



(12) 发明专利

(10) 授权公告号 CN 109344049 B

(45) 授权公告日 2020.10.16

(21) 申请号 201810942889.3

审查员 孙尧

(22) 申请日 2018.08.17

(65) 同一申请的已公布的文献号

申请公布号 CN 109344049 A

(43) 申请公布日 2019.02.15

(73) 专利权人 华为技术有限公司

地址 518129 广东省深圳市龙岗区坂田华为总部办公楼

(72) 发明人 刘新春 时金魁 许利杰

(74) 专利代理机构 北京三高永信知识产权代理

有限责任公司 11138

代理人 肖庆武

(51) Int. Cl.

G06F 11/36 (2006.01)

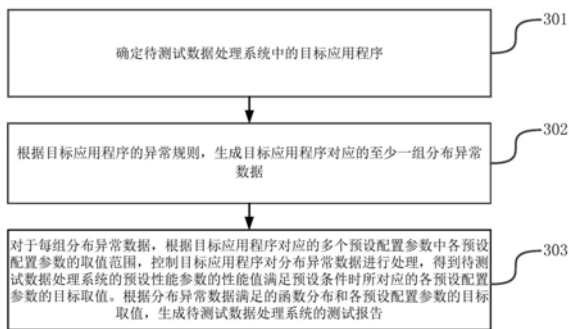
权利要求书4页 说明书12页 附图4页

(54) 发明名称

测试数据处理系统的方法和装置

(57) 摘要

本申请提供了一种测试数据处理系统的方法和装置,属于大数据处理领域。所述方法包括:在对待测试数据处理系统进行测试时,可以确定待测试数据处理系统中的目标应用程序,然后获取目标应用程序的异常规则,基于该异常规则,生成至少一组分布异常数据,然后对于每组分布异常数据,根据目标应用程序对应的多个预配置参数中各配置参数的取值范围,控制目标应用程序对分布异常数据进行处理,得到待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的各预配置参数的目标取值,然后基于分布异常数据的分布和各预配置参数的目标取值,生成测试报告。采用本申请,提供了一种测试数据处理系统的方法。



1. 一种测试数据处理系统的方法,其特征在于,所述方法包括:

确定待测试数据处理系统中的目标应用程序;

根据所述目标应用程序的异常原则,生成所述目标应用程序对应的至少一组分布异常数据;

对于每组分布异常数据,根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值;根据所述分布异常数据满足的函数分布和所述各预设配置参数的目标取值,生成所述待测试数据处理系统的测试报告。

2. 根据权利要求1所述的方法,其特征在于,所述目标应用程序为SQL应用程序、MLib应用程序和Graph应用程序中的一个或多个。

3. 根据权利要求1所述的方法,其特征在于,所述多个预设配置参数之间相互独立;

所述根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值,包括:

根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,确定所述各预设配置参数的取值;

根据所述各预设配置参数的取值和预设的贪心算法,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值。

4. 根据权利要求3所述的方法,其特征在于,所述预设条件为所述目标应用程序运行异常时所得到的性能值,或者所述贪心算法运行结束时,得到的最差性能值。

5. 根据权利要求1所述的方法,其特征在于,所述多个预设配置参数之间非相互独立;

所述根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值,包括:

根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围和所述各预设配置参数的数学函数,确定所述各预设配置参数的取值;

根据所述各预设配置参数的取值,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值;

将确定出的性能值中满足预设条件的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

6. 根据权利要求5所述的方法,其特征在于,所述根据所述各预设配置参数的取值,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值,包括:

根据所述各预设配置参数的取值和预设的统计次数,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值。

7. 根据权利要求5所述的方法,其特征在于,所述预设性能参数为一个;

所述将确定出的性能值中满足预设条件的性能值所对应的所述各预设配置参数的取

值,确定为各预设配置参数的目标取值,包括:

在确定出的性能值中选择最差性能值;

将所述最差性能值所对应的所述各预设配置参数的取值,确定为各预设配置参数的目标取值。

8. 根据权利要求5所述的方法,其特征在于,所述预设性能参数为多个;

所述将确定出的性能值中满足预设条件的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值,包括:

对确定出的性能值中各预设性能参数的性能值进行加权处理,将加权值最大的预设性能参数的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

9. 根据权利要求1所述的方法,其特征在于,所述根据所述分布异常数据满足的函数分布和所述各预设配置参数的目标取值,生成所述待测试数据处理系统的测试报告,包括:

根据所述分布异常数据满足的函数分布、所述各预设配置参数的目标取值和所述目标取值对应的所述待测试数据处理系统的预设性能参数的性能值,生成所述待测试数据处理系统的测试报告。

10. 一种测试数据处理系统的装置,其特征在于,所述装置包括:

确定模块,用于确定待测试数据处理系统中的目标应用程序;

数据生成模块,用于根据所述目标应用程序的异常原则,生成所述目标应用程序对应的至少一组分布异常数据;

处理模块,用于对于每组分布异常数据,根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值;根据所述分布异常数据满足的函数分布和所述各预设配置参数的目标取值,生成所述待测试数据处理系统的测试报告。

11. 根据权利要求10所述的装置,其特征在于,所述目标应用程序为SQL应用程序、MLib应用程序和Graph应用程序中的一个或多个。

12. 根据权利要求10所述的装置,其特征在于,所述多个预设配置参数之间相互独立;

所述处理模块,用于:

根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,确定所述各预设配置参数的取值;

根据所述各预设配置参数的取值和预设的贪心算法,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值。

13. 根据权利要求12所述的装置,其特征在于,所述预设条件为所述目标应用程序运行异常时所得到的性能值,或者所述贪心算法运行结束时,得到的最差性能值。

14. 根据权利要求10所述的装置,其特征在于,所述多个预设配置参数之间非相互独立;

所述处理模块,用于:

根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围和所

述各预设配置参数的数学函数,确定所述各预设配置参数的取值;

根据所述各预设配置参数的取值,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值;

将确定出的性能值中满足预设条件的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

15. 根据权利要求14所述的装置,其特征在于,所述处理模块,用于:

根据所述各预设配置参数的取值和预设的统计次数,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值。

16. 根据权利要求14所述的装置,其特征在于,所述预设性能参数为一个;

所述处理模块,用于:

在确定出的性能值中选择最差性能值;

将所述最差性能值所对应的所述各预设配置参数的取值,确定为各预设配置参数的目标取值。

17. 根据权利要求14所述的装置,其特征在于,所述预设性能参数为多个;

所述处理模块,用于:

对确定出的性能值中各预设性能参数的性能值进行加权处理,将加权值最大的预设性能参数的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

