



(12) 发明专利

(10) 授权公告号 CN 108733508 B

(45) 授权公告日 2022.03.11

(21) 申请号 201710250379.5

US 2004172512 A1,2004.09.02

(22) 申请日 2017.04.17

US 2010269168 A1,2010.10.21

(65) 同一申请的已公布的文献号

CN 102573053 A,2012.07.11

申请公布号 CN 108733508 A

CN 103744620 A,2014.04.23

(43) 申请公布日 2018.11.02

CN 102414673 A,2012.04.11

CN 105512163 A,2016.04.20

(73) 专利权人 伊姆西IP控股有限责任公司

US 2005044446 A1,2005.02.24

地址 美国马萨诸塞州

US 9424074 B1,2016.08.23

(72) 发明人 张宇霆 高雪东 林超 彭飞

陈志佳等.云训练中基于自适应副本策略的容错研究.《微电子学与计算机》.2016,第33卷(第02期),第39-43页.

代文豪

(74) 专利代理机构 北京市金杜律师事务所

Yang Liu等.Q_Learning based on active backup and memory mechanism.《Proceedings of 2004 International Conference on Machine Learning and Cybernetics》.2004,第271-275页.

11256

代理人 王茂华

(51) Int.Cl.

G06F 11/14 (2006.01)

审查员 苗文娟

(56) 对比文件

US 2016342481 A1,2016.11.24

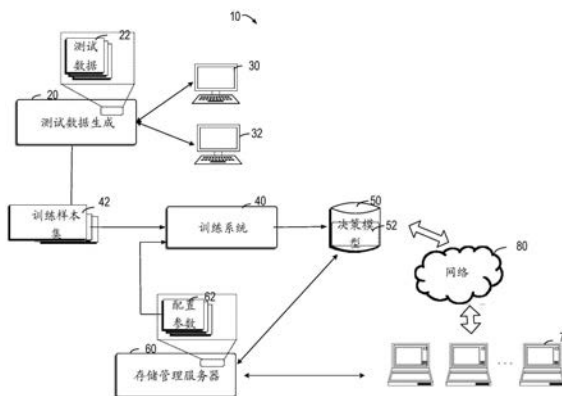
权利要求书3页 说明书13页 附图6页

(54) 发明名称

用于控制数据备份的方法和系统

(57) 摘要

本公开的实施例提供了用于控制数据备份的方法、系统以及计算机可读存储介质。该方法包括获得以多种备份方案执行数据备份的测试数据,所述多种所述备份方案与影响所述数据备份的多个因素的多组取值中的每组取值相关联;基于对所述测试数据的比较,生成训练样本集;将所述训练样本集转换成用于控制数据备份的决策模型;基于使用所述决策模型执行数据备份而获得的配置参数,来优化所述决策模型。



1. 一种用于控制数据备份的方法,包括:

获得以多种备份方案执行数据备份的测试数据,所述多种所述备份方案与影响所述数据备份的多个因素的多组预设取值中的每组预设取值相关联,其中获得所述测试数据包括:基于所述多个因素的每组预设取值,在第一数据存储系统和第二数据存储系统之间以所述多种备份方案中的每一个执行数据备份,来生成所述测试数据;

基于对所述测试数据的比较,生成训练样本集,其中生成训练样本集包括:对于所述每组预设取值,在所述第一数据存储系统和所述第二数据存储系统之间比较所述测试数据中不同备份方案执行所述数据备份所花费的时间;选择所述第一数据存储系统和所述第二数据存储系统之间执行所述数据备份所花费的时间低于阈值的备份方案;基于被选择的所述备份方案及其相关联的所述多个因素的所述每组预设取值,生成所述训练样本集的一条训练样本;

客户端将所述训练样本集转换成用于控制数据备份的决策模型并将所述决策模型存储到共享存储单元;

基于使用所述决策模型执行数据备份而获得的配置参数,来优化所述决策模型,其中使用所述决策模型执行数据备份而获得的配置参数如下获得:客户端使用存储在共享存储单元上的所述决策模型在线执行所述数据备份;存储管理服务器在线收集所述客户端的所述配置参数。

2. 根据权利要求1所述的方法,其中所述因素包括以下至少一个:待备份数据的大小、与上一次备份数据相比新增字节所占的比例、平均存储段的大小、最小存储段的大小、最大存储段的大小、存储段的总数量、待备份文件的数量、执行备份的设备与待备份的设备之间的网络带宽、执行备份的设备与待备份的设备之间的网络往返时延(RTT)。

3. 根据权利要求1所述的方法,其中将所述训练样本集转换成用于控制数据备份的决策模型包括:

基于快速决策树(VFDT)算法将所述训练样本集转换成以所述因素为属性的决策树。

4. 根据权利要求3所述的方法,其中所述决策树的内部节点为所述因素,叶节点表示用于控制数据备份的备份方案,从所述内部节点到所述叶节点的各分支为基于所述因素的取值范围的分类。

5. 根据权利要求1所述的方法,其中优化所述决策模型包括:

基于所获得的配置参数,进行增量的样本训练,以便优化所述决策模型;

基于所述优化后的决策模型进行数据备份的控制。

6. 根据权利要求5所述的方法,其中所述基于所获得的配置参数,进行增量的样本训练包括:

基于所述配置参数,生成增量的训练样本;

基于能够进行多变量系统分类并且支持新增数据的增量学习的机器学习算法,对所述增量的训练样本进行训练。

7. 根据权利要求6所述的方法,其中所述决策模型为以所述因素为属性的决策树;

所述对所述增量的训练样本进行训练包括:通过计算所述增量的训练样本的属性的信息增益,根据Hoeffding边界确定所述决策树中内部节点的分裂。

8. 一种用于控制数据备份的系统,包括:

存储单元,被配置为存储一个或多个程序和用于控制数据备份方案选择的决策模型;处理器,耦合至所述存储单元并且被配置为执行所述一个或多个程序使所述系统执行多个动作,所述动作包括:

获得以多种备份方案执行数据备份的测试数据,所述多种备份方案与影响所述数据备份的多个因素的多组预设取值中的每组预设取值相关联,其中获得所述测试数据包括:基于所述多个因素的每组预设取值,在第一数据存储系统和第二数据存储系统之间以所述多种备份方案中的每一个执行数据备份,来生成所述测试数据;

基于对所述测试数据的比较,生成训练样本集,其中生成训练样本集包括:对于所述每组预设取值,在所述第一数据存储系统和所述第二数据存储系统之间比较所述测试数据中不同备份方案执行所述数据备份所花费的时间;选择所述第一数据存储系统和所述第二数据存储系统之间执行所述数据备份所花费的时间低于阈值的备份方案;基于被选择的所述备份方案及其相关联的所述多个因素的所述每组预设取值,生成所述训练样本集的一条训练样本;

客户端将所述训练样本集转换成用于控制数据备份的决策模型并将所述决策模型存储到共享存储单元;

基于使用所述决策模型执行数据备份而获得的配置参数,来优化所述决策模型,

其中使用所述决策模型执行数据备份而获得的配置参数如下获得:客户端使用存储在共享存储单元上的所述决策模型在线执行所述数据备份;存储管理服务器在线收集所述客户端的所述配置参数。

9. 根据权利要求8所述的系统,其中所述因素包括以下至少一个:待备份数据的大小、与上一次备份数据相比新增字节所占的比例、平均存储段的大小、最小存储段的大小、最大存储段的大小、存储段的总数量、待备份文件的数量、执行备份的设备与待备份的设备之间的网络带宽、执行备份的设备与待备份的设备之间的网络往返时延(RTT)。

10. 根据权利要求8所述的系统,其中将所述训练样本集转换成用于控制数据备份的决策模型包括:

基于快速决策树(VFDT)算法将所述训练样本集转换成以所述因素为属性的决策树。

11. 根据权利要求10所述的系统,其中所述决策树的内部节点为所述因素,叶节点表示用于控制数据备份的备份方案,从所述内部节点到所述叶节点的分支为基于所述因素的取值范围的分类。

12. 根据权利要求8所述的系统,其中优化所述决策模型包括:

基于所获得的配置参数,进行增量的样本训练,以便优化所述决策模型;

基于所述优化后的决策模型进行数据备份的控制。

13. 根据权利要求12所述的系统,其中所述基于所获得的配置参数,进行增量的样本训练包括:

基于所述配置参数,生成增量的训练样本;

基于能够进行多变量系统分类并且支持新增数据的增量学习的机器学习算法,对所述增量的训练样本进行训练。

14. 根据权利要求13所述的系统,其中

所述决策模型为以所述因素为属性的决策树;

所述对所述增量的训练样本进行训练包括：

通过计算所述增量的训练样本的属性的信息增益，根据Hoeffding边界确定所述决策树中内部节点的分裂。

15. 一种计算机可读存储介质，其上存储了一个或多个计算机程序，当所述程序由处理器运行时执行权利要求1-7中任一所述的方法。

用于控制数据备份的方法和系统

技术领域

[0001] 本公开涉及数据备份领域,并且更具体地,涉及用于控制数据备份的方法和系统。

背景技术

[0002] 现代数据备份技术通常能够根据待备份的数据类型、数据大小、备份设备所在的网络环境以及带备份数据的放置位置来选择不同的备份方案。上述数据类型、数据大小、备份设备所在的网络环境以及带备份数据的放置位置等因素通常会影响备份方案的选择。选择不同的数据备份方案,完成数据备份所需的时间和通过备份设备间网络所发送的字节数会有很大的不同。因此,在特定的备份场景下如何控制数据备份,进而选择最为合适的备份方案在现有技术中是一个亟待解决的难题。

