



## (12) 发明专利申请

(10) 申请公布号 CN 118819927 A

(43) 申请公布日 2024. 10. 22

(21) 申请号 202410840617.8

(22) 申请日 2024.06.26

(71) 申请人 苏州元脑智能科技有限公司

地址 215000 江苏省苏州市吴中经济开发区郭巷街道官浦路1号9幢

(72) 发明人 苗永威

(74) 专利代理机构 北京润泽恒知识产权代理有限公司 11319

专利代理师 姜影

(51) Int. Cl.

G06F 11/07 (2006.01)

G06F 11/30 (2006.01)

G06F 1/20 (2006.01)

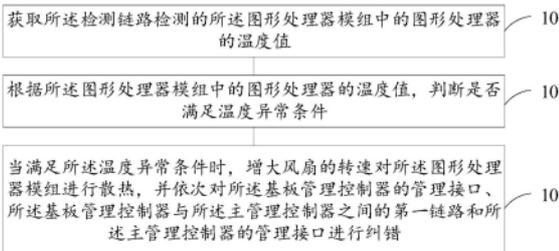
权利要求书2页 说明书14页 附图4页

## (54) 发明名称

一种服务器的检测链路的纠错方法、装置、设备及介质

## (57) 摘要

本发明实施例提供了一种服务器的检测链路的纠错方法、装置、设备及介质,服务器包括图形处理器模组,检测链路用于检测图形处理器的温度值,检测链路包括基板管理控制器的管理接口、基板管理控制器与主管理控制器之间的第一链路和主管理控制器的管理接口,该方法包括:获取检测链路检测的图形处理器模组中的图形处理器的温度值;根据图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件;当满足温度异常条件时,增大风扇的转速对图形处理器模组进行散热,并依次对检测链路进行纠错,从而可以在服务器的检测链路异常时,及时修复异常,避免失去对GPU温度的检测而导致GPU的温度长时间过热,提高了服务器的稳定性和计算效率。



1. 一种服务器的检测链路的纠错方法,其特征在于,所述服务器包括图形处理器模组,所述检测链路用于检测所述图形处理器模组中的图形处理器的温度值,所述检测链路包括基板管理控制器的管理接口、所述基板管理控制器与主管理控制器之间的第一链路和所述主管理控制器的管理接口,所述方法包括:

获取所述检测链路检测的所述图形处理器模组中的图形处理器的温度值;

根据所述图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件;

当满足所述温度异常条件时,增大风扇的转速对所述图形处理器模组进行散热,并依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错。

2. 根据权利要求1所述的方法,其特征在于,所述依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错,包括:

检测所述基板管理控制器的管理接口的工作状态;

若所述基板管理控制器的管理接口的工作状态异常,则对所述基板管理控制器的管理接口进行纠错;

若所述基板管理控制器的管理接口的工作状态正常,则检测所述基板管理控制器与所述主管理控制器之间的第一链路的工作状态;

若所述第一链路的工作状态异常,则对所述第一链路进行纠错;

若所述第一链路的工作状态正常,则检测所述主管理控制器的管理接口的工作状态;

若所述主管理控制器的管理接口的工作状态异常,则对所述主管理控制器的管理接口进行纠错;

在所述基板管理控制器的管理接口的工作状态,所述第一链路的工作状态,所述主管理控制器的管理接口的工作状态均正常后,获取所述图形处理器的温度值;

若所有图形处理器的温度值均正常,则恢复所述风扇的转速。

3. 根据权利要求2所述的方法,其特征在于,所述若所述基板管理控制器的管理接口的工作状态异常,则对所述基板管理控制器的管理接口进行纠错,包括:

若所述基板管理控制器的管理接口未激活或被占用,则确定所述基板管理控制器的管理接口的工作状态异常;

重新启用所述基板管理控制器的管理接口的权限,以对所述基板管理控制器的管理接口进行纠错。

4. 根据权利要求2所述的方法,其特征在于,所述若所述第一链路的工作状态异常,则对所述第一链路进行纠错,包括:

若所述第一链路的网络连接状态异常,则确定所述第一链路的工作状态异常;

重新建立所述第一链路的网络连接,以对所述第一链路进行纠错。

5. 根据权利要求2所述的方法,其特征在于,所述主管理控制器的管理接口包括第一主管理控制器接口,所述若所述主管理控制器的管理接口的工作状态异常,则对所述主管理控制器的管理接口进行纠错,包括:

若所述第一主管理控制器接口未激活或被占用,则确定所述第一主管理控制器接口的工作状态异常;

对所述主管理控制器对应的管理功能进行重置,以对所述主管理控制器的管理接口进行纠错。

6. 根据权利要求5所述的方法,其特征在于,所述主管理控制器的管理接口还包括第二主管理控制器接口,所述对所述主管理控制器的管理接口进行纠错,还包括:

若对所述主管理控制器对应的管理功能进行重置后,所述主管理控制器的管理接口的工作状态还存在异常,则通过I2C命令对所述主管理控制器进行重置,以对所述主管理控制器的管理接口进行纠错。

7. 根据权利要求1所述的方法,其特征在于,所述方法还包括:

将依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错的事件,记录到日志。

8. 一种服务器的检测链路的纠错装置,其特征在于,所述服务器包括图形处理器模组,所述检测链路用于检测所述图形处理器模组中的图形处理器的温度值,所述检测链路包括基板管理控制器的管理接口、所述基板管理控制器与主管理控制器之间的第一链路和所述主管理控制器的管理接口,所述装置包括:

模组检测温度获取模块,用于获取所述检测链路检测的所述图形处理器模组中的图形处理器的温度值;

温度异常条件判断模块,用于根据所述图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件;

异常检测链路纠错模块,用于当满足所述温度异常条件时,增大风扇的转速对所述图形处理器模组进行散热,并依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错。

9. 一种电子设备,其特征在于,包括:处理器、存储器及存储在所述存储器上并能够在所述处理器上运行的计算机程序,所述计算机程序被所述处理器执行时实现如权利要求1-7任一项所述的一种服务器的检测链路的纠错方法的步骤。

10. 一种计算机可读存储介质,其特征在于,所述计算机可读存储介质上存储计算机程序,所述计算机程序被处理器执行时实现如权利要求1-7任一项所述的一种服务器的检测链路的纠错方法的步骤。

## 一种服务器的检测链路的纠错方法、装置、设备及介质

### 技术领域

[0001] 本发明涉及电路检测技术领域,特别是涉及一种服务器的检测链路的纠错方法、装置、设备及介质。

### 背景技术

[0002] 图形处理器(GPU,Graphics Processing Unit)和GPU模组的更新换代带来了更高的浮点运行速度和显存带宽,GPU部件所带来的发热量也跟着水涨船高。GPU服务器通过检测链路实时获取GPU等高散热需求的元器件的温度,并根据温度执行对应的散热策略,以保证服务器能正常散热。当GPU服务器的检测链路异常时,就会失去对GPU实时温度状态的有效监测。

[0003] 然而,现有的GPU管理方法,缺少对检测链路异常情况的纠错机制。若检测链路的异常没有被及时修复,则服务器无法根据实时温度状态执行对应的散热策略,可能会使GPU的温度长时间过热。在高温环境下,GPU的稳定性和使用寿命会大大降低,当GPU温度超过设计规格后,不仅会出现降频降速,甚至会出现超温掉卡、算力应用报错等故障,给GPU模组造成不必要的损耗,也降低了服务器的稳定性和计算效率。

### 发明内容

[0004] 为了解决上述问题,本发明实施例公开了一种服务器的检测链路的纠错方法、装置、设备及介质。

[0005] 第一方面,本发明实施例提供了一种服务器的检测链路的纠错方法,所述服务器包括图形处理器模组,所述检测链路用于检测所述图形处理器模组中的图形处理器的温度值,所述检测链路包括基板管理控制器的管理接口、所述基板管理控制器与主管理控制器之间的第一链路和所述主管理控制器的管理接口,所述方法包括:

[0006] 获取所述检测链路检测的所述图形处理器模组中的图形处理器的温度值;

[0007] 根据所述图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件;

[0008] 当满足所述温度异常条件时,增大风扇的转速对所述图形处理器模组进行散热,并依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错。

[0009] 可选地,所述依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错,包括:

[0010] 检测所述基板管理控制器的管理接口的工作状态;

[0011] 若所述基板管理控制器的管理接口的工作状态异常,则对所述基板管理控制器的管理接口进行纠错;

[0012] 若所述基板管理控制器的管理接口的工作状态正常,则检测所述基板管理控制器与所述主管理控制器之间的第一链路的工作状态;

- [0013] 若所述第一链路的工作状态异常,则对所述第一链路进行纠错;
- [0014] 若所述第一链路的工作状态正常,则检测所述主管理控制器的管理接口的工作状态;
- [0015] 若所述主管理控制器的管理接口的工作状态异常,则对所述主管理控制器的管理接口进行纠错;
- [0016] 在所述基板管理控制器的管理接口的工作状态,所述第一链路的工作状态,所述主管理控制器的管理接口的工作状态均正常后,获取所述图形处理器的温度值;
- [0017] 若所有图形处理器的温度值均正常,则恢复所述风扇的转速。
- [0018] 可选地,所述若所述基板管理控制器的管理接口的工作状态异常,则对所述基板管理控制器的管理接口进行纠错,包括:
- [0019] 若所述基板管理控制器的管理接口未激活或被占用,则确定所述基板管理控制器的管理接口的工作状态异常;
- [0020] 重新启用所述基板管理控制器的管理接口的权限,以对所述基板管理控制器的管理接口进行纠错。
- [0021] 可选地,所述若所述第一链路的工作状态异常,则对所述第一链路进行纠错,包括:
- [0022] 若所述第一链路的网络连接状态异常,则确定所述第一链路的工作状态异常;
- [0023] 重新建立所述第一链路的网络连接,以对所述第一链路进行纠错。
- [0024] 可选地,所述主管理控制器的管理接口包括第一主管理控制器接口,所述若所述主管理控制器的管理接口的工作状态异常,则对所述主管理控制器的管理接口进行纠错,包括:
- [0025] 若所述第一主管理控制器接口未激活或被占用,则确定所述第一主管理控制器接口的工作状态异常;
- [0026] 对所述主管理控制器对应的管理功能进行重置,以对所述主管理控制器的管理接口进行纠错。
- [0027] 可选地,所述主管理控制器的管理接口还包括第二主管理控制器接口,所述对所述主管理控制器的管理接口进行纠错,还包括:
- [0028] 若对所述主管理控制器对应的管理功能进行重置后,所述主管理控制器的管理接口的工作状态还存在异常,则通过I2C命令对所述主管理控制器进行重置,以对所述主管理控制器的管理接口进行纠错。
- [0029] 可选地,所述方法还包括:
- [0030] 将依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错的事件,记录到日志。
- [0031] 第二方面,本发明实施例提供了一种服务器的检测链路的纠错装置,所述服务器包括图形处理器模组,所述检测链路用于检测所述图形处理器模组中的图形处理器的温度值,所述检测链路包括基板管理控制器的管理接口、所述基板管理控制器与主管理控制器之间的第一链路和所述主管理控制器的管理接口,所述装置包括:
- [0032] 模组检测温度获取模块,用于获取所述检测链路检测的所述图形处理器模组中的图形处理器的温度值;

- [0033] 温度异常条件判断模块,用于根据所述图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件;
- [0034] 异常检测链路纠错模块,用于当满足所述温度异常条件时,增大风扇的转速对所述图形处理器模组进行散热,并依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错。
- [0035] 可选地,所述异常检测链路纠错模块,包括:
- [0036] 第一检测子模块,用于检测所述基板管理控制器的管理接口的工作状态;
- [0037] 第一纠错子模块,用于若所述基板管理控制器的管理接口的工作状态异常,则对所述基板管理控制器的管理接口进行纠错;
- [0038] 第二检测子模块,用于若所述基板管理控制器的管理接口的工作状态正常,则检测所述基板管理控制器与所述主管理控制器之间的第一链路的工作状态;
- [0039] 第二纠错子模块,用于若所述第一链路的工作状态异常,则对所述第一链路进行纠错;
- [0040] 第三检测子模块,用于若所述第一链路的工作状态正常,则检测所述主管理控制器的管理接口的工作状态;
- [0041] 第三纠错子模块,用于若所述主管理控制器的管理接口的工作状态异常,则对所述主管理控制器的管理接口进行纠错;
- [0042] 温度获取子模块,用于在所述基板管理控制器的管理接口的工作状态,所述第一链路的工作状态,所述主管理控制器的管理接口的工作状态均正常后,获取所述图形处理器的温度值;
- [0043] 风扇调整子模块,用于若所有图形处理器的温度值均正常,则恢复所述风扇的转速。
- [0044] 可选地,所述第一纠错子模块,包括:
- [0045] 第一异常确定单元,用于若所述基板管理控制器的管理接口未激活或被占用,则确定所述基板管理控制器的管理接口的工作状态异常;
- [0046] 第一纠错单元,用于重新启用所述基板管理控制器的管理接口的权限,以对所述基板管理控制器的管理接口进行纠错。
- [0047] 可选地,所述第二纠错子模块,包括:
- [0048] 第二异常确定单元,用于若所述第一链路的网络连接状态异常,则确定所述第一链路的工作状态异常;
- [0049] 第二纠错单元,用于重新建立所述第一链路的网络连接,以对所述第一链路进行纠错。
- [0050] 可选地,所述主管理控制器的管理接口包括第一主管理控制器接口,所述第三纠错子模块,包括:
- [0051] 第三异常确定单元,用于若所述第一主管理控制器接口未激活或被占用,则确定所述第一主管理控制器接口的工作状态异常;
- [0052] 第三纠错单元,用于对所述主管理控制器对应的管理功能进行重置,以对所述主管理控制器的管理接口进行纠错。
- [0053] 可选地,所述主管理控制器的管理接口还包括第二主管理控制器接口,所述第三

纠错子模块,还包括:

[0054] 第四纠错单元,用于若对所述主管理控制器对应的管理功能进行重置后,所述主管理控制器的管理接口的工作状态还存在异常,则通过I2C命令对所述主管理控制器进行重置,以对所述主管理控制器的管理接口进行纠错。

[0055] 可选地,所述装置还包括:

[0056] 纠错事件记录日志模块,用于将依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错的事件,记录到日志。

[0057] 第三方面,本发明示出了一种电子设备,包括:存储器、处理器及存储在所述存储器上并可在所述处理器上运行的计算机程序,所述处理器执行所述计算机程序时实现上述的一种服务器的检测链路的纠错方法的步骤。

[0058] 第四方面,本发明示出了一种计算机可读存储介质,所述计算机可读存储介质上存储计算机程序,所述计算机程序被处理器执行时实现上述的一种服务器的检测链路的纠错方法的步骤。

[0059] 本发明实施例包括以下优点:

[0060] 本发明实施例可以通过获取检测链路检测的图形处理器模组中的图形处理器的温度值,以实时监测图形处理器的温度值;根据图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件,以便根据图形处理器的温度值即时判断检测链路的工作状态;当满足温度异常条件时,增大风扇的转速对图形处理器模组进行散热,并依次对基板管理控制器的管理接口、基板管理控制器与主管理控制器之间的第一链路和主管理控制器的管理接口进行纠错,从而可以在服务器的检测链路异常时,及时修复异常,避免失去对GPU实时温度状态的检测而导致GPU的温度长时间过热,提高了GPU模组服务器大规模计算的稳定性和计算效率。

## 附图说明

[0061] 为了更清楚地说明本发明实施例中的技术方案,下面将对实施例描述中所需要使用的附图作简单地介绍,显而易见地,下面描述中的附图仅仅是本发明的一些实施例,对于本领域普通技术人员来讲,在不付出创造性劳动的前提下,还可以根据这些附图获得其他的附图。