18. 根据权利要求10所述的装置,其特征在于,所述处理模块,用于:

根据所述分布异常数据满足的函数分布、所述各预设配置参数的目标取值和所述目标取值对应的所述待测试数据处理系统的预设性能参数的性能值,生成所述待测试数据处理系统的测试报告。

19. 一种测试数据处理系统的服务器,其特征在于,所述服务器包括处理器和存储器:

所述处理器,用于:

确定待测试数据处理系统中的目标应用程序;

根据所述目标应用程序的异常原则,生成所述目标应用程序对应的至少一组分布异常数据;

对于每组分布异常数据,根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值;根据所述分布异常数据满足的函数分布和所述各预设配置参数的目标取值,生成所述待测试数据处理系统的测试报告。

20. 根据权利要求19所述的服务器,其特征在于,所述多个预设配置参数之间相互独立;

所述处理器,用于:

根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,确定所述各预设配置参数的取值;

根据所述各预设配置参数的取值和预设的贪心算法,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条

件时所对应的所述各预设配置参数的目标取值。

21. 根据权利要求19所述的服务器,其特征在于,所述多个预设配置参数之间非相互独立;

所述处理器,用于:

根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围和所述各预设配置参数的数学函数,确定所述各预设配置参数的取值;

根据所述各预设配置参数的取值,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值;

将确定出的性能值中满足预设条件的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

22. 根据权利要求21所述的服务器,其特征在于,所述处理器,用于:

根据所述各预设配置参数的取值和预设的统计次数,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值。

23. 根据权利要求21所述的服务器,其特征在于,所述预设性能参数为一个;

所述处理器,用于:

在确定出的性能值中选择最差性能值;

将所述最差性能值所对应的所述各预设配置参数的取值,确定为各预设配置参数的目标取值。

24. 根据权利要求21所述的服务器,其特征在于,所述预设性能参数为多个;

所述处理器,用于:

对确定出的性能值中各预设性能参数的性能值进行加权处理,将加权值最大的预设性能参数的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

25. 根据权利要求19所述的服务器,其特征在于,所述处理器,用于:

根据所述分布异常数据满足的函数分布、所述各预设配置参数的目标取值和所述目标取值对应的所述待测试数据处理系统的预设性能参数的性能值,生成所述待测试数据处理系统的测试报告。

26. 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存储有指令,当所述计算机可读存储介质在服务器上运行时,使得所述服务器执行所述权利要求1-9中任一权利要求所述的方法。

测试数据处理系统的方法和装置

技术领域

[0001] 本申请涉及大数据处理领域,特别涉及一种测试数据处理系统的方法和装置。

背景技术

[0002] 近年来,随着互联网、电子商务以及物联网等技术的快速发展,数据的产生速度呈现爆炸性的增长方式的趋势。这些数据具有产生速度快、产生规模大等特点,为了挖掘这些海量数据背后隐藏的巨大商业价值,许多数据处理系统应用而生,例如,Hadoop、Spark、Storm、Flink等,这些数据处理系统分别偏重于不同的处理场景。数据处理系统在处理数据时,经常会出现内存溢出、I/O异常、任务无响应等错误,这些错误会直接导致数据处理系统的任务执行失败。

[0003] 相关技术中,一般是数据处理系统上线后,数据处理系统出现问题,才会对数据处理系统进行分析诊断。

[0004] 这样,由于只能在数据处理系统上线之后,发生问题才会解决问题,然而发生问题后再解决问题,会导致任务处理严重延迟,所以急需提供一种提前测试数据处理系统的方法。

发明内容

[0005] 为了解决相关技术的问题,本发明实施例提供了一种测试数据处理系统的方法和装置。所述技术方案如下:

[0006] 第一方面,提供了一种测试数据处理系统的方法,所述方法包括:

[0007] 确定待测试数据处理系统中的目标应用程序;根据所述目标应用程序的异常原则,生成所述目标应用程序对应的至少一组分布异常数据;对于每组分布异常数据,根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值;根据所述分布异常数据满足的函数分布和所述各预设配置参数的目标取值,生成所述待测试数据处理系统的测试报告。

[0008] 本发明实施例所示的方案,在对待测试数据处理系统进行测试时,可以确定待测试数据处理系统中的目标应用程序,然后获取目标应用程序的异常规则,基于该异常规则,生成至少一组分布异常数据,然后对于每组分布异常数据,根据目标应用程序对应的多个预设配置参数中各配置参数的取值范围,控制目标应用程序对分布异常数据进行处理,得到待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的各预设配置参数的目标取值,然后基于分布异常数据的分布和各预设配置参数的目标取值,生成测试报告。后续技术人员可以基于测试报告,对该数据处理系统进行修复,使其在运行分布异常数据时,也可以正常运行。

[0009] 在一种可能的实施方式中,所述目标应用程序为SQL应用程序、Mlib应用程序和

Graph应用程序中的一个或多个。

[0010] 在一种可能的实施方式中,所述多个预设配置参数之间相互独立;

[0011] 所述根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值,包括:根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,确定所述各预设配置参数的取值;根据所述各预设配置参数的取值和预设的贪心算法,控制所述目标程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值。所述预设条件为所述目标应用程序运行异常时所得到的性能值,或者所述贪心算法运行结束时,得到的最差性能值。

[0012] 本发明实施例所示的方案,多个预设配置参数相互独立,假设有 n 个预设配置参数,第 i 个预设配置参数的 m_i 个取值与预设性能参数的性能值正相关或负相关。

[0013] 服务器可以根据目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,确定各预设配置参数的取值为取值范围对应的端点值,也就是说每个预设配置参数仅有两个取值。例如,某个预设配置参数的取值范围为 $1\sim 10$,该预设配置参数的取值为 1 和 10 。

[0014] 然后服务器可以为目标应用程序选定一组各预设配置参数的取值,然后将分布异常数据输入的目标应用程序中,使目标应用程序处理分布异常数据,得到输出结果,并统计在处理分布异常数据过程中,预设性能参数的性能值,然后可以调整第一个预设配置参数的取值,其余预设配置参数的取值不变,然后重新使用目标应用程序处理异常分布数据,得到输出结果,并统计在处理分布异常数据的过程中,预设性能参数的性能值。如果第二次的性能值比上次差(如CPU占用率比上次的高),则将第一个预设配置参数的第二次的取值作为接下来的测试过程中的固定配置,也就是后续均使用第一个预设配置参数的第二次的取值,如果这次的性能值比上次的好,则将第一个预设配置参数上次的取值,确定为后续的固定配置。并调整第二个预设配置参数的取值,其余预设配置参数的取值不改变,重新使用目标应用程序处理异常分布数据,得到输出结果,并统计在处理分布异常数据的过程中,预设性能参数的性能值,如果这次的性能值比第二次的差,则将第二个预设配置参数这次的取值,确定为后续的固定配置,如果这次的性能值比上次的好,则将第二个预设配置参数上次的取值,确定为后续的固定配置。