发明内容

[0003] 本公开提供一种用于控制数据备份的方法,能够针对不同影响数据备份的因素选择合适的备份方案。

[0004] 在本公开的第一方面,提供了一种用于控制数据备份的方法。该方法包括:获得以多种备份方案执行数据备份的测试数据,该多种该备份方案与影响该数据备份的多个因素的多组取值中的每组取值相关联;基于对该测试数据的比较,生成训练样本集;将该训练样本集转换成用于控制数据备份的决策模型;基于使用该决策模型执行数据备份而获得的配置参数,来优化该决策模型。

[0005] 在本公开的第二方面,还提供一种用于控制数据备份的系统。该系统包括:存储单元,被配置为存储一个或多个程序和用于控制数据备份方案选择的决策模型;处理器,耦合至该存储单元并且被配置为执行该一个或多个程序使该设备执行多个动作,该动作包括:获得以多种备份方案执行数据备份的测试数据,该多个备份方案与影响该数据备份的多个因素的多组取值中的每组取值相关联;基于对该测试数据的比较,生成训练样本集;将该训练样本集转换成用于控制数据备份的决策模型;基于使用该决策模型执行数据备份而获得的配置参数,来优化该决策模型。

[0006] 此外,本公开还提供一种用于控制数据备份的方法,能够利用客户端的实际复制数据快速对数据备份的控制决策进行优化。

[0007] 在本公开的第三方面,提供了一种用于控制数据备份的方法。该方法包括:将用于控制数据备份方案选择的决策模型存储在共享存储单元上,该决策模型与影响该数据备份的多个因素的取值及其对应的备份方案相关联;客户端使用该决策模型执行数据备份而获得配置参数;基于该配置参数,进行增量的样本训练,以便优化该决策模型;基于该优化后的决策模型进行数据备份的控制。

[0008] 在本公开的第四方面,提供了一种用于控制数据备份的系统。该系统包括:存储单元,被配置为存储一个或多个程序;处理器,耦合至该存储单元并且被配置为执行该一个或多个程序使该设备执行多个动作,该动作包括:将用于控制数据备份方案选择的决策模型

存储在共享存储单元上,该决策模型与影响该数据备份的多个因素的取值及其对应的备份方案相关联;客户端使用该决策模型执行数据备份而获得配置参数;基于该配置参数,进行增量的样本训练,以便优化该决策模型;基于该优化后的决策模型进行数据备份的控制。

[0009] 根据本公开的第五方面,提供了一种计算机可读存储介质。该计算机可读存储介质上存储了一个或多个计算机程序,当该程序由处理器运行时执行本公开的第一或三方面用于控制数据备份的方法。

[0010] 提供发明内容部分是为了简化的形式来介绍对概念的选择,它们在下文的具体实施方式中将被进一步描述。发明内容部分无意标识本公开的关键特征或主要特征,也无意限制本公开的范围。

附图说明

[0011] 通过结合附图对本公开示例性实施例进行更详细的描述,本公开的上述以及其它目的、特征和优势将变得更加明显,其中,在本公开示例性实施例中,相同的参考标号通常代表相同部件。

[0012] 图1示出了能够实施本公开的多个实现的用于控制数据备份的系统10的框图。;

[0013] 图2示出了根据本公开的一个实现的用于控制数据备份的存储管理服务器、客户端的交互场景100的示意图;

[0014] 图3示意性示出了根据本公开的一个实现的用于控制数据备份的方法300的流程图;

[0015] 图4示意性示出了根据本公开的一个实现的用于测试数据生成方法400的流程图;

[0016] 图5示意性示出了根据本公开的一个实现的用于训练样本集生成方法500的流程图;

[0017] 图6示意性示出了根据快速决策树(VFDT)算法将训练样本集转换后的决策树600的示意图;

[0018] 图7示意性示出了根据本公开的一个实现的用于决策模型优化方法700的流程图;

[0019] 图8示意性示出了根据本公开的另一个实现的用于控制数据备份的方法800的流程图;

[0020] 图9示意性示出了适于用来实现本公开实施例的电子设备900的系统框图;

[0021] 图10示意性示出了根据本公开的一个实现的用于控制数据备份的装置1000的框图;

[0022] 图11示意性示出了根据本公开的另一个实现的用于控制数据备份的装置1100的框图。

具体实施方式

[0023] 现在将参照若干示例实现来论述本公开。应当理解,论述了这些实现仅是为了使本领域普通技术人员能够更好地理解且因此实现本公开,而不是暗示对本主题的范围的任何限制。

[0024] 如本文所使用的,术语“包括”及其变体要被解读为意味着“包括但不限于”的开放式术语。术语“基于”要被解读为“至少部分地基于”。术语“一个实现”和“一种实现”要被解

读为“至少一个实现”。术语“另一个实现”要被解读为“至少一个其他实现”。术语“第一”、“第二”等等可以指代不同的或相同的对象。下文还可能包括其他明确的和隐含的定义。

[0025] 传统上,控制数据备份的方案可以分为两大类。第一类方法是基于固定规则来选择备份方案,所述固定规则一般是来源于人们以往长期备份数据过程中积累的经验 and 直觉。使用固定规则来确定数据备份方法存在以下缺陷:一方面,固定规则中所考虑的情况通常不能覆盖所有备份场景;另一方面,对于有些情况,例如属性值为连续值的情形,人们的经验是有限的,甚至可能是不准确的。

[0026] 另一类方法是基于KNN匹配算法来选择备份方案,基于KNN匹配算法的控制数据备份的方法通常通过计算两个特征(属性)向量的相似性来确定最合适的备份方案。这种方法非常依赖于历史备份数据,因此在缺少历史备份数据情况下,基于KNN匹配算法的控制数据备份的方法可能导致确定不正确的类类型,进而错误选择备份方案。

[0027] 正如前文所描述的,传统的用于控制数据备份的方法通常是基于固定规则来做出关于备份方案的决策。在这些方法中,固定规则所考虑的影响数据备份的因素通常不能覆盖所有备份场景,另外,基于人们有限经验积累所概括的固定规则通常不够准确。

[0028] 近年来,机器学习算法预测模型被用于不同应用场景的优化决策过程,在这些决策过程中,通常基于特定的输入数据的进行机器学习与归纳,形成预测某一概率区间内的输出值,以及利用该输出值优化决策。因此,在推断中,这种基于机器学习算法预测模型的优化决策过程可以自动归纳数据与决策之间的内在联系,进而显著地避免人为确定决策规则的局限性,进一步提高决策效率和准确性。然而,一般的机器学习算法需要对所有样本数据进行反复训练,比较耗时并且计算也比较复杂,不适于集群系统中需要快速备份大量数据流的要求。另外,一般的机器学习算法不支持多变量的优化,不适于直接应用于控制数据备份方案决策这一受多个变量影响的应用场景。因此本公开针对如何选择合适的机器学习算法,以及怎样将机器学习算法具体结合到控制数据备份的决策场景提供了解决思路,以便降低对存储器的消耗,降低计算复杂度,提高控制数据备份的决策的效率及准确性等等。

[0029] 以下参考附图来说明本公开的基本原理和若干示例实现。

[0030] 图1示出了能够实施本公开的多个实现的用于控制数据备份的系统10的框图。应当理解,图1所示出的系统10仅仅是示例性的,而不应当构成对本公开所描述的实现的功能和范围的任何限制。如图1所示,系统10可以包括测试数据生成装置20、用于测试数据生成的第一数据存储系统30和第二数据存储系统32、训练系统40、共享存储设备50、存储管理服务器60、客户端70和网络80。

[0031] 在一些实施例中,共享存储设备50、存储管理服务器60、客户端70之间可以按照图2所示方式进行交互。图2示出了根据本公开的一个实现的用于控制数据备份的存储管理服务器、客户端的交互场景100的示意图。如图2所示,存储管理服务器60用于从多个客户端70捕获数据并且经由网络80(在图2中未示出)将数据存储于共享存储单元50中,例如,存储管理服务器60可以调用、保存并处理客户端70用于在线执行数据备份的配置参数62,如呼叫客户端主要数据(CALL HOME DATA)。客户端70也可直接从共享存储单元50中捕获数据,例如获取决策模型52。存储管理服务器60可以维护存储单元50中的决策模型52。存储管理服务器60例如可以包括一个或多个处理单元61和存储器68,存储器68上存储有一个或者多个程序63,处理单元61耦合至所述存储器运行一个或多个程序63用于执行客户端数据备份的控

制方法的各个操作或者动作。

[0032] 在一些实现中,客户端70可以被实现为各种用户终端或服务终端。服务终端可以是各种服务提供方提供的服务器、大型计算设备等。用户终端诸如是任意类型的移动终端、固定终端或便携式终端,包括移动手机、站点、单元、设备、多媒体计算机、多媒体平板、互联网节点、通信器、台式计算机、膝上型计算机、笔记本计算机、上网本计算机、平板计算机、个人通信系统(PCS)设备、个人导航设备、个人数字助理(PDA)、音频/视频播放器、数码相机/摄像机、定位设备、电视接收器、无线电广播接收器、电子书设备、游戏设备或者其任意组合,包括这些设备的配件和外设或者其任意组合。

