



(12) 发明专利申请

(10) 申请公布号 CN 104572306 A

(43) 申请公布日 2015. 04. 29

(21) 申请号 201510044203. 5

(22) 申请日 2015. 01. 28

(71) 申请人 中国石油集团川庆钻探工程有限公司地球物理勘探公司

地址 610213 四川省成都市双流县华阳镇华阳大道一段 216 号川庆地球物理勘探公司科技部

(72) 发明人 汤成兵 严飞 郭玲

(74) 专利代理机构 北京铭硕知识产权代理有限公司 11286

代理人 谭昌驰 张川绪

(51) Int. Cl.

G06F 9/50(2006. 01)

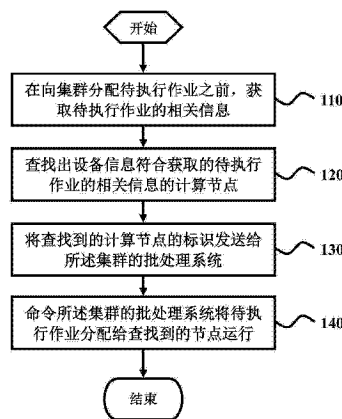
权利要求书2页 说明书6页 附图2页

(54) 发明名称

计算机集群的资源管理方法、资源管理器

(57) 摘要

本发明提供了一种计算机集群的资源管理方法、资源管理器,所述方法包括:在向所述集群的批处理系统管理的各个计算节点分配待执行作业之前,获取待执行作业的相关信息;从所述集群的批处理系统管理的各个计算节点中,查找出设备信息符合获取的待执行作业的相关信息的计算节点;将查找到的计算节点的标识发送给所述集群的批处理系统;命令所述集群的批处理系统将待执行作业分配给查找到的计算节点的标识所指示的计算机运行。



1. 一种计算机集群的资源管理方法,其特征在于,包括:

(A) 在向所述集群的批处理系统管理的各个计算节点分配待执行作业之前,获取待执行作业的相关信息;

(B) 从所述集群的批处理系统管理的各个计算节点中,查找设备信息符合获取的待执行作业的相关信息的计算节点;

(C) 将查找到的计算节点的标识发送给所述集群的批处理系统;

(D) 命令所述集群的批处理系统将待执行作业分配给查找到的计算节点的标识所指示的计算机运行。

2. 如权利要求 1 所述的方法,其特征在于,还包括:

如果没有查找出设备信息符合获取的待执行作业的相关信息的计算节点,则执行空闲节点的搜索步骤:从所述集群的批处理系统管理的各个计算节点中,查找空闲的计算节点;

按照获取的待执行作业的相关信息配置查找到的空闲的计算节点的计算环境;

当完成所述计算环境的配置时,将配置过的计算节点的标识发送给所述集群的批处理系统;

命令所述集群的批处理系统将待执行作业分配给配置过的计算节点的标识所指示的计算机运行。

3. 如权利要求 2 所述的方法,其特征在于,所述将配置过的计算节点的标识发送给所述集群的批处理系统的步骤包括:

检测配置过的计算节点的设备信息;

如果检测的设备信息符合获取的待执行作业的相关信息,则将配置过的计算节点的标识发送给所述集群的批处理系统;

如果检测的设备信息不符合获取的待执行作业的相关信息,则返回执行空闲节点的搜索步骤。

4. 如权利要求 2 或 3 所述的方法,其特征在于,还包括:

如果没有查找出空闲的计算节点,则在预定时间段之后,返回执行步骤 (B)。

5. 如权利要求 2 或 3 所述的方法,其特征在于,还包括:

当配置过的计算节点完成待执行作业或者配置过的计算节点的设备信息不符合获取的待执行作业的相关信息时,释放配置过的计算节点在配置所述计算环境时加载的资源,以便将配置过的计算节点恢复到原先的设备状态。

6. 一种计算机集群的资源管理器,其特征在于,包括:

作业获取单元,用于在向所述集群的批处理系统管理的各个计算节点分配待执行作业之前,获取待执行作业的相关信息;

