



(12) 发明专利申请

(10) 申请公布号 CN 115336236 A

(43) 申请公布日 2022. 11. 11

(21) 申请号 202080098258.3

(51) Int.Cl.

(22) 申请日 2020.03.31

H04L 45/00 (2022.01)

(85) PCT国际申请进入国家阶段日
2022.09.08

(86) PCT国际申请的申请数据
PCT/CN2020/082519 2020.03.31

(87) PCT国际申请的公布数据
W02021/195990 EN 2021.10.07

(71) 申请人 阿里巴巴集团控股有限公司
地址 英属开曼群岛大开曼资本大厦一座四
层847号邮箱

(72) 发明人 叶剑西 王绍创 冉仟元 冯飞
董建波

(74) 专利代理机构 北京同钧律师事务所 16037
专利代理师 柴海平 许怀远

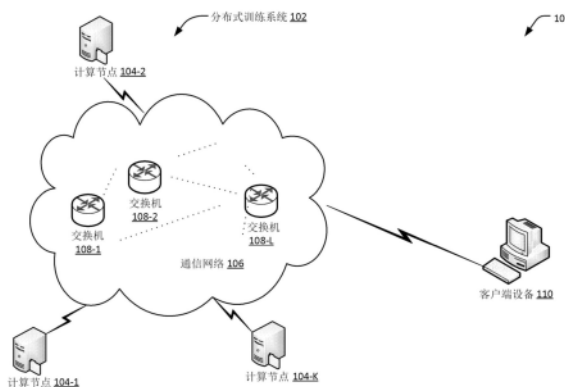
权利要求书3页 说明书25页 附图14页

(54) 发明名称

通过基于环的集群通信避免网络拥塞

(57) 摘要

在分布式训练中,为了避免网络拥塞,第一计算节点可以至少部分地基于与第一进程相关联的第一网络接口控制器和与第二进程相关联的第二网络接口控制器是否位于同一计算节点中或链接到同一片交换机来确定将数据从第一进程路由到第二进程的路由标识符,第一进程和第二进程属于连接网络拓扑下的多个不同节点的特定节点间环。然后,第一计算节点可以根据路由标识符将数据从第一进程路由到第二进程。



1. 一种由第一计算节点实施的方法,所述方法包括:

至少部分地基于与第一进程相关联的第一网络接口控制器和与第二进程相关联的第二网络接口控制器是否位于同一计算节点中或链接到同一片交换机来确定将数据从所述第一进程路由到所述第二进程的路由标识符,所述第一进程和所述第二进程属于网络拓扑下连接多个不同节点的特定节点间环;以及

根据所述路由标识符将所述数据从所述第一进程路由到所述第二进程。

2. 根据权利要求1所述的方法,其中与所述第一进程相关联的所述第一网络接口控制器被配置为向环形拓扑中的第二计算节点发送数据或从环形拓扑中的第二计算节点接收数据,所述第二计算节点不同于所述第一计算节点。

3. 根据权利要求1所述的方法,其中与所述第一进程相关联的所述第一网络接口控制器还与一个或多个进程相关联,其中从所述第一进程和所述一个或多个进程发送的数据通过所述第一网络接口控制器发送。

4. 根据权利要求1所述的方法,其中所述网络拓扑包括胖树拓扑。

5. 根据权利要求1所述的方法,还包括:响应于确定与所述第一进程相关联的所述网络接口控制器和与所述第二进程相关联的所述第二网络接口控制器位于所述同一计算节点中或链接到所述同一片交换机,将所述路由标识符设置为默认标识符。

6. 根据权利要求1所述的方法,还包括:响应于确定与所述第一进程相关联的所述第一网络接口控制器和与所述第二进程相关联的所述第二网络接口控制器位于不同的计算节点或链接到不同的片交换机,将所述路由标识符设置为等于与所述第一进程相关联的所述第一网络接口控制器的标识符。

7. 根据权利要求6所述的方法,其中根据所述路由标识符将所述数据从所述第一进程路由到所述第二进程包括:至少通过与所述第一网络接口控制器连接的片交换机和具有与所述第一网络接口控制器的所述标识符具有对应关系的标识符的汇聚交换机将所述数据从所述第一进程路由到所述第二进程,所述第一网络接口控制器与所述第一进程相关联。

8. 一个或多个机器可读介质,其存储有机器可读指令,当所述机器可读指令由第一计算节点执行时,使所述第一计算节点执行动作,所述动作包括:

至少部分地基于与第一进程相关联的第一网络接口控制器和与第二进程相关联的第二网络接口控制器是否位于同一计算节点中或链接到同一片交换机来确定将数据从所述第一进程路由到所述第二进程的路由标识符,所述第一进程和所述第二进程属于网络拓扑下连接多个不同节点的特定节点间环;以及

根据所述路由标识符将所述数据从所述第一进程路由到所述第二进程。

9. 根据权利要求8所述的一个或多个机器可读介质,其中与所述第一进程相关联的所述第一网络接口控制器被配置为向环形拓扑中的第二计算节点发送数据或从环形拓扑中的第二计算节点接收数据,所述第二计算节点不同于所述第一计算节点。

10. 根据权利要求8所述的一个或多个机器可读介质,其中与所述第一进程相关联的所述第一网络接口控制器还与一个或多个进程相关联,其中从所述第一进程和所述一个或多个进程发送的数据通过所述第一网络接口控制器发送。

11. 根据权利要求8所述的一个或多个机器可读介质,其中所述网络拓扑包括胖树拓扑。

12. 根据权利要求8所述的一个或多个机器可读介质,其中所述动作还包括:响应于确定与所述第一进程相关联的所述第一网络接口控制器和与所述第二进程相关联的所述第二网络接口控制器位于所述同一计算节点中或链接到所述同一片交换机,将所述路由标识符设置为默认标识符。

13. 根据权利要求8所述的一个或多个机器可读介质,其中所述动作还包括:响应于确定与所述第一进程相关联的所述第一网络接口控制器和与所述第二进程相关联的所述第二网络接口控制器位于不同的计算节点或链接到不同的片交换机,将所述路由标识符设置为等于与所述第一进程相关联的所述第一网络接口控制器的标识符。

14. 根据权利要求13所述的一个或多个机器可读介质,其中根据所述路由标识符将所述数据从所述第一进程路由到所述第二进程包括:至少通过与所述第一网络接口控制器连接的片交换机和具有与所述第一网络接口控制器的所述标识符具有对应关系的标识符的汇聚交换机将所述数据从所述第一进程路由到所述第二进程,所述第一网络接口控制器与所述第一进程相关联。

15. 一种第一计算节点,包括:

一个或多个处理单元;和

存储器,其存储机器可执行指令,当所述机器可执行指令由一个或多个处理单元执行时,使所述一个或多个处理单元执行动作,所述动作包括:

至少部分地基于与第一进程相关联的第一网络接口控制器和与第二进程相关联的第二网络接口控制器是否位于同一计算节点中或链接到同一片交换机来确定将数据从所述第一进程路由到所述第二进程的路由标识符,所述第一进程和所述第二进程属于网络拓扑下连接多个不同节点的特定节点间环;以及

根据所述路由标识符将所述数据从所述第一进程路由到所述第二进程。

16. 根据权利要求15所述的第一计算节点,其中与所述第一进程相关联的所述第一网络接口控制器被配置为向环形拓扑中的第二计算节点发送数据或从环形拓扑中的第二计算节点接收数据,所述第二计算节点不同于所述第一计算节点。

17. 根据权利要求15所述的第一计算节点,其中与所述第一进程相关联的所述第一网络接口控制器还与一个或多个进程相关联,其中从所述第一进程和所述一个或多个进程发送的数据通过所述第一网络接口控制器发送。

18. 根据权利要求15所述的第一计算节点,其中所述动作还包括:响应于确定与所述第一进程相关联的所述第一网络接口控制器和与所述第二进程相关联的所述第二网络接口控制器位于所述同一计算节点中或链接到所述同一片交换机,将所述路由标识符设置为默认标识符。

19. 根据权利要求15所述的第一计算节点,其中所述动作还包括:响应于确定与所述第一进程相关联的所述第一网络接口控制器和与所述第二进程相关联的所述第二网络接口控制器位于不同的计算节点或链接到不同的片交换机,将所述路由标识符设置为等于与所述第一进程相关联的所述第一网络接口控制器的标识符。

20. 根据权利要求19所述的第一计算节点,其中根据所述路由标识符将所述数据从所述第一进程路由到所述第二进程包括:至少通过与所述第一网络接口控制器连接的片交换机和具有与所述第一网络接口控制器的所述标识符具有对应关系的标识符的汇聚交换机

将所述数据从所述第一进程路由到所述第二进程,所述第一网络接口控制器与所述第一进程相关联。

通过基于环的集群通信避免网络拥塞

背景技术

[0001] 随着诸如深度神经网络 (Deep Neural Network, DNN) 的神经网络迅速发展, 各种应用领域 (例如, 计算机视觉、自然语言处理、语音识别等) 都得到了发展, 并且会从神经网络固有的多功能性和灵活性中受益。然而, 由于神经网络应用日益增加的复杂性和越来越严格的准确性要求, 神经网络模型的大小和训练模型所需的训练数据的大小也显著增加, 这将不可避免地导致训练时间越来越长, 从而对训练模型能够满足不断变化的应用环境的有效性和及时性产生不利影响。

[0002] 为了减少训练神经网络模型的时间, 可以使用一种采用并行训练的分布式训练系统。一般而言, 分布式训练系统可以包括分布在网络上的大量计算节点或服务器, 并且将计算任务的子集分配给计算节点或服务器, 用于采用并行训练来执行计算。然而, 分布式训练系统中的计算节点或服务器之间的数据通信造成了分布式训练系统中可能发生的训练时间的减少量的下限或瓶颈。当分布式训练系统包括计算节点或服务器内部的和之间的各种类型的异构连接或互连时, 尤其如此, 这些异构连接或互连在延迟、带宽、拓扑等方面表现出不同的特性。这种连接或互连的异构性增加了为分布式训练系统中的计算节点或服务器设计数据通信网络的难度和复杂性。

[0003] 此外, 由于过量的数据流通过分布式训练系统中的计算节点或服务器之间的特定网络交换机或连接, 可能引发网络拥塞, 这样可能会由于处理训练结果的延迟而导致训练时间延长。之所以有过量的数据流通过特定网络交换机或连接, 可能是因为计算节点或服务器之间发送的路由数据的路径选择失去控制。

附图说明

[0004] 参考附图进行详细描述。附图中, 附图标记的最左侧数字表示该附图标记第一次出现的附图。不同附图中使用相同的附图标记表示相似或相同的指代。

[0005] 图1示出了分布式训练系统可应用的示例环境。

[0006] 图2示出了更详细的示例性计算节点。

[0007] 图3A示出了将预设数量的节点互相连接的环形配置。

[0008] 图3B示出了将预设数量的节点互相连接的减半加倍配置。

[0009] 图4示出了示例性集群通信库的示意图。

[0010] 图5示出了示例性拓扑感知多阶段算法。

[0011] 图6示出了用于计算节点的节点内归约散布阶段的示例性基于环的算法。

[0012] 图7示出了用于计算节点的节点内归约散布阶段的示例性减半加倍算法。

[0013] 图8示出了节点间全局归约阶段的示例性减半加倍算法。

[0014] 图9示出了更详细的节点间全局归约阶段的示例性减半加倍算法。

[0015] 图10示出了示例性基于环的集群通信算法。

[0016] 图11示出了用并行或重叠的方式执行节点内归约散布阶段、节点间全局归约阶段、节点内全局聚集阶段的示例场景。

- [0017] 图12示出了示例性胖树 (fat-tree) 网络拓扑。
- [0018] 图13示出了使用第一拥塞避免方法的示例场景。
- [0019] 图14示出了使用第二拥塞避免方法的示例场景。
- [0020] 图15示出了示例性拓扑感知多阶段方法。
- [0021] 图16示出了第一示例性网络拥塞避免方法。
- [0022] 图17示出了第二示例性网络拥塞避免方法。
- [0023] 图18示出了分布式训练中基于混合架构的示例性并行方法。

具体实施方式

[0024] 概述

[0025] 如上所述,现有的分布式训练系统由于分布式训练系统中计算节点间的数据通信而对良好可扩展性造成了性能瓶颈。此外,由于网络结构的多样性(包括例如以太网、无限带宽(InfiniBand)、PCIe、NVLink、NVSwitch、QPI/UPI等)和网络特征(例如延迟、带宽和拓扑等)的高度差异,分布式训练系统通常不能很好地利用这种异构类型的连接或互连来执行计算节点中和计算节点之间的集群数据运算以及计算节点之间的数据传输。此外,由于对计算节点间发送的数据进行路由的路径选择可能失去控制,而发生网络拥塞,从而导致过量的数据流通过分布式训练系统中计算节点之间的特定网络交换机或连接,且导致由处理训练结果的延迟引发的训练时间延长。此外,现有的分布式训练系统未能区分用于集群运算的不同类型的底层结构的算法,因此导致较差的性能。

[0026] 本公开描述了示例性分布式训练系统。在实施中,示例性分布式训练系统可以采用结构感知的集群通信库,该库使得分布式训练系统能够线性扩展。在实施中,集群通信库可以至少部分地基于对底层结构和支持网络架构的分析来定制通信算法,以获得期望的或最大的效率。在实施中,分布式训练系统可以将基本运算分成多个子运算,每个子运算使用一种类型的结构。

[0027] 在实施中,示例性分布式训练系统可以实现混合算法,该混合算法允许多个算法在单个集群运算中共存,并且选择性地采用用于特定结构的算法来提高或最大化整个通信路径的效率。在实施中,分布式训练系统可以采用双进程并行算法,该算法启动两个并发进程,并流水线化(pipeline)节点内和节点间连接的使用,从而通过重叠节点内通信和节点间通信来提高通信效率。

[0028] 在实施中,示例性分布式训练系统可以采用基于探测的路由控制机制,该机制生成从连接到路径的映射,从而通过对集群运算中的参与者或过程重新排序并将分布式训练系统上的数据流映射到特定物理链路,来将连接分布或散布到通信网络中的不同汇聚或中间交换机,从而避免网络拥塞。

[0029] 本申请描述了多个不同的实施例和实施方式。以下部分描述了适用于实践各种实施方式的示例性框架。接下来,本申请描述了用于实现分布式训练系统的示例性系统、设备和过程。

[0030] 示例性环境

[0031] 图1示出了可用于实施分布式训练系统的示例性环境100。该环境100可以包括分布式训练系统102。在该示例中,分布式训练系统102可以包括多个计算节点或服务器104-

1、104-2、...、104-K(下文统称为计算节点104),其中K是大于1的正整数。在实施中,多个计算节点104可以通过通信网络106彼此通信数据。

[0032] 计算节点104可以实施为具有计算/处理和通信能力的任意多种计算设备,包括但不限于服务器、台式电脑、笔记本电脑或便携式电脑、手持设备、上网本、互联网设备、平板电脑、移动设备(如移动电话、个人数字助理、智能电话等)等,或其组合。

[0033] 通信网络106可以是无线或有线网络,或其组合。网络106可以是彼此互连并用作单个大型网络(例如,互联网或内联网)的独立网络的集合。这种独立网络的示例包括但不限于电话网络、电缆网络、局域网(Local Area Network,LAN)、广域网(Wide Area Network,WAN)和城域网(Metropolitan Area Network,MAN)。此外,独立网络可以是无线或有线网络,或其组合。有线网络可以包括电载波连接(例如通信电缆等)和/或光学载体或连接(例如光纤连接等)。无线网络可以包括例如WiFi网络、其他射频网络(例如蓝牙(Bluetooth®)、紫峰(Zigbee)等)等。在实施中,通信网络106可以包括用于提供计算节点104之间的连接的多个节点间互连或交换机108-1、108-2、...、108-L(下文统称为节点间交换机108),其中L是大于1的正整数。

[0034] 在实施中,环境100可以进一步包括客户端设备110。用户可以指示分布式训练系统102基于从客户端设备110发送给分布式训练系统102的数据,对特定的学习模型(例如深度神经网络模型)执行训练。

[0035] 示例性计算节点

[0036] 图2示出了更详细的计算节点104。在实施中,计算节点104可以包括但不限于一个或多个处理单元202、输入/输出(input/output,I/O)接口204和/或一个或多个网络接口206以及存储器208。在实施中,计算节点104可以进一步包括一个或多个节点内互连或交换机210。

[0037] 在实施中,处理单元202可以被配置成执行存储在存储器208中的和/或从输入/输出接口204和/或网络接口206接收的指令。在实施中,处理单元202可被实施为一个或多个硬件处理器,包括例如微处理器、专用指令集处理器、物理处理单元(Physics Processing Unit,PPU)、中央处理单元(Central Processing Unit,CPU)、图形处理单元、数字信号处理器、张量处理单元等。附加地或替代地,这里描述的功能可以至少部分地由一个或多个硬件逻辑组件来执行。例如,但不限于,可以使用的硬件逻辑组件的示例类型包括现场可编程门阵列(Field Programmable Gate Array,FPGA)、专用集成电路(Application-Specific Integrated Circuit,ASIC)、专用标准产品(Application-Specific Standard Product,ASSP)、片上系统(System-on-a-Chip System,SOC)、复杂可编程逻辑器件(Complex Programmable Logic Device,CPLD)等。

[0038] 存储器208可以包括易失性存储器形式的机器可读介质,例如随机存取存储器(Random Access Memory,RAM)和/或非易失性存储器,例如只读存储器(Read Only Memory,ROM)或闪存RAM。存储器208是机器可读介质的一种示例。

[0039] 机器可读介质可以包括易失性或非易失性类型、可移动或不可移动介质,其可以使用任何方法或技术来实现信息的存储。该信息可以包括机器可读指令、数据结构、程序模块或其他数据。机器可读介质的示例包括但不限于相变存储器(Phase-Change Memory,PRAM)、静态随机存取存储器(Static Random Access Memory,SRAM)、动态随机存取存储器

