

文章编号: 1007-4619(2005)04-0421-08

空间数据多服务器群集及高性能实现策略

严志民^{1,2}, 刘仁义¹, 刘南¹, 陆丽珍¹

(1. 浙江大学 浙江省 GIS 重点实验室, 浙江 杭州 310028; 2. 杭州市房产信息中心, 浙江 杭州 310006)

摘 要: 针对当前海量空间数据网络化管理的特殊要求, 分析了改善新一代地理信息系统 (GIS) 整体性能所需要解决的问题, 剖析并引入了群集技术, 充分发挥群集技术在高性能运算和高可靠性多服务器架构等方面的技术优势, 优化设计了 GIS 群集多 PC 服务器体系架构, 并结合空间索引分区等常规海量空间数据库管理策略, 开发了一套实用化的高效空间定位、多缓冲策略等基于群集的多服务器海量空间数据管理技术。Oracle 真正应用群集 (RAC) 数据库系统在图库管理中的应用, 验证了群集空间数据多服务器管理平台的整体性能优越性。

关键词: 计算机群集; 地理空间数据; 高性能策略; Oracle RAC

中图分类号: P208/TP392 **文献标识码:** A

1 引 言

在海量空间数据的管理和分析处理中, 当前流行的空间数据管理系统如 ArcSDE 等在管理机制和性能优化方面做了很多努力, 但是在对海量数据的上下下载、查询和分析等处理中, 频显 CPU、硬盘等资源不足和效率低下。随着基于网络的 GIS 服务需求的不断扩大, 海量空间数据的处理逐步从科学计算型过渡到了决策支持型, 要求设计开发出高性能的多服务器空间数据管理平台, 既支持海量空间数据处理、复杂业务逻辑处理, 服务器端又要具备高可伸缩性和高可用性; 面向海量空间数据的大型 GIS 应用更是对整个系统的稳定性、运算效率等性能方面提出了更高的要求。国内外正在进行研究开发的第四代 GIS 平台, 其突破重点仍主要集中在解决跨图层访问、多源空间数据一体化和网络环境下的异构空间数据库系统互操作上^[1]; 分布式计算主要建立基于 TCP/IP 协议的快速空间数据传输构件和高效空间数据搜索算法, 整体平台的性能并没有实质性的改善。

群集是通过高速网络互联并以单一系统模式加以管理的计算机组合, 具有良好的可伸缩性、高度的可用性、负载平衡和并行运算高效性^[2]等方面特

点, 实现了将 PC 机以及其他低档机组成虚拟超级机, 并达到性能并行超级计算机化。群集技术中的负载和故障动态转移机制、多节点处理器的分布式共享内存直接访问技术、低开销消息传递系统如消息传递接口, 以及基于群集的双机热备份技术等, 为解决大型关系型数据库在海量分布式空间数据管理中存在的性能问题, 实现优化并高效地管理多服务器海量地理空间信息提供了新的技术和思路。

国内将群集技术与 GIS 相结合的可参考研究和文献极少, 本文尝试将群集技术引入新一代 GIS 多服务器应用系统的设计中, 在数据库服务器层利用 Oracle 的真正应用群集 (Real Application Cluster RAC) 搭建群集数据库服务器, 并将群集技术融入多服务器的空间数据管理技术中, 拟解决大型关系型数据库在海量空间分布式数据管理中存在的性能问题, 优化并高效地实现海量空间信息的分布式管理。

2 群集空间数据管理服务器体系和 RAC

群集环境, 包含紧耦合的同构部件, 其具有常规分布式环境的一般特性, 如环境中所有节点作为一个逻辑整体统一对外提供透明数据存取等服务。群集系统的目标是充分发挥并行计算机的优势, 充分

收稿日期: 2003-09-27; 修订日期: 2004-05-18
 (c) 1994-2012 China Academic Journal Electronic Publishing House. All rights reserved. http://www.cnki.net