[0033] 在一些实现中,如图1所示,第一数据存储系统30、第二数据存储系统32和存储管理服务器60可以是服务器、大型计算设备等。存储管理服务器60可以包括存储器62、处理设备单元61等组成部分。共享存储设备50可以由各种存储介质来实现,包括但不限于易失性和非易失性介质、可拆卸和不可拆卸介质。网络80可以是任何能够提供设备之间的通信的有线或无线网络,例如因特网、局域网(WLAN)、内联网等。

[0034] 在一些实现中,测试数据生成装置20、训练系统40可以是专用的处理设备,也可以由相应的程序模块实现,例如由存储在存储管理服务器60中存储器68内的程序指令实现。

[0035] 如图1所示,测试数据生成装置20基于多个因素的每组预设取值,在用于测试数据生成的第一数据存储系统30和第二数据存储系统32之间利用多种备份方案,最好是所有备份方案执行真实的数据备份过程,用以产生测试数据。测试数据生成装置20进一步基于针对影响数据备份的多个因素的每组预设取值,比较不同备份方案执行数据备份的效果,例如比较执行备份数据所花费时间;选择所花费的时间小于阈值,或者选择所花费的时间最少的备份方案生成所述训练样本集42中的一条训练样本。

[0036] 如图1所示,训练系统40基于预定算法将所述训练样本集转换为用于控制数据备份的决策模型52,例如转换成以所述因素为属性的决策树。所述决策模型52被存储在共享存储设备50中。客户端70使用决策模型在线执行数据备份,存储管理服务器60通过网络80收集客户端70进行数据备份而获得的配置参数62。配置参数62至少包括关于备份方案、与所述备份方案相关联的因素及其取值、备份效果的信息,配置参数62例如是客户端调用主要数据(CALL HOME DATA)。训练系统40基于配置参数62生成增量的训练样本,利用预定算法对增量训练样本进行训练(而不需要对所有训练样本集进行重新学习),以便对存储在共享存储设备50中的决策模型52(例如决策树)进行优化。客户端70根据共享存储设备50中存储的经优化后的决策模型52进行数据备份的控制。

[0037] 图3示意性示出了根据本公开的一个实现的用于控制数据备份的方法300的流程图。方法300可以由存储管理服务器60、诸如存储管理服务器60中的控制器或处理设备来执行。

[0038] 在框302,针对影响数据备份的多个因素的多组取值中的每组取值,获得以多种备份方案执行数据备份的测试数据。在一些实施例中,针对影响数据备份的多个因素可以包括:待备份数据的大小、与上一次备份数据相比新增字节所占的比例(后续简称为“新字节所占比例”)、平均存储段的大小、最小存储段的大小、最大存储段的大小、存储段的总数量、待备份文件的数量、执行备份的设备与待备份的设备之间的网络带宽、执行备份的设备与待备份的设备之间的网络往返时延(RTT)。针对这些因素不同的取值对应选择的最为匹配

的备份方案可能不同。

[0039] 在一些实施例中,针对多个因素的多组取值中的每组取值,可以通过图4所示框402至框404的操作加以实现。图4示意性示出了根据本公开的一个实现的用于测试数据生成方法400的流程图。其中,在框402设定多个因素的取值范围和值增加间隔。例如,下表1示例性地给出了九个因素的取值范围和值增加间隔的数值。所谓的取值范围的含义是该因素可取值的最小值与最大值之间的范围,所谓的值增加间隔,是指每组预设取值相对于前一组预设取值增加的数值。例如,对于待复制数据大小这一因素,预设的取值范围例如为20KB~12TB,值增加间隔例如为100MB。其含义就是说,在生成测试数据的过程中,复制数据大小这一因素最小预设取值例如为20KB,最大预设取值例如为12TB。复制数据大小这一因素的每个取值例如以100MB为间隔增加。需要强调的是,这些数值仅仅是示例,无意以任何方式限制本公开的范围。

[0040] 表1

影响因素	取值范围	值增加间隔
待复制数据大小	20KB~12TB	100MB
新字节所占比例	1%~70%	2%
平均存储段大小	4KB~40MB	10KB
最小存储段大小	100B~40MB	10KB
最大存储段大小	100B~40MB	10KB
存储段总数量	$1\sim 2^{32}$	10
待备份文件的数量	$1\sim 2^{64}$	100
备份设备间的网络带宽	10Mbps~10Gbps	100Mbps
备份设备间的网络往返时延	1ms~200ms	20ms

[0042] 在设定多个因素的取值范围和取值增加间隔之后,预设多个因素的一组取值,即框404。例如下表2示例性地给出了九个因素的一组预设取值,表2中的数值仅仅是示例,无意以任何方式限制本公开的范围。例如在该组预设取值中,待复制数据大小这一因素的取值例如为200G,新字节所占比例这一因素的取值例如为1%,平均存储段大小这一因素的取值例如为4KB,最小存储段大小这一因素的取值例如为100B,最大存储段大小这一因素的取值例如为4MB,存储段总数量这一因素的取值例如为547,待备份文件的数量这一因素的取值例如为143,备份设备间的网络带宽这一因素的取值例如为10Mbps,备份设备间的网络往返时延这一因素的取值例如为1ms。需要强调的是,这些数值仅仅是示例,无意以任何方式限制本公开的范围。

[0043] 表2

影响因素	取值
待复制数据大小	200G
新字节所占比例	1%
平均存储段大小	4KB
最小存储段大小	100B
最大存储段大小	4MB
存储段总数量	547

待备份文件的数量	143
备份设备间的网络带宽	10Mbps
备份设备间的网络往返时延	1ms

[0045] 在一些实施例中,针对每组取值获得以多种备份方案执行数据备份的测试数据可以通过图4所示框406至框412的操作加以实现。其中,针对每组预设取值,例如表2中所示的这组预设取值,执行框406,即在第一数据存储系统30和第二数据存储系统32之间以多种备份方案(例如备份方案1至4)分别执行真实的数据备份,进而生成一组测试数据,例如包括执行备份方案1的测试数据、执行备份方案2的测试数据、执行备份方案3的测试数据和执行备份方案4的测试数据;然后判断这组预设取值是否超出预先设定的各因素的取值范围,即框408;如果没有超出预先设定的取值范围,则在九个因素的当前预设取值上增加预先设定的值增加间隔,以便形成下一组预设取值,即框410;如果已经超出预先设定的取值范围,输出记录的测试数据,即框412;针对形成的每一组预设取值,在第一数据存储系统30和第二数据存储系统32之间以多种备份方案(例如备份方案1至4)分别执行真实的数据备份,以便生成每一组测试数据。

[0046] 在框304,基于对所述测试数据的比较,生成训练样本集。在一些实施例中,针对基于对所述测试数据的比较生成训练样本集,可以通过图5所示框502至框506的方式加以实现。图5示意性示出了根据本公开的一个实现的用于训练样本集生成方法500的流程图。其中,框502例如是针对多个因素的每组预设取值,比较相应组测试数据中以不同备份方案执行数据备份对应生成的测试数据。例如,针对表2示例的预设取值,具体比较执行备份方案1的测试数据、执行备份方案2的测试数据、执行备份方案3的测试数据以及执行备份方案4的测试数据中的备份效果信息。例如表3示例性显示了对应执行备份方案1-4的测试数据中有关复制效果的部分信息,包括“执行数据备份所花费时间”和“备份设备间网络发送的字节数”。需要强调的是,这些数值仅仅是示例,无意以任何方式限制本公开的范围。

[0047] 表3

备份方案	执行数据备份所花费时间	备份设备间网络发送的字节数
1	0000h: 05m: 29s	1,498,638,048 (1.396 GB)
2	0000h: 03m: 14s	10,867,267,796 (10.12 GB)
3	0000h: 06m: 10s	5,949,861,888 (5.541 GB)
4	0000h: 02m: 47s	4,718,062,572 (4.394 GB)

[0048] 通过比较对应执行备份方案1-4的测试数据,执行框504,即选择执行数据备份所花费时间较少或低于阈值的备份方案;然后基于被选择的备份方案及其相关联的多个因素的预设取值,生成训练样本集的一条训练样本,即框506。例如,通过比较表3中对应执行备份方案1-4的四条测试数据中有关“执行数据备份所花费时间”的信息。可以看到备份方案4执行数据备份所花费的时间相对于备份方案1-3更少,因此选择备份方案4及其相关联的九个因素的所述预设取值,生成图1所示训练样本集42中的一条训练样本。下表4示例性显示该条训练样本,在该条训练样本中包括九个因素的一组预设取值以及与该组预设取值最为匹配的备份方案。由于针对多个因素的每组预设取值,都需要通过执行框502至框506,生成

如表4所示的一条对应的训练样本,直至针对多个因素的所有组预设取值,生成整个训练样本集42中的所有训练样本。需要强调的是,表4中的数值仅仅是示例,无意以任何方式限制本公开的范围。

[0050] 表4

[0051]	待复制数据大小	200G
	新字节所占比例	1%
	平均存储段大小	4KB
	最小存储段大小	100B
	最大存储段大小	4MB
	存储段总数量	547
	待备份文件的数量	143
	备份设备间的网络带宽	10Mbps
	备份设备间的网络往返时延	1ms
	备份方案	4