(Dynamic Random Access Memory, DRAM)、其他类型的随机存取存储器 (Random-Access Memory, RAM)、只读存储器 (Read-Only Memory, ROM)、电可擦除可编程只读存储器 (Electrically Erasable Programmable Read-Only Memory, EEPROM)、快速闪存或其他内部存储技术、光盘只读存储器 (Compact Disk Read-Only Memory, CD-ROM)、数字通用光盘 (Digital Versatile Disc, DVD) 或其他光存储器、盒式磁带、磁盘存储或其他磁存储设备、或任何其他可用于存储可被计算节点访问的信息的非传输介质。如这里所定义的, 机器可读介质不包括任何瞬态介质, 例如调制数据信号和载波。

[0040] 在实施中, 网络接口206可以被配置成通过通信网络106将计算节点104连接到其他计算节点。在实施中, 网络接口206可以通过网络接口控制器 (Network Interface Controller, NIC) 来建立, 该网络接口控制器可以使用硬件和软件来将计算节点104连接到通信网络106。在实施中, 每种类型的NIC可以使用不同类型的结构或连接器来连接到与通信网络106相关联的物理介质。在IEEE 802规范中可以找到各种类型的结构或连接器的示例, 并且可以包括例如以太网 (在802.3中定义)、令牌环网 (在802.5中定义)、无线网络 (在802.11中定义)、无限带宽 (InfiniBand) 等。

[0041] 在实施中, 节点内交换机210可以包括各种类型的互连或交换机, 其可以包括但不限于高速串行计算机扩展总线 (例如PCIe等)、串行多通道近距离通信链路 (例如Nolan, 其是基于有线的通信协议串行多通道近距离通信链路)、具有多个端口的交换机芯片 (例如NVSwitch等)、点对点处理器互连 (例如英特尔QPI/UPI等) 等。

[0042] 尽管在该示例中, 在计算节点104中仅描述了硬件组件, 但是在其他实例中, 异常检测系统110可以进一步包括其他硬件组件和/或其他软件组件, 例如执行存储在存储器208中的指令以执行各种操作的程序模块212, 以及用于存储接收的用于训练的数据、训练期间计算的中间和最终结果等的程序数据214。

[0043] 示例性集群通信算法

[0044] 图3A和3B示出了可以在分布式训练系统102中使用的示例性集群通信算法。在实施中, 集群通信算法可以包括, 但不限于, 基于环的通信算法、减半加倍通信算法等。

[0045] 图3A示出了环形配置, 该配置将预定数量的节点 (例如, N 个节点, 其中 N 是大于1的正整数) 与多个连接 (即, N 个连接) 互连, 并将数据 (例如, 数据包或消息) 划分成多个数据块 (即, N 个数据块) 用于传输, 且需要多个步骤 (在该示例中为 $N-1$ 个步骤) 的通信来完成集群运算。在每个步骤中, 节点可以从其相邻节点之一接收数据, 对所接收的数据进行特定操作以获得本地结果, 并且将所接收的数据转发给相邻节点中的另一个。在 $N-1$ 个步骤之后, 环中的每个节点具有来自环中其他节点的数据, 并且最终结果被散布到所有节点, 这需要另外的 $N-1$ 个步骤来广播各自的本地结果。对于每个节点, 转发的总数据大小是 $2S$, 其中 S 表示数据大小或消息大小。

[0046] 图3B示出了互连预定数量的节点 (例如, N 个节点, 其中 N 是大于1的正整数) 的减半加倍配置。在这种减半的配置中, 节点以成对的方式相互通信, 通信的每一步只需要 $N/2$ 个连接。在第一步中, 相邻节点被配对在一起, 向各自的配对节点发送一半的消息或数据, 并接收另一半的消息或数据进行处理。因此, 中间结果可以散布至配对节点。在随后的步骤中, 以增加或加倍的距离形成新的配对, 并且用于处理的数据大小减半。在 $\log_2 N$ 步的通信之后, 结果被散布在减半加倍配置中的所有节点中。然后将节点中的本地结果通过附加的

$\log_2 N$ 步通信广播到其他节点。

[0047] 示例性集群通信库

[0048] 图4示出了描绘可以由分布式训练系统102使用的示例集群通信库400的示意图。在实施中,集群通信库是一种设计成提供高性能、高可扩展性和强可用性的通信库,并且可以被配置成不仅对诸如全局归约 (Allreduce) 和全局聚集 (Allgather) 运算之类的标准集群运算提供支持,还对定制应用的其他自定义运算提供支持。在实施中,集群通信库400可以采用具有不同特性 (例如,延迟、带宽、拓扑) 的不同类型的互连或交换机,并提供一种机制来收集网络和计算节点中底层硬件的信息,从而可以基于一条或多条所收集的信息来开发拓扑感知算法的设计。

[0049] 在实施中,集群通信库400可以提供灵活性,以允许在单个运算中执行多个算法,并通过利用节点内通信和节点间通信之间的并行性来提高性能 (例如,通信和训练的性能等)。此外,集群通信库400可以在具有传统或新的映射算法的计算节点中利用多个NIC,并且通过连接的拓扑感知布置来消除网络拥塞。

[0050] 在实施中,集群通信库400可以包括软件栈402。在实施中,软件栈402可以包括多个组件,其中可以包括但不限于传输组件404、运算组件406、通信器组件408和库上下文组件410。在实施中,软件栈402可以用模块化方式来设计,以允许通用性和可扩展性。

[0051] 在实施中,传输组件404可以负责节点内和节点间通信中的对等 (Peer-to-Peer, P2P) 数据的转移或传输。作为示例而非限制,集群通信库400可以支持用于节点间通信的TCP (Transmission Control Protocol, 传输控制协议) 和RDMA (Remote Direct Memory Access, 远程直接存储器访问), 以及用于节点内通信的P2P结构, 例如PCIe (Peripheral Component Interconnect Express, 快速外围组件互连)、NVLink/NVSwitch和QPI/UPI (Quick Path Interconnect/Ultra Path Interconnect, 快速路径互连/超路径互连) 等。对于RDMA通信,传输组件404还可被配置成管理处理单元 (如图形处理单元 (Graphics Processing unit, GPU) 设备) 和主机存储器中的存储器区域 (Memory Region, MR) 和对应的存储器缓冲区。

[0052] 在实施中,运算组件406可以提供一组基本运算和各种网络算法。例如,基本运算可以被配置有集群通信库400所支持的算法。此外,运算组件406可允许用户基于这些基本运算定义新运算,以实现异构性感知运算,该运算可对每种类型的结构采用最优或更好的算法。

[0053] 在实施中,通信器组件408可以与软件进程相关联,并且可以被配置成在处理单元 (诸如GPU设备) 上执行操纵和处理。通信器组件408可以保存或记录关于其他对等组件的信息 (例如,排序ID、IP地址等), 并保持与对等组件的连接。在实施中,通信器组件408可以进一步收集节点内和节点间拓扑信息,并使用该信息来指导算法设计。在实施中,节点内信息可以包括但不限于互连的类型、处理单元的位置之间的距离、处理单元和网络接口控制器之间的距离等。在实施中,节点间信息可以包括,但不限于,例如,可用网络接口控制器的数量、集群或计算节点的拓扑、集群中计算节点的位置。

[0054] 在实施中,库上下文组件410可被配置为开放 (expose) 用于设置系统配置 (例如环境变量)、管理通信器组件408的一个或多个应用接口,并提供诸如日志记录等其他功能。

[0055] 此外,在一些情况下,集群通信库400可以进一步包括或提供用于拓扑感知设计、

测试和评估以及可用性改进的多个工具和实用程序412。作为示例而非限制，工具和实用程序412可包括用于传输组件404的性能测试工具以辅助算法设计和评估、用于确保系统可用性的基于探测的路由机制以及其他功能，例如可扩展到支持除GPU之外的设备的设备管理功能。

[0056] 用于集群通信的示例性拓扑感知多阶段算法

[0057] 在实施中，集群通信可以被定义为涉及一组处理单元或进程的通信，并且集群通信的运算可以由该组中包括的所有处理单元或进程一起执行。集群通信运算的实例可以包括但不限于全局归约运算、全局聚集运算、归约散布 (Reduce-Scatter) 运算等。在实施中，全局归约运算是分布式训练中集群通信的许多重要基础之一，并且涉及组中跨进程对数据执行归约。归约的示例可以包括但不限于求和运算、求平均值运算、求最大值运算、求最小值运算等。

[0058] 作为示例而非限制，这里以全局归约运算为示例来说明如何将集群运算划分成多个微运算或子运算。在实施中，分布式训练系统102可以采用拓扑感知的多阶段算法，该算法将全局归约运算分成多个微运算或子运算，并根据需要选择性地挑选一个或多个微运算或子运算，从而通过消除可能不需要的微运算或子运算来减少传输的数据量。在实施中，分布式训练系统102可以将集群通信算法从微运算或子运算中分离，并允许基于底层结构信息在算法和微运算或子运算之间进行独立或单独的匹配，从而通过减少传输数据量，最大化或优化带宽利用。

[0059] 图5示出了可以用于分布式训练系统102的示例性拓扑感知多阶段算法500。在实施中，拓扑感知多阶段算法500可以包括多个阶段，例如，节点内归约散布阶段502、节点间全局归约阶段504和节点内全局聚集阶段506。

[0060] 在实施中，分布式训练系统102可以首先将用于训练的待处理数据的各部分分配给多个计算节点204，使得多个计算节点204的每个计算节点104可以接收数据的相应部分。在实施中，每个计算节点104可以将数据的相应部分划分成多个数据片（例如，N个数据片，其中N是正整数），并将这多个数据片分配给包括在相应计算节点104中的多个本地处理单元或进程（例如，N个本地处理单元或进程）。

[0061] 在实施中，在节点内归约散布阶段502中，包括在每个计算节点104中的每个本地处理单元或进程可以将分配给它的数据片划分成多个数据块（例如，M个块）。然后包括在每个计算节点104中的本地处理单元或进程可以协作地执行节点内归约散布子运算，以根据特定的集群通信算法在多个步骤或迭代中获得相应计算节点104中的多个数据块的所有归约结果。在节点内归约散布阶段502结束时，计算节点104中包括的本地处理单元或进程可能具有不同数据块中的该计算节点104中包括的所有处理单元或进程的归约结果（或称为归约散布结果）。

[0062] 作为示例而非限制，以基于环的算法和减半加倍算法这两个示例性集群通信算法作为示例来描述特定的集群通信算法，以示出节点内归约散布阶段502中的特定机制或运算。然而，在该节点内归约散布阶段502中可以使用其他集群通信算法。例如，分布式训练系统102可以基于由集群通信库400收集的多个因素的信息，选择在节点内归约散布阶段502中使用的特定集群通信算法。在实施中，该多个因素可以包括但不限于计算节点中的处理单元（或进程）之间的互连类型、计算节点中的互连数量等。

[0063] 例如,在节点内归约散布阶段502中,分布式训练系统102可以为第一计算节点采用第一集群通信算法,为第二计算节点采用第二集群通信算法,其中第二计算节点具有与第一计算节点相同或不同的处理和连接能力,且第一集群通信算法可以与第二集群通信算法相同或不同。作为示例而非限制,对于使用NVSwitch或PCIe进行互连、并包括数量为2的幂的多个用于训练的处理单元或进程的计算节点,分布式训练系统102可采用减半加倍算法,而对使用NVLink或其他进行互连、并使用数量为非2的幂的多个用于训练的处理单元或进程的另一计算节点,分布式训练系统102采用基于环的算法,等等。

[0064] 图6示出了在节点内归约散布阶段502中用于计算节点的示例性基于环的算法600。出于简洁和描述的目的,示例性基于环的算法仅包括一个环的配置。然而,任何包括多于一个环的配置的基于环的算法都可被采用,例如,每个环处理数据块的一部分。

[0065] 在该示例中,描述的计算节点包括M个处理单元或进程(具有排序标识符或编号1、2、...、M),且分配给每个处理单元或进程的数据被分成M个数据块。在第一步,处理单元或进程(如P1)可以将其M个数据块中的一个发送给环中的下一个处理单元或进程(如P2),从环中的前一个处理单元或进程(如PM)接收另一个数据块,并用对应的本地数据块来归约所接收的数据块以获得部分归约的结果。在每个后续步骤中(如第k个步骤),处理单元或进程(如P1)可以向环中的下一个处理单元或进程(如P2)发送部分归约结果(在本示例中是由P1在第k-1个步骤获得的部分归约结果),从前一个处理单元或进程(如PM)接收部分归约结果(在本例中是由PM在第k-1个步骤获得的部分归约结果),并且用先前没有发送或用其他数据归约的另一个本地数据块来归约所接收的部分归约结果。

[0066] 如图6所示,每一步中,不同的数据块可以被计算节点中不同的处理单元或进程接收和归约或发送。此外,每个处理单元或进程可以在不同的步骤发送或接收并归约不同的数据块(或部分结果)。在节点内归约散布阶段502结束时(即,在M-1个步骤之后),每个处理单元或进程可以包括结果数据块,该结果数据块存储该计算节点中M个处理单元或进程的M个相应数据块的归约结果。例如,在M-1个步骤之后,“顶部位置”的P1的数据块存储对应于该“顶部位置”的M个处理单元或进程的所有数据块的归约结果,如图6所示。

[0067] 图7示出了在节点内归约散布阶段502中用于计算节点的示例性减半加倍算法700。在本示例中,描述的计算节点包括M个处理单元或进程(本示例中以M设定为8进行描述)。在第一步中,处理单元或进程(如P1)可以将分配给它的一半数据发送给附近的另一个处理单元或进程(如P2),并接收分配给该另一个处理单元或进程(如P2)的一半数据,并用分配给本处理单元或进程(如P1)的另一半数据来归约所接收的数据,以获得部分归约结果。在每个后续步骤中,处理单元或进程(如P1)可以将先前步骤中本地获得的部分归约结果的一半发送给位于离处理单元或进程(即P1)越来越远的位置的不同处理单元或进程,并用先前步骤中本地获得的部分归约结果的另一半来归约所接收的部分归约结果,以获得处理单元或进程(即P1)的新的部分归约结果。在节点内归约散布阶段502结束时(即,在 $\log_2 M$ 个步骤之后,即,在如图7所示的该示例中的3个步骤之后),每个处理单元或进程可以包括结果数据块,该结果数据块存储该计算节点中M个(在该示例中,如图7所示的8个)处理单元或进程的M个相应数据块的归约结果。例如,在 $\log_2 M$ 个步骤之后,“底部位置”的P1的数据块存储对应于如图7所示的该“底部位置”的M个(在该示例中,如图7所示为8个)处理单元或进程的所有数据块的归约结果。

[0068] 在实施中,在节点间全局归约阶段504中,节点间全局归约子运算是基于节点的(即,在不同计算节点之间),并且可以在不同计算节点中包括的处理单元(或进程)之间执行。在实施中,持有相同的归约结果(或归约散布结果)的数据块的不同计算节点的处理单元(或进程)形成同一组,并且在该同一组中相互传送各自的结果以执行节点间全局归约子运算。在节点间全局归约阶段504结束时,特定组中的每个计算节点的每个处理单元或进程可以拥有该同一组中所有处理单元或进程的归约结果的特定数据块,不同组的处理单元或进程拥有不同组中相应处理单元或进程的归约结果的不同数据块。

[0069] 在实施中,分布式训练系统102可以基于一个或多个选择标准来选择特定的集群通信算法,并且可以基于所选择的集群通信算法来实现节点间全局归约子运算。特定集群通信算法的示例可以包括但不限于基于环的算法(例如分级环算法、多环算法等)、减半加倍算法等。在实施中,一个或多个选择标准可以包括但不限于连接计算节点的通信网络(例如,通信网络206)的拓扑、通信网络中使用的交换机的数量、通信网络中使用的交换机的类型、通信网络的网络类型等。

[0070] 作为示例而非限制,以基于环的算法和减半加倍算法这两个示例性集群通信算法作为示例来描述特定的集群通信算法,以示出节点间全局归约阶段504中的特定机制或运算。然而,基于上述一个或多个选择标准,可以将其他集群通信算法用在该节点间全局归约阶段504中。

[0071] 图8和9出了节点间全局归约阶段504的示例性减半加倍算法。在该示例中,如图8所示,出于简洁和描述的目的,描述的分布式训练系统102包括多个计算节点(即,节点0、节点1、节点2、...节点N-1,为了进行说明,N在图8中被示为4),其中每个计算节点包括八个处理单元或进程,具有对应的排序编号(即,序位0、序位1、序位2、...序位M-1,为了进行说明,M在图8中被示为8)。如图8所示,在对应的计算节点中具有相同排序编号的处理单元或进程包括相同的归约结果(或归约散布结果)的数据块,并且其形成同一组。例如,在对应的计算节点中具有排序编号0的处理单元或进程包括各本地数据块中第一位置处的归约结果的数据块,并且其形成同一组(例如,组0)。在实施中,不同组中的处理单元或进程可以不相互通信。

[0072] 在实施中,节点间全局归约子运算可以在每个组中的处理单元(或进程)之间单独执行,使得组中的每个处理单元(或进程)可以获得同一组中所有处理单元(或进程)的同一数据块的所有归约结果。类似于上面针对节点内归约散布阶段描述的减半加倍算法的机制,每个组中的处理单元或进程可以与相应组中的其他处理单元或进程迭代地发送对应数据块的本地归约结果,从加倍或增加的距离中的其他处理单元或进程接收对应数据块各本地归约结果,并且用本地归约结果对所接收的归约结果执行归约运算。

[0073] 图9示出了对八个计算节点应用减半加倍算法的示例场景。在本示例中,如图8所示,利用减半加倍算法在节点间全局归约阶段504中执行步骤的数量为 $\log_2 N = \log_2 8 = 3$,其中N是计算节点的数量。在第一步中,第一计算节点(例如,节点0)中的某个组的第一处理单元或进程(例如,排序编号为0的处理单元或进程)可以向第二计算节点(例如,节点1)中的同一组的第二处理单元或进程发送其本地归约结果,从第二计算节点中的同一组的第二处理单元或进程接收本地归约结果,并对其本地归约结果和接收到的本地归约结果执行归约运算,以获得新的本地归约结果。