基金项目: 国家自然科学基金 (40271087), 国家 863 高科技研究项目 (2001AA135180), 浙江省自然科学基金 (401006)。

作者简介: 严志民 (1971—), 男, 浙江衢州, 工程师, 博士, 主要从事空间数据库基础和分布式地理信息系统研究和应用开发。

调动系统中所有资源并行地完成任 务,其与常规分布式的区别在于突出强调系统整体性能的提高。由于目前网络传输技术的水平,群集技术在局域网或高速广域网中的优势更明显。

群集技术与大型关系型数据库的结合相对成熟。至目前为止,Oracle⁹ⁱ和 IBM DB²等商用大型关系型数据库系统都提供了数据库的群集并行服务技术,这些系统中每个服务器一般都有自己的系统引导盘,可以独立运行,而数据存储在 RAID 阵列或网络存储系统中。当数据库群集中某个服务器出现硬件或软件故障时,其他服务器中数据库仍继续对外服务,故障解除后的节点可以重新方便地加入群集中恢复正常服务。群集数据库技术在分布并行计算、集中管理方面的优势,特别适合于 GIS 业这种有大数据量运算的空间数据中心。

RAC^[3]是 ORACLE 公司提供的空间数据库群集技术,它管理多个相互连接计算机的处理能力,将软件和硬件集合组成一个群集单元,每个组成部分的处理能力产生一强劲的计算环境。RAC 在可扩展性和高可用性特征的群集软件结构是一个突破,可以利用 RAC 传递高性能、增加吞吐量和高可用性等。在 RAC 环境中,所有活动的实例可以同时执行基于一个共享数据库的事务。RAC 协调每个实例对共享数据的访问,保证数据的一致性和数据的完

整性。RAC 充分利用了群集的优势,可以高效地分解和分发大任务,并提高了处理大工作负载和容纳快速增长的用户群的处理能力。利用 RAC,可以在不改变应用编码的前提下,扩展应用以适应增长的数据处理需求。加入资源如节点或存储器,RAC 延伸了这些资源单组成制约之外的处理能力。数据仓库应用是 RAC 访问只读数据的基本代表。并且,RAC 成功管理在线事务处理 (OLTP) 系统和混合系统,组合了可以读写和只读应用的特性。

RAC 应用了 Oracle 高速缓存熔合体系结构新技术,能迅速、有效地在群集的所有计算机上共享那些经常被访问的数据。RAC 支持多进程的并行访问,群集节点具有一致的数据文件和控制文件,却有不同的 SGA、日志文件和回滚段等。RAC 物理存储采用了外挂共享磁盘阵列 (如 RAID⁵) 能保持共享数据的一致性,进一步提高了数据库的稳定性、降低故障率。而 IBM DB² 的非共享存储系统过分依赖网络,其可应用性受到了当前网络带宽小、稳定性差等现状的极大限制。

本文设计的空间数据库群集利用 Oracle 的 RAC 搭建群集多服务器架构 (图 1),结合 Oracle Spatial 对地理空间数据特有的支持能力,并利用下文设计的双机互备等技术,拟提高大型海量空间数据库的可靠性、可扩展性以及并发处理能力等。

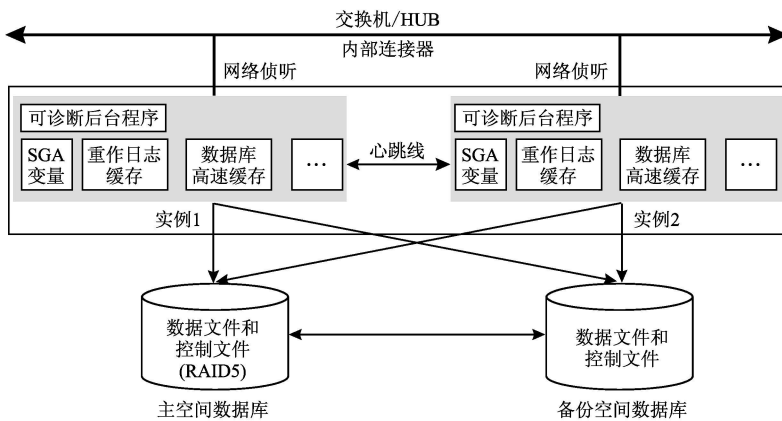


图 1 RAC 共享存储群集与互备空间数据库架构

Fig 1 Sharing storage cluster and mutual-backup spatial database structure of RAC