[0062] 图1是本发明实施例的一种服务器的结构框图;

[0063] 图2是本发明实施例的一种GPU温度获取异常的带外管理方法的逻辑图;

[0064] 图3是本发明实施例的一种服务器的检测链路的纠错方法的步骤流程图;

[0065] 图4是本发明实施例的另一种服务器的检测链路的纠错方法的步骤流程图;

[0066] 图5是本发明实施例的一种服务器的检测链路的纠错方法的逻辑图;

[0067] 图6是本发明实施例的一种服务器的检测链路的纠错装置的结构框图;

[0068] 图7是本发明实施例的一种电子设备的结构框图;

[0069] 图8是本发明实施例的一种计算机可读存储介质的结构框图。

## 具体实施方式

[0070] 目前人工智能已成为热点产业并渗透到不同行业不同领域,随着人工智能算法的飞速发展,联动人工智能的算力需求飞速增长。GPU服务器作为智能计算的核心载体,在处理并行计算密集型任务时具有显著优势,出色的图形处理能力和高性能计算能力提供极致计算性能,通过将应用程序计算密集部分的工作负载转移到GPU上,同时仍由CPU (Central Processing Unit,中央处理器) 运行其余程序代码,GPU服务器能够大幅提升应用程序的运行速度,有效解放计算压力,提升产品的计算处理效率与竞争力。

[0071] 人工智能、复杂模拟和海量数据集需要多个具有极快互连速度的GPU和完全加速的软件堆栈。GPU模组可以实现以多个GPU、各GPU间高速互联为整体的GPU计算单元,进而实现高带宽低延迟的强劲性能和极致扩展。GPU模组服务器更是成为面向超大规模数据中心和智算中心的核心资源及算力“发动机”,GPU模组服务器的稳定性对智算中心的算力保障息息相关。

[0072] 参照图1,示出了本发明实施例提供的一种服务器的结构框图。服务器的主板上设置有BMC(Baseboard Management Controller,基板管理控制器)模块,服务器的GPU模组上设置有HMC(Host Management Controller,主管理控制器)模块,BMC模块和HMC模块之间通过物理链路通信连接,HMC模块与GPU模组上的GPU也通过物理链路通信连接。其中,通信协议可以为Redfish协议或I2C协议,Redfish是一种开放的管理接口和协议,I2C是一种串行通信协议。BMC模块可以使用redfish带外管理协议(OOB,Out of Band)对GPU模组中各GPU、switch等组件进行带外管理监控。

[0073] 参照图2,示出了本发明实施例的一种GPU温度获取异常的带外管理方法的逻辑图。在该方法中,GPU模组服务器通过带外管理方法来控制GPU的温度的主要过程为:通过BMC模块物理链路连接到GPU模组的主管理控制器模块,使用带外管理协议对GPU模组中各GPU组件进行带外管理监控。当BMC或HMC的redfish接口或HMC与BMC之间建链的链路无响应时,就会导致BMC模块获取GPU带外温度返回值异常,当BMC轮询3次所有GPU温度返回值均异常时,BMC自动触发异常风扇调控策略以保证GPU的散热。

[0074] 图2中的带外管理方法缺少对BMC的接口、HMC的接口以及BMC与HMC之间的链路异常的纠错机制,BMC模块不会主动对以上异常现象进行干预修复。而且,若不是所有GPU的带外温度获取异常(即仅有个别GPU带外温度获取异常),则BMC模块并不会触发异常风扇转速调控。此时,风扇模组的实时转速或实时转速的提升速度可能无法满足GPU温度继续上升所带来的散热需求,从而会出现GPU超温掉卡或服务器宕机的现象。GPU模组的超温掉卡问题必须对服务器整系统做交流断电开关机操作,来恢复所有GPU的状态。

[0075] 针对上述存在的问题,本发明提出了一种服务器的检测链路的纠错方法,目的是在服务器的检测链路异常时,及时修复异常,避免失去对GPU实时温度状态的检测而导致GPU的温度长时间过热。为了实现这一目标,本发明实施例通过检测链路实时检测GPU模组中的GPU的温度值,当温度值返回异常满足温度异常条件时,增大风扇的转速对GPU模组进行散热,并依次对BMC的管理接口、BMC与HMC之间的第一链路和HMC的管理接口进行纠错,从而可以提高GPU模组服务器的稳定性和计算效率。

[0076] 为使本发明的上述目的、特征和优点能够更加明显易懂,下面结合附图和具体实施方式对本发明作进一步详细的说明。

[0077] 参照图3,示出了本发明实施例的一种服务器的检测链路的纠错方法的步骤流程图,所述服务器包括图形处理器模组,所述检测链路用于检测所述图形处理器模组中的图形处理器的温度值,所述检测链路包括基板管理控制器的管理接口、所述基板管理控制器与主管理控制器之间的第一链路和所述主管理控制器的管理接口,所述方法具体可以包括如下步骤:

[0078] 步骤101,获取所述检测链路检测的所述图形处理器模组中的图形处理器的温度值;

[0079] 在本发明实施例中,服务器包括GPU模组,GPU模组中包括多个GPU,服务器可以通过检测链路检测GPU模组中的所有GPU的温度值,并将温度值的检测结果返回至BMC,从而避免失去对GPU实时温度状态的检测而导致GPU的温度长时间过热。其中,检测链路包括BMC的管理接口、BMC与HMC之间的第一链路和HMC的管理接口,BMC的管理接口可以为redfish接口,HMC的管理接口可以为redfish接口,BMC与HMC之间的第一链路可以为redfish链路或I2C链路。

[0080] 步骤102,根据所述图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件;

[0081] 在本发明实施例中,BMC可以根据GPU模组中的每个GPU的温度值,判断GPU模组是否满足温度异常条件,从而即时判断服务器的检测链路是否异常。具体的,BMC可以实时监测GPU模组中每个GPU的温度值,并在GPU模组中至少一个GPU的温度值异常时,记录异常次数。当异常次数达到预设次数阈值时,确定GPU模组满足温度异常条件,其中,预设次数阈值可以为2。本领域技术人员可以根据本发明的思想,设置预设次数阈值为其他恰当的值,本发明对此不做限制。

[0082] 步骤103,当满足所述温度异常条件时,增大风扇的转速对所述图形处理器模组进行散热,并依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错。

[0083] 在本发明实施例中,当满足温度异常条件时,可以确定服务器的检测链路异常,MBC可以将风扇的转速增大到100%对GPU模组进行散热,并依次对BMC的管理接口、BMC与HMC之间的第一链路和HMC的管理接口进行纠错,从而可以在服务器的检测链路异常时,及时修复异常,避免失去对GPU实时温度状态的检测而导致GPU的温度长时间过热。

[0084] 在一种实施例中,所述依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错的步骤,可以包括如下子步骤:

[0085] 子步骤S11,检测所述基板管理控制器的管理接口的工作状态;

[0086] 子步骤S12,若所述基板管理控制器的管理接口的工作状态异常,则对所述基板管理控制器的管理接口进行纠错;

[0087] 在一种实施例中,所述若所述基板管理控制器的管理接口的工作状态异常,则对所述基板管理控制器的管理接口进行纠错,包括:若所述基板管理控制器的管理接口未激活或被占用,则确定所述基板管理控制器的管理接口的工作状态异常;重新启用所述基板管理控制器的管理接口的权限,以对所述基板管理控制器的管理接口进行纠错。

[0088] 在本发明实施例中,当服务器的检测链路异常时,BMC可以先对BMC的redfish接口

进行诊断排查,查看BMC模块对应的redfish接口状态是否为Enabled(激活),以及是否被I2C的SMBPBI权限占用。若发现BMC的redfish接口未激活或被占用,则确定BMC的redfish接口的工作状态异常,此时可以对BMC的redfish接口的权限进行修复,即重新启用BMC redfish接口权限,以对BMC的redfish接口进行纠错,从而确保Redfish接口的稳定性和可靠性,提高了BMC的管理功能可用性。