[0074] 在第二步中,第一计算节点(例如,节点0)中的第一处理单元或进程(例如,排序编号为0的处理单元或进程)可以将其新的本地归约结果发送给第三计算节点(即,本例中的节点2)中的同一组的第三处理单元或进程(例如,排序编号为0),从第一计算节点中的同一组的第三处理单元或进程接收本地归约结果,并对其新的本地归约结果和所接收的本地归约结果执行归约运算,以获得另一个新的本地归约结果。

[0075] 在第三步(或最后一步),对第一处理单元或进程执行相同的运算,但是此时是对第四计算节点(即,本例中的节点4)中的同一组的第四处理单元或进程执行相同的运算。

[0076] 在节点间全局归约阶段504结束时,特定组中的每个计算节点的每个处理单元或进程可以拥有该同一组中所有处理单元或进程的归约结果的特定数据块,不同组的处理单元或进程拥有不同组中相应处理单元或进程的归约结果的不同数据块。

[0077] 与减半加倍算法类似,节点间全局归约子运算可以利用基于环的算法在多个计算节点(如N个计算节点)的每个组中的处理单元(或进程)之间单独执行,使得组中的每个处理单元(或进程)可以获得同一组中所有处理单元(或进程)的同一数据块的所有归约结果。类似于以上针对节点内归约散布阶段描述的基于环的算法的机制,计算节点中的每个组中的处理单元或进程可以迭代地将对应数据块的本地归约结果发送给下一个计算节点中的相应组中的处理单元或进程,从前一个计算节点中的相应组中的处理单元或进程接收对应数据块的本地归约结果,并利用其本地归约结果对所接收的归约结果执行归约运算。在节点间全局归约阶段504结束时(即在N-1步之后),特定组中的每个计算节点的每个处理单元或进程可以拥有该同一组中所有处理单元或进程的归约结果的特定数据块,不同组的处理单元或进程拥有不同组中相应处理单元或进程的归约结果的不同数据块。

[0078] 在实施中,类似于节点内归约散布阶段502,在节点内全局聚集阶段506中,全局聚集子运算可以跨分布式训练系统102的多个计算节点的每个计算节点中的本地处理单元或进程来执行,以在同一计算节点中向彼此本地广播在节点间全局归约阶段504中获得的相应归约结果。在节点内全局聚集阶段506结束时,分布式训练系统102的每个计算节点中的每个处理单元或过程可以具有分布在多个计算节点中的整个数据的归约结果。

[0079] 作为示例而非限制,本文使用基于环的算法来说明如何广播由分布式训练系统102的计算节点中的本地处理单元或过程(在节点间全局归约阶段504中)获得的归约结果。然而,分布式训练系统102可以对不同的计算节点采用不同的或相同的集群通信算法(例如减半加倍算法等)。例如,分布式训练系统102可以基于与每个单独的计算节点相关联的多个因素,为不同的计算节点采用不同或相同的集群通信算法。在实施中,该多个因素可以包括但不限于计算节点中的处理单元(或进程)之间的互连类型、计算节点中的互连数量等。

[0080] 图10示出了示例性的基于环的集群通信算法1000,用于在分布式训练系统102的计算节点内相互广播处理单元或过程的个体归约结果。如图10所示,在第一步中,计算节点中的M个处理单元或进程中的每个处理单元或进程(例如,P1)可以根据环形配置将其在节点间全局归约阶段504中获得的归约结果发送给两个相邻处理单元或进程中的一个(例如,本例中的P2),并从两个相邻处理单元或进程中的另一个(例如,本例中的PM)接收归约结果。在每个后续步骤中,每个处理单元或进程(例如,P1)可以根据环形配置向两个相邻处理单元或进程中的一个(例如,本例中的P2)发送新接收的归约结果,并从两个相邻处理单元或进程中的另一个(例如,本例中的PM)接收另一个归约结果。在节点内全局聚集阶段506结

束时(即,在M-1个步骤之后),计算节点中的每个处理单元或进程可具有计算节点中所有处理单元或进程的归约结果的归约结果。

[0081] 示例性并行算法

[0082] 在实施中,分布式训练系统102可以执行包括在拓扑感知多阶段算法中的多个阶段,即,依次为节点内归约散布阶段502、节点间全局归约阶段504和节点内全局聚集阶段506等。在实施中,可替代地,分布式训练系统102可以对节点内归约散布阶段502、节点间全局归约阶段504和节点内全局聚集阶段506中的一些进行部分或实质重叠,并且并行地执行这些阶段的一些部分。

[0083] 例如,由于节点内归约散布阶段502和节点内全局聚集阶段506涉及节点内数据通信或传输(即,计算节点内的数据通信或传输),并且节点间全局归约阶段504涉及节点间数据通信或传输(即,计算节点之间的数据通信或传输),在实施中,分布式训练系统102可以允许节点内归约散布阶段502和节点间全局归约阶段504的至少一部分并行执行,以及节点间全局归约阶段504和节点内全局聚集阶段506的一部分,从而提高节点内和节点间链路(或连接)的利用率,并避免节点内链路在节点间链路被使用时空闲,反之亦然。

[0084] 图11示出了用并行或重叠的方式执行节点内归约散布阶段、节点间全局归约阶段、节点内全局聚集阶段的示例场景。如图11所示,计算节点的处理单元或进程可以将数据块划分为多个块(在本示例中是如图11所示的4个块),并将这些块分发至至少两个并发线程(如第一线程1102和第二线程1104)。以这种方式,处理单元或进程可以流水线化节点内和节点间子运算,以便由至少两个并发线程(在该示例中,第一线程1102和第二线程1104)执行。

[0085] 作为示例而非限制,第一线程1102可以对第一数据块(例如,数据块1106)执行节点间全局归约子运算(即,节点间全局归约阶段504中的运算),而第二线程1104对第二数据块(例如,数据块1108)执行节点内归约散布子运算(即,节点内归约散布阶段502中的运算)。例如,第一线程1102可以对第三数据块(例如,数据块1110)执行节点内全局聚集子运算(即,节点内全局聚集阶段506中的运算),而第二线程1104对第四数据块(例如,数据块1112)执行节点间全局归约子运算。

[0086] 作为示例而非限制,分布式神经网络训练中涉及的另一个运算可以进一步用作示例。在实施中,分布式训练系统102可以将分布式神经网络训练中涉及的全局聚集运算分成多个子运算,即节点间全局聚集子运算、节点内全局聚集子运算和数据复制子运算。在实施中,节点间全局聚集子运算可以类似于如上所述的节点间全局归约子运算,但进行的是广播数据(例如,归约的结果)而不是归约运算(例如,用本地结果归约接收的结果),而节点间全局聚集子运算可以类似于或等同于如上所述的节点间全局聚集子运算。在实施中,数据复制子运算可以包括复制结果数据(例如,最终的归约结果)作为输出参数的运算。

[0087] 在实施中,计算节点的处理单元或进程可以将数据块划分成多个块(例如,四个块),并将这些块分发到至少两个并发线程(例如,第一线程和第二线程),并且流水线化节点内和节点间子运算以供至少两个并发线程执行。

[0088] 例如,第一线程可以对第一数据块执行节点间全局聚集子运算,而第二线程对第二数据块执行节点内全局聚集子运算。此外,第一线程可以对第三数据块执行数据复制子运算,而第二线程对第四数据块执行节点间全局聚集子运算。

[0089] 示例性拥塞避免方法

[0090] 在实施中,由于分布式训练系统102中的多个计算节点之间的数据传输,在通信网络206中的一些交换机或链路处可能发生数据或流量拥塞。为了避免拥塞,分布式训练系统102可以采用预定的拥塞避免策略来在通信网络206中的各种交换机或链路之间分发或转移数据流量,从而避免过量的数据在训练期间(例如,节点间全局归约子运算或阶段,或者节点间全局聚集子运算或阶段)通过通信网络206中的某个交换机或链路。

[0091] 在实施中,分布式训练系统102可以采用第一拥塞避免方法,该方法包括环生成策略,以及随后的网络流的路由管理。附加地或替代地,分布式训练系统102可以采用第二拥塞避免方法,该方法包括对节点标识重新排序的策略,以及随后的网络流的路由管理。取决于通信网络206的网络拓扑的类型以及多个计算节点204的处理和通信能力等,分布式训练系统102可以选择第一拥塞避免方法或第二拥塞避免方法中的一个或多个,用于在分布式训练系统102中的多个计算节点的全部或部分之间路由数据流。此外,分布式训练系统102可以选择性地组合第一拥塞避免方法和第二拥塞避免方法的部分,以实现新的拥塞避免方法。在实施中,第一拥塞避免方法和第二拥塞避免方法都可以旨在以节点间数据流彼此没有或有很少冲突的方式为节点间数据流的每个方向指定专用网络路径。

[0092] 在实施中,分布式训练系统102可以预先获得或建立通信连接和路由路径(例如,物理链路)之间的映射关系。在实施中,可以创建表格、链表等形式的连接路径数据结构并用于存储映射关系的信息。在实施中,分布式训练系统102可以基于连接路径数据结构选择性地或策略性地使用特定路径来建立任何两个计算节点之间的连接。

[0093] 在实施中,分布式训练系统102可以通过使分布式训练系统102的每个计算节点能够通过变化探测数据包的源/目的地端口向其他计算节点发送探测数据包,来确定通信连接和路由路径之间的映射关系,以穷尽分布式训练系统102的计算节点之间可能的通信连接。显然,分布式训练系统102可以采用其他方法来探索通信连接和路由路径之间的映射关系,在此不做限定。

[0094] 作为示例而非限制,第一计算节点可以向第二计算节点发送多个探测数据包,每个探测数据包具有源端口和目的端口的不同组合,而源地址和目的地址分别是第一计算节点的地址和第二计算节点的地址。每个探测数据包可以记录相应探测数据包所经过的交换机,因此当相应探测数据包返回到第一计算节点时,第一计算节点可以知道用于映射的相应探测数据包的整个路由路径。据此,可以在第一计算节点和第二计算节点之间建立连接路径数据结构(例如,连接路径数据结构)。类似地,可以相应地建立分布式训练系统102中的其他计算节点对的通信连接和路由路径之间的映射关系(以及连接路径数据结构)。

[0095] 出于简洁和说明的目的,示例性网络拓扑,即胖树网络(或者特别是全网状拓扑中的双层Clos网络架构)在此被用作与分布式训练系统102相关联的通信网络206的示例网络拓扑。然而,这里描述的示例性拥塞避免策略也可以适用于其他网络拓扑。

[0096] 图12示出了示例性胖树网络拓扑1200。在该示例中,示例性胖树网络拓扑是全网状拓扑中的双层Clos网络架构。一层对应于直接连接到计算节点1204的一层片交换机1202,每个片交换机1202连接到一个或多个计算节点1204。在实施中,计算节点1204可以包括连接到片交换机1202的一个或多个端口(例如,四个端口)的一个或多个网络接口控制器(例如,四个网络接口控制器)。在实施中,每个计算节点1204的网络接口控制器的数量可以

相同也可以不同。另一层对应于连接到一个或多个片交换机1202的一层汇聚交换机1206(或称为主干交换机1206)。

[0097] 在实施中,如果包括在不同计算节点中的两个处理单元或进程连接在同一片交换机下,则在这两个处理单元或进程之间传输的数据包将通过该同一片交换机,而不通过任何汇聚交换机。可替换地,如果包括在不同计算节点中的两个处理单元或进程在不同的片交换机下连接,则在这两个处理单元或进程之间传输的数据包将通过汇聚交换机之一。使用如上所述的连接路径数据结构,通过在数据包中设置源端口和目的端口的适当组合,可以使在两个处理单元或进程之间传输的数据包流经指定的汇聚交换机。在实施中,第一拥塞避免方法和/或第二拥塞避免方法的路由管理可以旨在使得从同一片交换机到不同目的地片交换机的数据流能够通过不同的汇聚交换机,和/或使得从不同源片交换机到同一目的地片交换机的数据流能够通过不同的汇聚交换机,从而避免数据流之间的冲突,并且使得汇聚交换机处没有网络拥塞。

[0098] 在实施中,如前面的描述中所述,第一拥塞避免方法可以包括环生成策略,以及随后的网络流的路由管理。第一拥塞避免方法可以支持各种基于环的算法,包括但不限于环算法、环分块算法、多环算法、分层环算法、涉及多个分层环的算法和节点感知环算法等。

[0099] 在实施中,环生成的策略可以包括环生成的拓扑感知策略。作为示例和限制,环生成的拓扑感知策略可以包括多个规则来建立处理单元或进程的环或环配置。在实施中,计算节点中的处理单元或进程可以通过网络接口控制器向/从另一计算节点中的处理单元或进程发送/接收数据。在实施中,计算节点中的处理单元或进程可以与单个网络接口控制器或多个网络接口控制器相关联,以向其他计算节点中的处理单元或进程传输数据。附加地或可替换地,多个处理单元或进程可以与单个网络接口控制器相关联,并且使用该网络接口控制器来向其他计算节点中的处理单元或进程传输数据。

[0100] 在实施中,多个规则可以包括但不限于第一计算节点中的处理单元或进程选择相邻处理单元或进程的优先级、第一计算节点中的网络接口控制器发送或接收数据的条件、第一计算节点中的网络接口控制器向/从第二计算节点中的网络接口控制器路由数据的条件等。

[0101] 在实施中,第一计算节点中的处理单元或进程选择相邻处理单元或进程的优先级可以包括,按照优先级的降序,选择第一计算节点中的处理单元或进程并使用进程间通信(如果适用的话),选择连接到与第一计算节点所连的片交换机相同的片交换机的第二计算节点中的处理单元或进程,在连接到与第一计算节点所连的片交换机不同的片交换机的第三计算节点中选择处理单元或进程,其中第一计算节点不同于第二计算节点和第三计算节点。

[0102] 在实施中,第一计算节点中的网络接口控制器发送或接收数据的条件可以包括,例如,网络接口控制器能够仅向第二计算节点中的网络接口控制器发送数据,和/或网络接口控制器能够仅从第三计算节点中的网络接口控制器接收数据,其中第一计算节点不同于第二计算节点和第三计算节点,并且第二计算节点可以与第三计算节点相同或不同。

[0103] 在实施中,第一计算节点中的网络接口控制器向/从第二计算节点中的网络接口控制器路由数据的条件可以包括,例如,如果数据是通过第一计算节点中的网络接口控制器发送的,则将由属于多个环的处理单元或进程发送的数据路由到第二计算节点中的网络

接口控制器。在实施中,第一计算节点中的网络接口控制器向/从第二计算节点中的网络接口控制器路由数据的条件可以进一步包括,如果数据是由属于多个环的处理单元或进程通过第二计算节点中的网络接口控制器发送的,则通过第一计算节点中的网络接口控制器接收数据。

[0104] 在实施中,第一拥塞避免方法的路由管理可以将网络接口控制器(Network Interface Controller, NIC)标识符分配给连接或链接到同一片交换机的每个网络接口控制器。第一拥塞避免方法的路由管理还可以向通信网络206中的每个汇聚交换机分配汇聚标识符。对于某个环中的处理单元或进程,路由管理可以确定用于路由来自该处理单元或进程的数据包的路由标识符。

[0105] 例如,如果处理单元或进程的网络接口控制器和环中的下一个处理单元或进程的网络接口控制器位于相同的计算节点中,或者直接连接或链接到相同的片交换机,则路由标识符可以被确定为默认值或默认标识符。该默认路由标识符指示数据或在计算节点内路由,或通过片交换机路由,而不通过通信网络中的任何汇聚交换机。否则,路由标识符可被确定为等于该处理单元或进程的NIC标识符或其他预定义值。基于路由标识符和汇聚标识符之间的映射关系,可以基于所确定的路由标识符来确定汇聚标识符。在实施中,例如,路由标识符和汇聚标识符之间的映射关系可以使用基于探测的路由机制(例如,如前面的描述中所描述的在计算节点之间发送探测数据包)来预先确定。

[0106] 换句话说,包括在同一计算节点中的处理单元(或进程)或同一片交换机的网络接口控制器之间的数据流将不经过通信网络中的任何汇聚交换机。另一方面,不同计算节点中包括的处理单元(或进程)和不同片交换机的网络接口控制器之间的数据流将基于预定的映射关系通过指定的汇聚交换机,从而实现数据流的路由控制和管理,并将数据流分发到不同的汇聚交换机,以避免网络拥塞。

[0107] 图13示出了使用第一拥塞避免方法的示例场景。在该示例中,生成了包含八个计算节点(节点0、节点1、…、节点7)的四个节点间环(或环配置,R0、R1、R2和R3),并且每个环使用不同的汇聚交换机来发送和接收数据(例如,在节点间全局归约阶段504期间)。因此,这四个环不存在冲突。此外,任何环的每个片交换机只有一个数据流进入,一个数据流离开,从而避免了网络拥塞的发生。

[0108] 在实施中,如上所述,第二拥塞避免方法可以包括关于节点标识的重新排序的策略,以及随后的网络流的路由管理。在实施中,为了最小化通信成本,第二拥塞避免方法可以根据基于多个规则连接计算节点和处理单元(或进程)的网络拓扑来重新排序计算节点和处理单元(或进程)的标识符。

[0109] 在实施中,多个规则可以包括例如通过各片交换机对计算节点进行分组。例如,连接到同一片交换机的计算节点(例如,具有链接到同一片交换机的网络接口控制器的计算节点)被组成一个组,并且每个计算节点被分配有节点标识符。由于计算节点连接到同一片交换机,这些计算节点彼此(物理上)相邻。

[0110] 在实施中,多个规则还可以包括使用相同的顺序序列将排序标识符(或排序编号)分配给计算节点中的每个处理单元或进程。例如,第一计算节点中的k个处理单元(或进程)可以被分配有排序标识符0、1、…、k-1,而第二计算节点中的k个处理单元(或进程)可以被分配有排序标识符k、k+1、…、2k-1等,对于其他计算节点也是如此。可以根据处理单元(或