3 群集多服务器空间数据管理技术

3.1 群集扩展的常规分布式空间数据库管理

3.1.1 透明访问机制和命名机制

前面提到群集系统中不同层次的群集服务器都

是作为一个整体透明对外提供服务的,进程具体连接到那个服务器实例,对客户端来说都是透明的。GIS 业务深入到空间对象管理层次,客户端更关心的是面向几何对象的操作。多个群集空间库体现了大型系统的广义分布性,新型群集多服务器数据库对象访问透明机制的实现借鉴了普通网络环境的命

名机制^[4],结合多层多级特点建立一套树状命名机制,海量空间多级数据库相对复杂一些,需要深入考虑几何要素集 DATASET和表对象 FEATURECLASS 层次的全局命名唯一性。

多层多服务器数据库系统采用自底向上的设计方法,当全局概念模型建立成功之后,在全局模型之上为需要访问全局数据的用户定义全局视图,全局用户使用全局统一的结构化空间数据查询语言透明访问全局数据库。另外,利用别名或同义词方式可以更好地隐藏对象的实际地址,使访问更便捷和安全。

3.1.2 数据存储策略

群集要求数据高度集中存放,在全局分布式总体存放原则的基础上,对于局域网内部的不同层次级别的不同应用专题的空间数据都可以考虑集中存储在统一的数据中心集中管理。

但目前多服务器数据库系统受网络速度限制,其多层架构服务器之间应尽量减少数据传输。在 GIS 应用中,现有的空间数据库内容丰富而且数据量从 GB、TB 乃至 PB 日渐庞大,物理存放机制直接影响着整个系统的性能。集中和分布如何取舍是目

前需要协调的关键问题。

总体存放的原则有^[5]:局部或本地的数据本地存放;使用频率最多使用原则,调用最多的地方为数据库存放地;面向主题或应用专题存放。

单点存放需要综合考虑服务器性能:各结点 CPU、内存、磁盘容量;每一结点需传递的事务量;每一结点使用的数据量;网络的性能与可靠性;若结点间连接不通后的访问规则等。

3.1.3 远程内嵌过程和替代触发器策略

群集强调的是局部整体性和高效性,分布式空间应用程序要处理远程跨多群集或非群集数据库的大量数据,本文利用远程数据库内嵌对象提高应用程序的性能。在远程数据库上执行内嵌的过程或者函数,仅把执行结果返回应用程序,降低网络负担,改善远程数据操作的性能。另外,本文结合远程内嵌过程和替代触发器策略解决了大型数据库如 Oracle 在跨服务器远程空间对象管理中存在的许多问题,据 ORACLE 公司确认的管理空间数据的几个 BUG,发挥了重要作用。

图 2 展示了一个基本跨服务器远程对象处理的实现策略。

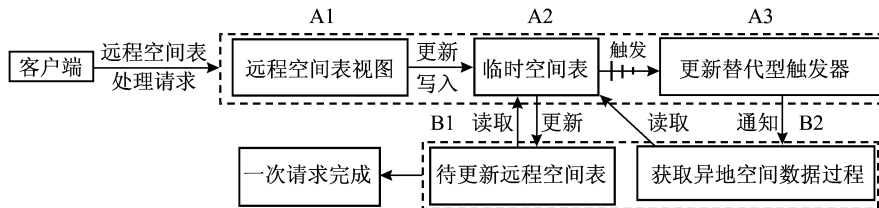


图 2 大型分布式数据库原子级的空间数据处理

Fig 2 An atom of spatial data processing in large-scale distributed database

3.1.4 元数据库管理机制

空间数据库特有的空间性,使得元数据库包含更多内容,如空间数据格式、数据精度质量等诸多信息。群集支撑下的空间元数据的管理在局域网络内部相对一般空间库管理要简单的多,但寓于广域网络中多级分布的集群环境,涉及了空间数据的异地、多源、异构等,本文引进空间目录树技术,结合轻量型目录访问技术(LDAP)管理多级元数据,空间数据目录记录