[0089] 具体的,对BMC的redfish接口的权限进行修复的过程为:将带外管理权限从redfish协议切换至I2C协议,再释放I2C协议,并重新切换至redfish协议。对BMC的redfish接口的权限进行修复的代码如下:

```
[0090] ***
[0091] #切换Out-of-band privilege至I2C
[0092] echo"----Get SMBPBI fencing privilege command for BMC--Opcode a3H,
arg1 01H----"
[0093] i2cset-y 11 0x60 0x5c 0x04 0xa3 0x01 0x00 0x80 i
[0094] i2ctransfer-y 11w1@0x60 0x5c r5
[0095] i2ctransfer-y 11w1@0x60 0x5d r5
[0096] #释放Out-of-band privilege的I2C并切换至redfish
[0097] echo"-----release SMBPBI fencing privilege command for BMC--Opcode
a3H,arg1 00H-----"
[0098] i2cset-y 11 0x60 0x5c 0x04 0xa3 0x00 0x00 0x80 i
[0099] i2ctransfer-y 11w1@0x60 0x5c r5
[0100] i2ctransfer-y 11w1@0x60 0x5d r5
[0101] ***
```

[0102] 子步骤S13,若所述基板管理控制器的管理接口的工作状态正常,则检测所述基板管理控制器与所述主管理控制器之间的第一链路的工作状态;

[0103] 子步骤S14,若所述第一链路的工作状态异常,则对所述第一链路进行纠错;

[0104] 在一种实施例中,所述若所述第一链路的工作状态异常,则对所述第一链路进行纠错,包括:若所述第一链路的网络连接状态异常,则确定所述第一链路的工作状态异常;重新建立所述第一链路的网络连接,以对所述第一链路进行纠错。

[0105] 在本发明实施例中,若BMC的redfish接口的工作状态正常,则可以继续检测BMC与HMC之间的第一链路的工作状态。具体的,BMC模块可以对HMC模块自定义的默认IP(默认IP:192.168.31.1)进行ping通尝试,并查看HMC的启动进度状态,若IP无法ping通,则说明第一链路的网络连接状态异常,即第一链路的USB枚举链路建链出现问题,此时可以重新建立第一链路的网络连接,以对第一链路进行纠错,从而确保BMC模块可以和HMC模块进行交互。重新建立第一链路的网络连接的代码如下:

```
[0106] ***
[0107] ping 192.168.31.1
[0108] i2cdump-y 11 0x54
[0109] no size specified(using byte-data access)
[0110] ***
```

```
[0111] ###USB重新枚举链路建链
[0112] cd/sys/bus/platform/drivers/ehci-platform/
[0113] /sys/bus/platform/drivers/ehci-platform#ls le6a3000.usb bind
[0114] /sys/bus/platform/drivers/ehci-platform#echo le6a3000.usb>unbind
[0115] /sys/bus/platform/drivers/ehci-platform#echo le6a3000.usb>bind
[0116] sleep 30
[0117] ping 192.168.31.1
[0118] ***
```

[0119] 子步骤S15,若所述第一链路的工作状态正常,则检测所述主管理控制器的管理接口的工作状态;

[0120] 子步骤S16,若所述主管理控制器的管理接口的工作状态异常,则对所述主管理控制器的管理接口进行纠错;

[0121] 在一种实施例中,所述主管理控制器的管理接口包括第一主管理控制器接口和第二主管理控制器接口,所述若所述主管理控制器的管理接口的工作状态异常,则对所述主管理控制器的管理接口进行纠错,包括:

[0122] 若所述第一主管理控制器接口未激活或被占用,则确定所述第一主管理控制器接口的工作状态异常;

[0123] 对所述主管理控制器对应的管理功能进行重置,以对所述主管理控制器的管理接口进行纠错;

[0124] 若对所述主管理控制器对应的管理功能进行重置后,所述主管理控制器的管理接口的工作状态还存在异常,则通过I2C命令对所述主管理控制器进行重置,以对所述主管理控制器的管理接口进行纠错。

[0125] 在本发明实施例中,在完成HMC模块到BMC模块的USB重新建链后,BMC可以对HMC的redfish接口进行诊断排查。其中,HMC的redfish接口包括第一HMC接口和第二HMC接口,当第一HMC接口和第二HMC接口都返回状态enable的返回值时,说明HMC的redfish接口的工作状态正常。

[0126] 若第一HMC接口未激活或被占用,则BMC可以对HMC对应的redfish功能相关软件栈做重置初始化,以对第一HMC接口进行纠错,从而可以解决由于HMC配置错误或损坏导致的问题,并清除软件栈的错误状态,通过重置初始化,可以简化第一HMC接口的维护过程,减少故障诊断和修复的时间。若BMC对HMC对应的redfish功能相关软件栈进行重置后,HMC的redfish接口的工作状态还存在异常,则可以通过I2C命令对HMC模块进行重置,即通过I2C命令对HMC模块各redfish接口进行恢复出厂设置,在恢复出厂设置后,可以根据当前的需求和最佳实践重新配置Redfish接口,确保HMC的所有Redfish接口配置一致,便于管理和维护。对第一HMC接口和第二HMC接口进行修复的代码如下:

```
[0127] ***
[0128] ###reset HMC redfish interface
[0129] curl--insecure-uroot:OpenBmc-XPOSThttp://192.168.31.1/redfish/v1/Managers/HGX_BMC_0/Actions/Manager.ResetToDefaults-d'{"ResetToDefaultsType":"ResetAll"}
```

[0130] \*\*\*

[0131] ###HMC factory reset

[0132] i2cset-f-y 13 0x54 0x00 0x0f

[0133] \*\*\*

[0134] 需要说明的是,HMC模块中类似集成一个微型特制版的linux系统,对HMC模块的管理接口进行修复有两种层级:第一,reset HMC redfish interface,通过redfish命令做HMC的redfish硬件重置,以实现针对HMC集成的linux系统中redfish功能相关软件栈的重置初始化;第二,HMC factory reset,通过I2C命令对HMC进行恢复出厂设置,以实现HMC集成的linux系统的重置初始化。其中,第二层级的时长长于第一层级的时长。

[0135] 子步骤S17,在所述基板管理控制器的管理接口的工作状态,所述第一链路的工作状态,所述主管理控制器的管理接口的工作状态均正常后,获取所述图形处理器的温度值;

[0136] 子步骤S18,若所有图形处理器的温度值均正常,则恢复所述风扇的转速。

[0137] 在本发明实施例中,在BMC的redfish接口的工作状态,第一链路的工作状态,以及HMC的redfish接口的工作状态均正常后,BMC可以重新获取GPU的温度值,以确定所有GPU的温度值是否均正常。BMC模块可以再次对HMC模块的redfish接口进行检测并重新对GPU的温度轮询进行异常判断。重新检测的代码如下:

[0138] \*\*\*

[0139] curl--insecure-uroot:OpenBmc-XGEThttp://192.168.31.1/redfish/v1/Managers/HGX\_BMC\_0

[0140] curl--insecure-uroot:OpenBmc-XGEThttp://192.168.31.1/redfish/v1/TelemetryService/MetricReports/HGX\_PlatformEnvironmentMetrics\_0

[0141] \*\*\*

[0142] 在确定所有GPU的温度值均正常后,即当BMC模块可以获取并确认所有GPU带外温度均正常可用,以及对当前GPU模组服务器整机状态满足度做进一步判断检查后,可以恢复风扇的转速为检测链路异常前的值,从而降低服务器的噪音。