[0015] 依上述类推,直到目标应用程序异常时,使用的一组预设配置参数的取值为目标取值。或者,使用贪心算法,将预设配置参数的组合都全部使用时,将最终得到的最差性能值所使用的一组预设配置参数的取值,确定为目标取值。

[0016] 在一种可能的实施方式中,所述多个预设配置参数之间非相互独立;

[0017] 所述根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值,包括:根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围和所述各预设配置参数的数学函数,确定所述各预设配置参数的取值;根据所述各预设配置参数

的取值,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值;将确定出的性能值中满足预设条件的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

[0018] 本发明实施例所示的方案,服务器可以获取多个预设配置参数中各预设配置参数的取值范围和各预设配置参数的函数表达式,然后使用该取值范围和函数表达式,确定各预设配置参数的所有取值。然后服务器可以将不同组合的各预设配置参数的取值,作为目标应用程序的配置参数,处理分布异常数据,并统计每次得到的预设性能参数的性能值。然后服务器将统计到的性能值中满足预设条件的性能值所使用的各预设配置参数的取值,确定为各预设配置参数对应的目标取值。

[0019] 在一种可能的实施方式中,根据所述各预设配置参数的取值和预设的统计次数,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值。

[0020] 在一种可能的实施方式中,所述预设性能参数为一个;在确定出的性能值中选择最差性能值;将所述最差性能值所对应的所述各预设配置参数的取值,确定为各预设配置参数的目标取值。

[0021] 在一种可能的实施方式中,所述预设性能参数为多个;对确定出的性能值中各预设性能参数的性能值进行加权处理,将加权值最大的预设性能参数的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

[0022] 本发明实施例所示的方案,服务器对于一组预设性能参数的取值,服务器使用加权的方式,进行加权处理,在加权处理后,选择加权值最大的一组预设性能参数的性能值所对应的各预设配置参数的取值,确定为各预设配置参数的目标取值。

[0023] 在一种可能的实施方式中,所述根据所述分布异常数据满足的函数分布和所述各预设配置参数的目标取值,生成所述待测试数据处理系统的测试报告,包括:根据所述分布异常数据满足的函数分布、所述各预设配置参数的目标取值和所述目标取值对应的所述待测试数据处理系统的预设性能参数的性能值,生成所述待测试数据处理系统的测试报告。

[0024] 本发明实施例所示的方案,服务器可以获取预设的统计次数,如果各预设配置参数的取值的组合的数目小于或等于预设的统计次数,则可以使用所有的配置参数的组合,进行测试,如果各预设配置参数的取值的组合的数目大于预设的统计次数,则可以在所有的配置参数的组合中,选取预设的统计次数个组合,进行测试,后续在这些测试结果中,选取出目标取值。

[0025] 第二方面,提供了一种测试数据处理系统的服务器,该服务器包括处理器和存储器,所述处理器通过执行指令来实现上述第一方面所提供的测试数据处理系统的方法。

[0026] 第三方面,提供了一种测试数据处理系统的装置,该装置包括一个或多个模块,该一个或多个模块通过执行指令来实现上述第一方面所提供的测试数据处理系统的方法。

[0027] 第四方面,提供了一种计算机可读存储介质,计算机可读存储介质存储有指令,当计算机可读存储介质在服务器上运行时,使得服务器执行上述第一方面所提供的测试数据处理系统的方法。

[0028] 第五方面,提供了一种包含指令的计算机程序产品,当其在服务器上运行时,使得服务器执行上述第一方面所提供的测试数据处理系统的方法。

[0029] 本发明实施例提供的技术方案带来的有益效果至少包括：

[0030] 本发明实施例中，在对待测试数据处理系统进行测试时，可以确定待测试数据处理系统中的目标应用程序，然后获取目标应用程序的异常规则，基于该异常规则，生成至少一组分布异常数据，然后对于每组分布异常数据，根据目标应用程序对应的多个预设配置参数中各配置参数的取值范围，控制目标应用程序对分布异常数据进行处理，得到待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的各预设配置参数的目标取值，然后基于分布异常数据的分布和各预设配置参数的目标取值，生成测试报告。这样，由于待测试数据处理系统在未上线前，就有运行分布异常数据的测试，并且得到测试报告，提供了一种测试数据处理系统的方法，后续技术人员可以基于测试报告对该数据处理系统进行修复，尽可能的防止数据处理系统上线后，运行分布异常数据时，出现错误，也提高了数据处理系统的可靠性。

附图说明

[0031] 图1是本发明实施例提供的一种数据处理系统的结构示意图；

[0032] 图2是本发明实施例提供的一种服务器的结构示意图；

[0033] 图3是本发明实施例提供的一种测试数据处理系统的方法示意图；

[0034] 图4是本发明实施例提供的一种生成分布异常数据的示意图；

[0035] 图5是本发明实施例提供的一种生成分布异常数据的示意图；

[0036] 图6是本发明实施例提供的一种生成分布异常数据的示意图；

[0037] 图7是本发明实施例提供的一种测试数据处理系统的装置结构示意图。

具体实施方式

[0038] 为使本申请的目的、技术方案和优点更加清楚，下面将结合附图对本申请实施方式作进一步地详细描述。

[0039] 为了便于对本发明实施例的理解，下面首先介绍本发明实施例涉及的系统架构、以及所涉及名词的概念。

[0040] 本发明实施例可以适用于数据处理系统，该数据处理系统可以是大数据处理系统，如图1所示，如面向Spark大数据处理系统（也可以称为是Spark集群），该系统可以部署在多个计算机节点上，使用多个计算机节点并行处理大规模数据（可以简称为大数据）。

[0041] 异常：非正常现象，数据处理系统运行过程中可能出现的内存溢出、IO失败、任务运行超时等错误，这些错误都会导致数据处理系统异常。

[0042] 异常数据：当某组数据出现以下（数据量大、数据倾斜、数据稀疏、数据维度高和数据分布异常）一种或多种情况时，将该组数据称为异常数据。

[0043] 在进行实施前，首先介绍一下本发明实施例的应用场景，在数据处理系统中通常会包括不同的应用程序，每个应用程序处理的数据不一样。数据处理系统中的应用程序在处理数据时，通常会出现内存溢出、I/O异常、任务无响应的错误，所以一般在数据处理系统上线前，需要对数据处理系统进行检测，发现是否存在潜在的问题，在数据处理系统正式上线前，尽可能的将可能出现的问题解决掉。