空间数据的地点、信息类型或来源以及各级服务器所能操作远程数据库服务器结点间的逻辑关系等。

本文为多服务器群集数据库系统的实现设计了一套专用空间数据表。

表 1 用来管理空间几何数据,据此建立空间数据库逻辑表达关系。FID 在整个系统唯一,包含数据库识别信息。此处‘共享级别’是专门用来适应多级服务器数据库设计引入的。

表 1 特征元数据表 (SMDB Feature)

Table 1 Metadata of features

特征 ID	数据集 ID	特征类型	特征名	所有者	...	别名	共享级别
Unique	1	Simple

表 2 记录整个多服务器群集数据库系统中的所有服务器结点信息,以便配置网络服务。并快速方

便地定位多服务器中库结点以及记录所在的具体层次。

表 2 多服务器数据库系统结点元数据表 (SMDB D is=DBNode)

Table 2 Metadata of nodes in multi-servers database

结点 ID	主机名	数据库唯一标识 SID	...	本地命名	网络服务名	上级主结点	结点级别
Unique	121.11.0.1	GIS863	...	GIS863	Dell-a	Solaris-A	1
...		

群集空间库中,考虑到多个实例并发的存在,应数据访问一致性的要求则必须保证实例调用中获取的元数据信息严格保持一致。全库备份应该同时对空间元数据进行备份。锁机制应用帮助共享存储设备的群集系统数据一致性的实现,而对于广度非共享存储器的数据库群集难度较大。

内存以及整体性能状态,以保证高效地为客户提供服务。在群集高性能并发机制的基础上,本文在各级群集服务器的主控服务器端设计了一套与空间数据相关的信息表(表 3 和表 4),从不同群集层次的主控信息中实现高效空间定位。

3.2 群集特色的服务器性能优化

利用上述并发主控信息表可以达到快速定位被处理的空间对象,其粒度可达到几何要素级,即数据表的行。这样实现在服务器的最外层就能直接定位判断要素级的占用和锁定状态,以指导群集服务器实例快速地重新分配负载,合理调度。主控信息字段的刷新和传输相比所指物理空间数据的调度效率,仍然是可行和有效的,能达到进一步性能的提高。

3.2.1 高效空间定位机制

群集架构中,同一群集中一般有一个独立的主控节点或类似于主控功能的节点。它统一调度群集中多个实例,即时探测并监控各个实例中的 CPU、

表 3 并发应用服务主控信息表

Table 3 Main control information of concurrent application services

应用服务器	群集实例号	CPU %	Memory %	进程数	进程号	进程相关数据库标识	进程占用
200.0.0.10	1	36	20	8	1	racdb	1055346
...							

表 4 并发数据库服务主控信息表

Table 4 Main control information of concurrent database services

数据库服务器	群集实例号	CPU %	Memory %	进程数	进程号	线程数	进程占用	进程占用要素类	进程占用要素
200.0.0.16	1	21	38	8	1	4	1055346	2456732	8654312
...									

3.2.2 多缓冲区策略

多客户并发访问需要一套先进的调度策略和缓存机制^[6],否则频繁的服务启动和数据直接读取会大大影响整个系统的性能。针对 GIS 大型应用中的三层应用体系结构中,本文对各应用服务层次设计了高速缓存机制(图 3)以进一步优化系统整体性能,从而大幅度地提高海量空间数据的处理分析能力。

系统和数据的可靠性。

双机互备的好处是,备份机不是简单通过硬件连接在本地的备份,它可以是通过软件进行的远程网络动态全复制,为主数据库分担客户负载支持,另外备份数据库还能保持自己特有的数据,并对外提供与主服务器不相关的其他服务,充分利用了系统资源。当然,最重要的是即使出现火灾、水淹、线路故障造成的系统损坏和逻辑损坏等的严重物理故障也能很好地防止。