第一查找单元,用于从所述集群的批处理系统管理的各个计算节点中,查找设备信息符合获取的待执行作业的相关信息的计算节点;

作业分配单元,用于将查找到的计算节点的标识发送给所述集群的批处理系统,并命令所述集群的批处理系统将待执行作业分配给查找到的计算节点的标识所指示的计算机运行。

7. 如权利要求 6 所述的资源管理器,其特征在于,还包括:

第二查找单元,用于如果没有查找出设备信息符合获取的待执行作业的相关信息的计算节点,则执行空闲节点的搜索步骤:从所述集群的批处理系统管理的各个计算节点中,查找空闲的计算节点;

环境配置单元,用于按照获取的待执行作业的相关信息配置查找到的空闲的计算节点的计算环境;

所述作业分配单元在完成所述计算环境的配置时,将配置过的计算节点的标识发送给所述集群的批处理系统,并命令所述集群的批处理系统将待执行作业分配给配置过的计算节点的标识所指示的计算机运行。

8. 如权利要求 7 所述的资源管理器,其特征在于,所述作业分配单元包括:

信息检测单元,用于检测配置过的计算节点的设备信息;

第一处理单元,用于如果检测的设备信息符合获取的待执行作业的相关信息,则将配置过的计算节点的标识发送给所述集群的批处理系统;

第二处理单元,用于如果检测的设备信息不符合获取的待执行作业的相关信息,则返回所述第二查找单元中执行空闲节点的搜索步骤。

9. 如权利要求 7 或 8 所述的资源管理器,其特征在于,还包括:

第三处理单元,用于如果没有查找出空闲的计算节点,则在预定时间段之后,并返回所述第一查找单元执行步骤 (B)。

10. 如权利要求 7 或 8 所述的资源管理器,其特征在于,还包括:

资源回收单元,用于当配置过的计算节点完成待执行作业或者配置过的计算节点的设备信息不符合获取的待执行作业的相关信息时,释放配置过的计算节点在配置所述计算环境时加载的资源,以便将查找到的计算节点恢复到原先的设备状态。

计算机集群的资源管理方法、资源管理器

技术领域

[0001] 本申请属于高性能计算领域,特别涉及一种在计算机集群中实现动态调度管理的技术。

背景技术

[0002] 随着现代计算机技术的发展,地震勘探、预测模型的构造和模拟、工业设计以及自动化等各个领域的计算系统规模越来越大,这些大型计算系统中包含了各种应用系统、各种计算类型及海量计算资源。

[0003] 为管理协调这些计算资源,通常采用 PBS(portable batch system) 资源管理器(即,批处理系统)解决大型计算系统中计算作业的批量提交问题,目前有三个版本:OpenPBS、Torque 以及 PBS pro,其中 Torque 为 PBS 的开源版本。目前大部分计算应用软件都通过 PBS 完成批量作业的提交及资源管理。

[0004] 然而,现有的批处理系统在实际应用中,通常存在以下问题:

[0005] 1、各应用软件的批处理系统间无法通信,导致各应用软件的计算资源边界模糊,容易引起资源竞争失败。

[0006] 2、批处理系统策略预置,灵活性差。批处理系统所需要的计算节点、应用软件、环境配置、存储等必须事先配置好,一旦不满足提交条件,则会导致批处理作业提交失败,中断批处理作业提交进程。

[0007] 3、批处理系统没有作业运行条件检查功能,即使提交的作业进入计算实体,仍然有可能因为运行条件不满足,从而导致批处理作业运行失败,造成作业重复提交,引起资源浪费。

[0008] 因此,如何管理协调这些计算资源,最大限度的利用资源,在这类大规模计算系统中显得尤其重要。

发明内容

[0009] 本发明的目的在于提供一种计算机集群的资源管理方法、资源管理器,以解决上述问题。

[0010] 根据本发明的一方面,提供一种计算机集群的资源管理方法,所述方法包括:(A) 在向所述集群的批处理系统管理的各个计算节点分配待执行作业之前,获取待执行作业的相关信息;(B) 从所述集群的批处理系统管理的各个计算节点中,查找设备信息符合获取的待执行作业的相关信息的计算节点;(C) 将查找到的计算节点的标识发送给所述集群的批处理系统;(D) 命令所述集群的批处理系统将待执行作业分配给查找到的计算节点的标识所指示的计算机运行。