[0044] 本发明实施例提供了一种数据处理系统的方法，该方法的执行主体可以是服务

器。

[0045] 图2示出了本发明实施例中服务器的结构框图,该服务器至少可以包括接收器201、处理器202、存储器203和发射器204。其中,接收器201可以用于实现数据的接收,具体可以用于数据的接收,发射器204可以用于数据的发送,具体可以用于处理结果的发送,存储器203可以用于存储软件程序以及模块,处理器202通过运行存储在存储器203中的软件程序以及模块,从而执行各种功能应用以及数据处理。存储器203主要包括存储程序区和存储数据区,其中,存储程序区可存储操作系统、至少一个功能所需的应用程序等;存储数据区可存储根据服务器的使用所创建的数据等。此外,存储器203可以包括高速随机存取存储器,还可以包括非易失性存储器,例如至少一个磁盘存储器件、闪存器件、或其他易失性固态存储器件。相应地,存储器203还可以包括存储器控制器,以提供处理器202、接收器201和发射器204对存储器203的访问。处理器202是服务器的控制中心,利用各种接口和线路连接整个服务器的各个部分,通过运行或执行存储在存储器203内的软件程序和/或模块,以及调用存储在存储器203内的数据,执行服务器的各种功能和处理数据,从而对服务器进行整体监控。

[0046] 可选的,处理器202可包括一个或多个处理核心;优选的,处理器202可集成应用处理器和调制解调处理器,其中,应用处理器主要处理操作系统、用户界面和应用程序等,调制解调处理器主要处理无线通信。可以理解的是,上述调制解调处理器也可以不集成到处理器202中。

[0047] 本发明实施例,提供了一种测试数据处理系统的方法,如图3所示,该方法的执行步骤流程可以如下:

[0048] 步骤301,确定待测试数据处理系统中的目标应用程序。

[0049] 其中,待测试数据处理系统为任一数据处理系统,用于进行大数据处理。目标应用程序为数据处理系统中常用的用于处理大数据的应用程序,一般是由技术人员预设,例如,SQL (Structure Query language,结构化查询语言) 应用程序、Machine Learning应用程序(后续可以简称为MLib应用程序)和Graph应用程序等。

[0050] 在实施中,服务器中存储有待测试数据处理系统中的目标应用程序的列表,技术人员可以通过发送指令的方式,通知服务器确定待测试数据处理系统中的目标应用程序。

[0051] 步骤302,根据目标应用程序的异常规则,生成目标应用程序对应的至少一组分布异常数据。

[0052] 其中,不同的应用程序有不同的异常规则,异常规则可以预设,并且与应用程序对应存储在服务器中。

[0053] 在实施中,服务器在确定出目标应用程序后,可以获取目标应用程序的异常规则,然后基于目标应用程序的异常规则,生成目标应用程序对应的至少一组分布异常数据,如果是多组分布异常数据,每组分布异常数据所满足的分布不相同。

[0054] 需要说明的是,如果目标应用程序是多个,服务器中可以存储有应用程序与异常原则的对应关系,可以用于查询目标应用程序对应的异常原则。

[0055] 需要说明的是,异常规则的确定方式可以如下:

[0056] 目标应用程序是SQL应用程序,SQL应用程序中常用的基本操作语句在处理key/value对时,计算复杂度与key的分布相关,当key分布不均匀时,会影响SQL应用程序查询的

复杂度,所以SQL应用程序生成需要考虑数据倾斜带来的key值分布不均匀的问题,所以异常规则为数据量大、数据倾斜。

[0057] 目标应用程序是Graph应用程序,Graph应用程序通常需要迭代计算,当出现数据倾斜时,会出现单个定点压力过大,容易出现内存溢出等问题,又如TriangleCount在提供的数据中有重复的边的情况下,会使计算结果不正确,所以异常规则为数据量大、数据稀疏、数据分布异常等

[0058] 目标应用程序是MLib应用程序,MLib应用程序中有些要进行迭代计算,有些需要生成宽度优先树,例如,Logistics Regression及K-means等应用程序的特征是迭代计算,且以矩阵作为输入数据,因此需要考虑矩阵特征对应用程序带来的影响,同时,矩阵特征会在迭代中表现的更为明显。又如,Random Forest等应用程序在计算过程中,需要保存宽度优先树,当数据维度过高时,每个节点存储的信息也会相应增多,容易出现内存溢出等问题。因此,MLib应用程序的数据生成需要考虑数据维度、数据稀疏性等带来的内存占用问题,所以异常规则为数据量大、数据稀疏、数据维度高、数据分布异常等。

[0059] 步骤303,对于每组分布异常数据,根据目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制目标应用程序对分布异常数据进行处理,得到待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的各预设配置参数的目标取值。根据分布异常数据满足的函数分布和各预设配置参数的目标取值,生成待测试数据处理系统的测试报告。

[0060] 其中,预设配置参数可以预设,并且与目标应用程序对应存储至服务器中,每个预设配置参数都对应取值范围。例如,对于SQL应用程序,预设配置参数可以为处理数据的属性个数(处理数据为涉及年龄、性别、身高的数据,属性可以是年龄、性别、身高,预设配置参数的数目为3),对于Graph应用程序,预设配置参数可以为顶点个数等,对于MLib应用程序,预设配置参数可以为各神经网络层的参数。预设性能参数指中央处理器(Central Processing Unit,CPU)占用率、内存占用率等。预设条件可以预设,并且存储至服务器中。

[0061] 在实施中,对于至少一组分布异常数据中的每组分布异常数据,服务器在针对目标应用程序,生成分布异常数据时,服务器可以获取对应目标应用程序存储的多个预设配置参数,并获取每个预设配置参数的取值范围,然后根据各预设配置参数的取值范围,确定各预设配置参数的多个取值,将各预设配置参数的多个取值,组成多组取值,每组取值中包括各预设配置参数的一个取值。在这多组取值下,将分布异常数据输入到目标应用程序,进行处理,并统计待测试数据处理系统的预设性能参数的性能值,确定预设性能参数的性能值中,满足预设条件的性能值,确定得到该性能值所对应的各预设配置参数的取值,分别为各预设配置参数的目标取值。服务器在得到目标取值后,可以确定分布异常数据满足的函数分布(可选的,可以从函数分布对应的属性信息中获取),然后获取测试报告的模板,将函数分布和目标取值填入测试报告的模板中,生成一个测试报告,在该测试报告中,函数分布与目标取值相对应,也就是说明,在这种函数分布的异常数据中,各预设配置参数的取值分别为多少。