3.2.3 空间数据库的双机互备^[7]

双机空间数据库群集系统对数据库系统的可用性有了极大的提高,但仍存在双机同时出现故障的可能。为此,本文借鉴以往数据库的备份机制,为系统设计了磁盘阵列和双机互备方式(图 1),以保证

本文采用异步复制技术的多主复制方案支持全表在主数据库和备用数据库间的对称复制,允许主数据库和备用数据库对主表都有更新操作的权利。主数据库上和备用数据库上复制表的更新都会被传

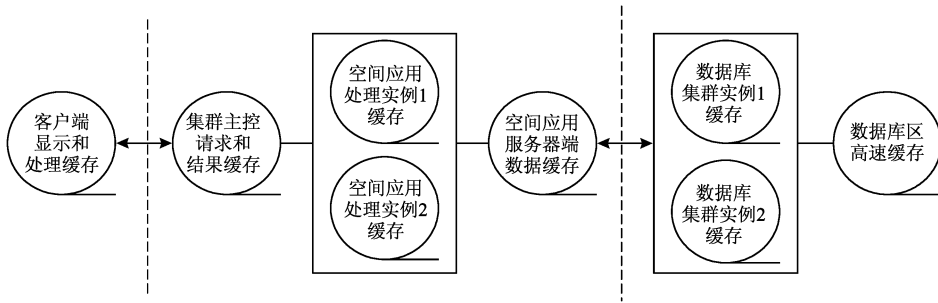


图 3 多级缓存结构示意图

Fig 3 Multi-levels caches architecture

播到对方。在传播数据变化时,如果其中的一个远端数据库系统没有准备好,传播变化的延迟远程过程调用就会保存在其远端系统的本地队列中,等到系统准备好以后再执行。当 Oracle 主服务器因磁盘阵列发生硬件故障或其他原因而造成系统停机和服务中断时,系统事务可以立刻转移到备用数据库上,当主数据库恢复后,备用数据库上的数据的变化将被复制到主数据库。这种机制大大提高了系统的可靠性。

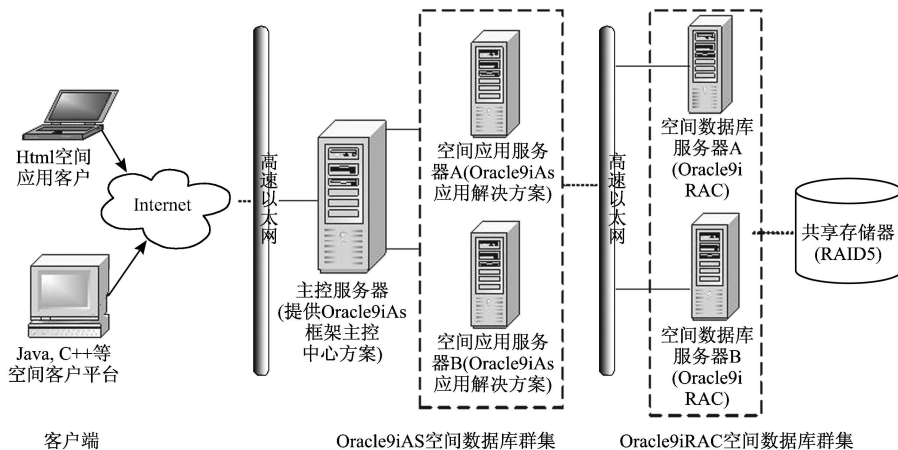
4 群集空间数据管理实例

空间数据是 GIS 技术的基础和研究主要对象, GIS 应用服务是空间数据管理的目的。结合前文的群集多服务器空间数据库管理技术与 Oracle 的 RAC 等技术,实例中结合应用了 Oracle 的 Oracle9iAs 应用服务器群集技术,设计基于群集技术的大型多服务器 GIS 应用系统,系统结构见图 4。大型多服务器空间库需考虑实际应用中空间库的多