[0052] 在框306,将所述训练样本集转换成用于控制数据备份的决策模型。在一些实施例中,可以基于多种预定的机器学习算法来实现训练样本集到所述决策模型的转换。该预定的机器学习算法需要满足两方面特点:一是,能够用监督学习解决多变量系统分类问题;二是,对于新增的样本数据,支持增量学习。这是因为,训练样本集42中的每一条训练样本,如表4所示,都含有多个因素及其关联的备份方案,也就是说,从训练样本集42到决策模型52的转换过程需要涉及多个变量的分类与学习过程,因此,所采用的机器学习算法应当能够用监督学习解决多变量系统分类问题。在一些实施例中,例如,可以采用快速决策树(VFDT)算法将所述训练样本集转换成以多个因素为属性的决策树;在一些实施例中,例如可以利用人工神经网络算法(Learning++)来实现决策模型的构建;除此之外,还可以采用具备在线递归算法的增量支持向量机(SVM)算法来将所述训练样本集转换成用于控制数据备份的决策模型。

[0053] 在一些实施例中,例如采用快速决策树(VFDT)算法将所述训练样本集转换成的以多个因素为属性的决策树结构。表5示意性示出了训练样本集中的部分训练样本。表5中的数值仅仅是示例,无意以任何方式限制本公开的范围。

[0054] 图6示意性示出了根据快速决策树(VFDT)算法将训练样本集转换后的决策树600的示意图。其中该决策树600的内部节点为因素,叶节点表示用于控制数据备份的备份方案,从所述内部节点到所述叶节点的分支为基于所述因素的取值范围的分类。如图6所示,决策树600中的内部节点602为待备份数据大小这一因素,内部节点602的不同的分支,即分支一604、分支二606和分支三608,对应基于待备份数据大小这一因素的不同取值或取值范围的分类,例如分支一604为待备份数据大小这一因素取值例如为100M的分类,分支二606为待备份数据大小这一因素取值例如为200M的分类,分支三608为待备份数据大小这一因素取值例如为2G的分类。叶节点630为与分支二606匹配的备份方案3。内部节点610为备份设备间的网络带宽这一因素,内部节点610的不同的分支,即子分支一612、子分支二614和子分支三616,为基于备份设备间的网络带宽这一因素的不同取值或取值范围的分类。例如,子分支一612为网络带宽这一因素取值例如为10Mbps的分类,子分支二614为网络带宽

这一因素取值例如为100Mbps的分类,子分支三616为网络带宽这一因素取值例如为200Mbps的分类。叶节点618为与子分支一612匹配的备份方案1,叶节点620为与子分支二614匹配的备份方案3,叶节点622为与子分支三616匹配的备份方案3。决策树600中的内部节点640为存储段总数量这一因素,内部节点640的不同的分支,即子分支四642和子分支五644,对应基于存储段总数量这一因素的不同取值或取值范围的分类,例如子分支四642为存储段总数量这一因素取值例如为1的分类,子分支五644为存储段总数量这一因素取值例如为4的分类。叶节点646为与分支四642匹配的备份方案4,叶节点648为与子分支五644匹配的备份方案2。

[0055] 如图6所示,当使用决策树600决策备份方案的时候,如果当前备份数据的任务所涉及的待备份数据大小为2G,并且用于备份数据设备之间的网络带宽为10Mbps,根据决策树600的决策会推荐选择备份方案1来进行数据备份;如果当前待备份数据大小为100MB,并且用于备份数据的存储段总数量为1,根据决策树600的决策会推荐选择备份方案4来进行数据备份。需要强调的是,前述数值仅仅是示例,无意以任何方式限制本公开的范围。

[0056] 表5

待复制数据大小	新字节所占比例	平均存储段大小	最小存储段大小	最大存储段大小	存储段总数量	待备份文件的数量	备份设备间的网络带宽	备份设备间的网络往返时延	备份方案
100M	1%	4KB	100B	100B	1	1	10Mbps	1ms	4
100M	2%	8KB	10KB	1MB	4	5	100Mbps	10ms	2
200M	3%	12KB	30B	2MB	8	23	200Mbps	20ms	3
2G	3%	12KB	30B	2MB	8	23	200Mbps	20ms	3
2G	3%	12KB	30B	2MB	8	23	10Mbps	20ms	1
2G	3%	12KB	30B	2MB	8	23	1Gbps	20ms	2

[0057] 在框308,基于使用所述决策模型执行数据备份而获得的配置参数,来优化所述决策模型。在一些实施例中,可以通过图7所示框702至框712的操作加以实现。图7示意性示出了根据本公开的一个实现的用于决策模型优化方法700的流程图。其中,如图7所示,将决策模型52存储在共享存储设备50中,即框702;客户端70使用所述决策模型52执行数据备份进而获得的配置参数,即框704;存储服务器60获取客户端的配置参数,例如,通过调用客户端70的主要数据(CALL HOME DATA),即框706;基于客户端的配置参数生成增量的训练样本,即框708;在一些实施例中,可以采用与图5所示的框502至框506相同的方式生成该增量的训练样本。然后,基于前面所述的预定的机器学习算法来对增量的训练样本进行增量训练,以便优化该决策模型52,即框710;之后,基于优化后的决策模型52进行数据备份的控制,即框712。

[0059] 在一些实施例中,可以采用前面提及的、满足两方面特点的所述的预定的机器学习

习算法来对增量的训练样本进行增量训练。该两方面特点包括：一是，能够用监督学习解决多变量系统分类问题；二是，对于新增的样本数据，支持增量学习。由于仅对增量训练样本进行增量训练，而不需要对所有训练样本集进行重新学习，因此降低对存储器的消耗，同时也降低计算复杂度，进而提高控制数据备份的决策的效率。。在一些实施例中，例如，可以采用Hoeffding决策树更新算法对新增的样本进行增量训练，计算样本属性的信息增益，并根据Hoeffding边界决定决策树分裂哪个节点，进而实现决策模型的优化；在一些实施例中，例如可以利用神经网络算法 (Learning++) 来实现决策模型的优化。神经网络算法 (Learning++) 是在2002年发布的，其是神经网络 (NN) 模式分类器的增量训练算法，该算法在后续的增量学习会话中不需要访问以前使用的数据，但同时也不会忘记先前获得的知识或模型，因此神经网络算法 (Learning++) 适于对新增的样本进行训练以便优化已有决策模型。除此之外，具备在线递归算法的增量支持向量机 (SVM) 算法也适于对新增的样本进行训练以便优化已有决策模型。

[0060] 图8示意性示出了根据本公开的另一个实现的用于控制数据备份的方法800的流程图。方法800可以由存储管理服务器60、诸如存储管理服务器60中的控制器或处理设备来执行。

[0061] 在框802，将用于控制数据备份方案选择的决策模型存储在共享存储单元上，该决策模型与影响所述数据备份的多个因素的取值及其对应的备份方案相关联。在一些实施例中，将决策模型52存储在共享存储设备50中。该决策模型52可以采用图3中的框302至框306的操作来实现，也可以采用其他操作实现。其中，决策模型52用于控制数据备份方案选择并且与影响数据备份的多个因素的取值及其对应的备份方案相关联，也就是能够根据多个因素的取值的不同，对应选择匹配的备份方案。

[0062] 在框804，客户端使用所述决策模型执行数据备份而获得配置参数。在一些实施例中，可以通过图7所示的框702至框706的方式加以实现。其中，如图7所示，将决策模型52存储在共享存储设备50中，即框702；客户端70使用所述决策模型52执行数据备份进而获得的配置参数，即框704；存储服务器60获取客户端的配置参数，例如，通过调用客户端70的主要数据 (CALL HOME DATA)，即框706。

[0063] 在框806，基于配置参数，进行增量的样本训练，以便优化决策模型。在一些实施例中，可以通过图7所示的框708至框712的操作加以实现。其中，基于客户端的配置参数生成增量的训练样本，即框708；在一些实施例中，可以采用与图5所示的框502至框506相同的操作生成该增量的训练样本。然后，基于前面所述的预定的机器学习算法来对增量的训练样本进行增量训练，以便优化该决策模型52，即框710；之后，基于优化后的决策模型52进行数据备份的控制，即框712。在一些实施例中，例如采用可以采用Hoeffding决策树更新算法对增量的训练样本进行增量训练，以便优化该决策模型。表6具体示意性示出了Hoeffding决策树更新算法处理逻辑的实现。需要强调的是，后续处理逻辑的实现仅仅是示例，无意以任何方式限制本公开的范围。其中S是一系列的训练样本，训练样本包括两部分，一部分训练样本来源于基于对所述测试数据的比较而生成训练样本集，另一部分训练样本是基于客户端的配置参数而生成新增的训练样本；X是一组离散的属性，也就是影响数据备份的一组因素；G(.)是分割评估函数； δ 减为去任何给定的选择正确属性的期望概率；n为独立观察次数；R为随机变量范围 X_i 的实值范围；i为属性索引；j为每个属性 X_i 的索引值。

[0064] 表6

[0065]

让 HT 为一棵树，其仅有单叶 l_1 （根节点）让 $G_1(X_i)$ 为通过预测 S 中最常见的类来获得的 \bar{G}

对于每个类 y_k

对于每个属性 X_i 的每个值 X_{ij}

令 $n_{ijk}(l_1) = 0$

对于 S 中的每个样本 (x, y_k)

使用 HT 将 (x, y) 分拣成叶

对于 x 中的每个 X_{ij} ，使得 $X_i \in X_l$

增加 $n_{ijk}(l)$

以到 l 为止看到的样本中主要的类来标签 l

如果到 l 为止看到的样本中，全都不是相同的类，则

使用计数 $n_{ijk}(l)$ 为每个属性 X_i 计算 $G_1(X_i)$

让 X_a 为最高 G_1 的属性

让 X_b 为第二高 G_1 的属性

使用公式计算 ε :