[0011] 优选地,所述方法还包括:如果没有查找出设备信息符合获取的待执行作业的相关信息的计算节点,则执行空闲节点的搜索步骤:从所述集群的批处理系统管理的各个计算节点中,查找空闲的计算节点;按照获取的待执行作业的相关信息配置查找到的空闲的

计算节点的计算环境；当完成所述计算环境的配置时，将配置过的计算节点的标识发送给所述集群的批处理系统；命令所述集群的批处理系统将待执行作业分配给配置过的计算节点的标识所指示的计算机运行。

[0012] 优选地，所述将配置过的计算节点的标识发送给所述集群的批处理系统的步骤包括：检测配置过的计算节点的设备信息；如果检测的设备信息符合获取的待执行作业的相关信息，则将配置过的计算节点的标识发送给所述集群的批处理系统；如果检测的设备信息不符合获取的待执行作业的相关信息，则返回执行空闲节点的搜索步骤。

[0013] 优选地，所述方法还包括：如果没有查找出空闲的计算节点，则在预定时间段之后，返回执行步骤 (B)。

[0014] 优选地，所述方法还包括：当配置过的计算节点完成待执行作业或者配置过的计算节点的设备信息不符合获取的待执行作业的相关信息时，释放配置过的计算节点在配置所述计算环境时加载的资源，以便将配置过的计算节点恢复到原先的设备状态。

[0015] 根据本发明的另一方面，提供一种计算机集群的资源管理器，所述资源管理器包括：作业获取单元，用于在向所述集群的批处理系统管理的各个计算节点分配待执行作业之前，获取待执行作业的相关信息；第一查找单元，用于从所述集群的批处理系统管理的各个计算节点中，查找设备信息符合获取的待执行作业的相关信息的计算节点；作业分配单元，用于将查找到的计算节点的标识发送给所述集群的批处理系统，并命令所述集群的批处理系统将待执行作业分配给查找到的计算节点的标识所指示的计算机运行。

[0016] 优选地，所述资源管理器还包括：第二查找单元，用于如果没有查找出设备信息符合获取的待执行作业的相关信息的计算节点，则执行空闲节点的搜索步骤；从所述集群的批处理系统管理的各个计算节点中，查找空闲的计算节点；环境配置单元，用于按照获取的待执行作业的相关信息配置查找到的空闲的计算节点的计算环境；所述作业分配单元在完成所述计算环境的配置时，将配置过的计算节点的标识发送给所述集群的批处理系统，并命令所述集群的批处理系统将待执行作业分配给配置过的计算节点的标识所指示的计算机运行。

[0017] 优选地，所述作业分配单元包括：信息检测单元，用于检测配置过的计算节点的设备信息；第一处理单元，用于如果检测的设备信息符合获取的待执行作业的相关信息，则将配置过的计算节点的标识发送给所述集群的批处理系统；第二处理单元，用于如果检测的设备信息不符合获取的待执行作业的相关信息，则返回所述第二查找单元中执行空闲节点的搜索步骤。

[0018] 优选地，所述资源管理器还包括：第三处理单元，用于如果没有查找出空闲的计算节点，则在预定时间段之后，并返回所述第一查找单元执行步骤 (B)。

[0019] 优选地，所述资源管理器还包括：资源回收单元，用于当配置过的计算节点完成待执行作业或者配置过的计算节点的设备信息不符合获取的待执行作业的相关信息时，释放配置过的计算节点在配置所述计算环境时加载的资源，以便将查找到的计算节点恢复到原先的设备状态。

[0020] 与现有技术相比，本发明不仅提高了作业在批处理过程的灵活性、消除了集群中计算资源的孤岛，而且进一步提升了大型计算系统的易用性，提高资源利用率。

附图说明

[0021] 通过下面结合附图进行的描述,本发明的上述和其他目的和特点将会变得更加清楚,其中:

[0022] 图 1 是示出根据本发明的示例性实施例的计算机集群的资源管理方法的流程图;

[0023] 图 2 是示出根据本发明的示例性实施例的计算机集群的资源管理器的结构框图;

[0024] 图 3 是示出根据本发明的示例性实施例的在 Torque 批处理系统中实现动态调度的示意图。

具体实施方式

[0025] 以下,将参照附图来详细说明本发明的实施例。

[0026] 图 1 示出了本发明一种计算机集群的资源管理方法的优选实施例的流程图。

[0027] 参照图 1,在步骤 110 中,在向所述集群的批处理系统管理的各个计算节点分配待执行作业之前,获取待执行作业的相关信息。

[0028] 通常,在待执行作业中包含有待执行作业所需计算节点数量、运行作业所需的应用软件名称、运行时对设备内存及缓冲区的要求、文件系统等各种参数信息,因此,在本发明中,可通过获取待执行作业的相关信息,分析待执行作业的需求,从而选择出合适的计算节点来执行作业,为筛选出合适的计算节点,在本发明一个可选的实施例中,可选择待执行作业的应用软件的种类、内存需求量、许可需求的数量、缓冲区大小、文件系统中的多个信息作为筛选计算节点的条件。

[0029] 在步骤 120 中,从所述集群的批处理系统管理的各个计算节点中,查找设备信息符合获取的待执行作业的相关信息的计算节点。

[0030] 在本发明中,计算节点的设备信息主要包括计算节点的资源信息和状态信息,例如,节点上安装的各种应用程序、内存容量、缓冲区的大小以及运行的状况等。

[0031] 在步骤 130 中,将查找到的计算节点的标识发送给所述集群的批处理系统,并在步骤 140 中,命令所述集群的批处理系统将待执行作业分配给查找到的计算节点的标识所指示的计算机运行。换句话讲,就是由所述集群的批处理系统来完成待执行作业的分配、监控以及提交等。

[0032] 这里,所述的标识可以是设备的名称、IP 地址以及编号等能够唯一区别集群中各个计算机设备的标识信息。优选地,可利用集群的批处理系统提供预定的作业命令接口,将查找到的标识作为参数传递给所述集群的批处理系统。以下示出了 PBS 批处理系统中提供的 qsub 命令接口:

```
[0033] qsub[-a date_time][-A account_string][-e path][-h][-I][-j join]
[-k keep][-l resource_list][-m mail_options][-n Node_allocation_Method[-L
v1,[v2,[v3,[v4]]]]][-M user_list][-N name][-o path][-p priority][-q pool][-r
y|n][-u user_list][-v variable_list][-V][script]
```

[0034] 具体实施时,可通过 qsub 命令接口中的 resource_list 参数将查找到的计算节点的标识信息传递给批处理系统,批处理系统在接收到 qsub 命令后,可通过解析 resource_list 参数获取到查找到的计算节点。

[0035] 需要说明的是,本发明包括但不限于 qsub 命令接口来实现信息的传递,也可利用

系统提供的其他通信接口,将节点信息发送给批处理系统。

[0036] 在上述实施过程中,可能会出现集群中现有的计算资源都不符合计算需求,但是集群中仍然存在空闲的计算节点,只是这种空闲节点不能运行待执行作业。为使得空闲的计算节点的资源也能得到利用,在图 1 所示的实施例中,所述方法还包括:如果没有查找出设备信息符合获取的待执行作业的相关信息的计算节点,则从所述集群的批处理系统管理的各个计算节点中,查找空闲的计算节点;按照获取的待执行作业的相关信息配置查找到的空闲的计算节点的计算环境;当完成所述计算环境的配置时,将配置过的计算节点的标识发送给所述集群的批处理系统;命令所述集群的批处理系统将待执行作业分配给配置过的计算节点的标识所指示的计算机运行。