级层次性,本应用采用了大比例尺地形图 800 幅,总数据量 8 G 的矢量数据。

在图库管理系统的实施中,综合应用了上文提到的群集多服务器空间数据库管理的策略和机制,如多应用实例并发服务;多服务器群集服务、高效空间定位机制、多缓冲区策略等、过程策略等,另外还利用了 ORACLE 提供的群集相关技术和同类服务器等相关策略与技术。图 4 是应用实例的系统结构,有服务器 5 台,其中数据库服务器和应用服务器各 2 台,系统运行环境中的服务器都为主频 2.4GHz 内存 512M,单 CPU 的 PC 机。

在该系统中,空间数据的应用与配置的物理服务器无关,所有应用通过逻辑应用名称由群集服务器统一进行调度和任务分派。在一个空间数据服务站中,至少有一个主群集服务器^[8],可以同时配置一个或多个空间数据管理服务器,每一空间数据管理服务器可以配置一个或多个地空间分析应用,同一应用也可以配置于不同的物理服务器上。用户访问群集服务器,具体执行处理任务的空间服务器



(C)1994-2021 China Academic Journal Electronic Publishing House. All rights reserved. http://www.cnki.net
图 4 Oracle9i 搭建的群集多服务器空间信息系统

Fig 4 Multi-servers spatial information system in Oracle9i clustering

对用户来讲可以是透明的,用户不用关心数据或处理由哪一台服务器来完成,从而提高了性能,同时方便开发。系统中采用多实例服务器处理并发用户请求,网络应答传输由群集支撑软件来完成,本文的自主空间数据管理组件(zSDM)提供数据管理和数据处理服务。对于每台服务器,可以利用业务复制和业务分割方法来配置空间数据库管理服务、地图管理服务、Web服务等多种服务,每一个应用可以启动多个实例,每一个实例的数据请求处理是同步的,而数据传输是异步的。

作者基于上述研究成果,在 Windows 平台下,利用 Oracle9 iAS 和 Java 基于 J2EE 设计开发了一套空间数据管理服务包,并在先前开发的大型多级分布式空间数据管理系统中拓展了群集功能。经过程序调用配置以后,客户端用户就可以透明高效地对操作更新或显示远程空间对象。群集空间库在整个系统中只是局部的一个节点,对用户来说与单点数据库无异。