[0066]

$$\varepsilon = \sqrt{\frac{R^2 \ln\left(\frac{1}{\delta}\right)}{2n}}$$

如果 $G_1(X_a) - G_1(X_b) > \varepsilon$ ，那么

在 X_a 分割的内部节点替换 l

对于分裂的每个分支

添加新叶 l_m ，并且让 $X = X - \{X_a\}$

让 $G_m(X_a)$ 为通过预测在 l_m 最常见的类所获得的 \bar{G}

对于每个类 y_k 和每个属性 X_i 的每个值 X_{ij}

令 $n_{ijk}(l_m) = 0$

返回 HT

[0067] 在框808，基于所述优化后的决策模型进行数据备份的控制。在一些实施例中，将优化后的据决策模型存储在共享存储设备50上，客户端70可直接从共享存储单元50中获取优化后的据决策模型52，使用该优化后的决策模型52控制数据的备份。在一些实施例中，也可以由存储管理服务器60从共享存储单元50中获取优化后的据决策模型52，并使用该优化后的决策模型52控制客户端70的数据的备份。在一些实施例中，还可以利用在线数据的备份过程中产生的数据对决策模型52反复优化和应用。

[0068] 图9示意性示出了适于用来实现本公开实施例的电子设备900的框图。设备900可以用于实现图1的存储管理服务器60中的一个或多个主机。如图所示，设备900包括中央处理单元(CPU)901，其可以根据存储在只读存储器(ROM)902中的计算机程序指令或者从存储单元608加载到随机访问存储器(RAM)903中的计算机程序指令，来执行各种适当的动作和

处理。在RAM903中,还可存储设备900操作所需的各种程序和数据。CPU 901、ROM 902以及RAM903通过总线804彼此相连。输入/输出(I/O)接口905也连接至总线904。

[0069] 设备900中的多个部件连接至I/O接口905,包括:输入单元906,例如键盘、鼠标等;输出单元907,例如各种类型的显示器、扬声器等;存储单元908,例如磁盘、光盘等;以及通信单元909,例如网卡、调制解调器、无线通信收发机等。通信单元909允许设备900通过诸如因特网的计算机网络和/或各种电信网络与其他设备交换信息/数据。

[0070] 处理单元901执行上文所描述的各个方法和处理,例如执行用于控制数据备份的方法300和800。例如,在一些实施例中,方法300和方法800可被实现为计算机软件程序,其被存储于机器可读介质,例如存储单元908。在一些实施例中,计算机程序的部分或者全部可以经由ROM 902和/或通信单元909而被载入和/或安装到设备900上。当计算机程序加载到RAM 903并由CPU 901执行时,可以执行上文描述的方法300和方法800的一个或多个操作。备选地,在其他实施例中,CPU 901可以通过其他任何适当的方式(例如,借助于固件)而被配置为执行方法300和方法800的一个或多个动作。

[0071] 图10示意性示出了根据本公开的一个实现的用于控制数据备份的装置1000的框图。在一些实施例中,该装置1000包括多个模块,每个模块与方法300中的框对应一致。在一些实施例中,该装置还可以是一种计算机程序产品。所述计算机程序产品被存储在计算机存储介质中并且包括机器可执行指令,所述机器可执行指令在设备中运行时使得所述设备执行方法300中的各个框或者动作。如图10所示,该装置1000包括:测试数据生成模块1002,其用于获得以多种备份方案执行数据备份的测试数据,所述多种所述备份方案与影响所述数据备份的多个因素的多组取值中的每组取值相关联;训练样本集生成模块1004,其用于基于对所述测试数据的比较生成训练样本集;决策模型转换模块1006,其用于将所述训练样本集转换成用于控制数据备份的决策模型;以及决策模型优化模块1008,其基于使用所述决策模型执行数据备份而获得的配置参数,来优化所述决策模型。

[0072] 图11示意性示出了根据本公开的一个实现的用于控制数据备份的装置1100的框图。在一些实施例中,该装置1100包括多个模块,每个模块与方法800中的动作对应一致。在一些实施例中,该装置还可以是一种计算机程序产品。所述计算机程序产品被存储在计算机存储介质中并且包括机器可执行指令,所述机器可执行指令在设备中运行时使得所述设备执行方法800中的各个动作。如图11所示,该装置1100包括:决策模型存储模块1102,其将用于控制数据备份方案选择的决策模型存储在共享存储单元上,所述决策模型与影响所述数据备份的多个因素的取值及其对应的备份方案相关联;配置参数获取模块1104,其客户端使用所述决策模型执行数据备份而获得配置参数;决策模型优化模块1106,其基于所述配置参数,进行增量的样本训练,以便优化所述决策模型;以及决策模型控制模块1108,其基于所述优化后的决策模型进行数据备份的控制。尽管已经采用特定于结构特征和/或方法逻辑动作的语言描述了本主题,但是应当理解所附权利要求书中所限定的主题未必局限于上面描述的特定特征或动作。相反,上面所描述的特定特征和动作仅仅是实现权利要求书的示例形式。

[0073] 本公开可以是方法、装置、系统和/或计算机程序产品。计算机程序产品可以包括计算机可读存储介质,其上载有用于执行本公开的各个方面的计算机可读程序指令。

[0074] 计算机可读存储介质可以是保持和存储由指令执行设备使用的指令的有形

设备。计算机可读存储介质例如可以是一一但不限于一一电存储设备、磁存储设备、光存储设备、电磁存储设备、半导体存储设备或者上述的任意合适的组合。计算机可读存储介质的更具体的例子(非穷举的列表)包括:便携式计算机盘、硬盘、随机存取存储器(RAM)、只读存储器(ROM)、可擦式可编程只读存储器(EPR0M或闪存)、静态随机存取存储器(SRAM)、便携式压缩盘只读存储器(CD-ROM)、数字多功能盘(DVD)、记忆棒、软盘、机械编码设备、例如其上存储有指令的打孔卡或凹槽内凸起结构、以及上述的任意合适的组合。这里所使用的计算机可读存储介质不被解释为瞬时信号本身,诸如无线电波或者其他自由传播的电磁波、通过波导或其他传输媒介传播的电磁波(例如,通过光纤电缆的光脉冲)、或者通过电线传输的电信号。

[0075] 这里所描述的计算机可读程序指令可以从计算机可读存储介质下载到各个计算/处理设备,或者通过网络、例如因特网、局域网、广域网和/或无线网下载到外部计算机或外部存储设备。网络可以包括铜传输电缆、光纤传输、无线传输、路由器、防火墙、交换机、网关计算机和/或边缘服务器。每个计算/处理设备中的网络适配卡或者网络接口从网络接收计算机可读程序指令,并转发该计算机可读程序指令,以供存储在各个计算/处理设备中的计算机可读存储介质中。

[0076] 用于执行本公开操作的计算机程序指令可以是汇编指令、指令集架构(ISA)指令、机器指令、机器相关指令、微代码、固件指令、状态设置数据、或者以一种或多种编程语言的任意组合编写的源代码或目标代码,所述编程语言包括面向对象的编程语言—诸如Smalltalk、C++等,以及常规的过程式编程语言—诸如“C”语言或类似的编程语言。计算机可读程序指令可以完全地在用户计算机上执行、部分地在用户计算机上执行、作为一个独立的软件包执行、部分在用户计算机上部分在远程计算机上执行、或者完全在远程计算机或服务器上执行。在涉及远程计算机的情形中,远程计算机可以通过任意种类的网络—包括局域网(LAN)或广域网(WAN)—连接到用户计算机,或者,可以连接到外部计算机(例如利用因特网服务提供商来通过因特网连接)。在一些实施例中,通过利用计算机可读程序指令的状态信息来个性化定制电子电路,例如可编程逻辑电路、现场可编程门阵列(FPGA)或可编程逻辑阵列(PLA),该电子电路可以执行计算机可读程序指令,从而实现本公开的各个方面。

[0077] 这里参照根据本公开实施例的方法、装置(系统)和计算机程序产品的流程图和/或框图描述了本公开的各个方面。应当理解,流程图和/或框图的每个方框以及流程图和/或框图中各方框的组合,都可以由计算机可读程序指令实现。

[0078] 这些计算机可读程序指令可以提供给通用计算机、专用计算机或其它可编程数据处理装置的处理单元,从而生产出一种机器,使得这些指令在通过计算机或其它可编程数据处理装置的处理单元执行时,产生了实现流程图和/或框图中的一个或多个方框中规定的功能/动作的装置。也可以把这些计算机可读程序指令存储在计算机可读存储介质中,这些指令使得计算机、可编程数据处理装置和/或其他设备以特定方式工作,从而,存储有指令的计算机可读介质则包括一个制品,其包括实现流程图和/或框图中的一个或多个方框中规定的功能/动作的各个方面的指令。

[0079] 也可以把计算机可读程序指令加载到计算机、其它可编程数据处理装置、或其它设备上,使得在计算机、其它可编程数据处理装置或其它设备上执行一系列操作步骤,以产