[0062] 这样,如果目标应用程序是一个,分布异常数据为三组,则一般可以得到三份测试报告。

[0063] 可选的,对于目标应用程序,多个预设配置参数之间可以相互独立,相应的步骤

303中确定目标取值的处理可以如下：

[0064] 根据目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围，确定各预设配置参数的取值；根据各预设配置参数的取值和预设的贪心算法，控制目标程序对分布异常数据进行处理，得到待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的各预设配置参数的目标取值。

[0065] 在实施中，多个预设配置参数相互独立，假设有 n 个预设配置参数，第 i 个预设配置参数的 m_i 个取值与预设性能参数的性能值正相关或负相关。

[0066] 服务器可以根据目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围，确定各预设配置参数的取值为取值范围对应的端点值，也就是说每个预设配置参数仅有两个取值。例如，某个预设配置参数的取值范围为1~10，该预设配置参数的取值为1和10。

[0067] 然后服务器可以为目标应用程序选定一组各预设配置参数的取值，然后将分布异常数据输入的目标应用程序中，使目标应用程序处理分布异常数据，得到输出结果，并统计在处理分布异常数据过程中，预设性能参数的性能值，然后可以调整第一个预设配置参数的取值，其余预设配置参数的取值不变，然后重新使用目标应用程序处理异常分布数据，得到输出结果，并统计在处理分布异常数据的过程中，预设性能参数的性能值。如果第二次的性能值比上次差（如CPU占用率比上次的高），则将第一个预设配置参数的第二次的取值作为接下来的测试过程中的固定配置，也就是后续均使用第一个预设配置参数的第二次的取值，如果这次的性能值比上次的好，则将第一个预设配置参数上次的取值，确定为后续的固定配置。并调整第二个预设配置参数的取值，其余预设配置参数的取值不改变，重新使用目标应用程序处理异常分布数据，得到输出结果，并统计在处理分布异常数据的过程中，预设性能参数的性能值，如果这次的性能值比第二次的差，则将第二个预设配置参数这次的取值，确定为后续的固定配置，如果这次的性能值比上次的好，则将第二个预设配置参数上次的取值，确定为后续的固定配置（该过程可以称为是贪心算法）。

[0068] 依上述类推，直到目标应用程序异常时，使用的一组预设配置参数的取值为目标取值。或者，使用贪心算法，将预设配置参数的组合都全部使用时，将最终得到的最差性能值所使用的一组预设配置参数的取值，确定为目标取值。

[0069] 例如，预设性能参数为CPU占用率，一共有3个预设配置参数，分别为A/B/C，A的取值范围为1~10，B的取值范围为2~10，C的取值范围为1~20，那么A的取值为1和10，B的取值为2和10，C的取值为1和20。A、B、C的取值分别为1、2、1，得到预设性能参数的性能值为60%，然后调整A的取值为10，A、B、C的取值分别为10、2、1，得到预设性能参数的性能值为70%，则下次调整B的取值为10，A、B、C的取值分别为10、10、1，得到预设性能参数的性能值为80%，由于70%小于80%，所以调整C的取值为20，A、B、C的取值分别为10、10、20，得到预设性能参数的性能值为50%，由于50%小于80%，所以A、B、C的目标取值为10、10、1。

[0070] 再例如，预设性能参数为CPU占用率，一共有3个预设配置参数，分别为A/B/C，A的取值范围为1~10，B的取值范围为2~10，C的取值范围为1~20，那么A的取值为1和10，B的取值为2和10，C的取值为1和20。A、B、C的取值分别为1、2、1，得到预设性能参数的性能值为60%，然后调整A的取值为10，A、B、C的取值分别为10、2、1，得到预设性能参数的性能值为70%，则下次调整B的取值为10，A、B、C的取值分别为10、10、1，在A、B、C的取值分别为10、10、

1时,目标应用程序运行异常,A/B/C的目标取值分别为10、10、1。

[0071] 另外,如果预设性能参数为多个,可以采用加权的方式,确定最差的性能值,最差的性能值对应的加权值最大。例如,预设性能参数为CPU占用率和内存占用率,加权表达式为: $Y=a*x+b*y$,其中,a、b是加权系数,可以是0.6、0.4等,x为CPU占用率,y为内存占用率。

[0072] 可选的,对于目标应用程序,多个预设配置参数之间可以非相互独立,相应的步骤303中确定目标取值的处理可以如下:

[0073] 根据目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围和各预设配置参数的函数表达式,确定各预设配置参数的取值;根据各预设配置参数的取值,确定目标应用程序在对分布异常数据进行处理时,待测试数据处理系统的预设性能参数的性能值;将得到确定出的性能值中最差的性能值所使用的各预设配置参数的取值,确定为各预设配置参数的目标取值。

[0074] 其中,各预设配置参数的函数表达式一般不相同,可以是不同预设配置参数对应的函数表达式中的系数不相同。

[0075] 在实施中,服务器可以获取多个预设配置参数中各预设配置参数的取值范围和各预设配置参数的函数表达式,然后使用该取值范围和函数表达式,确定各预设配置参数的所有取值。例如,有3个预设配置参数A\B\C,取值范围分别是0~150,假设A的函数表达式 $F(n)=0.5x2^n+1$,那么A可以取1.5、2、3、5、9、17、33、65、129,B的函数表达式 $F(n)=3*2^n+3$,那么B可以取6、9、15、27、51、99、C的函数 $F(n)=9*2^n+3$,C可以取12、21、39、75,这样,3个配置参数的取值组合的数目为 $9*6*4$,n为大于或等于零的正整数。

[0076] 然后服务器可以将不同组合的各预设配置参数的取值,作为目标应用程序的配置参数,处理分布异常数据,并统计每次得到的预设性能参数的性能值。然后服务器将统计到的性能值中满足预设条件的性能值所使用的各预设配置参数的取值,确定为各预设配置参数对应的目标取值。

[0077] 可选的,预设性能参数为一个时,可以是内存占用率或CPU占用率,在确定目标取值时,方式可以如下:

[0078] 在确定出的性能值中选择最差性能值;将最差性能值所对应的各预设配置参数的取值,确定为各预设配置参数的目标取值。

[0079] 其中,最差性能值用于表示内存占用率最高或者CPU占用率最高。

[0080] 在实施中,服务器可以在确定出的性能值中,选择最差性能值,然后确定得到最差性能值所使用的各预设配置参数的取值,将各预设配置参数的取值,确定为目标取值。