[0037] 具体实施时,可根据获取的待执行作业的相关信息中的应用软件的种类配置所述空闲的计算节点的计算环境(即,各种应用软件的安装),从而使得查找到的空闲的计算节点能够运行待执行作业。此外,计算环境的配置操作还包括:各种应用软件的挂接、配置文件的更改、后台进程的启动等操作。

[0038] 尽管被配置的计算节点的计算环境发生了变化,具备了运行作业的基本条件,但是,为进一步确保被配置的计算节点能够更好地运行待执行作业,还需要考虑运行作业所需的内存、缓冲区等其他相关条件,以便查找到更为合适的节点来运行待执行作业。在搜索空闲节点的实施例中,可检测配置过的计算节点的设备信息;如果检测的设备信息符合获取的待执行作业的相关信息,则将配置过的计算节点的标识发送给所述集群的批处理系统;如果检测的设备信息不符合获取的待执行作业的相关信息,则继续从所述集群的批处理系统管理的各个计算节点中,查找空闲的计算节点。换句话讲,就是通过检测被配置的计算节点的设备信息,来进一步确定该计算节点是否为执行作业所需的节点。

[0039] 由于集群中的节点资源在预定时间段之后会发生变化,例如,当有其他的作业完成之后,会将计算节点释放回集群中。因此,在搜索空闲节点的实施例中,所述方法还包括:如果没有查找出空闲的计算节点,则在预定时间段之后,并返回执行步骤 120。换言之,就是将待执行任务挂起,等待预定时间段之后,返回步骤 120 重新查找符合获取的待执行作业的相关信息的计算节点。

[0040] 此外,为合理地利用集群中的各个节点,在配置计算节点的实施例中,所述方法还包括:当配置过的计算节点完成待执行作业或者配置过的计算节点的设备信息不符合获取的待执行作业的相关信息时,释放配置过的计算节点在配置所述计算环境时加载的资源,以便将配置过的计算节点恢复到原先的设备状态。

[0041] 由此可见,上述实施过程是以待执行作业的需求为中心,计算机集群的计算资源(例如,计算节点、网络、存储、本地缓存区、许可服务器、后台进程等)为对象,通过对待执行作业的需求和计算机集群的计算资源进行对比分析,找出最佳资源配比,从而实现计算机集群的计算资源的动态分配与回收。

[0042] 图 2 示出了本发明一种计算机集群的资源管理器的优选实施例的结构框图。

[0043] 参照图 2,该资源管理器至少包括作业获取单元 201、第一查找单元 202 以及作业分配单元 203。

[0044] 其中,作业获取单元 201 在向所述集群的批处理系统管理的各个计算节点分配待执行作业之前,获取待执行作业的相关信息;第一查找单元 202 从所述集群的批处理系统

管理的各个计算节点中,查找出设备信息符合获取的待执行作业的相关信息的计算节点;作业分配单元 203 将查找到的计算节点的标识发送给所述集群的批处理系统,并命令所述集群的批处理系统将待执行作业分配给查找到的计算节点的标识所指示的计算机运行。

[0045] 通常,在待执行作业中包含有待执行作业所需计算节点数量、运行作业所需的应用软件名称、运行时对设备内存及缓冲区的要求、文件系统等各种参数信息,因此,在本发明中,可通过获取待执行作业的相关信息,分析待执行作业的需求,从而选择出合适的计算节点来执行作业,为筛选出合适的计算节点,在本发明一个可选的实施例中,可选择待执行作业的应用软件的种类、内存需求量、许可需求的数量、缓冲区大小、文件系统中的某一个或多个信息作为筛选计算节点的条件。

[0046] 在上述实施过程中,可能会出现集群中现有的计算资源都不符合计算需求,但是集群中仍然存在空闲的计算节点,只是这种空闲节点不能运行待执行作业。为使得空闲的计算节点的资源也能得到利用,在图 2 所示的实施例中,所述资源管理器还包括:第二查找单元(图中未示出),用于如果没有查找出设备信息符合获取的待执行作业的相关信息的计算节点,则执行空闲节点的搜索步骤:从所述集群的批处理系统管理的各个计算节点中,查找空闲的计算节点;环境配置单元(图中未示出),用于按照获取的待执行作业的相关信息配置查找到的空闲的计算节点的计算环境;所述作业分配单元 203 在完成所述计算环境的配置时,将配置过的计算节点的标识发送给所述集群的批处理系统,并命令所述集群的批处理系统将待执行作业分配给配置过的计算节点的标识所指示的计算机运行。