生计算机实现的过程,从而使得在计算机、其它可编程数据处理装置、或其它设备上执行的指令实现流程图和/或框图中的一个或多个方框中规定的功能/动作。

[0080] 附图中的流程图和框图显示了根据本公开的多个实施例的系统、方法和计算机程序产品的可能实现的体系架构、功能和操作。在这点上,流程图或框图中的每个方框可以代表一个模块、程序段或指令的一部分,所述模块、程序段或指令的一部分包含一个或多个用于实现规定的逻辑功能的可执行指令。在有些作为替换的实现中,方框中所标注的功能也可以以不同于附图中所标注的顺序发生。例如,两个连续的方框实际上可以基本并行地执行,它们有时也可以按相反的顺序执行,这依所涉及的功能而定。也要注意的,框图和/或流程图中的每个方框、以及框图和/或流程图中的方框的组合,可以用执行规定的功能或动作的专用的基于硬件的系统来实现,或者可以用专用硬件与计算机指令的组合来实现。

[0081] 以上已经描述了本公开的各实施例,上述说明是示例性的,并非穷尽性的,并且也不限于所披露的各实施例。在不偏离所说明的各实施例的范围和精神的情况下,对于本技术领域的普通技术人员来说许多修改和变更都是显而易见的。本文中所用术语的选择,旨在最好地解释各实施例的原理、实际应用或对市场中的技术改进,或者使本技术领域的其它普通技术人员能理解本文披露的各实施例。以上所述仅为本公开的可选实施例,并不用于限制本公开,对于本领域的技术人员来说,本公开可以有各种更改和变化。凡在本公开的精神和原则之内,所作的任何修改、等效替换、改进等,均应包含在本公开的保护范围之内。