群集大型多服务器空间数据库配置管理系统运行界面一览见图 5。

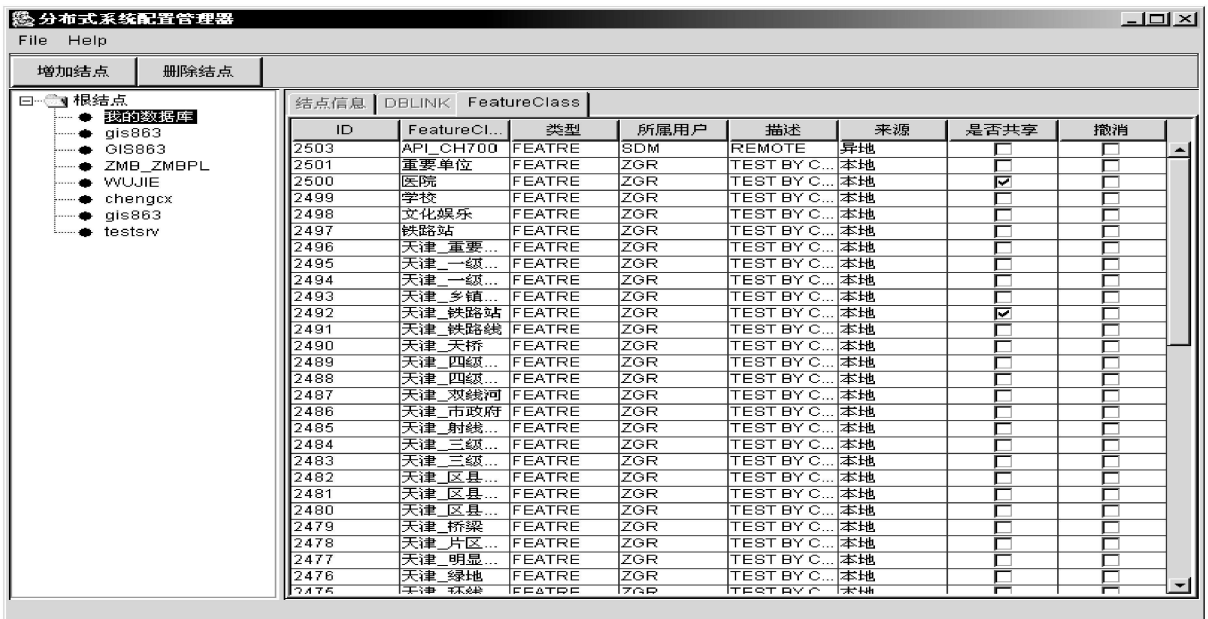


图 5 群集大型多服务器空间数据库管理系统

Fig 5 A large-scale multi-servers-cluster spatial database manage system

示例数据为图库中的一个矢量面数据表,有记录 1980257 条,物理数据量 286MB。

4.1 可用性测试

群集环境可用性测试试验以 RAC 群集为例,本文先保持两个数据库实例正常启动状态,在客户端正常连接到实例 RACDB1 的情况下,突然中断其中数据库实例 RACDB1,此时再查看客户端连接状态,可见客户端连接已经从实例 1 自动转移到了实例 2。

4.2 并行运算性能测试

表 5 为群集高性能测试结果表,执行的对空间表的全表下载指令

测试中的表格并行度和指令并行度都设为 4(并行度是从进程的数量),由线程监控程序可见,在表格

本身设置并行选项和指令级设立并行时,同样有多个线程被启用。海量表格的访问效率,会因服务器 CPU 个数和并行度大小配合而不同,一般原则上并行度是 CPU 个数的两倍,并非并行度越大越快。

空间数据库的管理需要在应用中体现。空间查询与分析包括求交、相邻、覆盖等多种关系,表 6 为启动应用服务器 iAS 后,结合数据库层 RAC 群集的空间表的空间包含操作性能测试表。

表 5 和表 6 的效率提升是与不采用群集技术的相同软硬件和基本配置环境下独立服务器数据管理性能的对比结果,由测试表数据表明群集环境下的空间数据处理效率等可量化性能的提升均值在 20% 以上。海量空间表格的操作相对复杂,运算量很大,大型系统的三层结构中多缓存和多层次的群集机制设计不同程度地提高了系统的性能。如果群

表 5 两层体系全表下载群集性能测试数据表(√:运行,×:放弃)

Table 5 Performance test data in 2-tiers cluster architecture

测试方案	表并行	指令并行	无 RAC 单服务器 (m : s m s)	RAC 单实例 (m : s m s)	RAC 双实例 (m : s m s)	性能优化率 /%
1	×	×	19:22.10	18:35.17	16:37.21	14.19
2	√	×	18:30.18	18:15.51	14:07.11	23.70
3	×	√	18:35.27	18:46.43	14:21.36	22.77
4	√	√	18:20.08	18:22.56	14:01.01	23.55

(注:表并行指物理表启用永久并行查询;指令并行指单个指令临时指定并行)

表 6 三层体系求空间包含关系的群集性能测试数据表(√:运行,×:放弃)

Table 6 Performance test data in 3-tiers cluster architecture

测试方案	无群集独立服务器 (m : s m s)	IAS 服务器单实例	IAS 服务器双实例	RAC 单实例	RAC 双实例	耗时 (m : s m s)	效率提升 /%
1	32:07.00	×	√	√	×	23:33.12	26.7
2	32:07.00	×	√	×	√	21:44.09	32.3
3	32:07.00	√	×	√	×	26:30.03	17.5
4	32:07.00	√	×	×	√	23:52.07	25.7

