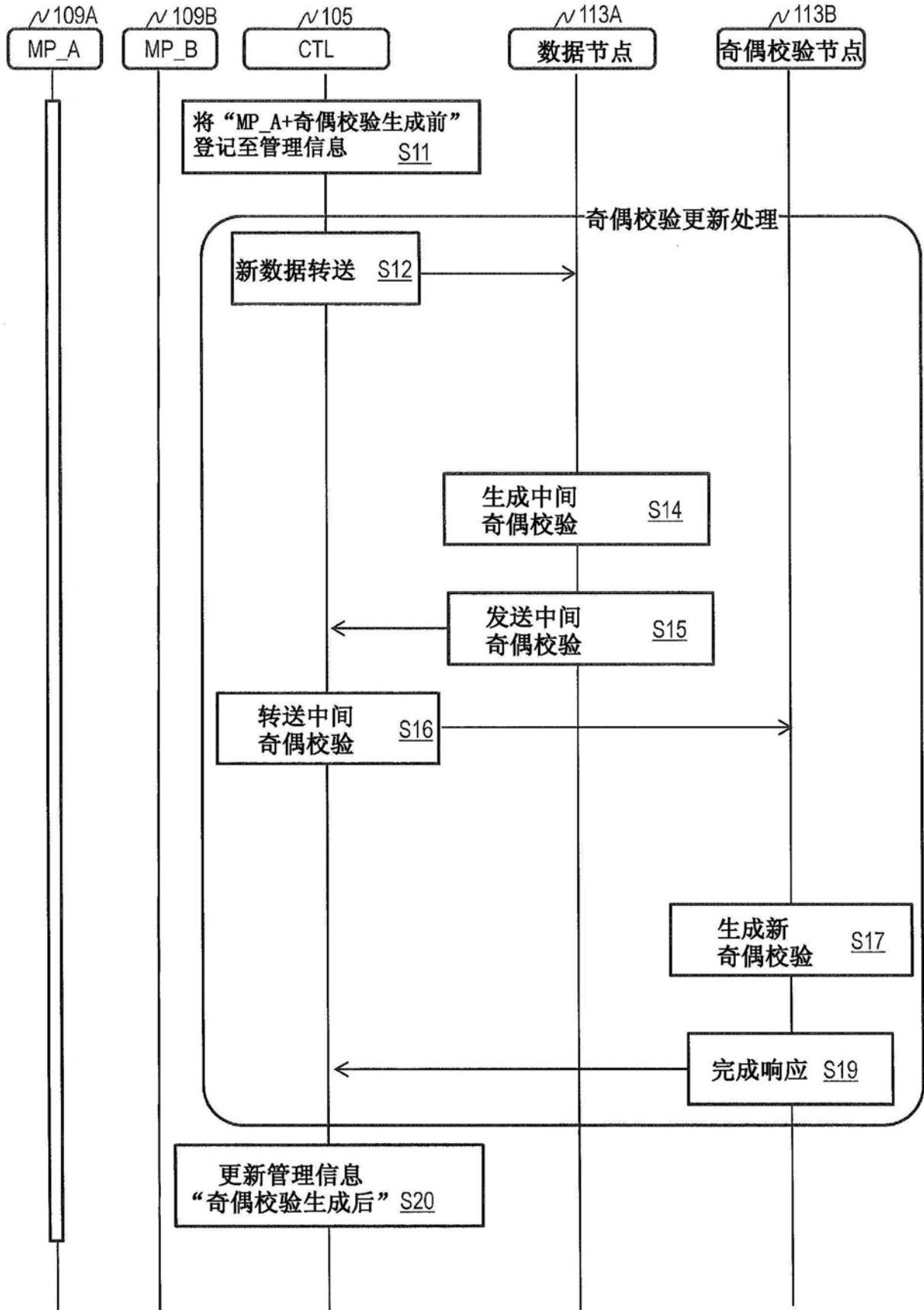


图14