[0143] 在本发明实施例中,基板管理控制器可以通过获取检测链路检测的图形处理器模组中的图形处理器的温度值,以实时监测图形处理器的温度值;根据图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件,以便根据图形处理器的温度值即时判断检测链路的工作状态;当满足温度异常条件时,增大风扇的转速对图形处理器模组进行散热,并依次对基板管理控制器的管理接口、基板管理控制器与主管理控制器之间的第一链路和主管理控制器的管理接口进行纠错,从而可以在服务器的检测链路异常时,及时修复异常,避免失去对GPU实时温度状态的检测而导致GPU的温度长时间过热,提高了GPU模组服务器大规模计算的稳定性和计算效率。

[0144] 参照图4,示出了本发明实施例的另一种服务器的检测链路的纠错方法的步骤流程图,所述服务器包括图形处理器模组,所述检测链路用于检测所述图形处理器模组中的图形处理器的温度值,所述检测链路包括基板管理控制器的管理接口、所述基板管理控制器与主管理控制器之间的第一链路和所述主管理控制器的管理接口,所述方法具体可以包括如下步骤:

[0145] 步骤201,获取所述检测链路检测的所述图形处理器模组中的图形处理器的温度

值;

[0146] 对于步骤201而言,由于其与步骤101相同,相关之处参见步骤101的说明即可。

[0147] 步骤202,根据所述图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件;

[0148] 对于步骤202而言,由于其与步骤102相同,相关之处参见步骤102的说明即可。

[0149] 步骤203,当满足所述温度异常条件时,增大风扇的转速对所述图形处理器模组进行散热,并依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错;

[0150] 对于步骤203而言,由于其与步骤103相同,相关之处参见步骤103的说明即可。

[0151] 步骤204,将依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错的事件,记录到日志。

[0152] 在本发明实施例中,在依次对BMC的redfish接口、BMC与HMC之间的第一链路和HMC的redfish接口进行纠错的过程中,可以将每次纠错对应的事件记录到日志,从而可以对服务器的检测链路进行长期监控,通过分析日志还可以减少未来的错误和故障。此外,服务器还可以根据纠错事件的日志,分别统计BMC的redfish接口、HMC的redfish接口和第一链路出现异常的概率,并根据概率设置对上述各组件进行纠错的优先级,其中,出现异常的概率越高的组件的纠错优先级越高,从而可以在检测链路出现异常时,根据检测链路各组件的纠错优先级对检测链路进行纠错,提高了对检测链路纠错的效率。

[0153] 记录日志的代码如下:

[0154] \*\*\*

[0155] ###selftest dump log

[0156] curl--insecure-uroot:OpenBmc-XPOSThttp://\$HMC\_IP/redfish/v1/Systems/HGX\_Baseboard\_0/LogServices/Dump/Actions/LogService.CollectDiagnosticData-d'{"DiagnosticDataType":"OEM","OEMDiagnosticDataType":"DiagnosticType=SelfTest"}">>\$00B\_Redfish\_log 2>&1

[0157] ###FPGAdump log

[0158] curl--insecure-uroot:OpenBmc-XPOSThttp://\$HMC\_IP/redfish/v1/Systems/HGX\_Baseboard\_0/LogServices/Dump/Actions/LogService.CollectDiagnosticData-d'{"DiagnosticDataType":"OEM","OEMDiagnosticDataType":"DiagnosticType=FPGA"}">>\$00B\_Redfish\_log 2>&1

[0159] ###EROT dump log

[0160] curl--insecure-uroot:OpenBmc-XPOSThttp://\$HMC\_IP/redfish/v1/Systems/HGX\_Baseboard\_0/LogServices/Dump/Actions/LogService.CollectDiagnosticData-d'{"DiagnosticDataType":"OEM","OEMDiagnosticDataType":"DiagnosticType=EROT"}">>\$00B\_Redfish\_log 2>&1

[0161] ###HMC dump log

[0162] curl--insecure-uroot:OpenBmc-XPOSThttp://\$HMC\_IP/redfish/v1/Managers/HGX\_BMC\_0/LogServices/Dump/Actions/LogService.CollectDiagnosticData-d'{"

DiagnosticDataType\":"\Manager\}"'>>\$OOB\_Redfish\_log 2>&l

[0163] \*\*\*

[0164] 本发明实施例通过对GPU模组服务器中BMC模块和GPU模组的带外管理功能的优化,实现GPU模组带外管理接口和带外管理redfish链路的实时监听和自动诊断纠错,确保GPU模组服务器整系统在GPU温度等高散热需求组件的带外监控的稳定性和风扇调控的及时性和准确性,实现了智算中心GPU模组服务器整体带外管理监控的优化和散热策略调控的优化,降低了GPU模组服务器因过热导致的不必要损耗和寿命衰减,提升了GPU模组服务器算力计算的可持续性及其算力能效。

[0165] 参照图5,示出了本发明实施例提供的一种服务器的检测链路的纠错方法的逻辑图,为了使本领域技术人员能够更好地理解本发明实施例,下面通过图5对本发明实施例加以说明:

[0166] 1) BMC通过HMC对GPU模组进行监控;

[0167] 2) 判断BMC进行2次温度轮询是否出现GPU温度异常;

[0168] 3) 若BMC进行2次温度轮询都没有出现GPU温度异常,则服务器风扇正常运转;

[0169] 4) 若BMC进行2次温度轮询出现了GPU温度异常,则BMC触发异常诊断纠错;

[0170] 5) 服务器风扇转速调整至100%;

[0171] 6) 判断BMC的redfish接口的工作状态是否异常;

[0172] 7) 若BMC的redfish接口的工作状态异常,则重启BMC的redfish接口权限,并记录日志;

[0173] 8) 若BMC的redfish接口的工作状态正常,则判断BMC模块对HMC模块的IP进行ping通是否异常;

[0174] 9) 若BMC模块对HMC模块的IP进行ping通异常,则重新建立第一链路的网络连接,并记录日志;

[0175] 10) 若BMC模块对HMC模块的IP进行ping通正常,判断HMC的redfish接口的工作状态是否正常;

[0176] 11) 若HMC的redfish接口的工作状态异常,则通过redfish命令对HMC模块进行重置,并记录日志;

[0177] 12) 对HMC模块进行重置后,再次判断HMC的redfish接口的工作状态是否正常;

[0178] 13) 若对HMC模块进行重置后,HMC的redfish接口的工作状态还异常,则通过I2C命令对HMC模块进行恢复出厂设置,并记录日志;

[0179] 14) 若对HMC模块进行重置后,HMC的redfish接口的工作状态正常,则再次判断BMC进行2次温度轮询是否都出现GPU温度异常;

[0180] 15) 若HMC的redfish接口的工作状态正常,则再次判断BMC进行2次温度轮询是否出现GPU温度异常;

[0181] 16) 若BMC再次进行2次温度轮询出现了GPU温度异常,则重新判断BMC的redfish接口的工作状态是否异常;

[0182] 17) 若BMC再次进行2次温度轮询没有出现GPU温度异常,则恢复服务器风扇的转速。

[0183] 本发明实施例通过获取检测链路检测的图形处理器模组中的图形处理器的温度

值,以实时监测图形处理器的温度值;根据图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件,以便根据图形处理器的温度值即时判断检测链路的工作状态;当满足温度异常条件时,增大风扇的转速对图形处理器模组进行散热,并依次对基板管理控制器的管理接口、基板管理控制器与主管理控制器之间的第一链路和主管理控制器的管理接口进行纠错,从而可以在服务器的检测链路异常时,及时修复异常,避免失去对GPU实时温度状态的检测而导致GPU的温度长时间过热,提高了GPU模组服务器大规模计算的稳定性和计算效率。

[0184] 需要说明的是,对于方法实施例,为了简单描述,故将其都表述为一系列的动作组合,但是本领域技术人员应该知悉,本发明实施例并不受所描述的动作顺序的限制,因为依据本发明实施例,某些步骤可以采用其他顺序或者同时进行。其次,本领域技术人员也应该知悉,说明书中所描述的实施例均属于优选实施例,所涉及的动作并不一定是本发明实施例所必须的。