集中服务器数量越多, CPU 数越多, 内存越大, 理论上应该是性能成倍增加, 但事实上的多主机调度、并行事务的协调、网络速度稳定性和读写的一致性处理等对资源的消耗使整体性能大受影响。

5 结 论

本文提出的基于群集的多服务器空间服务器架构、空间数据库管理的效率机制、透明访问机制以及多缓冲机制等, 开拓了新一代 GIS 技术快速发展的崭新思路; 本文基于群集技术的研究以及应用尝试, 初步证明该技术与 GIS 技术的结合是可行的, 从实验的结果表明群集技术确实能为海量空间数据库的管理提供更好的可用性、可靠性以及负载均衡与高性能的运算。基于目前的网络性能, 在高速网络内部能极大地提高大型多服务器地理信息的整体性能。群集技术引入 GIS 研究是势在必行, 它能解决当今大型地理信息系统发展面临的几个瓶颈问题。群集技术在性能和价格等方面的优势, 使大型 GIS 的平民化成为可能。

参 考 文 献 (References)

[1] Chen B, Fang Y. Large-scale Distributed GIS Techniques and

Development [J]. Journal of Image and Graphics China, 2001, 6A (9): 862-864. [陈斌, 方裕. 大型分布式地理信息系统的技术与发展 [J]. 中国图象图形学报, 2001, 6A (9): 862-864.]

[2] Rajkumar Buyya. High Performance Cluster Computing Architectures and Systems (Volume 1). Beijing: Publishing House of People's Post and Telecommunication, 2002.

[3] Mark Bauer. Oracle9i Real Application Clusters Concepts Release 2 (9.2) [DB/OL]. <http://otn.oracle.com/OracleCorporationOnlineDocumentation>, 2002.

[4] Doreen L. Galli. Distributed Operating Systems Concepts and Practice [M]. Beijing: China Machine Press, 2003.

[5] Managing a Distributed Database [DB/OL]. <http://otn.oracle.com/pls/db901/homepage>, Oracle9i Database Online Documentation Release 9.0.1.

[6] Feisi Technique Research&developer Center. Oracle9i Application Server in Detail [M]. Publishing House of Electronics industry, 2003. [飞思科技产品研发中心. Oracle9i 应用服务器详解 [M]. 北京: 电子工业出版社, 2003.]

[7] Tan Z K, Kuang Z J. Oracle9i Database Administrator's Advanced Techniques Guides [M]. China Railway Publishing House, 2003. [谈竹奎, 况志军. Oracle9i 数据库管理员高级技术指南 [M]. 北京: 中国铁道出版社, 2003.]

[8] Oracle9i Application Server (An Oracle White Paper), <http://otn.oracle.com/OracleCorporationOnlineDocumentation>, 2002.

Multi-servers Clustering of Spatial Data and Realization Strategies in High Performance

YAN Zhimin^{1,2}, LIU Renyi¹, LIU Nan¹, LU Lizheng¹

(1. Zhejiang Provincial Keylab of GIS, Zhejiang University, Hangzhou 310028, China;

2. Information Center, Hangzhou Administration Bureau of Housing Property, Hangzhou 310006, China)

Abstract Analysis for the need that should be resolved in developing new geographic information system more deeper study and analysis for cluster techniques and Oracle Real Application Cluster (RAC) were taken. To meet special requirement for current management of massive spatial data on net, an optimal multi-level PC servers architecture in GIS with clustering techniques was designed to take full advantages of high performance computing, high reliability and other cluster features. Integrated the reliability, scalability, flexibility and other features of RAC with general spatial database management strategies such as spatial index and partition, remote in-line procedures, mutual-backup of two computers and other strategies, a new type of distributed data management techniques such as high spatial locating, multi-caches strategies cluster based were developed. An application with RAC Servers and Oracle⁹ⁱ application servers (iAs) was tested, which showed the high performance of whole spatial database managing platform, such as high usability, load balance, high data downloading and high efficient spatial data query and analysis.

Key words computers cluster; massive spatial data; high performance strategy; Oracle RAC