进程)使用的相应网络接口控制器来对计算节点中的处理单元(或进程)进行排序,并且使用相同网络接口控制器的处理单元(或进程)彼此(物理上)相邻。

[0111] 在这种情况下,在前 $\log_2 L$ 个步骤中,计算节点中的处理单元(或进程)之间的数据流可以被限制为通过相应的片交换机,其具有比汇聚交换机更优的延迟,因此不会产生网络拥塞。在实施中, L 是如在上文描述中的节点感知的减半加倍算法的每片交换机的计算节点数量。在实施中,对于传统的减半加倍算法, L 是每片交换机的计算节点数和每计算节点的处理单元(或进程)数的乘积。

[0112] 在实施中,第二拥塞避免方法的路由管理可以包括确定从具有第一节点标识符的第一计算节点中具有第一排序标识符的第一处理单元(或进程)发送给具有第二节点标识符的第二计算节点中具有第二排序标识符的第二处理单元(或进程)的数据流或数据包的数据流或数据包的汇聚标识符,其中第一计算节点可以与第二计算节点相同或不同。

[0113] 在实施中,汇聚标识符可以至少部分地基于排序标识符、节点标识符、每计算节点的网络接口控制器的数量以及每片交换机处的计算节点的最大数量中的至少一些来确定。作为示例而非限制,汇聚标识符可被确定为发送数据流或数据包的第一处理单元(或进程)的第一排序标识符+(具有第一处理单元(或进程)的第一计算节点的第一节点标识符)%每个片交换机处的计算节点的最大数量 \times 每计算节点的网络接口控制器的数量,其中%表示模数运算符。显然,只要能获得一致的结果,汇聚标识符的其他计算方法也是适用的。例如,汇聚标识符可以基于汇聚标识符与排序标识符和节点标识符的组合之间的预设映射关系来确定等。

[0114] 在实施中,第二拥塞避免方法的路由管理可以包括预先将汇聚标识符分配给与分布式训练系统102相关联的通信网络206中的每个汇聚交换机。如果第一处理单元(或进程)和第二处理单元(或进程)链接到同一片交换机或同一片交换机下(例如,通过各自的网络控制器),则数据流或数据包将通过该片交换机,而不需要通过通信网络206中的任何汇聚交换机。如果第一处理单元(或进程)和第二处理单元(或进程)没有链接到同一片交换机或不在同一片交换机下,则由第一处理单元(或进程)发送给第二处理单元(或进程)的数据流或数据包将通过具有所确定的汇聚标识符的汇聚交换机。在本示例中,所描述的包括在每个计算节点中的网络接口控制器的数量是4个。

[0115] 图14示出了使用第二拥塞避免方法的示例场景。在本示例中,所有计算节点包括相同数量的处理单元(或进程)和相同数量的网络接口控制器,每个网络接口控制器具有相同数量的待关联的处理单元(或进程)。此外,链接到片交换机的网络接口控制器的数量少于网络中汇聚交换机的数量。在本示例中,每计算节点的网络接口控制器的数量为4个,每个片交换机处的计算节点的数量最多为2个。在实施中,对于节点感知的减半加倍算法,同一片交换机下的计算节点的数量可以是2的幂,并且同一片交换机下的计算节点中包括的网络接口控制器的数量可以是2的幂,对于传统的减半加倍算法,使用同一网络接口控制器的处理单元(或进程)的数量可以是2的幂。在本示例中,所描述的包括在每个计算节点中的网络接口控制器的数量为4个。

[0116] 在本示例中,在节点感知的减半加倍算法中的节点间全局归约阶段,计算节点(节点0、节点2、节点4和节点6)的处理单元(或进程)将使用具有汇聚标识符的汇聚交换机(例如,A1、A2、A3和A4),并且计算节点(节点1、节点3、节点5和节点7)的处理单元(或进程)将使

用具有汇聚标识符的汇聚交换机(例如,A5、A6、A7和A8)。据此,计算节点之间的数据流中不存在冲突,从而避免了网络中任何汇聚交换机处的网络拥塞。

[0117] 在实施中,在节点间全局归约阶段的每个步骤,处理单元(或进程)可以与新的处理单元(或进程)进行数据通信。在实施中,可以进行同步以确保由使用网络接口控制器的处理单元(或进程)在当前步骤进行的数据流不与由使用相同网络接口控制器的相邻处理单元(或进程)在先前步骤进行的数据流重叠,以避免发生微突发流(incast),从而避免发生微突发流(incast)拥塞。

[0118] 示例性方法

[0119] 图15示出了示例性拓扑感知多阶段方法的示意图。图16示出了第一示例性网络拥塞避免方法的示意图。图17示出了第二示例性网络拥塞避免方法的示意图。图18示出了分布式训练中基于混合架构的示例性并行方法示意图。图15-18的方法可以(但不一定)在图1所示的环境中利用图2所示的计算节点在图3-14所示的方法和场景的辅助下实施。为了便于描述,参考图1-14来描述方法1500-1800。然而,方法1500-1800可以替代性地在其他环境和/或利用其他系统来实施。

[0120] 方法1500-1800在机器可执行指令的一般背景下描述。通常,机器可执行指令可以包括执行特定功能或实现特定抽象数据类型的例程、程序、对象、组件、数据结构、过程、模块、函数等。此外,将每个示例方法表示为逻辑流程图中的框的集合,该逻辑流程图表示可以在硬件、软件、固件或其组合中实现的一系列操作。描述该方法的顺序不意图被理解为进行限制,并且任何数量的所描述的方法框可以以任何顺序组合以实现该方法或替代方法。此外,在不脱离本文所述主题的精神和范围的情况下,可以从该方法中省略单独的框。在软件的背景下,框代表当由一个或多个处理器执行时执行所述操作的计算机指令。在硬件的背景下,一些或所有的框可以代表执行所述操作的专用集成电路(ASIC)或其他物理组件。

[0121] 参考图15,在框1502中,第一计算节点(例如,计算节点104)可以根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行归约散布子运算。

[0122] 在实施中,在执行归约散布子运算之前,第一计算节点可以至少部分地基于第一计算节点中的第一处理单元集合之间的节点内连接的类型或带宽来选择第一集群通信算法。在实施中,第一集群通信算法可以包括但不限于基于环的算法或减半加倍算法。

[0123] 在实施中,根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行归约散布子运算可以包括将数据划分成多个数据块;将多个数据块分配给第一处理单元集合;根据第一集群通信算法,在第一处理单元集合的第一处理单元处从第一处理单元集合的第二处理单元接收数据块;以及在第一处理单元处用本地数据块归约所接收的数据块。

[0124] 在框1504中,第一计算节点可以根据第二集群通信算法在第一计算节点中的第一处理单元集合和第二计算节点中的第二处理单元集合之间执行全局归约子运算。

[0125] 在实施中,在执行全局归约子运算之前,第一计算节点可以至少部分地基于第一计算节点和其他计算节点之间的节点间连接的类型或带宽,和/或第一计算节点和其他计算节点的连接拓扑来选择第二集群通信算法。在实施中,第一集群通信算法可以包括但不限于基于环的算法,或者减半加倍算法(比如节点感知减半加倍算法)等。

[0126] 在实施中,根据第二集群算法在第一计算节点中的第一处理单元集合和第二计算

节点中的第二处理单元集合之间执行全局归约子运算可以包括：第一处理单元集合接收第二计算节点中的第二处理单元集合根据第二集群算法所获得的归约散布结果的各部分，第一处理单元集合的每个处理单元与第二处理单元集合的相应处理单元形成组，并从相应处理单元接收归约散布结果的相应部分；第一处理单元集合通过在第一处理单元集合之间执行归约散布子运算后获得的归约散布结果的对应本地部分对归约散布结果的各部分执行归约。

[0127] 在框1506中，第一计算节点可以根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行全局聚集子运算。

[0128] 在实施中，根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行全局聚集子运算可以包括：根据第一集群通信算法，在第一处理单元集合的第一处理单元处从第一处理单元集合的第二处理单元接收数据块；以及在第一处理单元处用本地数据块归约所接收的数据块。

[0129] 参考图16，在框1602中，第一计算节点（例如，计算节点104）或第一进程可以至少部分地基于与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器是否位于同一计算节点中或链接到同一片交换机来确定将数据从第一进程路由到第二进程的路由标识符。

[0130] 在实施中，第一进程和第二进程可以属于特定网络拓扑下连接多个不同节点的特定节点间环。作为示例而非限制，特定网络拓扑可以包括胖树拓扑。

[0131] 在实施中，与第一进程相关联的网络接口控制器被配置为仅向环形拓扑中的第二计算节点发送数据或从第二计算节点接收数据，第二计算节点不同于第一计算节点。

[0132] 在实施中，与第一进程相关联的网络接口控制器还与一个或多个进程相关联，其中从第一进程和一个或多个进程发送的所有数据都通过网络接口控制器发送。

[0133] 在实施中，响应于确定与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器位于同一计算节点中或链接到同一片交换机，可以将路由标识符设置或确定为默认标识符。

[0134] 在实施中，响应于确定与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器位于不同的计算节点中或链接到不同的片交换机，可以将路由标识符设置或确定为等于与第一进程相关联的网络接口控制器的标识符。

[0135] 在框1604中，第一计算节点或第一进程可以根据路由标识符将数据从第一进程路由到第二进程。

[0136] 在实施中，根据路由标识符将数据从第一进程路由到第二进程可以包括至少通过与网络接口控制器连接的片交换机和具有与网络接口控制器的标识符具有对应关系的标识符的汇聚交换机将数据从第一进程路由到第二进程，该网络接口控制器与第一进程相关联。

[0137] 参考图17，在框1702中，第一计算节点（例如，计算节点104）或第一进程可以根据节点感知减半加倍算法确定用于从第一进程向第二进程发送数据包的汇聚标识符，第一进程和第二进程属于在特定网络拓扑下连接至不同片交换机的不同节点。

[0138] 在实施中，第一计算节点可以为定向到连接至不同片交换机的计算节点的数据包分配不同的汇聚标识符，以使得能够通过不同的汇聚交换机将数据包路由到连接至不同片

交换机的节点。

[0139] 在实施中,第一计算节点可以至少部分地基于预定对应关系,分配与汇聚标识符相关联的汇聚交换机对应的源端口和目的端口。在实施中,对应关系可以记录多个汇聚交换机的汇聚标识符与对应的源端口和目的端口对之间的关系。在实施中,特定的网络拓扑可以包括胖树拓扑。

[0140] 在框1704中,第一计算节点可以通过与汇聚标识符对应的汇聚交换机从第一进程向第二进程发送数据包。

[0141] 在实施中,第一计算节点还可以通过分配给各数据包的多个不同汇聚标识符对应的多个不同的汇聚交换机,从第一计算节点包括的第一进程集合向第二计算节点包括的第二进程集合发送各数据包。

[0142] 在实施中,第一计算节点还可以通过分配给各数据包的多个不同汇聚标识符对应的多个不同的汇聚交换机,由第一计算节点包括的第一进程集合从第二计算节点包括的第二进程集合接收数据包。

[0143] 参考图18,在框1802中,第一计算节点(例如,计算节点104)或处理单元可以将分配给该处理单元的数据块划分成多个数据段,该多个数据段至少包括第一数据段和第二数据段。

[0144] 在框1804中,第一计算节点或者处理单元可以将多个数据段分配给多个线程,该多个线程至少包括第一线程和第二线程。

[0145] 在框1806中,第一计算节点或处理单元可以使用第一线程对第一数据段的一部分执行节点内子运算,且并行地使用第二线程对第二数据段的一部分执行节点间子运算。

[0146] 在实施中,使用第一线程对第一数据段的一部分执行节点内子运算可以包括通过节点内连接在第一计算节点中包括的处理单元和另一处理单元之间传输上述第一数据段的一部分。

[0147] 在实施中,使用第二线程对第二数据段的一部分执行节点间子运算可以包括通过节点间连接在处理单元和不同于第一计算节点的第二计算节点中包括的另一处理单元之间传输上述第二数据段的一部分。

[0148] 在实施中,节点内子运算可以包括在第一计算节点内执行的归约散布子运算或全局聚集子运算,并且节点间子运算可以包括在第一计算节点和不同于第一计算节点的第二计算节点间执行的全局归约子运算。

[0149] 在实施中,节点内子运算可以包括在第一计算节点内执行的全局聚集子运算或复制子运算,节点间子运算可以包括在第一计算节点和不同于第一计算节点的第二计算节点间执行的全局聚集子运算。

[0150] 在实施中,第一计算节点或处理单元可以使用第一线程对上述第一数据段的一部分执行另一个节点间子运算,并且并行地使用第二线程对上述第二数据段的一部分执行另一个节点内子运算。

[0151] 在实施中,使用第一线程对第一数据段的该部分执行节点内子运算,以及并行地使用第二线程对第二数据段的该部分执行节点间子运算,使得使用节点内连接将第一数据段的一部分传输至第一计算节点包括的另一处理单元,以及并发地使用节点间连接将第二数据段的一部分传输至与第一计算节点不同的第二计算节点包括的另一处理单元。

[0152] 虽然上述描述的方法框是按一定顺序执行的,但在一些实施中,一些或全部方法框可以以其他顺序或并行地执行。

[0153] 总结

[0154] 虽然是用专用于结构特征和/或方法动作的语言描述了实施方式,但是应该理解,权利要求不一定限于所描述的特定特征或动作。相反地,这些具体特征和动作是作为实施所要求保护的主题的示例形式公开的。附加地或替代地,一些或所有操作可以由一个或多个ASICs、FPGA或其他硬件来实现。

[0155] 本公开可以进一步利用下述条款来理解:

[0156] 条款1:一种由第一计算节点实施的方法,所述方法包括:根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行归约散布子运算;根据第二集群通信算法在第一计算节点中的第一处理单元集合和第二计算节点中的第二处理单元集合之间执行全局归约子运算;并根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行全局聚集子运算。

[0157] 条款2:根据条款1所述的方法,还包括:至少部分基于第一计算节点中的第一处理单元集合间的节点内连接的类型或带宽选择第一集群通信算法。

[0158] 条款3:根据条款1所述的方法,还包括:至少部分基于第一计算节点和其他计算节点之间的节点间连接的类型或带宽,和/或,第一计算节点和其他计算节点的连接拓扑来选择第二集群通信算法。

[0159] 条款4:根据条款1所述的方法,其中第一集群通信算法包括基于环的算法,或减半加倍算法。

[0160] 条款5:根据条款1所述的方法,其中根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行归约散布子运算包括:将数据划分成多个数据块;将多个数据块分配给第一处理单元集合;根据第一集群通信算法,在第一处理单元集合的第一处理单元处从第一处理单元集合的第二处理单元接收数据块;以及在第一处理单元处用本地数据块归约所接收的数据块。

[0161] 条款6:根据条款1所述的方法,其中根据第二集群算法在第一计算节点中的第一处理单元集合和第二计算节点中的第二处理单元集合之间执行全局归约子运算包括:第一处理单元集合接收第二计算节点中的第二处理单元集合根据第二集群算法所获得的归约散布结果的各部分,其中第一处理单元集合的每个处理单元与第二处理单元集合的相应处理单元形成组,并从相应处理单元接收归约散布结果的相应部分;第一处理单元集合通过在第一处理单元集合之间执行归约散布子运算后获得的归约散布结果的对应本地部分对归约散布结果的各部分执行归约。

[0162] 条款7:根据条款1所述的方法,其中根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行全局聚集子运算包括:根据第一集群通信算法,在第一处理单元集合的第一处理单元处从第一处理单元集合的第二处理单元接收数据块;以及在第一处理单元处用本地数据块归约所接收的数据块。

[0163] 条款8:一个或多个机器可读介质,其存储有机器可读指令,当机器可读指令由第一计算节点执行时,使第一计算节点执行动作,所述动作包括:根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行归约散布子运算;根据第二集群通信算法在第

一计算节点中的第一处理单元集合和第二计算节点中的第二处理单元集合之间执行全局归约子运算;并根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行全局聚集子运算。

[0164] 条款9:根据条款8所述的一个或多个机器可读介质,所述动作还包括:至少部分基于第一计算节点中的第一处理单元集合间的节点内连接的类型或带宽选择第一集群通信算法。

[0165] 条款10:根据条款8所述的一个或多个机器可读介质,所述动作还包括:至少部分基于第一计算节点和其他计算节点之间的节点间连接的类型或带宽,和/或,第一计算节点和其他计算节点的连接拓扑来选择第二集群通信算法。

[0166] 条款11:根据条款8所述的一个或多个机器可读介质,其中第一集群通信算法包括基于环的算法,或减半加倍算法。

[0167] 条款12:根据条款8所述的一个或多个机器可读介质,其中根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行归约散布子运算包括:将数据划分成多个数据块;将多个数据块分配给第一处理单元集合;根据第一集群通信算法,在第一处理单元集合的第一处理单元处从第一处理单元集合的第二处理单元接收数据块;以及在第一处理单元处用本地数据块归约所接收的数据块。

[0168] 条款13:根据条款8所述的一个或多个机器可读介质,其中根据第二集群算法在第一计算节点中的第一处理单元集合和第二计算节点中的第二处理单元集合之间执行全局归约子运算包括:第一处理单元集合接收第二计算节点中的第二处理单元集合根据第二集群算法所获得的归约散布结果的各部分,第一处理单元集合的每个处理单元与第二处理单元集合的相应处理单元形成组,并从相应处理单元接收归约散布结果的相应部分;第一处理单元集合通过在第一处理单元集合之间执行归约散布子运算后获得的归约散布结果的对应本地部分对归约散布结果的各部分执行归约。

[0169] 条款14:根据条款8所述的一个或多个及其可读介质,其中根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行全局聚集子运算包括:根据第一集群通信算法,在第一处理单元集合的第一处理单元处从第一处理单元集合的第二处理单元接收数据块;以及在第一处理单元处用本地数据块归约所接收的数据块。

[0170] 条款15:一种第一计算节点,包括:第一处理单元集合;以及存储器,其存储有机器可执行指令,当机器可执行指令由第一处理单元集合执行时,使第一处理单元集合执行动作,所述动作包括:根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行归约散布子运算;根据第二集群通信算法在第一计算节点中的第一处理单元集合和第二计算节点中的第二处理单元集合之间执行全局归约子运算;并根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行全局聚集子运算。