[0185] 参照图6,示出了本发明实施例提供的一种服务器的检测链路的纠错装置的结构框图,所述服务器包括图形处理器模组,所述检测链路用于检测所述图形处理器模组中的图形处理器的温度值,所述检测链路包括基板管理控制器的管理接口、所述基板管理控制器与主管理控制器之间的第一链路和所述主管理控制器的管理接口,具体可以包括如下模块:

[0186] 模组检测温度获取模块301,用于获取所述检测链路检测的所述图形处理器模组中的图形处理器的温度值;

[0187] 温度异常条件判断模块302,用于根据所述图形处理器模组中的图形处理器的温度值,判断是否满足温度异常条件;

[0188] 异常检测链路纠错模块303,用于当满足所述温度异常条件时,增大风扇的转速对所述图形处理器模组进行散热,并依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错。

[0189] 在本发明实施例中,所述异常检测链路纠错模块303,包括:

[0190] 第一检测子模块,用于检测所述基板管理控制器的管理接口的工作状态;

[0191] 第一纠错子模块,用于若所述基板管理控制器的管理接口的工作状态异常,则对所述基板管理控制器的管理接口进行纠错;

[0192] 第二检测子模块,用于若所述基板管理控制器的管理接口的工作状态正常,则检测所述基板管理控制器与所述主管理控制器之间的第一链路的工作状态;

[0193] 第二纠错子模块,用于若所述第一链路的工作状态异常,则对所述第一链路进行纠错;

[0194] 第三检测子模块,用于若所述第一链路的工作状态正常,则检测所述主管理控制器的管理接口的工作状态;

[0195] 第三纠错子模块,用于若所述主管理控制器的管理接口的工作状态异常,则对所述主管理控制器的管理接口进行纠错;

[0196] 温度获取子模块,用于在所述基板管理控制器的管理接口的工作状态,所述第一链路的工作状态,所述主管理控制器的管理接口的工作状态均正常后,获取所述图形处理器的温度值;

[0197] 风扇调整子模块,用于若所有图形处理器的温度值均正常,则恢复所述风扇的转速。

[0198] 在本发明实施例中,所述第一纠错子模块,包括:

[0199] 第一异常确定单元,用于若所述基板管理控制器的管理接口未激活或被占用,则确定所述基板管理控制器的管理接口的工作状态异常;

[0200] 第一纠错单元,用于重新启用所述基板管理控制器的管理接口的权限,以对所述基板管理控制器的管理接口进行纠错。

[0201] 在本发明实施例中,所述第二纠错子模块,包括:

[0202] 第二异常确定单元,用于若所述第一链路的网络连接状态异常,则确定所述第一链路的工作状态异常;

[0203] 第二纠错单元,用于重新建立所述第一链路的网络连接,以对所述第一链路进行纠错。

[0204] 在本发明实施例中,所述主管理控制器的管理接口包括第一主管理控制器接口,所述第三纠错子模块,包括:

[0205] 第三异常确定单元,用于若所述第一主管理控制器接口未激活或被占用,则确定所述第一主管理控制器接口的工作状态异常;

[0206] 第三纠错单元,用于对所述主管理控制器对应的管理功能进行重置,以对所述主管理控制器的管理接口进行纠错。

[0207] 在本发明实施例中,所述主管理控制器的管理接口还包括第二主管理控制器接口,所述第三纠错子模块,还包括:

[0208] 第四纠错单元,用于若对所述主管理控制器对应的管理功能进行重置后,所述主管理控制器的管理接口的工作状态还存在异常,则通过I2C命令对所述主管理控制器进行重置,以对所述主管理控制器的管理接口进行纠错。

[0209] 在本发明实施例中,所述装置还包括:

[0210] 纠错事件记录日志模块,用于将依次对所述基板管理控制器的管理接口、所述基板管理控制器与所述主管理控制器之间的第一链路和所述主管理控制器的管理接口进行纠错的事件,记录到日志。

[0211] 对于装置实施例而言,由于其与方法实施例基本相似,所以描述的比较简单,相关之处参见方法实施例的部分说明即可。

[0212] 参照图7,示出了本发明实施例提供的一种电子设备的结构框图。本发明实施例还提供了一种电子设备,包括:存储器、处理器及存储在所述存储器上并可在所述处理器上运行的计算机程序,所述处理器执行所述计算机程序时实现上述的一种服务器的检测链路的纠错方法的步骤。

[0213] 参照图8,示出了本发明实施例提供的一种计算机可读存储介质的结构框图。本发明实施例还提供了一种计算机可读存储介质,所述计算机可读存储介质上存储计算机程序,所述计算机程序被处理器执行时实现上述的一种服务器的检测链路的纠错方法的步骤。

[0214] 本说明书中的各个实施例均采用递进的方式描述,每个实施例重点说明的都是与其他实施例的不同之处,各个实施例之间相同相似的部分互相参见即可。

[0215] 本领域内的技术人员应明白,本发明实施例的实施例可提供为方法、装置、或计算机程序产品。因此,本发明实施例可采用完全硬件实施例、完全软件实施例、或结合软件和硬件方面的实施例的形式。而且,本发明实施例可采用在一个或多个其中包含有计算机可用程序代码的机器可读介质(包括但不限于磁盘存储器、CD-ROM、光学存储器等)上实施的计算机程序产品的形式。

[0216] 本发明实施例是参照根据本发明实施例的方法、终端设备(系统)、和计算机程序产品的流程图和/或方框图来描述的。应理解可由计算机程序指令实现流程图和/或方框图中的每一流程和/或方框、以及流程图和/或方框图中的流程和/或方框的结合。可提供这些计算机程序指令到通用计算机、专用计算机、嵌入式处理机或其他可编程数据处理终端设备的处理器以产生一个机器,使得通过计算机或其他可编程数据处理终端设备的处理器执行的指令产生用于实现在流程图一个流程或多个流程和/或方框图一个方框或多个方框中指定的功能的装置。

[0217] 这些计算机程序指令也可存储在能引导计算机或其他可编程数据处理终端设备以特定方式工作的计算机可读存储器中,使得存储在该计算机可读存储器中的指令产生包括指令装置的制造品,该指令装置实现在流程图一个流程或多个流程和/或方框图一个方框或多个方框中指定的功能。

[0218] 这些计算机程序指令也可装载到计算机或其他可编程数据处理终端设备上,使得在计算机或其他可编程终端设备上执行一系列操作步骤以产生计算机实现的处理,从而在计算机或其他可编程终端设备上执行的指令提供用于实现在流程图一个流程或多个流程和/或方框图一个方框或多个方框中指定的功能的步骤。

[0219] 尽管已描述了本发明实施例的优选实施例,但本领域内的技术人员一旦得知了基本创造性概念,则可对这些实施例做出另外的变更和修改。所以,所附权利要求意欲解释为包括优选实施例以及落入本发明实施例范围的所有变更和修改。

[0220] 最后,还需要说明的是,在本文中,诸如第一和第二等之类的关系术语仅仅用来将一个实体或者操作与另一个实体或操作区分开来,而不一定要求或者暗示这些实体或操作之间存在任何这种实际的关系或者顺序。而且,术语“包括”、“包含”或者其任何其他变体意在涵盖非排他性的包含,从而使得包括一系列要素的过程、方法、物品或者终端设备不仅包括那些要素,而且还包括没有明确列出的其他要素,或者是还包括为这种过程、方法、物品或者终端设备所固有的要素。在没有更多限制的情况下,由语句“包括一个……”限定的要素,并不排除在包括所述要素的过程、方法、物品或者终端设备中还存在另外的相同要素。

[0221] 以上对本发明所提供的一种服务器的检测链路的纠错方法和一种服务器的检测链路的纠错装置,进行了详细介绍,本文中应用了具体个例对本发明的原理及实施方式进行了阐述,以上实施例的说明只是用于帮助理解本发明的方法及其核心思想;同时,对于本领域的一般技术人员,依据本发明的思想,在具体实施方式及应用范围上均会有改变之处,综上所述,本说明书内容不应理解为对本发明的限制。