[0081] 例如,预设性能参数为内存占用率,可以在确定出的内存占用率中,最大的内存占用率为95%,将95%所使用的各预设配置参数的取值,确定为目标取值。

[0082] 可选的,预设性能参数为多个时,可以是内存占用率和CPU占用率,在确定目标取值时,方式可以如下:

[0083] 对确定出的性能值中各预设性能参数的性能值进行加权处理,将加权值最大的预设性能参数的性能值所对应的各预设配置参数的取值,确定为各预设配置参数的目标取值。

[0084] 在实施中,服务器对于一组预设性能参数的取值,服务器使用加权的方式,进行加权处理,在加权处理后,选择加权值最大的一组预设性能参数的性能值所对应的各预设配

置参数的取值,确定为各预设配置参数的目标取值。例如,加权表达式为: $Y=a*x+b*y$,其中, a 、 b 是加权系数,可以是0.6、0.4等, x 为CPU占用率, y 为内存占用率。

[0085] 需要说明的是,如果加权值相同的有多个,可以记录多组目标取值。

[0086] 可选的,为了降低计算量,还可以控制使用配置参数的组合次数,相应的处理可以如下:

[0087] 根据各预设配置参数的取值和预设的统计次数,确定目标应用程序在对分布异常数据进行处理时,待测试数据处理系统的预设性能参数的性能值。

[0088] 其中,预设的统计次数可以预设,并且存储在服务器中,如200次,这样,可以测试的最大次数为200次。

[0089] 在实施中,服务器可以获取预设的统计次数,如果各预设配置参数的取值的组合的数目小于或等于预设的统计次数,则可以使用所有的配置参数的组合,进行测试,如果各预设配置参数的取值的组合的数目大于预设的统计次数,则可以在所有的配置参数的组合中,选取预设的统计次数个组合,进行测试,后续在这些测试结果中,选取出目标取值。详细过程已经在前面叙述,此处不再赘述。

[0090] 需要说明的是,在这种方式下,虽然不一定能得到最佳的目标取值,但是可以得到相对较好的目标取值,但是可以节约处理资源。

[0091] 可选的,在测试报告中,还可以包括性能值,相应的步骤204的处理可以如下:

[0092] 根据分布异常数据满足的函数分布、各预设配置参数的目标取值和目标取值对应的待测试数据处理系统的预设性能参数的性能值,生成待测试数据处理系统的测试报告。

[0093] 在实施中,服务器在得到目标取值后,可以确定分布异常数据满足的函数分布,并且获取使用目标取值时,得到的待测试数据处理系统的预设性能参数的性能值,然后获取测试报告的模板,将函数分布、目标取值和性能值填入测试报告的模板中,生成一个测试报告。

[0094] 另外,在对分布异常数据进行处理时,还可以记录处理时长,在生成测试报告时,将处理时长也填入测试报告模板中。

[0095] 需要说明的是,服务器在控制将分布异常数据输入到目标应用程序时,同时可以记录开始时刻,在目标应用程序处理完分布异常数据时,可以记录结束时刻,将结束时刻减去开始时刻,即为目标应用程序处理分布异常数据的处理时长。

[0096] 这样,在得到待测试数据处理系统的测试报告后,服务器可以将测试报告发送至技术人员所使用的终端,技术人员可以查看测试报告,对该数据处理系统,进行修复,使其既可以在分布异常数据下使用,也可以在分布正常数据下使用。

[0097] 还需要说明的是,上述提到的CPU占用率为在目标应用程序在处理分布异常数据过程中的最大CPU占用率,内存占用率为目标应用程序在处理分布异常数据过程中的最大内存占用率。

[0098] 另外,本发明实施例中,还给出了针对不同的目标应用程序,分布异常数据的生成方式:

[0099] 针对SQL应用程序,如图4所示,可以生成满足异常分布的分布异常数据,例如,可以生成满足Zif分布、泊松分布和高斯分布的数据。可以是生成倾斜数据(例如,单个key多次出现等),数据大小异常(一行特别长,value值过大)等。

[0100] 针对Graph应用程序,如图5所示,可以使用泊松分布生成顶点离散的图,还可以使用Zipf分布生成顶点度异常的稀疏图。

[0101] 针对MLib应用程序,如图6所示,可以随机合成不同维度、稀疏度、异常分布(如高斯分布、伽马分布、泊松分布、指数分布、Zipf分布及其混合)等的分布异常数据。可以是获取预先存储的随机数据,然后确定维度、实例数据设置,并且确定分布的类型(分布的类型有高斯分布、泊松分布等),将随机数据、维度和分布的类型输入到数据生成应用程序中,即可输出分布异常数据。

[0102] 本发明实施例中,在对待测试数据处理系统进行测试时,可以确定待测试数据处理系统中的目标应用程序,然后获取目标应用程序的异常规则,基于该异常规则,生成至少一组分布异常数据,然后对于每组分布异常数据,根据目标应用程序对应的多个预设配置参数中各配置参数的取值范围,控制目标应用程序对分布异常数据进行处理,得到待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的各预设配置参数的目标取值,然后基于分布异常数据的分布和各预设配置参数的目标取值,生成测试报告。这样,由于待测试数据处理系统在未上线前,就有运行分布异常数据的测试,并且得到测试报告,提供了一种测试数据处理系统的方法,后续技术人员可以基于测试报告对该数据处理系统进行修复,尽可能的防止数据处理系统上线后,运行分布异常数据时,出现错误,也提高了数据处理系统的可靠性。

[0103] 图7是本发明实施例提供的测试数据处理系统的装置的结构图。该装置可以通过软件、硬件或者两者的结合实现成为服务器中的部分或者全部。本发明实施例提供的服务器可以实现本发明实施例图3所述的流程,该装置包括:确定模块710、数据生成模块720和处理模块730,其中:

[0104] 确定模块710,用于确定待测试数据处理系统中的目标应用程序;具体可以实现上述步骤301中的确定功能,以及其它隐含步骤;

[0105] 数据生成模块720,用于根据所述目标应用程序的异常原则,生成所述目标应用程序对应的至少一组分布异常数据;具体可以实现上述步骤302中的数据生成,以及其它隐含步骤;

[0106] 处理模块730,用于对于每组分布异常数据,根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,控制所述目标应用程序对所述分布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值;根据所述分布异常数据满足的函数分布和所述各预设配置参数的目标取值,生成所述待测试数据处理系统的测试报告。具体可以实现上述步骤303中的处理功能,以及其它隐含步骤。