[0171] 条款16:根据条款15所述的第一计算节点,所述动作还包括:至少部分基于第一计算节点的第一处理单元集合间的节点内连接的类型或带宽选择第一集群通信;至少部分基于第一计算节点和其他计算节点之间的节点间连接的类型或带宽,和/或,第一计算节点和其他计算节点的连接拓扑来选择第二集群通信算法。

[0172] 条款17:根据条款15所述的第一计算节点,其中第一集群通信算法包括基于环的算法,或减半加倍算法。

[0173] 条款18:根据条款15所述的第一计算节点,其中根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行归约散布子运算包括:将数据划分成多个数据块;将多个数据块分配给第一处理单元集合;根据第一集群通信算法,在第一处理单元集合的第一处理单元处从第一处理单元集合的第二处理单元接收数据块;以及在第一处理单元处用本地数据块归约所接收的数据块。

[0174] 条款19:根据条款15所述的第一计算节点,其中根据第二集群算法在第一计算节点中的第一处理单元集合和第二计算节点中的第二处理单元集合之间执行全局归约子运算包括:第一处理单元集合接收第二计算节点中的第二处理单元集合根据第二集群算法所获得的归约散布结果的各部分,第一处理单元集合的每个处理单元与第二处理单元集合的相应处理单元形成组,并从相应处理单元接收归约散布结果的相应部分;第二处理单元集合通过在第一处理单元集合之间执行归约散布子运算后获得的归约散布结果的相应本地部分对归约散布结果的各部分执行归约。

[0175] 条款20:根据条款15所述的第一计算节点,其中根据第一集群通信算法在第一计算节点中的第一处理单元集合之间执行全局聚集子运算包括:根据第一集群通信算法,在第一处理单元集合的第一处理单元处从第一处理单元集合的第二处理单元接收数据块;以及在第一处理单元处用本地数据块归约所接收的数据块。

[0176] 条款21:一种由第一计算节点实施的方法,所述方法包括:至少部分地基于与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器是否位于同一计算节点中或链接到同一片交换机来确定将数据从第一进程路由到第二进程的路由标识符,第一进程和第二进程属于特定网络拓扑下连接多个不同节点的特定节点间环;以及根据路由标识符将数据从第一进程路由到第二进程。

[0177] 条款22:根据条款21所述的方法,其中与第一进程相关联的网络接口控制器被配置为仅向环形拓扑中的第二计算节点发送数据或从第二计算节点接收数据,第二计算节点不同于第一计算节点。

[0178] 条款23:根据条款21所述的方法,其中与第一进程相关联的网络接口控制器还与一个或多个进程相关联,其中从第一进程和一个或多个进程发送的所有数据都通过网络接口控制器发送。

[0179] 条款24:根据条款21所述的方法,其中特定网络拓扑包括胖树拓扑。

[0180] 条款25:根据条款21所述的方法,还包括:响应于确定与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器位于同一计算节点中或链接到同一片交换机,将路由标识符设置为默认标识符。

[0181] 条款26:根据条款21所述的方法,还包括:响应于确定与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器位于不同的计算节点或链接到不同的片交换机,将路由标识符设置为等于与第一进程相关联的网络接口控制器的标识符。

[0182] 条款27:根据条款26所述的方法,其中根据路由标识符将数据从第一进程路由到第二进程包括:至少通过与网络接口控制器连接的片交换机和具有与网络接口控制器的标识符具有对应关系的标识符的汇聚交换机将数据从第一进程路由到第二进程,该网络接口控制器与第一进程相关联。

[0183] 条款28:一个或多个机器可读介质,其存储有机器可读指令,当机器可读指令由第

一计算节点执行时,使第一计算节点执行动作,所述动作包括:至少部分地基于与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器是否位于同一计算节点中或链接到同一片交换机来确定将数据从第一进程路由到第二进程的路由标识符,第一进程和第二进程属于特定网络拓扑下连接多个不同节点的特定节点间环;以及根据路由标识符将数据从第一进程路由到第二进程。

[0184] 条款29:根据条款28所述的一个或多个机器可读介质,其中与第一进程相关联的网络接口控制器被配置为仅向环形拓扑中的第二计算节点发送数据或从第二计算节点接收数据,第二计算节点不同于第一计算节点。

[0185] 条款30:根据条款28所述的一个或多个机器可读介质,其中与第一进程相关联的网络接口控制器还与一个或多个进程相关联,其中从第一进程和一个或多个进程发送的所有数据都通过网络接口控制器发送。

[0186] 条款31:根据条款28所述的一个或多个机器可读介质,其中特定网络拓扑包括胖树拓扑。

[0187] 条款32:根据条款28所述的一个或多个机器可读介质,其中所述动作还包括:响应于确定与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器位于同一计算节点中或链接到同一片交换机,将路由标识符设置为默认标识符。

[0188] 条款33:根据条款28所述的一个或多个机器可读介质,其中所述动作还包括:响应于确定与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器位于不同的计算节点或链接到不同的片交换机,将路由标识符设置为等于与第一进程相关联的网络接口控制器的标识符。

[0189] 条款34:根据条款33所述的一个或多个机器可读介质,其中根据路由标识符将数据从第一进程路由到第二进程包括:至少通过与网络接口控制器连接的片交换机和具有与网络接口控制器的标识符具有对应关系的标识符的汇聚交换机将数据从第一进程路由到第二进程,该网络接口控制器与第一进程相关联。

[0190] 条款35:一种第一计算节点,包括:一个或多个处理单元;以及存储器,其存储机器可执行指令,当机器可执行指令由一个或多个处理单元执行时,使一个或多个处理单元执行动作,所述动作包括:至少部分地基于与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器是否位于同一计算节点中或链接到同一片交换机来确定将数据从第一进程路由到第二进程的路由标识符,第一进程和第二进程属于特定网络拓扑下连接多个不同节点的特定节点间环;以及根据路由标识符将数据从第一进程路由到第二进程。

[0191] 条款36:根据条款35所述的第一计算节点,其中与第一进程相关联的网络接口控制器被配置为仅向环形拓扑中的第二计算节点发送数据或从第二计算节点接收数据,第二计算节点不同于第一计算节点。

[0192] 条款37:根据条款35所述的第一计算节点,其中与第一进程相关联的网络接口控制器还与一个或多个进程相关联,其中从第一进程和一个或多个进程发送的所有数据都通过网络接口控制器发送。

[0193] 条款38:根据条款35所述的第一计算节点,其中所述动作还包括:响应于确定与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器位于同一计算节

点中或链接到同一片交换机,将路由标识符设置为默认标识符。

[0194] 条款39:根据条款35所述的第一计算节点,其中所述动作还包括:响应于确定与第一进程相关联的网络接口控制器和与第二进程相关联的网络接口控制器位于不同的计算节点或链接到不同的片交换机,将路由标识符设置为等于与第一进程相关联的网络接口控制器的标识符。

[0195] 条款40:根据条款39所述的第一计算节点,其中根据路由标识符将数据从第一进程路由到第二进程包括:至少通过与网络接口控制器连接的片交换机和具有与网络接口控制器的标识符具有对应关系的标识符的汇聚交换机将数据从第一进程路由到第二进程,该网络接口控制器与第一进程相关联。

[0196] 条款41:一种由第一计算节点实施的方法,所述方法包括:

[0197] 根据节点感知减半加倍算法确定用于从第一进程向第二进程发送数据包的汇聚标识符,第一进程和第二进程属于在特定网络拓扑下连接至不同片交换机的不同节点;以及通过与汇聚标识符对应的汇聚交换机从第一进程向第二进程发送数据包。

[0198] 条款42:根据条款41所述的方法,还包括:为定向到连接至不同片交换机的节点的数据包分配不同的汇聚标识符,以使得能够通过不同的汇聚交换机将数据包路由到连接至不同片交换机的节点。

[0199] 条款43:根据条款41所述的方法,还包括:至少部分地基于预定对应关系,分配与汇聚标识符相关联的汇聚交换机对应的源端口和目的端口。

[0200] 条款44:根据条款43所述的方法,其中对应关系记录多个汇聚交换机的汇聚标识符与对应的源端口和目的端口对之间的关系。

[0201] 条款45:根据条款41所述的方法,其中特定网络拓扑包括胖树拓扑。

[0202] 条款46:根据条款41所述的方法,还包括:通过分配给各数据包的多个不同汇聚标识符对应的多个不同的汇聚交换机,从第一计算节点包括的第一进程集合向第二计算节点包括的第二进程集合发送各数据包。

[0203] 条款47:根据条款41所述的方法,还包括:通过分配给各数据包的多个不同汇聚标识符对应的多个不同的汇聚交换机,由第一计算节点包括的第一进程集合从第二计算节点包括的第二进程集合接收数据包。

[0204] 条款48:一个或多个机器可读介质,其存储有机器可读指令,当机器可读指令由第一计算节点执行时,使第一计算节点执行动作,所述动作包括:根据节点感知减半加倍算法确定用于从第一进程向第二进程发送数据包的汇聚标识符,第一进程和第二进程属于在特定网络拓扑下连接至不同片交换机的不同节点;以及通过与汇聚标识符对应的汇聚交换机从第一进程向第二进程发送数据包。

[0205] 条款49:根据条款48所述的一个或多个机器可读介质,其中所述动作还包括:为定向到连接至不同片交换机的节点的数据包分配不同的汇聚标识符,以使得能够通过不同的汇聚交换机将数据包路由到连接至不同片交换机的节点。

[0206] 条款50:根据条款48所述的一个或多个机器可读介质,其中所述动作还包括:至少部分地基于预定对应关系,分配与汇聚标识符相关联的汇聚交换机对应的源端口和目的端口。

[0207] 条款51:根据条款50所述的一个或多个机器可读介质,其中对应关系记录多个汇

聚交换机的汇聚标识符与对应的源端口和目的端口对之间的关系。

[0208] 条款52:根据条款48所述的一个或多个机器可读介质,其中特定网络拓扑包括胖树拓扑。

[0209] 条款53:根据条款48所述的一个或多个机器可读介质,其中所述动作还包括:通过分配给各数据包的多个不同汇聚标识符对应的多个不同的汇聚交换机,从第一计算节点包括的第一进程集合向第二计算节点包括的第二进程集合发送各数据包。

[0210] 条款54:根据条款48所述的一个或多个机器可读介质,其中所述动作还包括:通过分配给各数据包的多个不同汇聚标识符对应的多个不同的汇聚交换机,由第一计算节点包括的第一进程集合从第二计算节点包括的第二进程集合接收数据包。

[0211] 条款55:一种第一计算节点,包括:一个或多个处理单元;以及存储器,其存储有机器可执行指令,当机器可执行指令由一个或多个处理单元执行时,使一个或多个处理单元执行动作,所述动作包括:根据节点感知减半加倍算法确定用于从第一进程向第二进程发送数据包的汇聚标识符,第一进程和第二进程属于在特定网络拓扑下连接至不同片交换机的不同节点;以及通过与汇聚标识符对应的汇聚交换机从第一进程向第二进程发送数据包。

[0212] 条款56:根据条款55所述的第一计算节点,其中所述动作还包括:为定向到连接至不同片交换机的节点的数据包分配不同的汇聚标识符,以使得能够通过不同的汇聚交换机将数据包路由到连接至不同片交换机的节点。

[0213] 条款57:根据条款55所述的第一计算节点,其中所述动作还包括:至少部分地基于预定对应关系,分配与汇聚标识符相关联的汇聚交换机对应的源端口和目的端口。

[0214] 条款58:根据条款57所述的第一计算节点,其中对应关系记录多个汇聚交换机的汇聚标识符与对应的源端口和目的端口对之间的关系。

[0215] 条款59:根据条款55所述的第一计算节点,其中所述动作还包括:通过分配给各数据包的多个不同汇聚标识符对应的多个不同的汇聚交换机,从第一计算节点包括的第一进程集合向第二计算节点包括的第二进程集合发送各数据包。

[0216] 条款60:根据条款55所述的第一计算节点,其中所述动作还包括:通过分配给各数据包的多个不同汇聚标识符对应的多个不同的汇聚交换机,由第一计算节点包括的第一进程集合从第二计算节点包括的第二进程集合接收数据包。

[0217] 条款61:一种由第一计算节点实施的方法,该方法包括:将分配给处理单元的数据包划分成多个数据段,该多个数据段至少包括第一数据段和第二数据段;将多个数据段分配给多个线程,该多个线程至少包括第一线程和第二线程;使用第一线程对第一数据段的一部分执行节点内子运算,且并行地使用第二线程对第二数据段的一部分执行节点间子运算。

[0218] 条款62:根据条款61所述的方法,其中使用第一线程对第一数据段的一部分执行节点内子运算包括:通过节点内连接在第一计算节点中包括的处理单元和另一处理单元之间传输第一数据段的该部分。

[0219] 条款63:根据条款61所述的方法,其中使用第二线程对第二数据段的一部分执行节点间子运算包括:通过节点间连接在处理单元和不同于第一计算节点的第二计算节点中包括的另一处理单元之间传输第二数据段的该部分。

[0220] 条款64:根据条款61所述的方法,其中节点内子运算包括:在第一计算节点内执行的归约散布子运算或全局聚集子运算,并且节点间子运算包括在第一计算节点和不同于第一计算节点的第二计算节点间执行的全局归约子运算。

[0221] 条款65:根据条款61所述的方法,其中节点内子运算包括在第一计算节点内执行的全局聚集子运算或复制子运算,节点间子运算包括在第一计算节点和不同于第一计算节点的第二计算节点间执行的全局聚集子运算。

[0222] 条款66:根据条款66所述的方法,还包括:使用第一线程对第一数据段的一部分执行另一个节点间子运算,以及并行地使用第二线程对第二数据段的一部分执行另一个节点内子运算。

[0223] 条款67:根据条款61所述的方法,其中使用第一线程对第一数据段的该部分执行节点内子运算,以及并行地使用第二线程对第二数据段的该部分执行节点间子运算,使得使用节点内连接来将第一数据段的一部分传输至第一计算节点包括的另一处理单元,以及并发地使用节点间连接将第二数据段的一部分传输至与第一计算节点不同的第二计算节点包括的另一处理单元。

[0224] 条款68:一个或多个机器可读介质,其存储有机器可读指令,当机器可读指令由第一计算节点执行时,使第一计算节点执行动作,所述动作包括:将分配给处理单元的数据包划分成多个数据段,该多个数据段至少包括第一数据段和第二数据段;将多个数据段分配给多个线程,该多个线程至少包括第一线程和第二线程;使用第一线程对第一数据段的一部分执行节点内子运算,且并行地使用第二线程对第二数据段的一部分执行节点间子运算。

[0225] 条款69:根据条款68所述的一个或多个机器可读介质,其中使用第一线程对第一数据段的部分执行节点内子运算包括通过节点内连接在第一计算节点中包括的另一处理单元和处理单元之间传输第一数据段的该部分。

[0226] 条款70:根据条款68所述的一个或多个机器可读介质,其中使用第二线程对第二数据段的部分执行节点间子运算包括:通过节点间连接在处理单元和不同于第一计算节点的第二计算节点中包括的另一处理单元之间传输第二数据段的该部分。

[0227] 条款71:根据条款68所述的一个或多个机器可读介质,其中节点内子运算包括在第一计算节点内执行的归约散布子运算或全局聚集子运算,并且节点间子运算包括在第一计算节点和不同于第一计算节点的第二计算节点间执行的全局归约子运算。

[0228] 条款72:根据条款68所述的一个或多个机器可读介质,其中节点内子运算包括在第一计算节点内执行的全局聚集子运算或复制子运算,并且节点间子运算包括在第一计算节点和不同于第一计算节点的第二计算节点间执行的全局聚集子运算。

[0229] 条款73:根据条款68所述的一个或多个机器可读介质,所述动作还包括:使用第一线程对第一数据段的部分执行另一个节点间子运算,以及并行地使用第二线程对第二数据段的部分执行另一个节点内子运算。

[0230] 条款74:根据条款68的一个或多个机器可读介质,其中使用第一线程对第一数据段的一部分执行节点内子运算,以及并行地使用第二线程对第二数据段的一部分执行节点间子运算,使得使用节点内连接来将第一数据段的一部分传输至第一计算节点包括的另一处理单元,以及并发地使用节点间连接将第二数据段的一部分传输至与第一计算节点不同

的第二计算节点包括的另一处理单元。

[0231] 条款75:一种第一计算节点,包括:一个或多个处理单元;以及存储器,其存储有机器可执行指令,当机器可执行指令由一个或多个处理单元执行时,使一个或多个处理单元执行动作,所述动作包括:将分配给处理单元的数据包划分成多个数据段,该多个数据段至少包括第一数据段和第二数据段;将多个数据段分配给多个线程,该多个线程至少包括第一线程和第二线程;使用第一线程对第一数据段的一部分执行节点内子运算,且并行地使用第二线程对第二数据段的一部分执行节点间子运算。

[0232] 条款76:根据条款75所述的第一计算节点,其中使用第一线程对第一数据段的一部分执行节点内子运算包括:通过节点内连接在第一计算节点中包括的另一处理单元和处理单元之间传输第一数据段的该部分。

[0233] 条款77:根据条款75所述的第一计算节点,其中使用第二线程对第二数据段的部分执行节点间子运算包括:通过节点间连接在处理单元和不同于第一计算节点的第二计算节点中包括的另一处理单元之间传输第二数据段的该部分。

[0234] 条款78:根据条款75所述的第一计算节点,其中节点内子运算包括在第一计算节点内执行的归约散布子运算或全局聚集子运算,并且节点间子运算包括在第一计算节点和不同于第一计算节点的第二计算节点间执行的全局归约子运算。

[0235] 条款79:根据条款75所述的第一计算节点,其中节点内子运算包括在第一计算节点内执行的全局聚集子运算或复制子运算,并且节点间子运算包括在第一计算节点和不同于第一计算节点的第二计算节点间执行的全局聚集子运算。

[0236] 条款80:根据条款75所述的第一计算节点,其中使用第一线程对第一数据段的该部分执行节点内子运算,以及并行地使用第二线程对第二数据段的一部分执行节点间子运算,使得使用节点内连接来将第一数据段的一部分传输至第一计算节点包括的另一处理单元,以及并发地使用节点间连接将第二数据段的一部分传输至与第一计算节点不同的第二计算节点包括的另一处理单元。