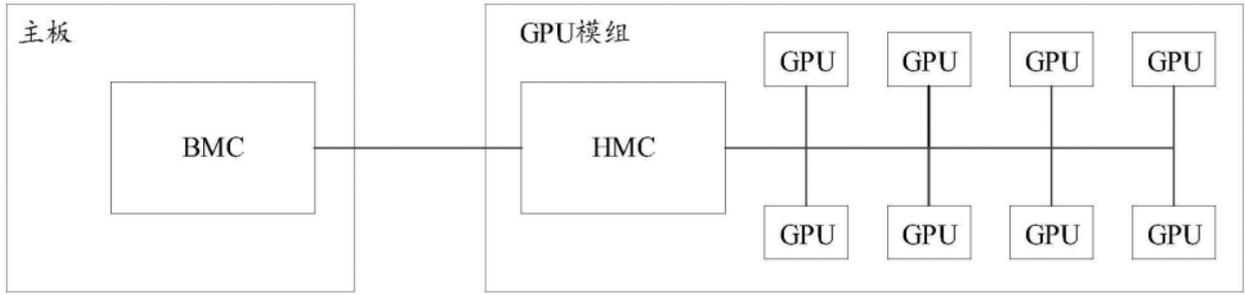


图1

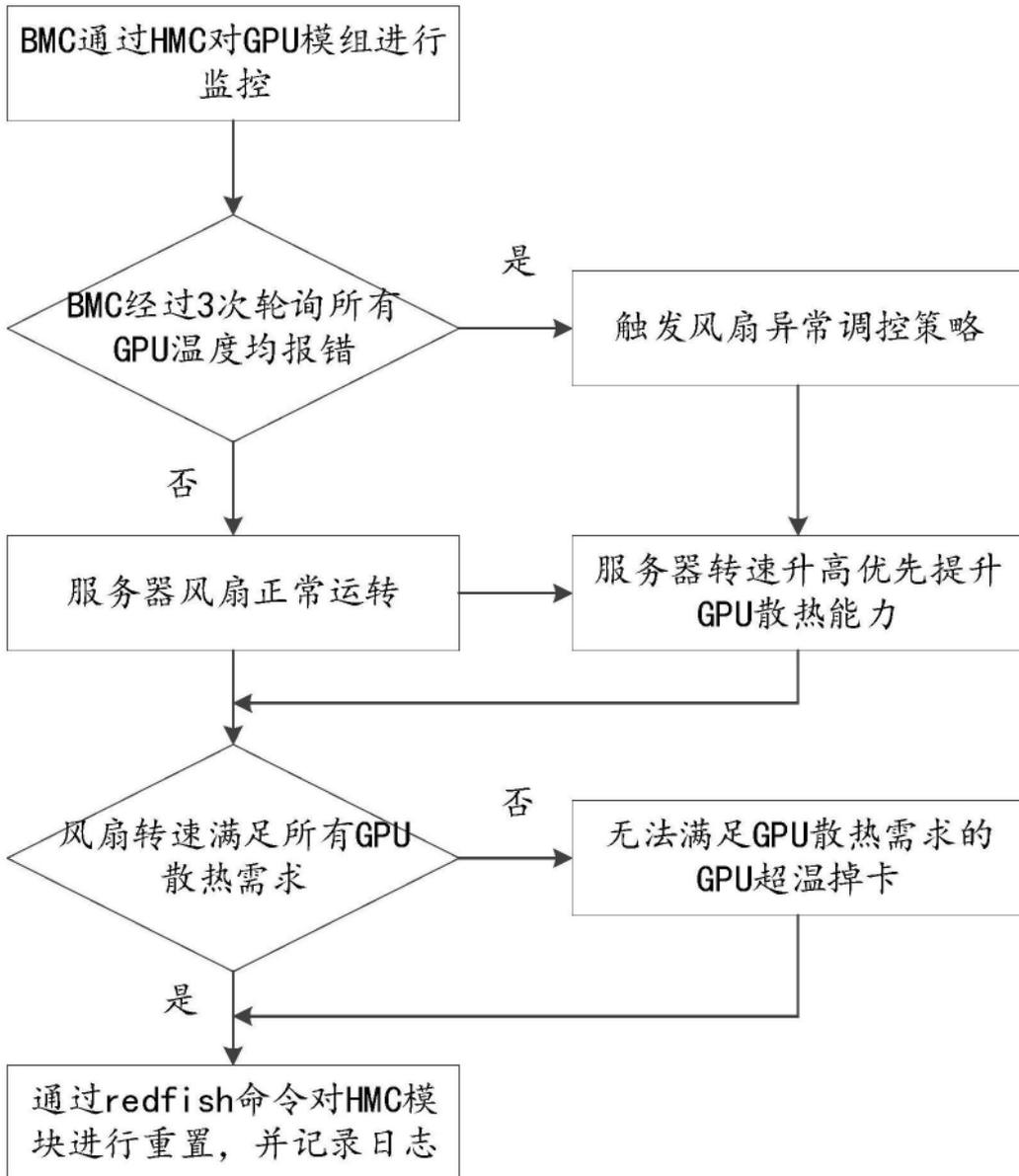


图2

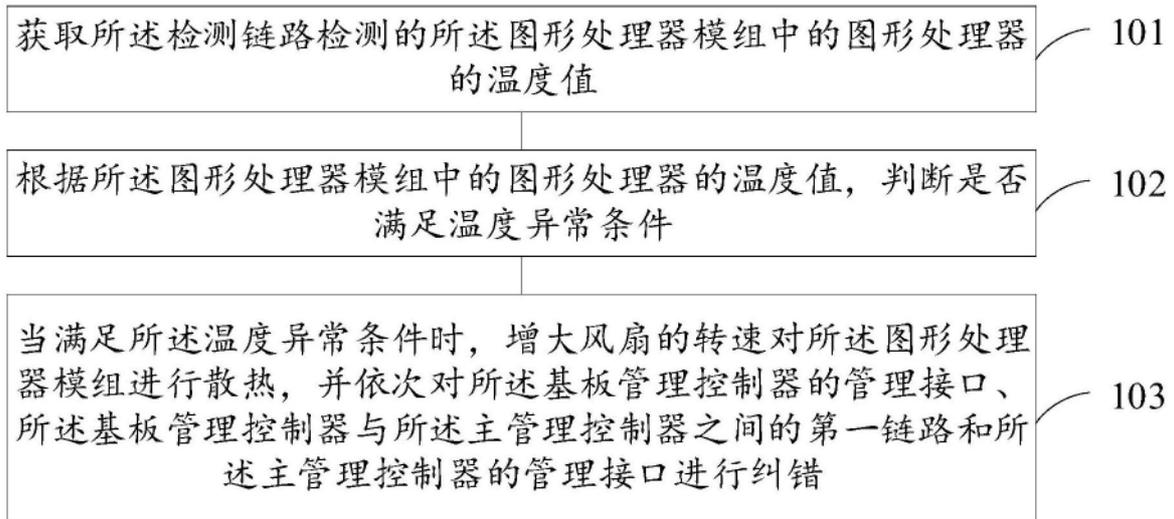


图3

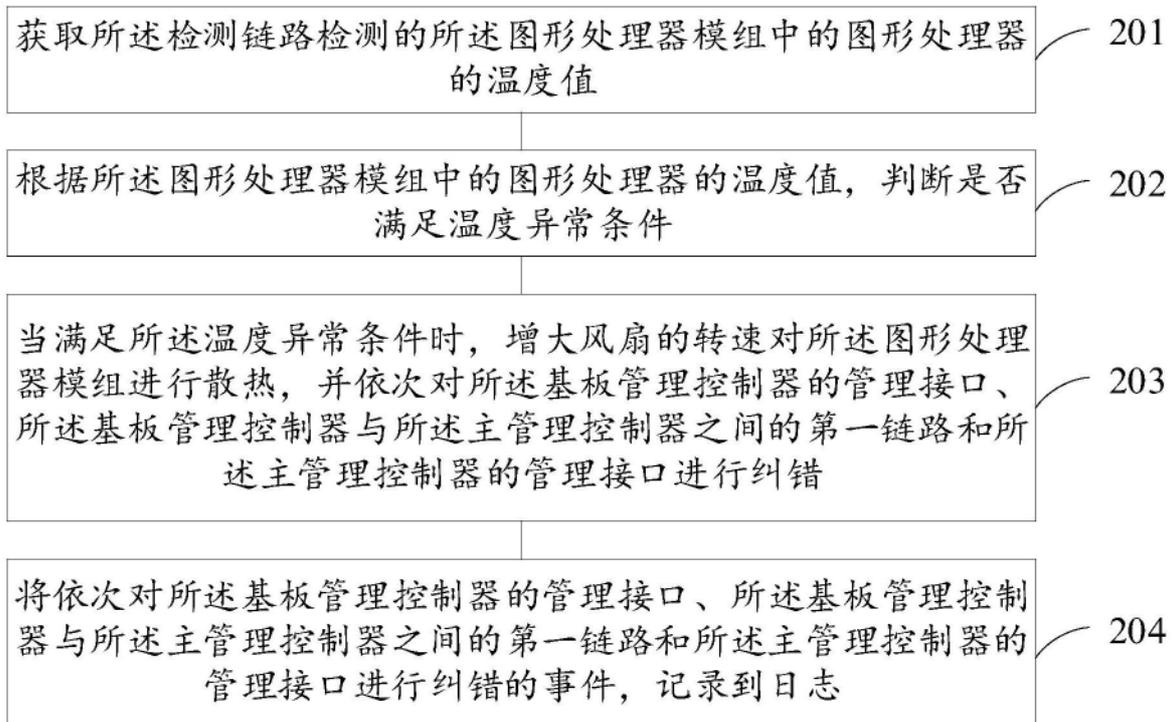