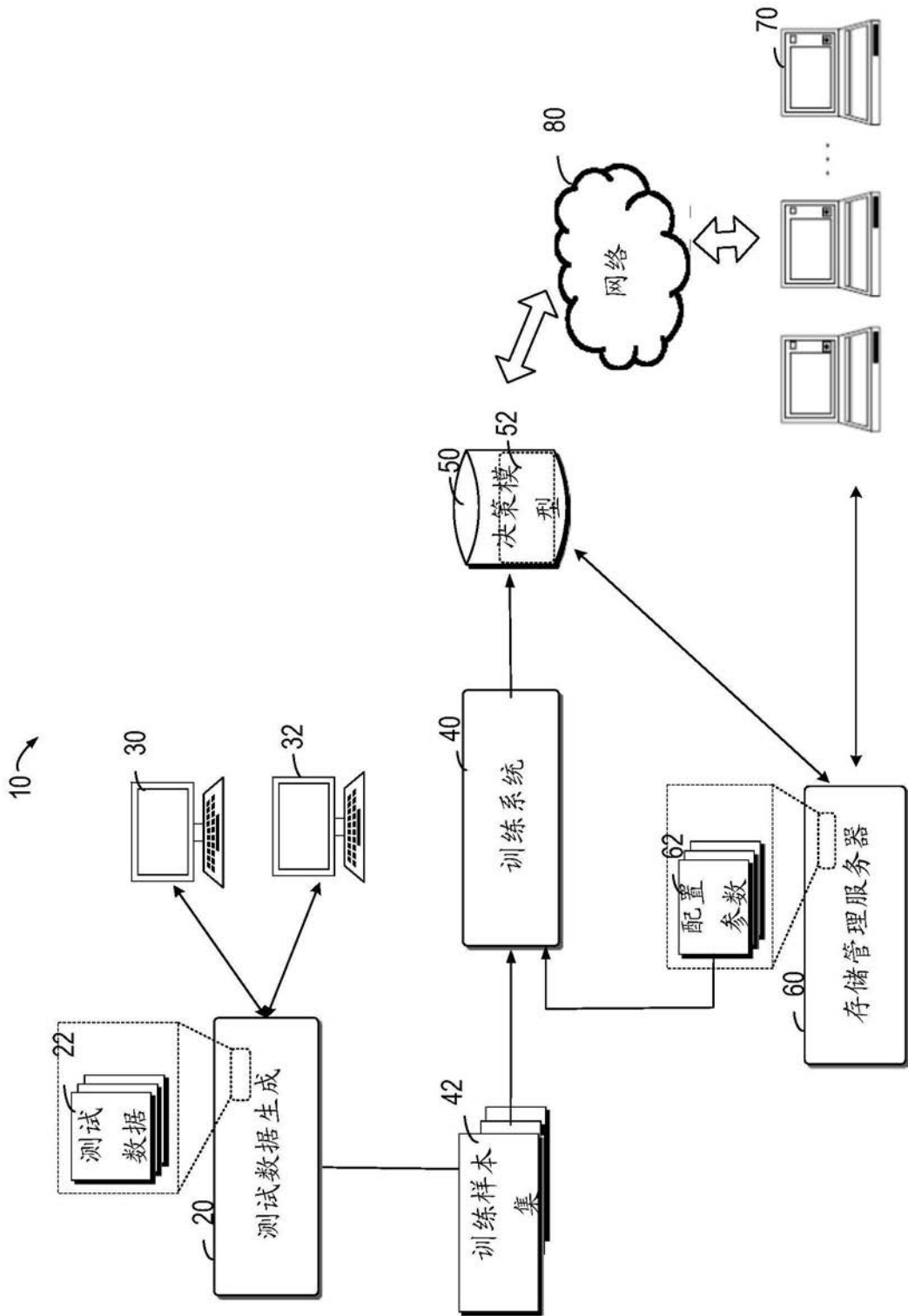


图1

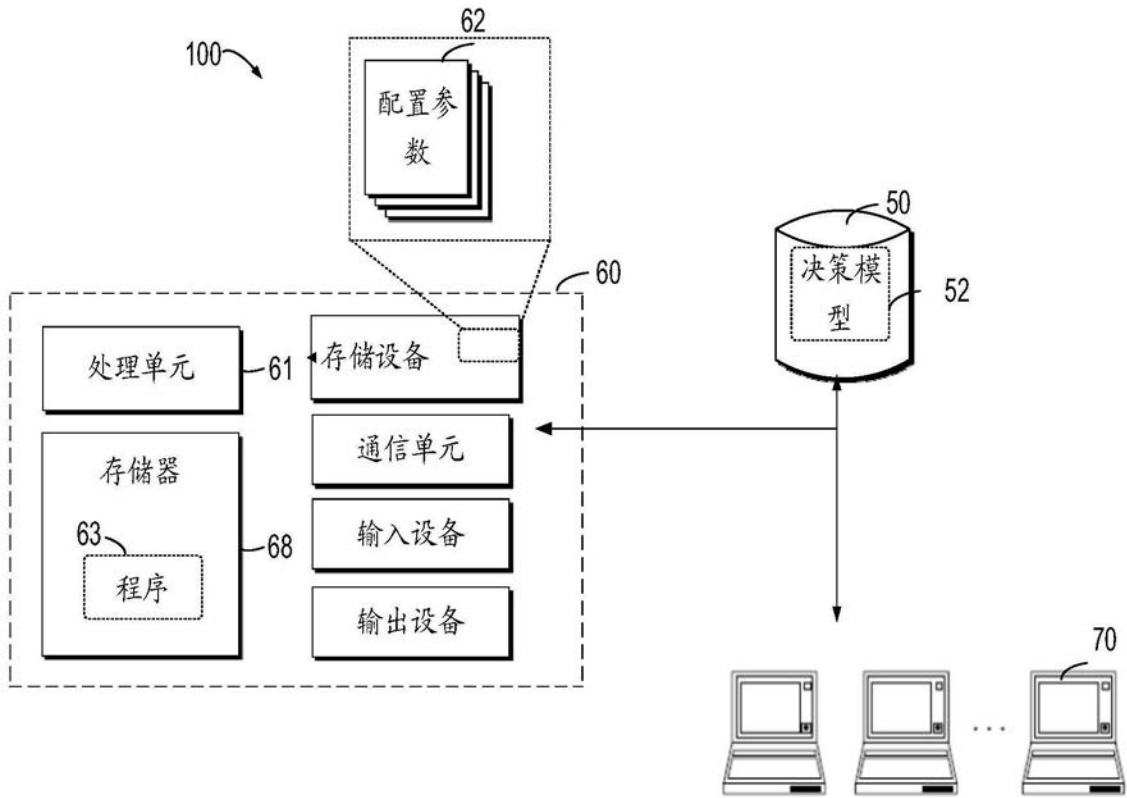


图2

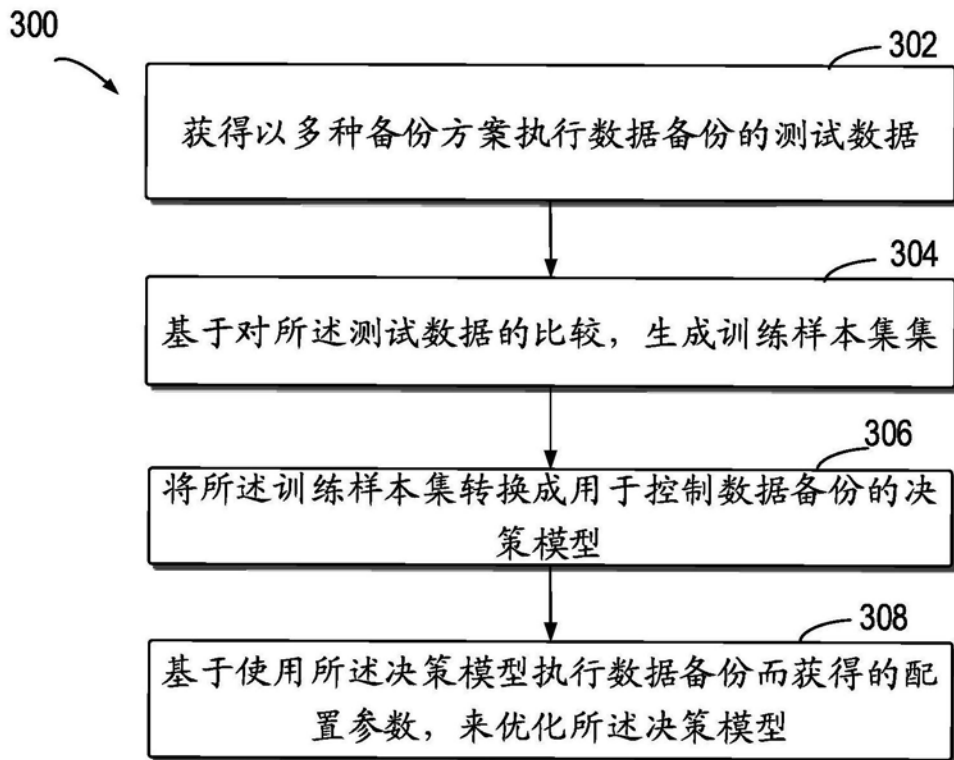


图3

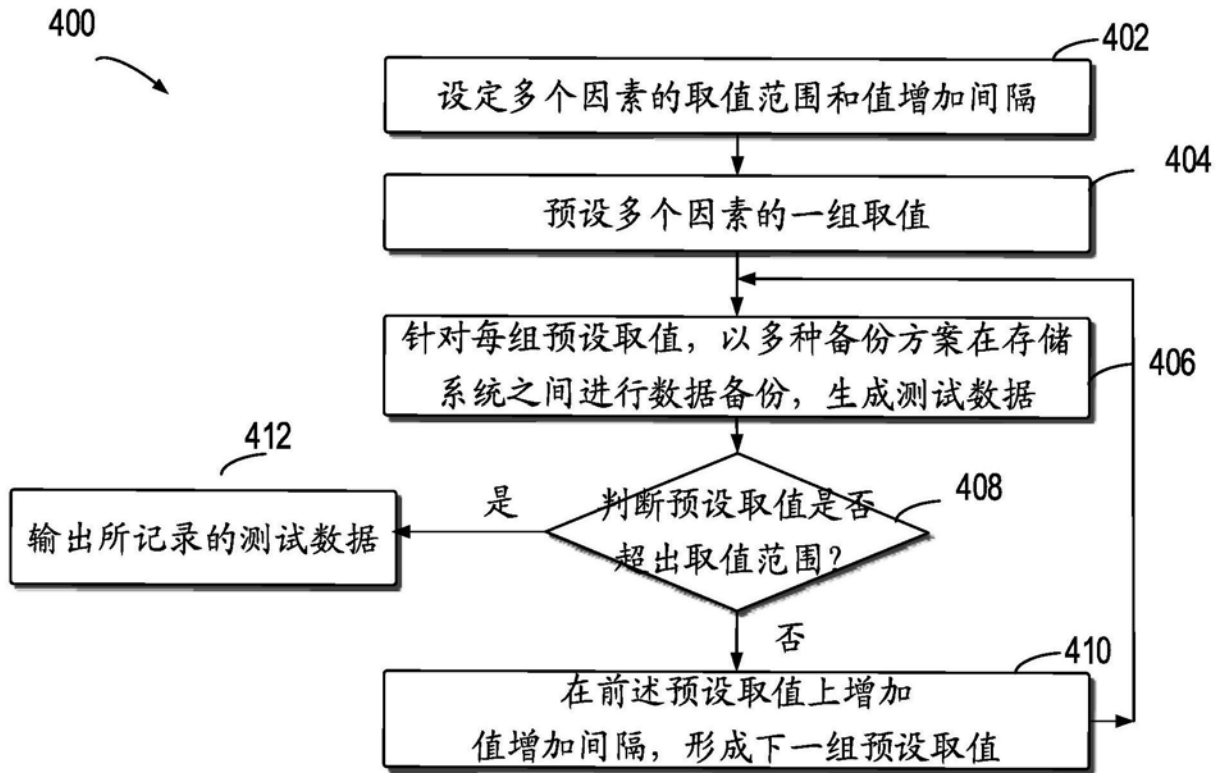


图4

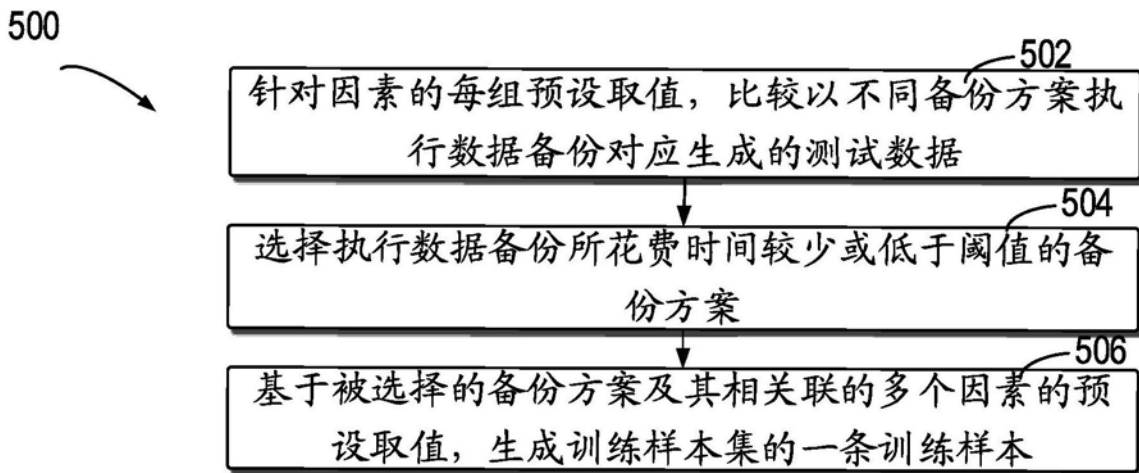


图5

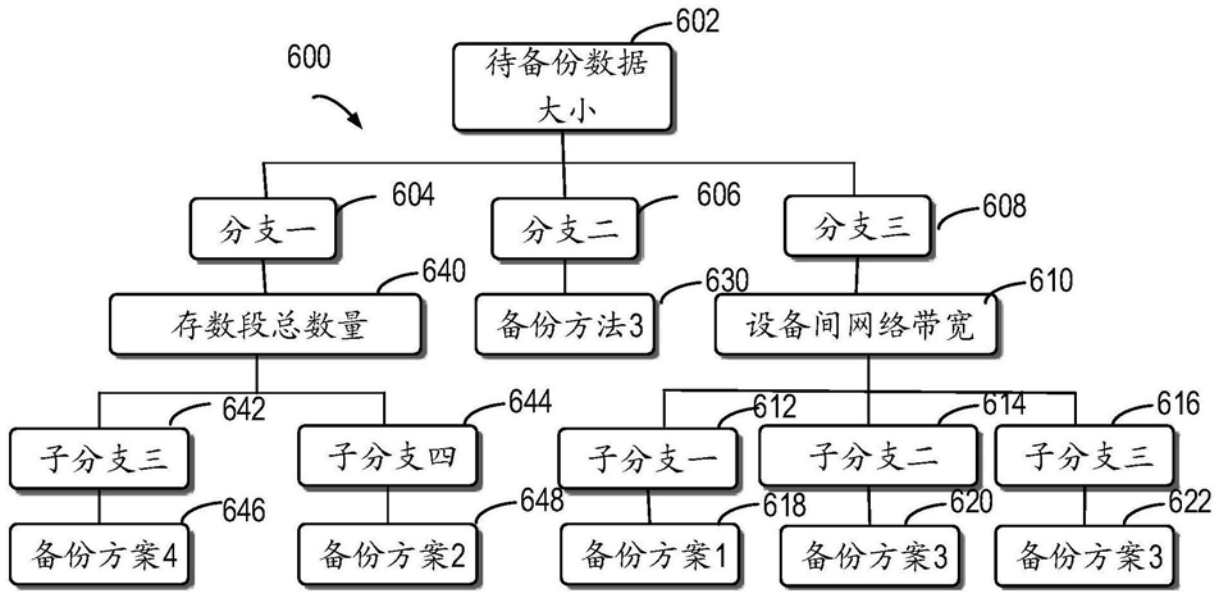