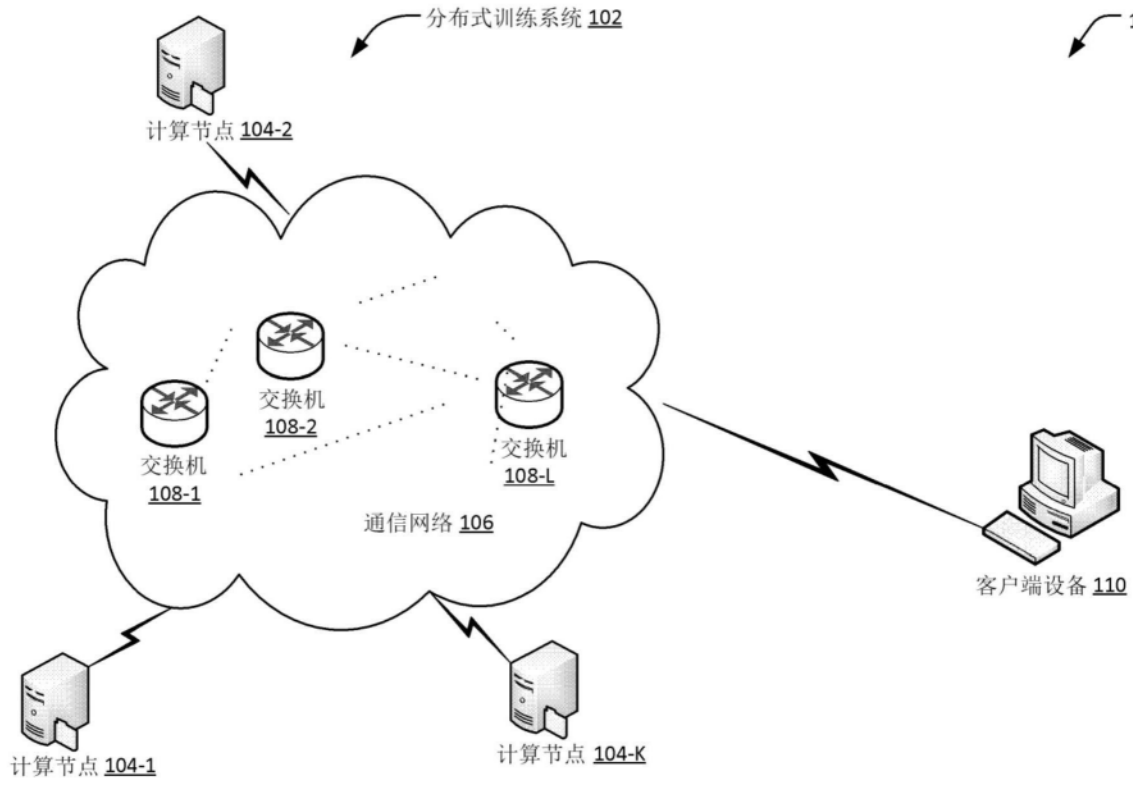


图1

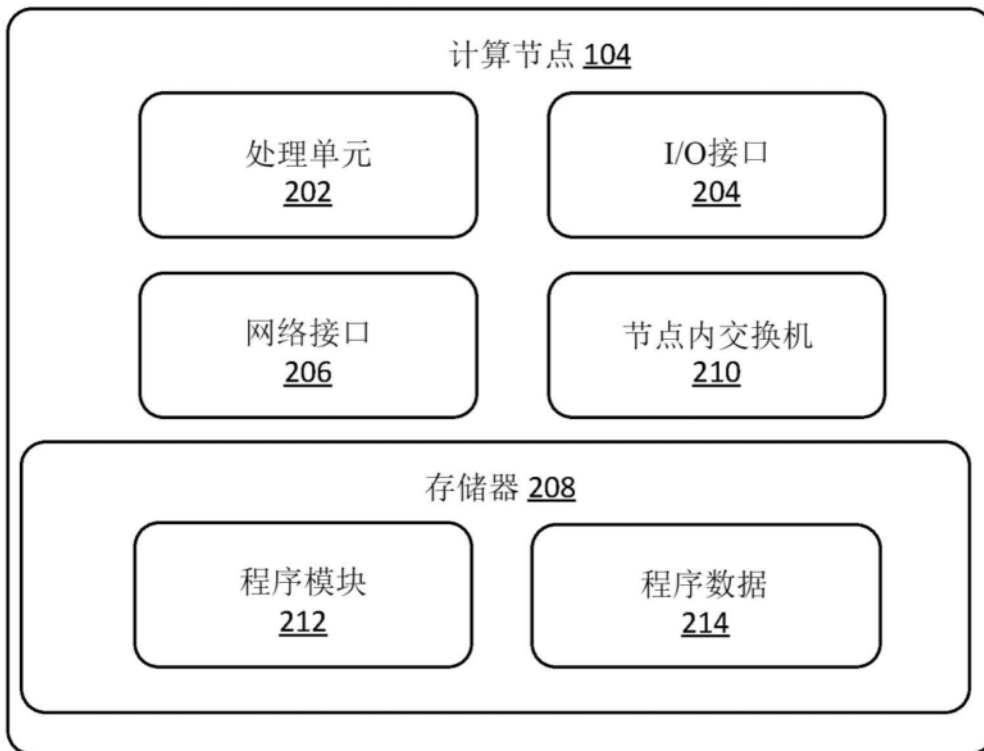


图2

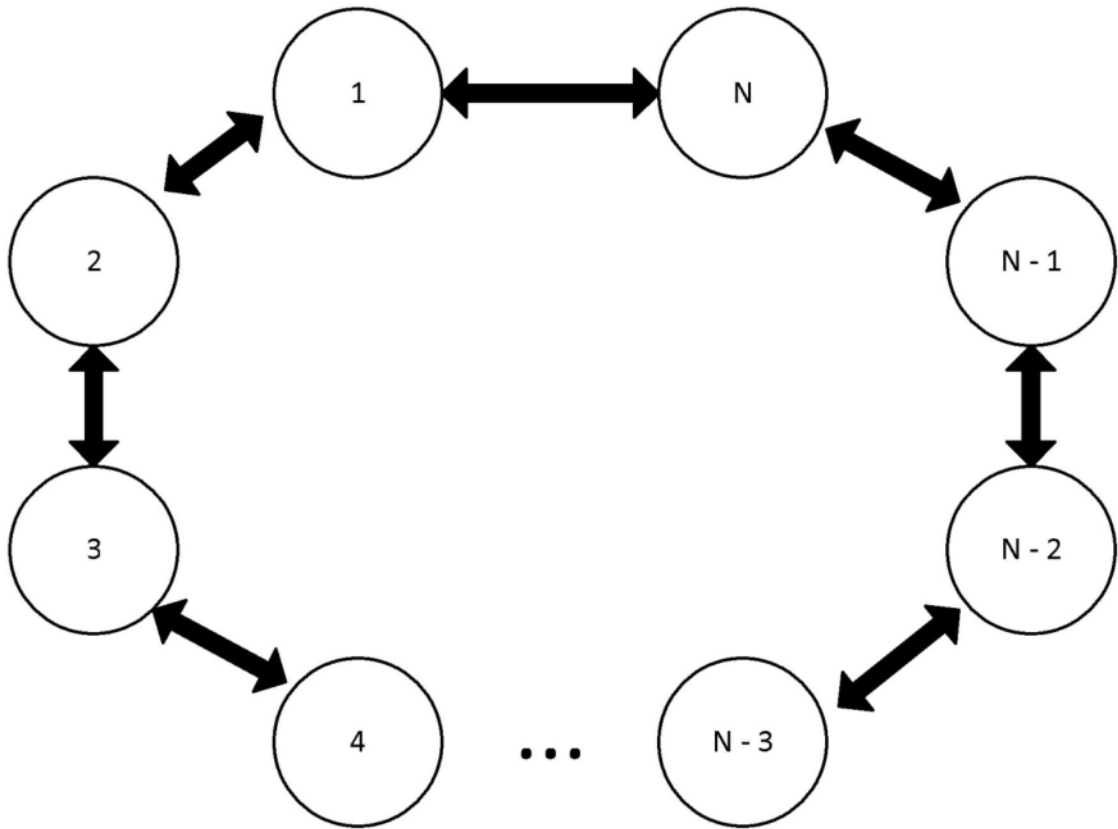


图3A

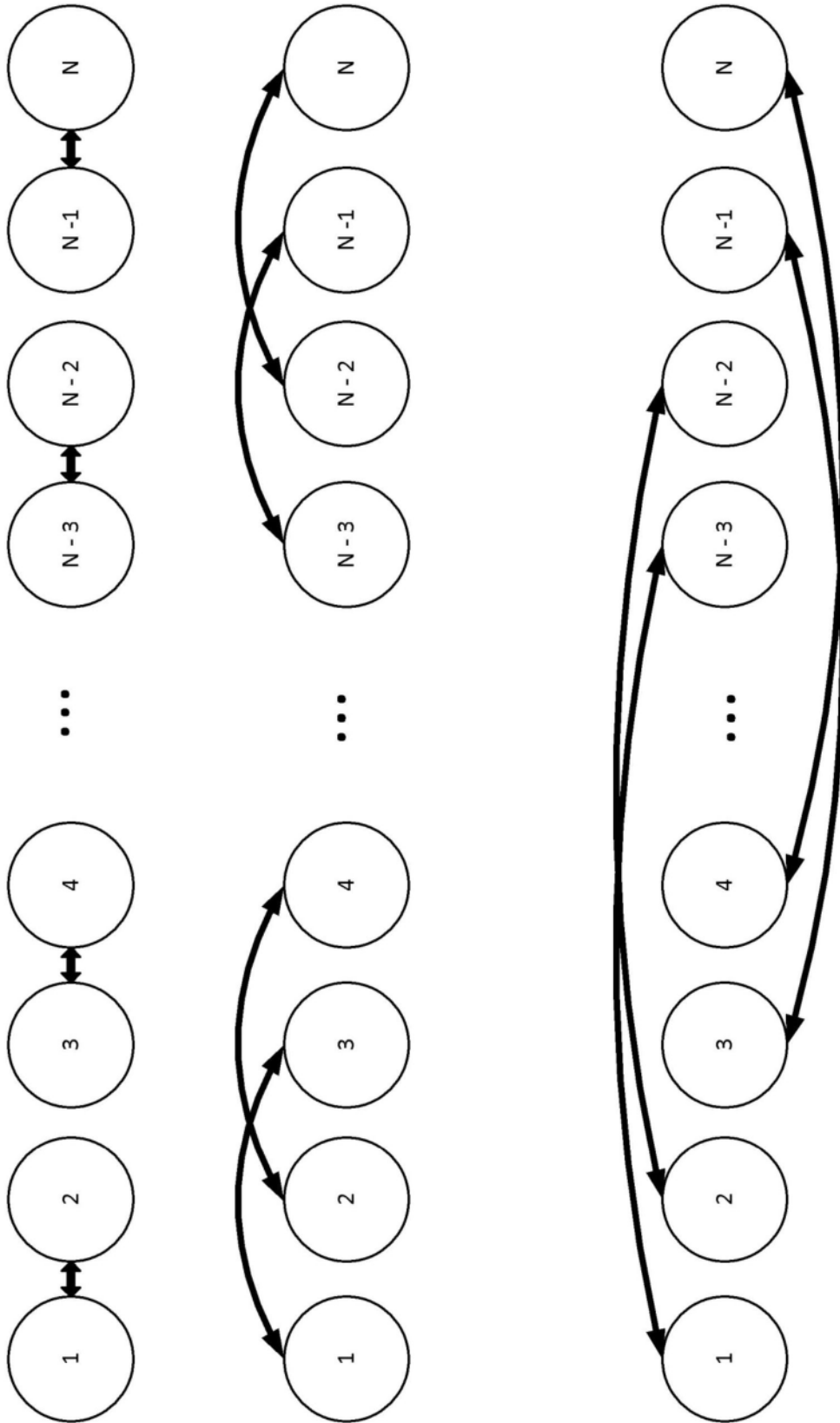


图3B

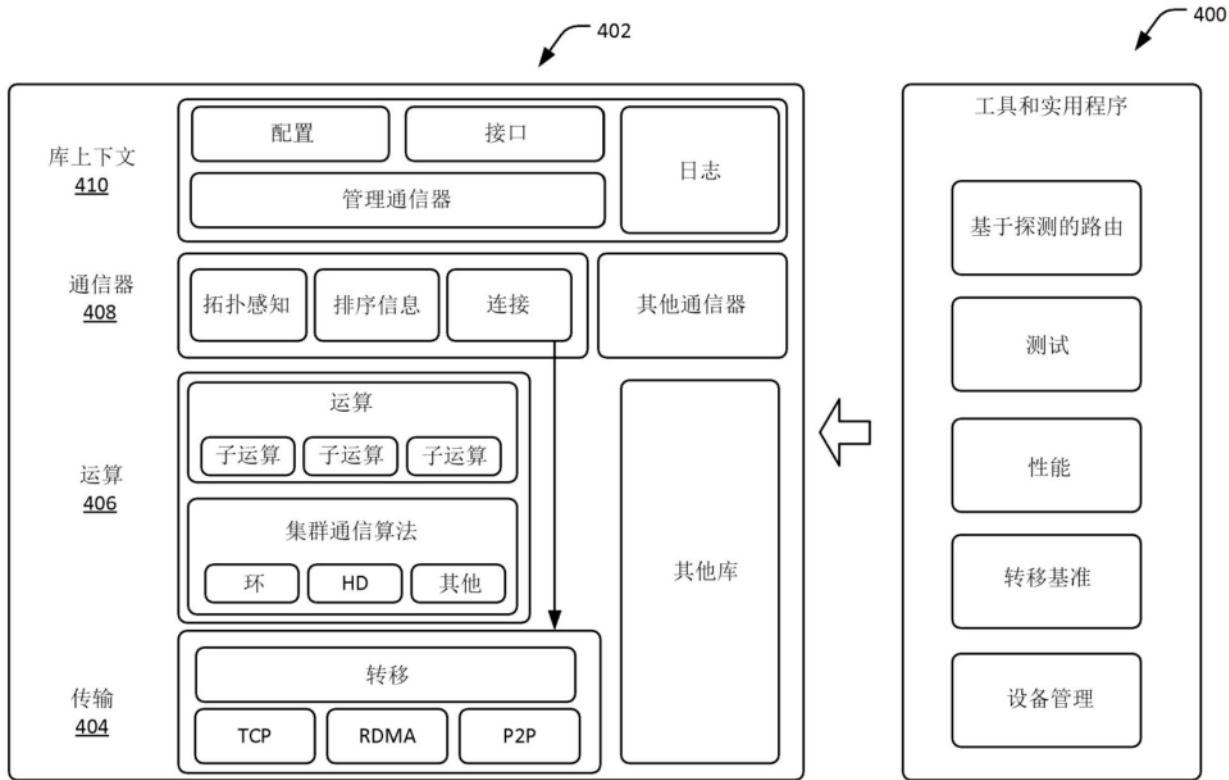


图4

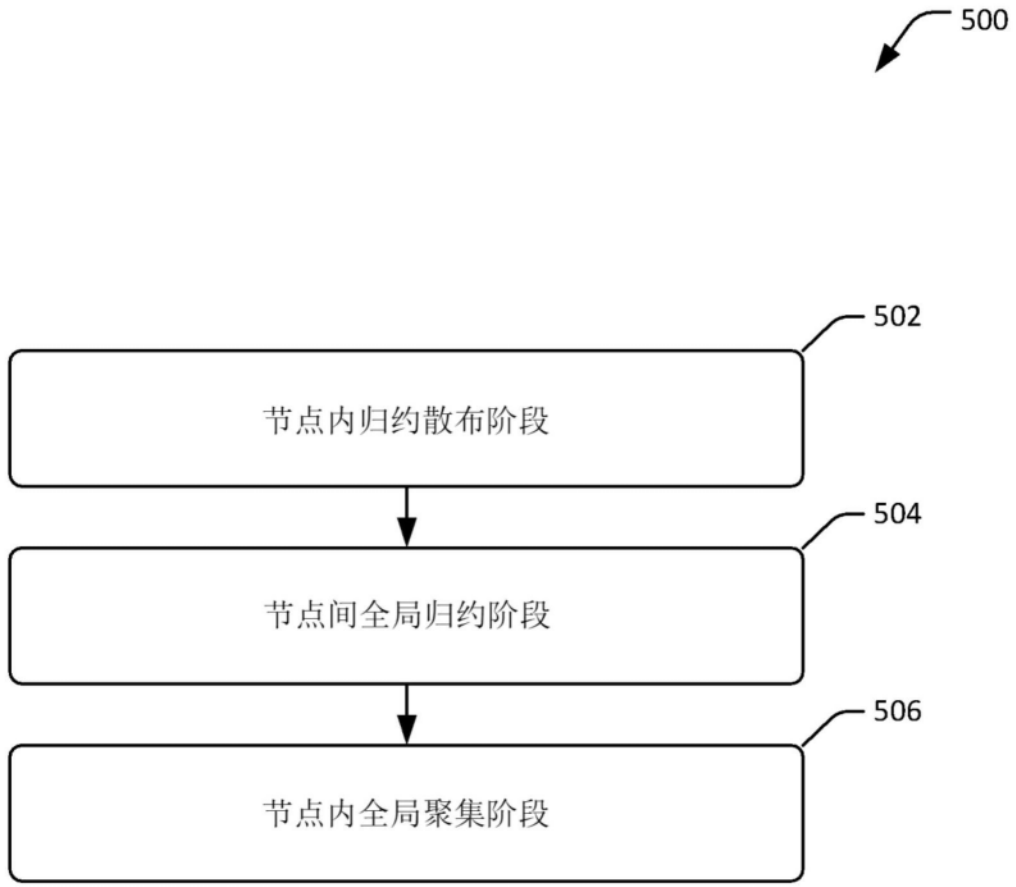


图5