[0047] 具体实施时,环境配置单元可根据获取的待执行作业的相关信息中的应用软件的种类配置所述空闲的计算节点的计算环境(即,各种应用软件的安装),从而使得查找到的空闲的计算节点能够运行待执行作业。此外,计算环境的配置操作还包括:各种应用软件的挂接、配置文件的更改、后台进程的启动等操作

[0048] 尽管被配置的计算节点的计算环境发生了变化,具备了运行作业的基本条件,但是,为确保被配置的计算节点能够更好地运行待执行作业,还需要考虑运行作业所需的内存、缓冲区等其他相关条件,以便查找到更为合适的节点来运行待执行作业。在搜索空闲节点的实施例中,所述作业分配单元 203 包括:信息检测单元(图中未示出),用于检测配置过的计算节点的设备信息;第一处理单元(图中未示出),用于如果检测的设备信息符合获取的待执行作业的相关信息,则将配置过的计算节点的标识发送给所述集群的批处理系统;第二处理单元(图中未示出),用于如果检测的设备信息不符合获取的待执行作业的相关信息,则返回第二查找单元中执行空闲节点的搜索步骤。

[0049] 由于集群中的节点资源在预定时间段之后会发生变化,例如,当有其他的作业完成之后,会将计算节点释放会集群中。因此,在搜索空闲节点的实施例中,还包括:第三处理单元(图中未示出),用于如果没有查找出空闲的计算节点,则在预定时间段之后,并返回所述第一查找单元执行步骤(B)。

[0050] 此外,为合理地利用集群中的各个节点,在配置计算节点的实施例中,所述资源管理器还包括:资源回收单元(图中未示出),用于当配置过的计算节点完成待执行作业或者配置过的计算节点的设备信息不符合获取的待执行作业的相关信息时,释放配置过的计算节点在配置所述计算环境时加载的资源,以便将配置过的计算节点恢复到原先的设备状态。

[0051] 以下结合 Torque 批处理系统,对上述实施过程作进一步的说明。

[0052] 图 3 是示出了本发明的示例性实施例的在 Torque 批处理系统中实现动态调度的示意图。图中所示 301 为计算机集群的主节点,主节点上部署有 Torque 批处理系统的资源与作业服务器模块 (PBS-SERVER) 和作业调度器模块 (PBS-SCHED);图中所示 302 为计算机集群的各个计算节点,各个计算节点上部署有 Torque 批处理系统的作业执行模块 (pbs-mom);图中所示 303 为用户提交作业脚、请求的资源 (job scripts),图中所示 304 为本发明所述的资源管理器。

[0053] 通常情况下, Torque 批处理系统作业服务器模块 (PBS-SERVER) 和作业调度器模块 (PBS-SCHED) 是按照管理员设定的调度策略将用户提交的作业脚、请求的资源 (job scripts),通过作业执行模块 (pbs-mom) 分配给管理员指定的计算节点来运行。由于批处理系统中的调度策略是管理员设定的,灵活性差。这意味着 Torque 批处系统不能按照待执行作业的需求实现动态分配和回收。

[0054] 为了在 Torque 批处理系统中实现计算机集群的计算资源的动态分配与回收,本发明提供了如图 3 所示的资源管理器 304。

[0055] 参照图 3,资源管理器 304 的工作流程如下:

[0056] (1) 资源管理器 304 的 Read_param 模块在 Torque 批处理系统向其管理的各个计算节点分配待执行作业之前,从用户提交的作业脚、请求的资源 (job scripts) 中,获取待执行作业的各项指标的参数信息。