[0107] 可选的,所述目标应用程序为SQL应用程序、MLib应用程序和Graph应用程序中的一个或多个。

[0108] 可选的,所述多个预设配置参数之间相互独立;

[0109] 所述处理模块730,用于:

[0110] 根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围,确定所述各预设配置参数的取值;

[0111] 根据所述各预设配置参数的取值和预设的贪心算法,控制所述目标程序对所述分

布异常数据进行处理,得到所述待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的所述各预设配置参数的目标取值。

[0112] 可选的,所述预设条件为所述目标应用程序运行异常时所得到的性能值,或者所述贪心算法运行结束时,得到的最差性能值。

[0113] 可选的,所述多个预设配置参数之间非相互独立;

[0114] 所述处理模块730,用于:

[0115] 根据所述目标应用程序对应的多个预设配置参数中各预设配置参数的取值范围和所述各预设配置参数的数学函数,确定所述各预设配置参数的取值;

[0116] 根据所述各预设配置参数的取值,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值;

[0117] 将确定出的性能值中满足预设条件的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

[0118] 可选的,所述处理模块730,用于:

[0119] 根据所述各预设配置参数的取值和预设的统计次数,确定所述目标应用程序在对所述分布异常数据进行处理时,所述待测试数据处理系统的预设性能参数的性能值。

[0120] 可选的,所述预设性能参数为一个;

[0121] 所述处理模块730,用于:

[0122] 在确定出的性能值中选择最差性能值;

[0123] 将所述最差性能值所对应的所述各预设配置参数的取值,确定为各预设配置参数的目标取值。

[0124] 可选的,所述预设性能参数为多个;

[0125] 所述处理模块730,用于:

[0126] 对确定出的性能值中各预设性能参数的性能值进行加权处理,将加权值最大的预设性能参数的性能值所对应的所述各预设配置参数的取值,确定为所述各预设配置参数的目标取值。

[0127] 可选的,所述处理模块730,用于:

[0128] 根据所述分布异常数据满足的函数分布、所述各预设配置参数的目标取值和所述目标取值对应的所述待测试数据处理系统的预设性能参数的性能值,生成所述待测试数据处理系统的测试报告。

[0129] 本发明实施例中,在对待测试数据处理系统进行测试时,可以确定待测试数据处理系统中的目标应用程序,然后获取目标应用程序的异常规则,基于该异常规则,生成至少一组分布异常数据,然后对于每组分布异常数据,根据目标应用程序对应的多个预设配置参数中各配置参数的取值范围,控制目标应用程序对分布异常数据进行处理,得到待测试数据处理系统的预设性能参数的性能值满足预设条件时所对应的各预设配置参数的目标取值,然后基于分布异常数据的分布和各预设配置参数的目标取值,生成测试报告。这样,由于待测试数据处理系统在未上线前,就有运行分布异常数据的测试,并且得到测试报告,提供了一种测试数据处理系统的方法,后续技术人员可以基于测试报告对该数据处理系统进行修复,尽可能的防止数据处理系统上线后,运行分布异常数据时,出现错误。

[0130] 需要说明的是:上述实施例提供的测试数据处理系统的装置在测试数据处理系统

时,仅以上述各功能模块的划分进行举例说明,实际应用中,可以根据需要而将上述功能分配由不同的功能模块完成,即将装置的内部结构划分成不同的功能模块,以完成以上描述的全部或者部分功能。另外,上述实施例提供的测试数据处理系统的装置与测试数据处理系统的方法实施例属于同一构思,其具体实现过程详见方法实施例,这里不再赘述。

[0131] 在上述实施例中,可以全部或部分地通过软件、硬件、固件或者其任意组合来实现,当使用软件实现时,可以全部或部分地以计算机程序产品的形式实现。所述计算机程序产品包括一个或多个计算机指令,在服务器或终端上加载和执行所述计算机程序指令时,全部或部分地产生按照本发明实施例所述的流程或功能。所述计算机指令可以存储在计算机可读存储介质中,或者从一个计算机可读存储介质向另一个计算机可读存储介质传输,例如,所述计算机指令可以从一个网站站点、计算机、服务器或数据中心通过有线(例如同轴光缆、光纤、数字用户线)或无线(例如红外、无线、微波等)方式向另一个网站站点、计算机、服务器或数据中心进行传输。所述计算机可读存储介质可以是服务器或终端能够存取的任何可用介质或者是包含一个或多个可用介质集成的服务器、数据中心等数据存储设备。所述可用介质可以是磁性介质(如软盘、硬盘和磁带等),也可以是光介质(如数字视盘(Digital Video Disk,DVD)等),或者半导体介质(如固态硬盘等)。