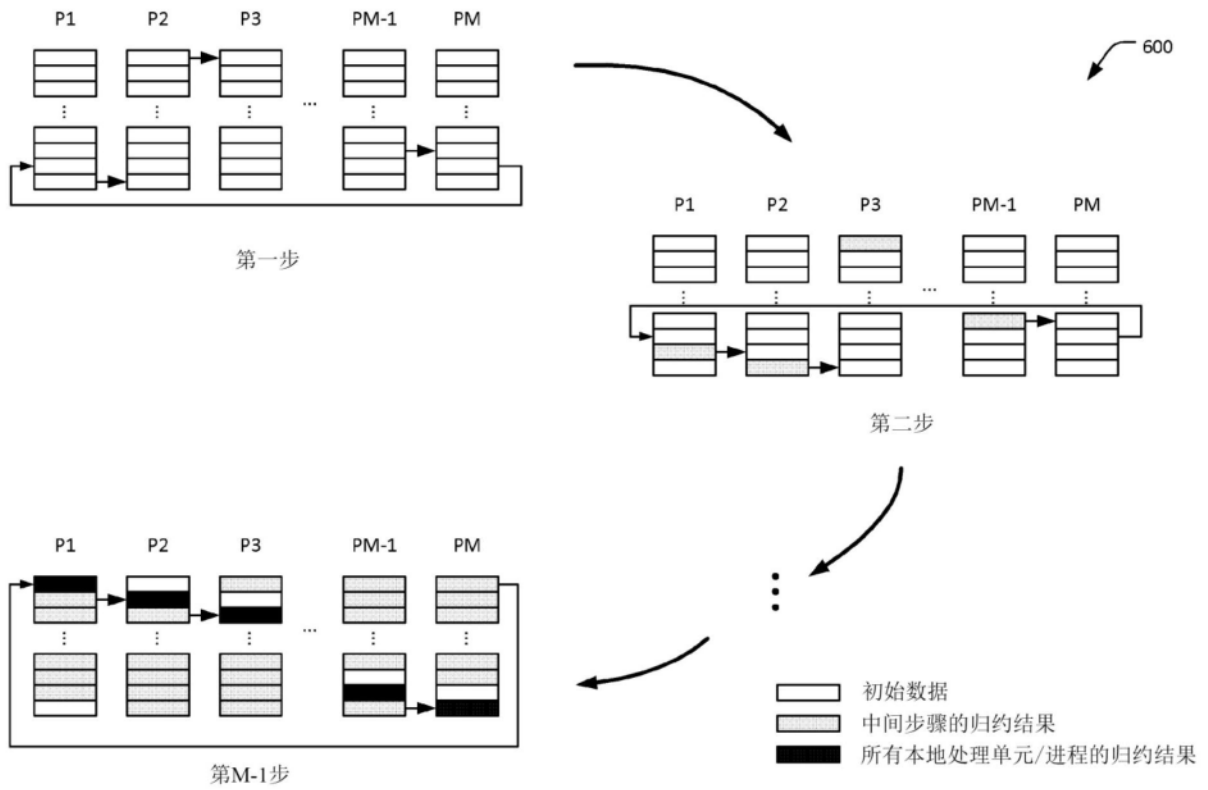


图6

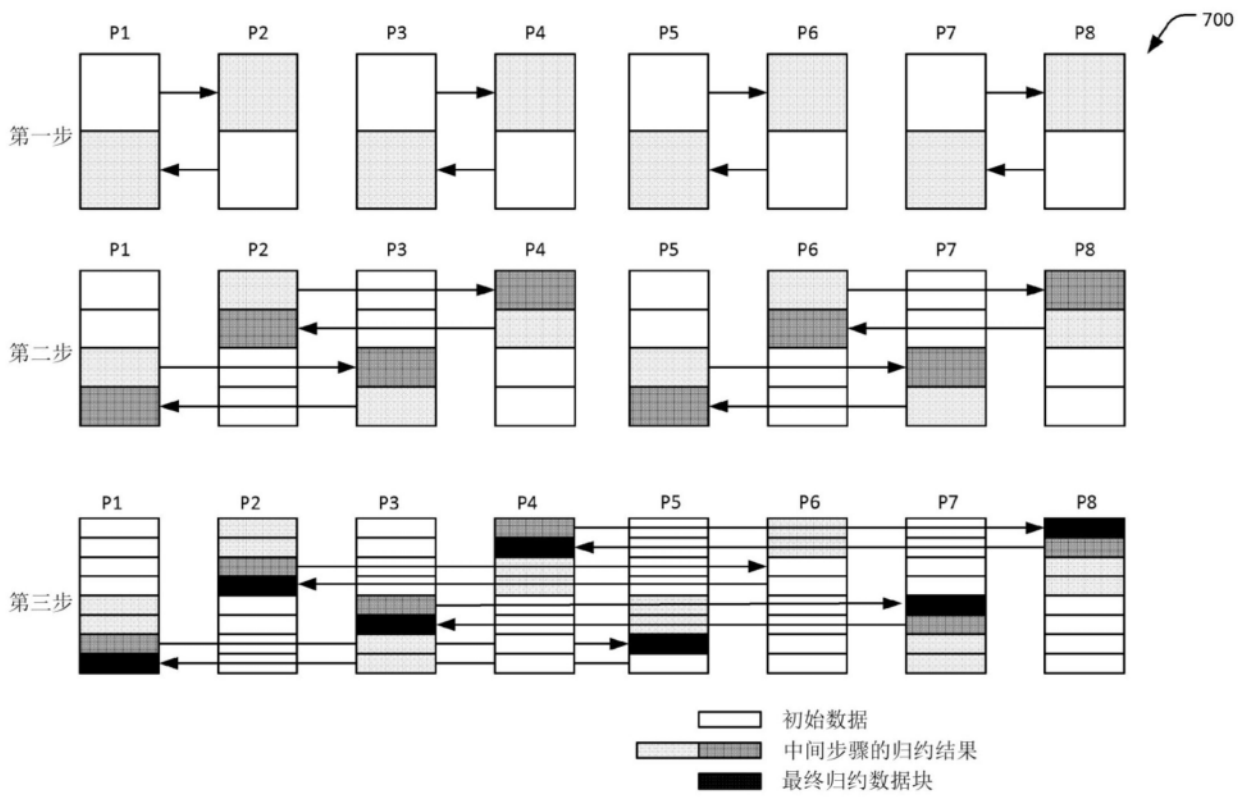


图7

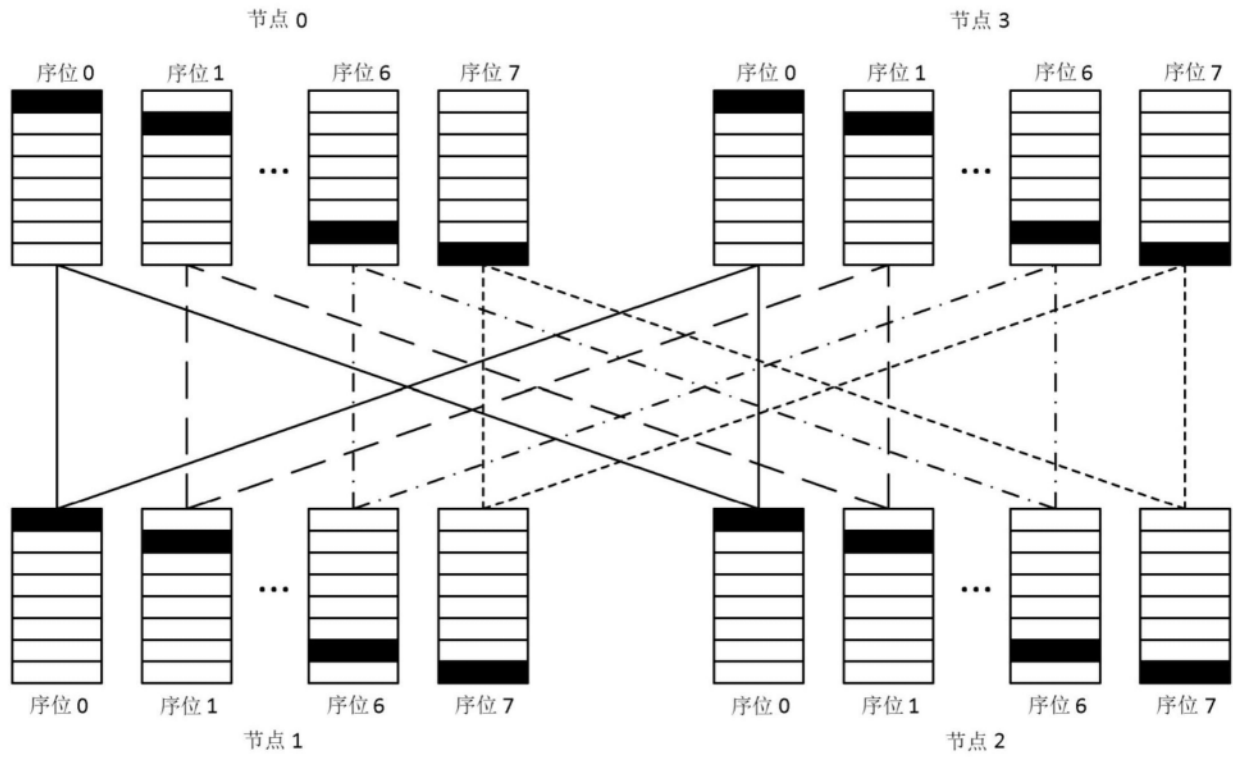


图8

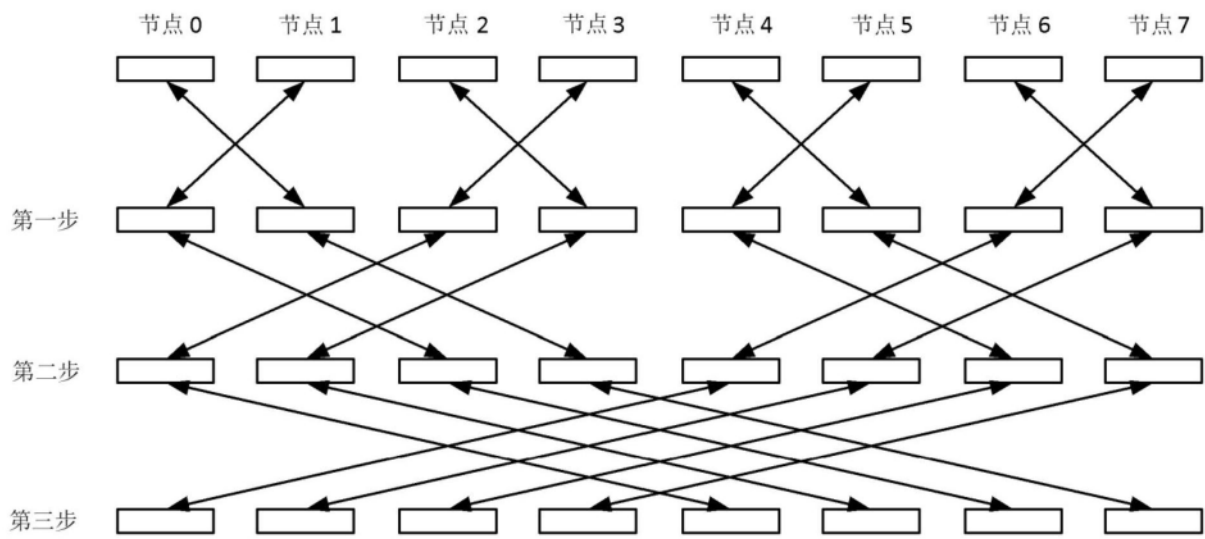


图9

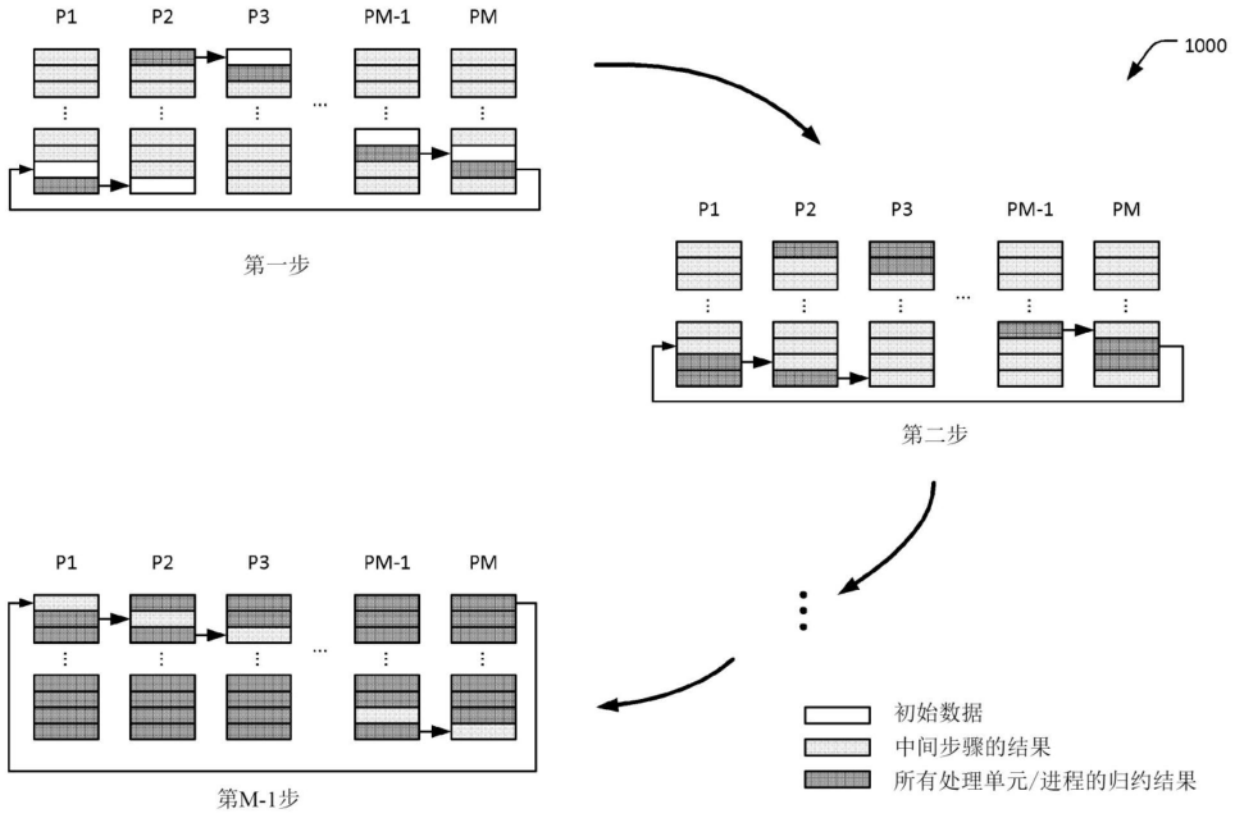


图10

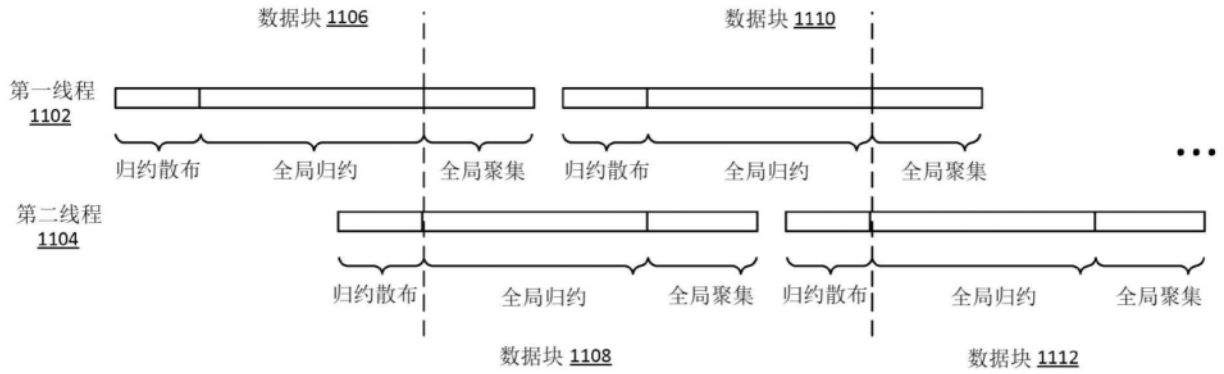


图11

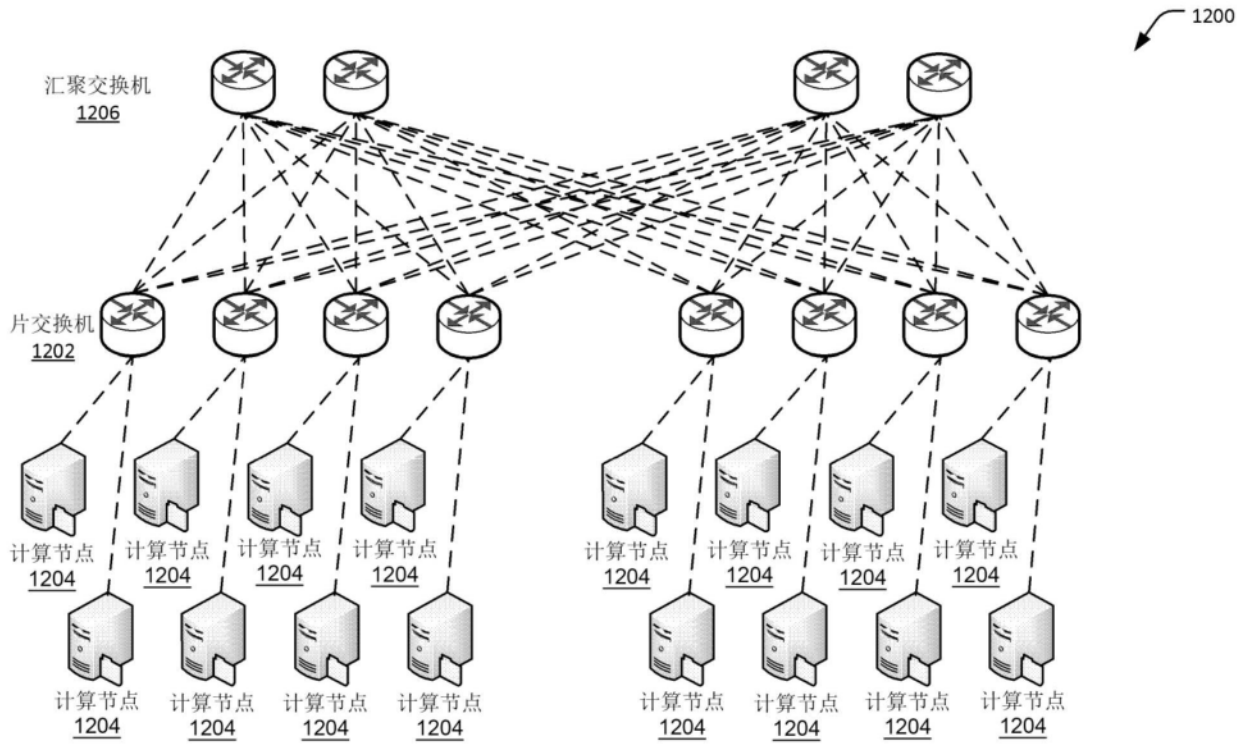