[0057] (2) 资源管理器 304 的 Node_status_check 模块扫描并检测集群中的各个计算节点的设备信息,从而查找出设备信息满足获取的待执行作业的相关信息的计算节点。

[0058] (3) 如果 Node_status_check 模块没有查找出设备信息满足获取的待执行作业的相关信息的计算节点,则查找出集群中空闲的计算节点,通过 App_config 模块按照获取的待执行作业的相关信息配置所述空闲的计算节点的计算环境。

[0059] (4) 为确保配置后的节点能全面满足作业运行的需要,可再次使用资源管理器 304 的 Node_status_check 模块对配置过的计算节点的设备信息进行检测,以便确定该计算节点能否满足作业运行的需要。

[0060] (5) 资源管理器 304 的 Scripts_startup 模块利用 Torque 批处理系统的 qsub 作业命令接口,将查找到的满足获取的待执行作业的相关信息的计算节点的标识信息发送给 Torque 批处理系统,并命令 Torque 批处理系统将待执行作业分配给查找到的计算节点的标识所指示的计算机运行。

[0061] (6) 为了将配置过的计算节点恢复到配置之前的设备状态,资源管理器 304 的 res_recycle 模块在配置过的计算节点完成待执行作业或者配置过的计算节点的设备信息不符合获取的待执行作业的相关信息时,释放配置过的计算节点在配置计算环境时加载的各种资源。

[0062] 与现有技术相比,本发明不仅提高了作业在批处理过程的灵活性、消除了集群中计算资源的孤岛,而且进一步提升了大型计算系统的易用性,提高资源利用率。

[0063] 尽管已参照优选实施例为和描述了本发明,但本领域技术人员应该理解,在不脱离由权利要求限定的本发明的精神和范围的情况下,可以对这些实施例进行各种修改和变换。