[0132] 以上所述仅为本申请的一个实施例,并不用以限制本申请,凡在本申请的精神和原则之内,所作的任何修改、等同替换、改进等,均应包含在本申请的保护范围之内。

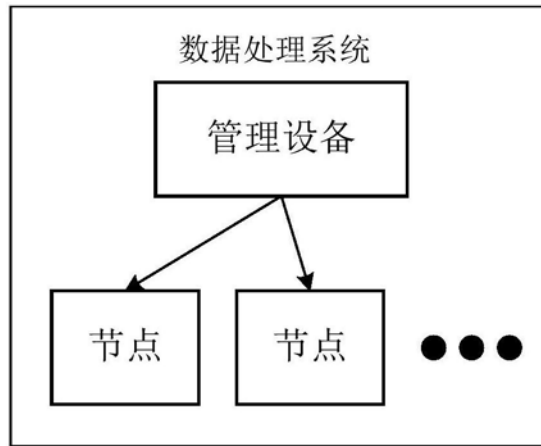


图1

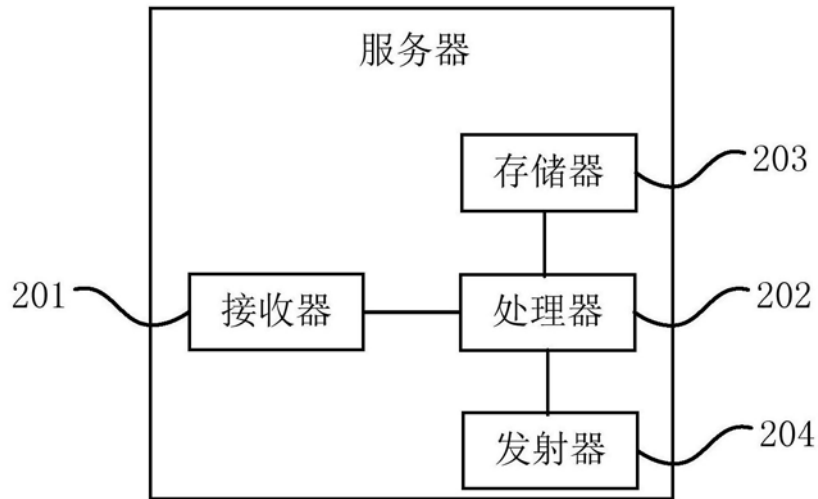


图2

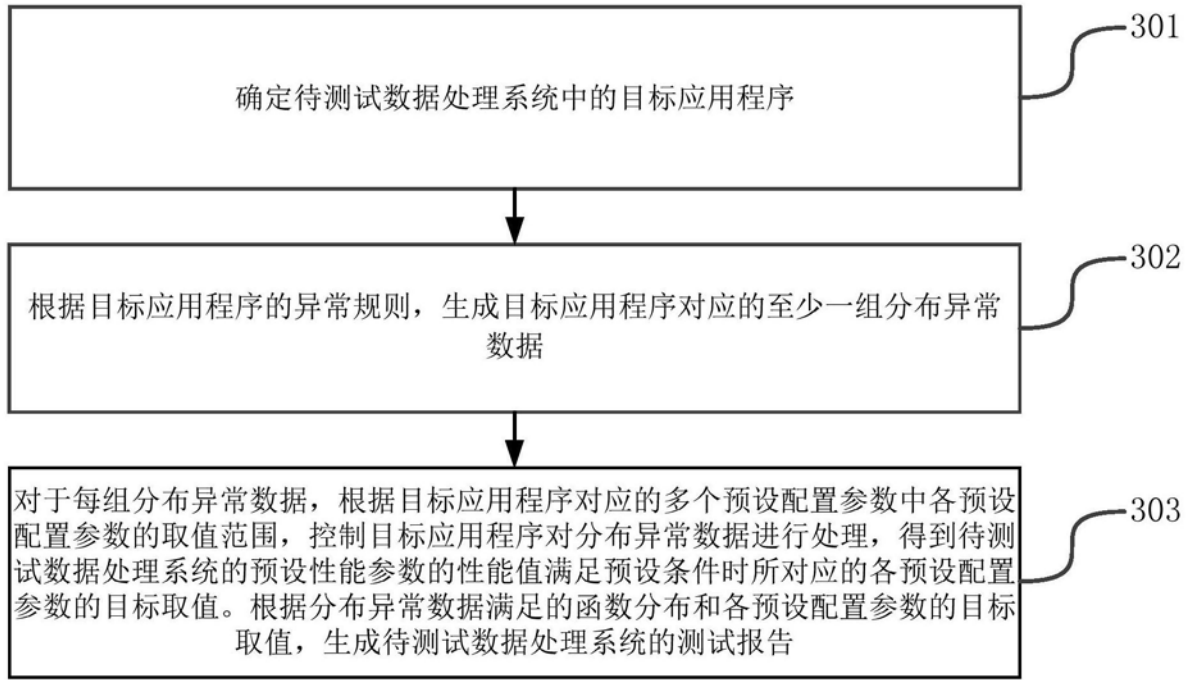


图3

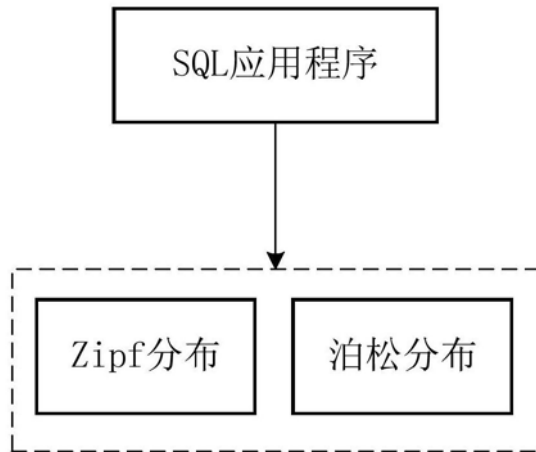


图4

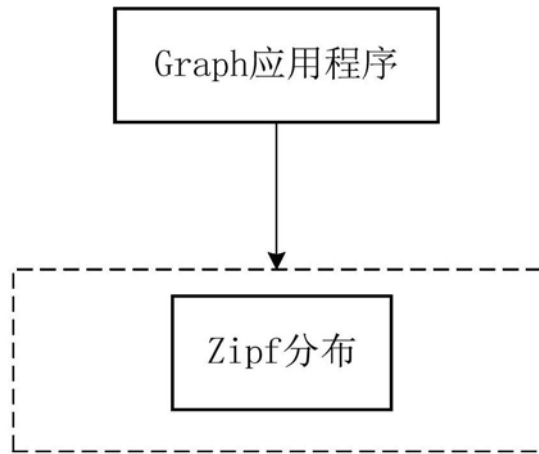


图5

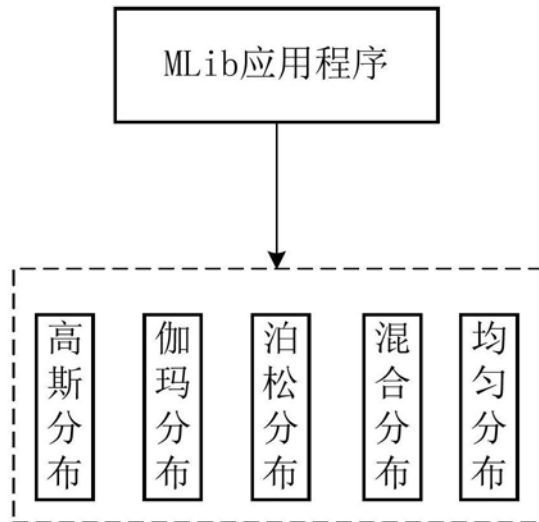


图6

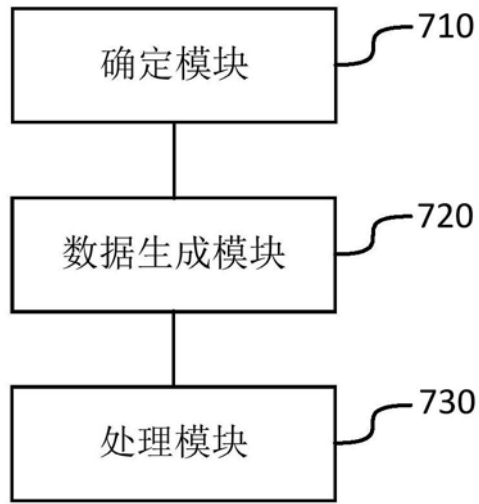


图7