图12

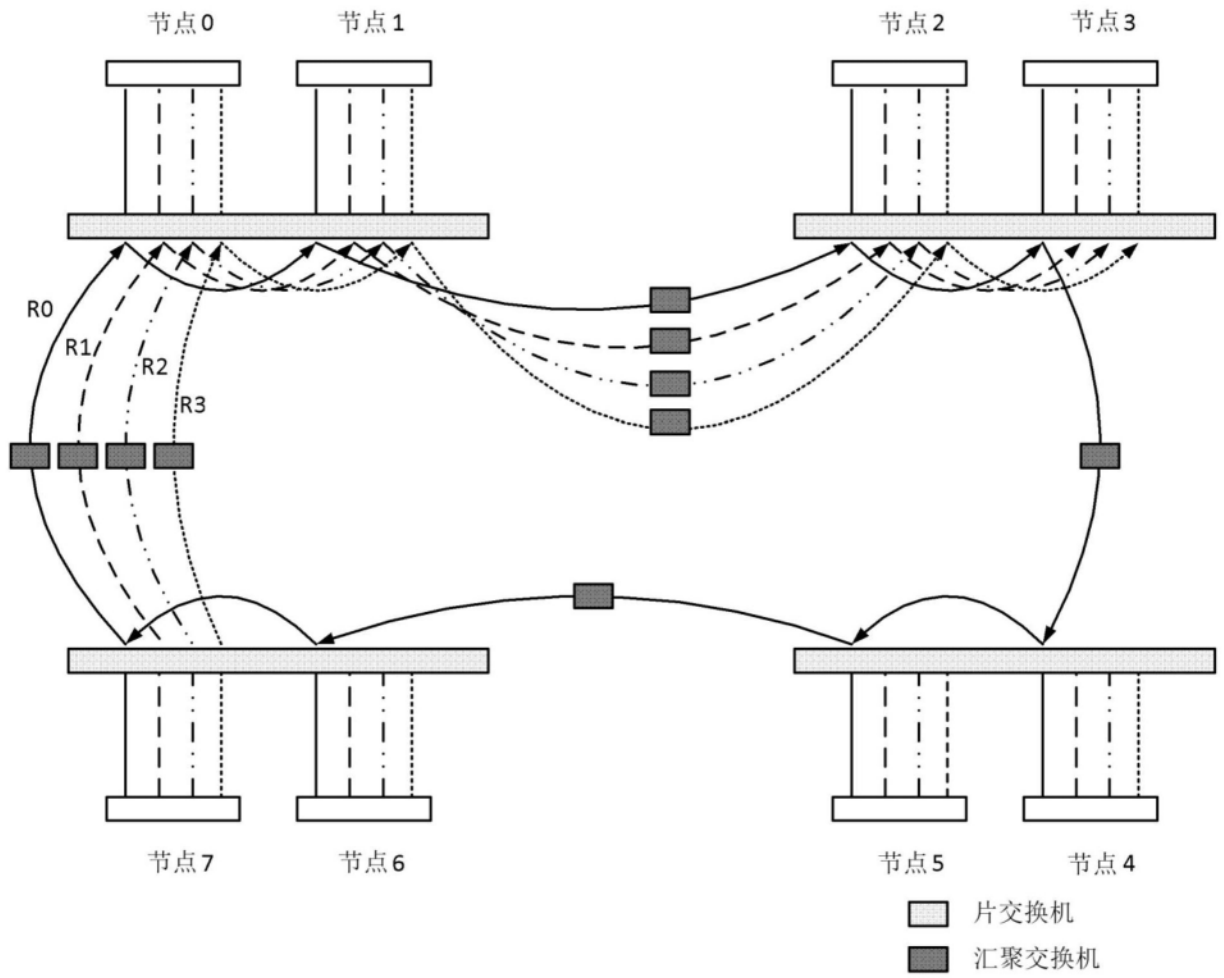


图13

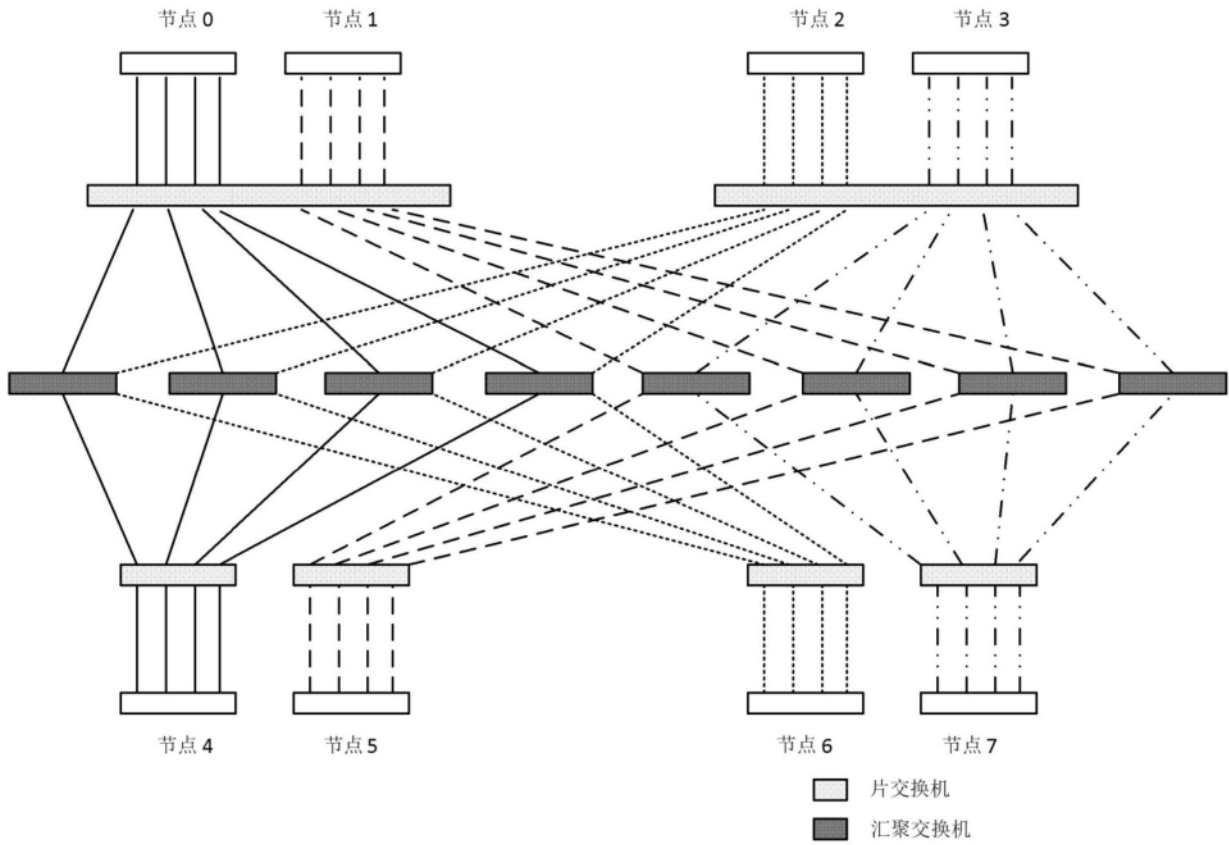


图14

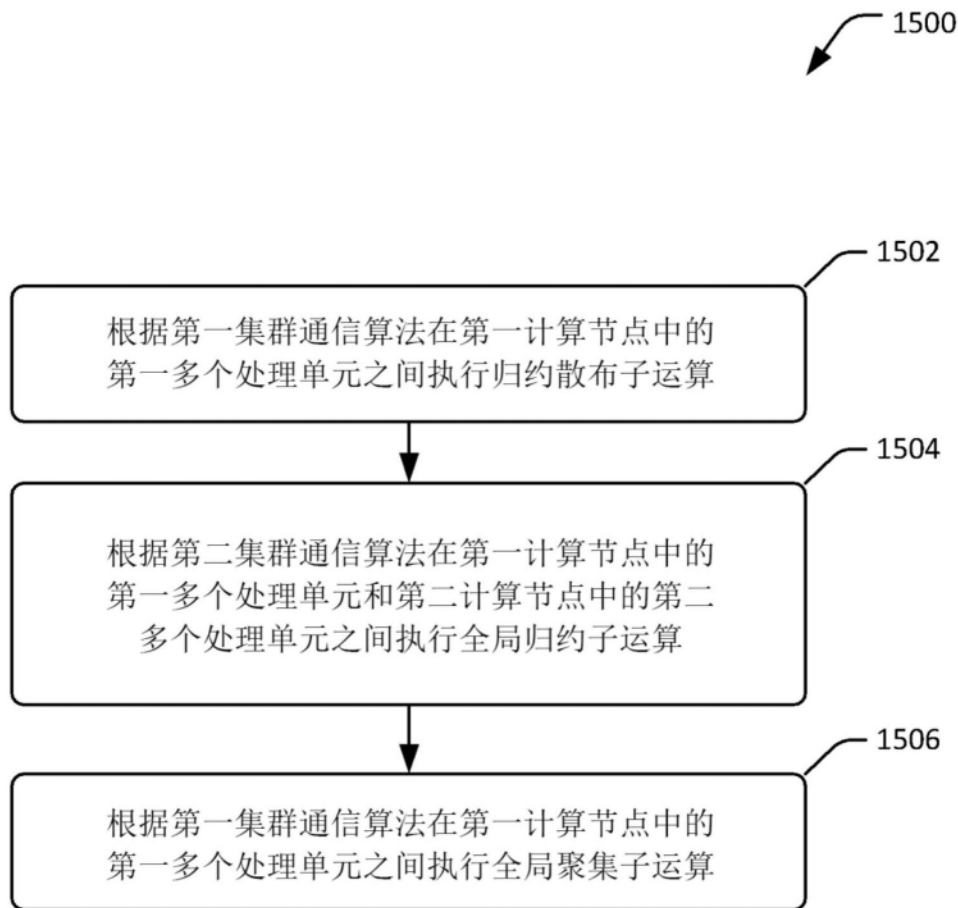


图15

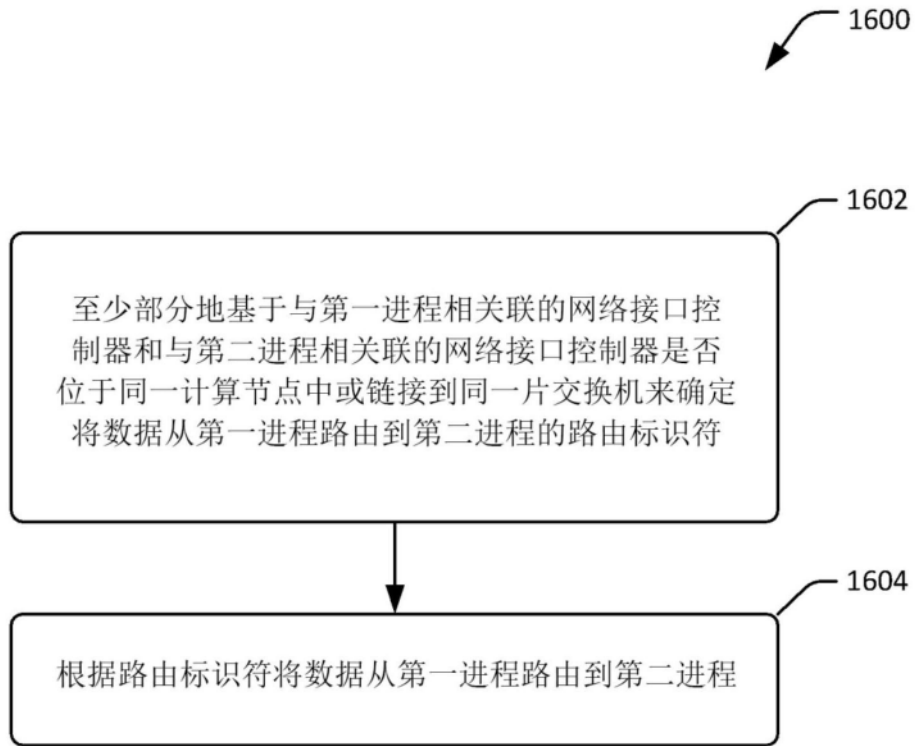


图16

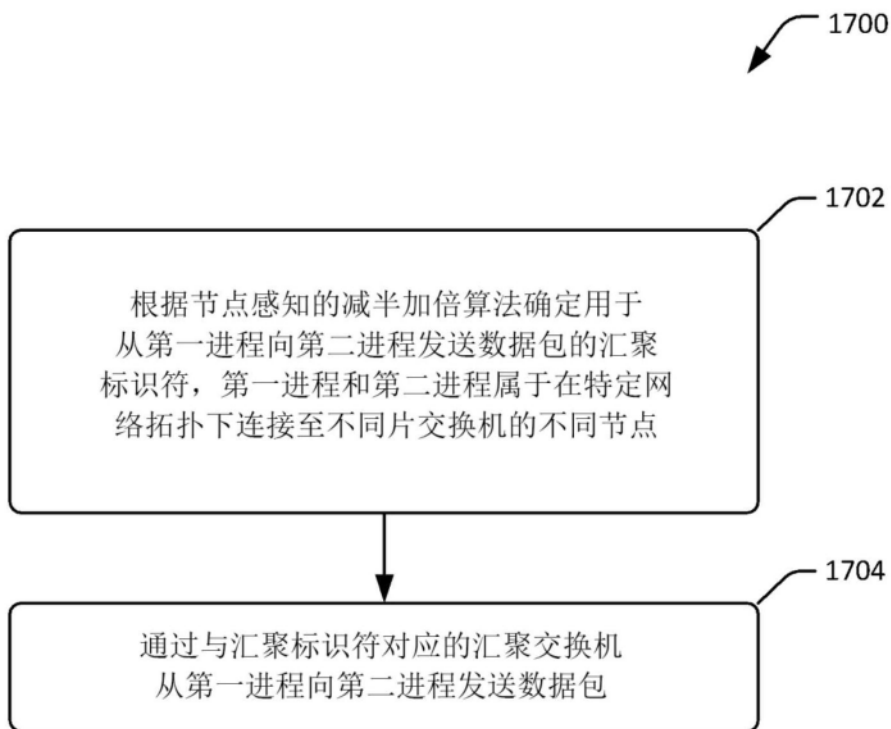


图17

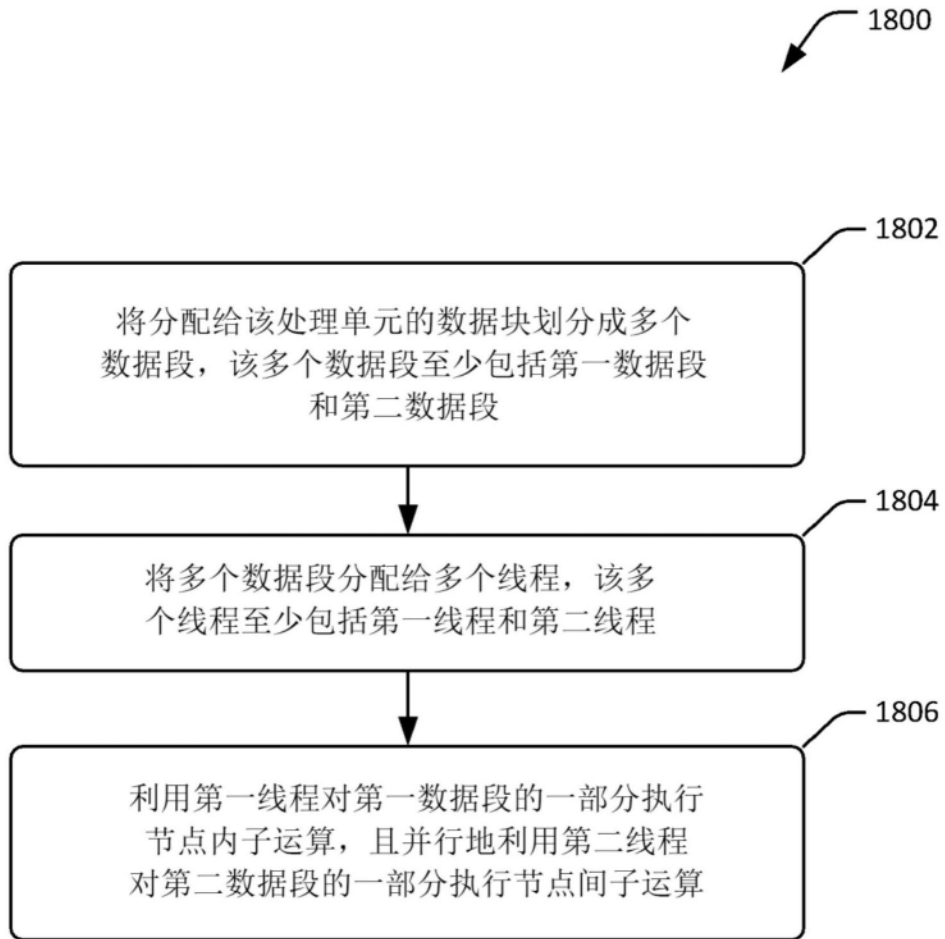


图18