图6

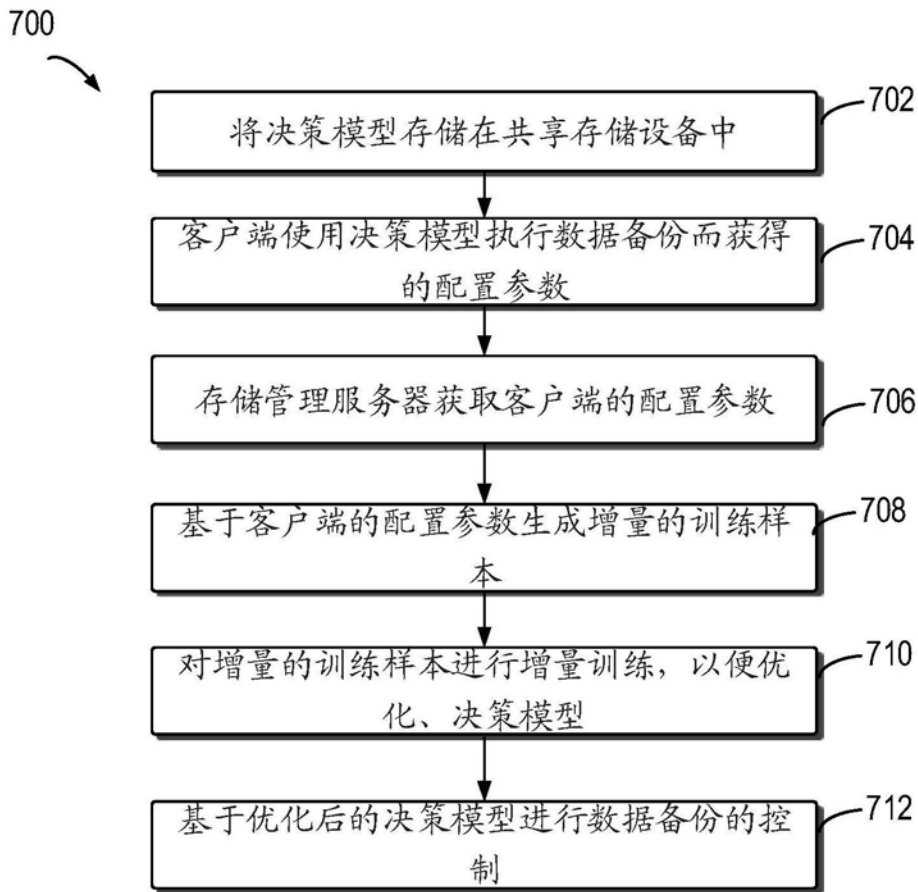


图7

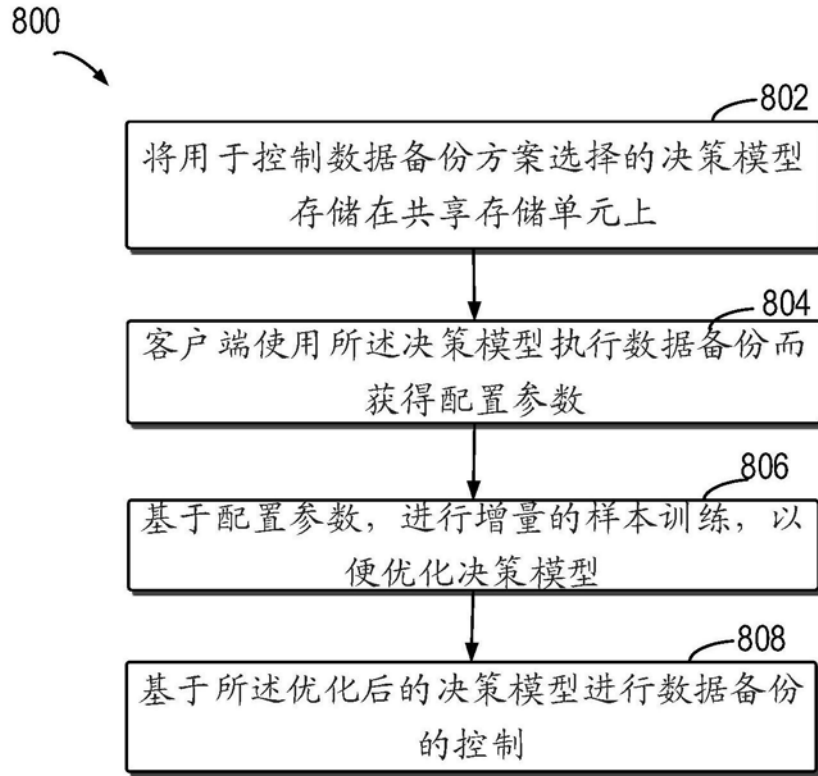


图8

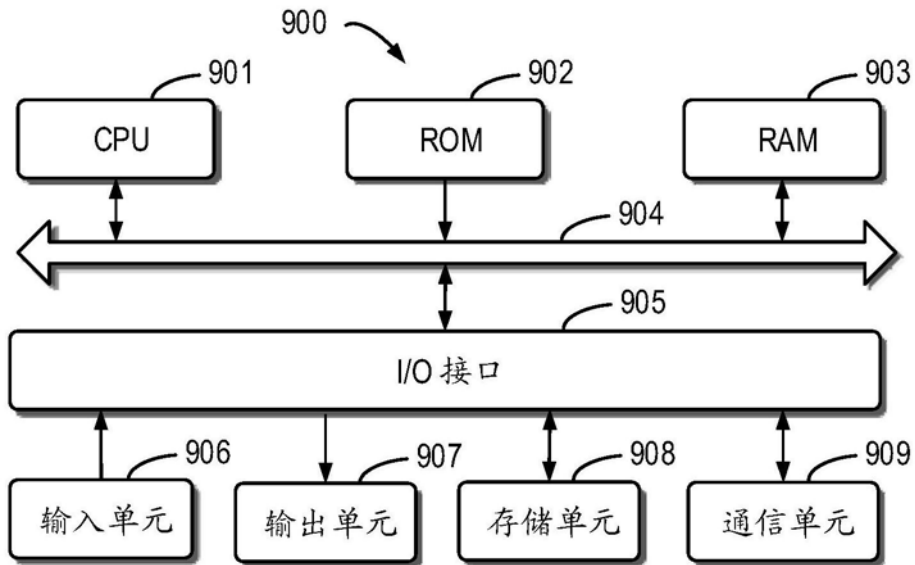


图9

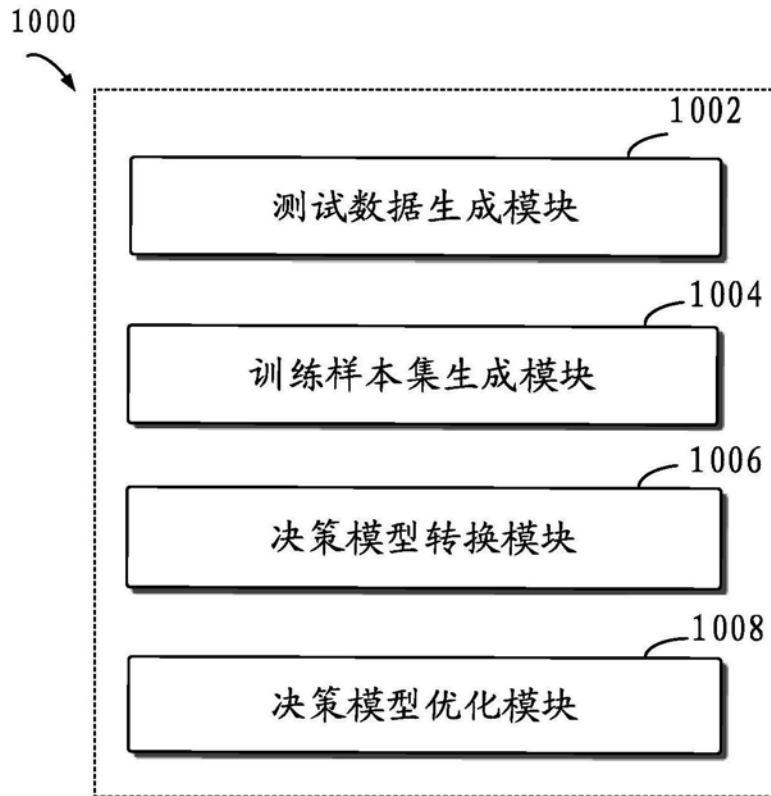


图10

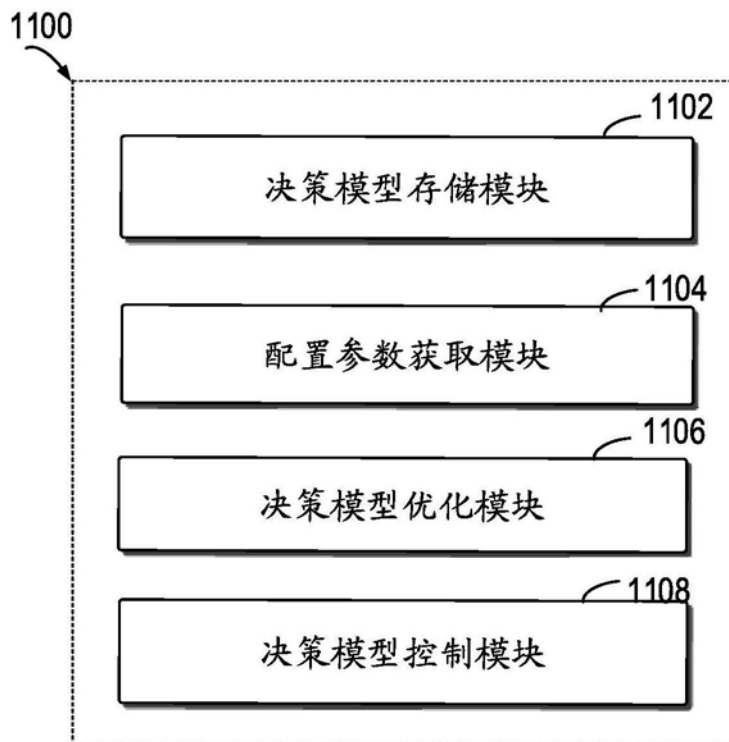


图11