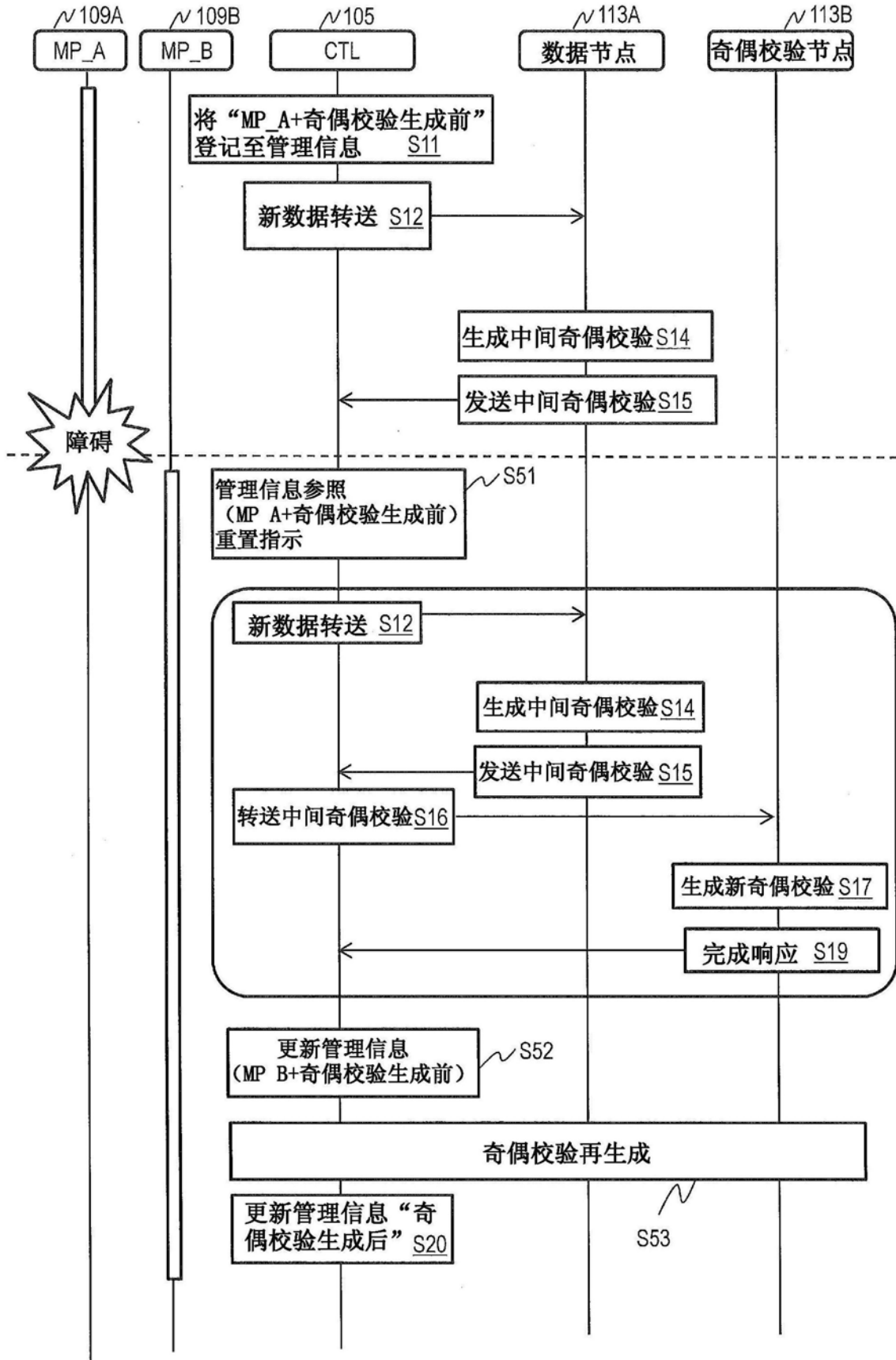


图15