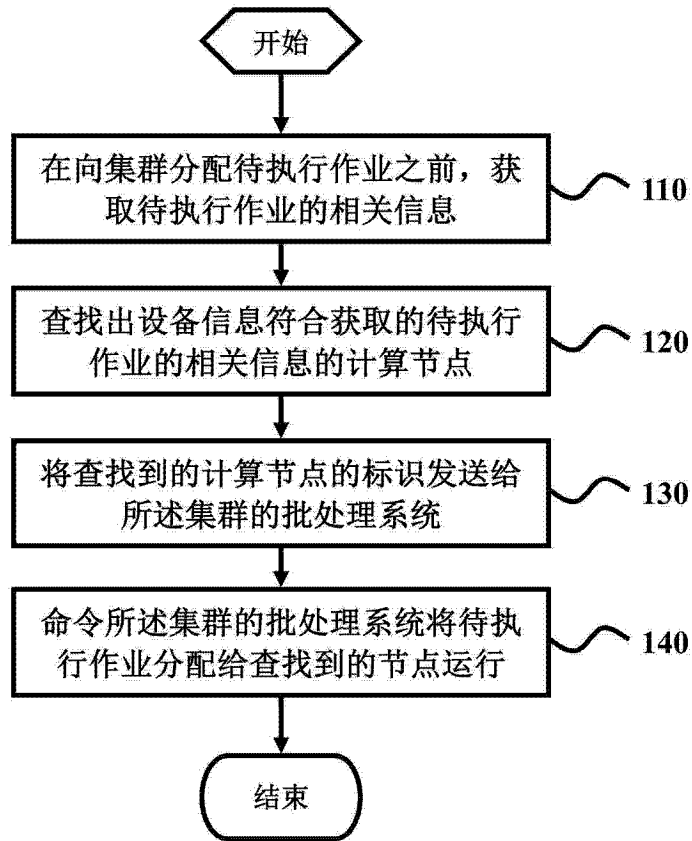


图 1

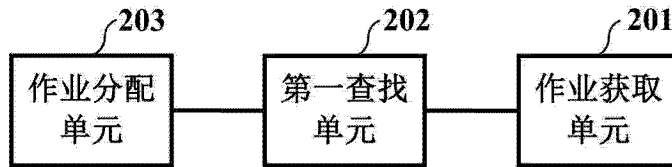


图 2

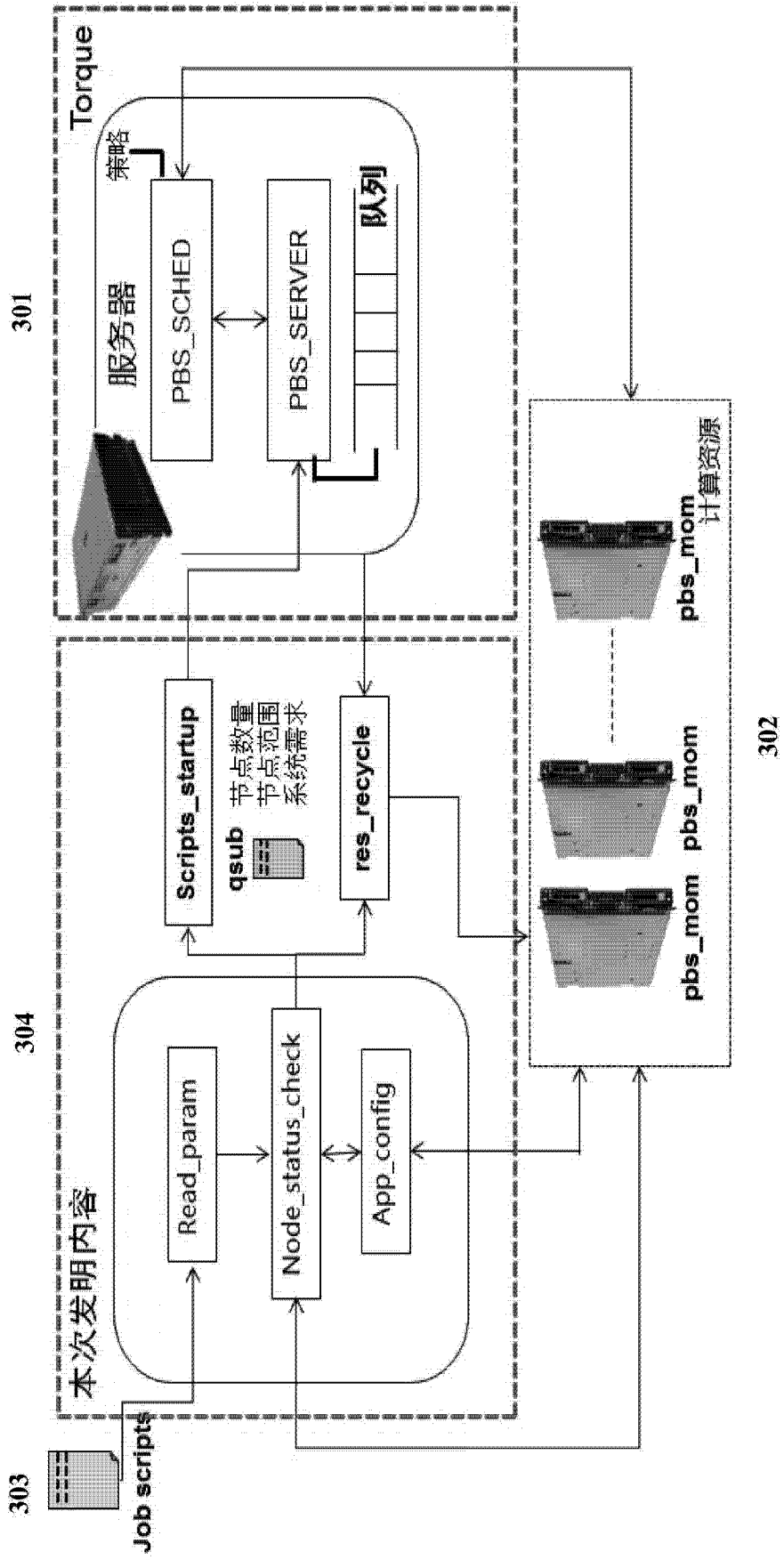


图 3