图4

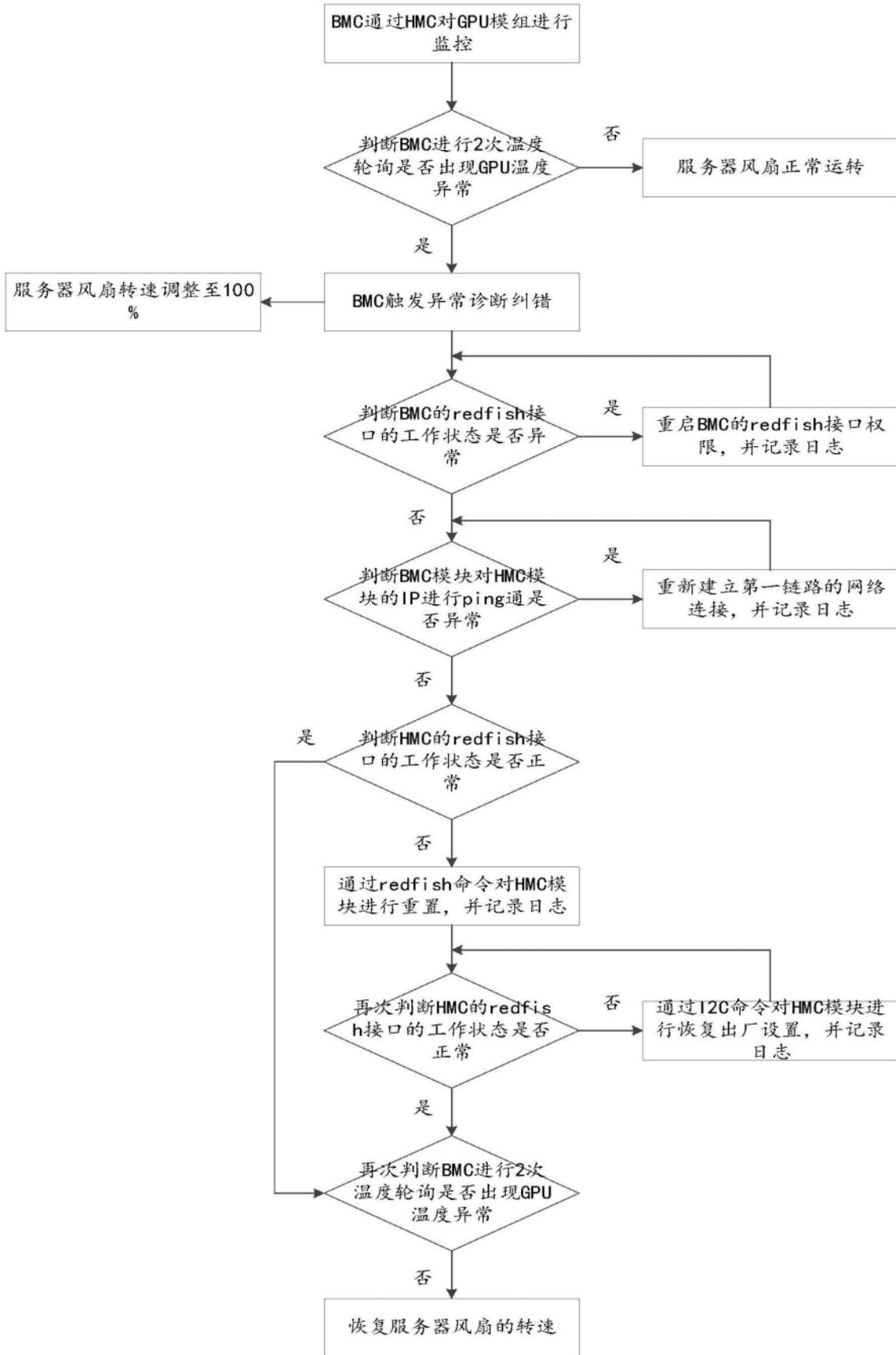


图5



图6

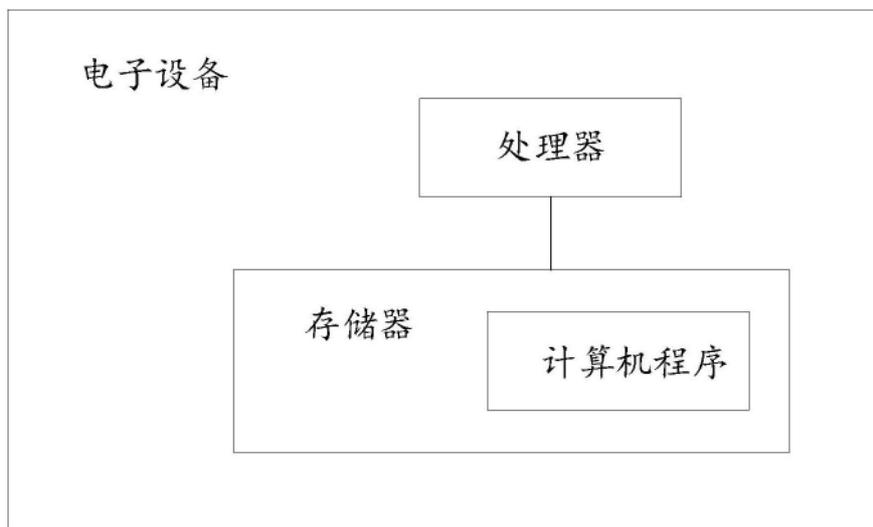


图7

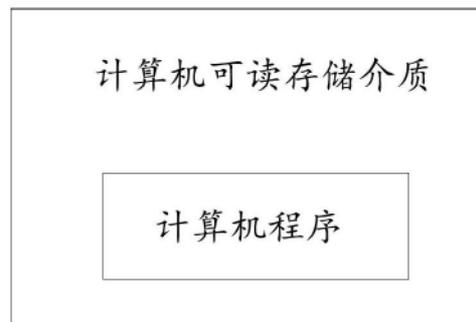


图8