

地球科学数据共享与数据网格技术

孙九林¹, 李 爽^{2,3}

(1. 中国科学院地理科学与资源研究所, 北京 100101; 2. 聊城大学地理系, 山东聊城 252059;
3. 河南大学环境与规划学院, 河南开封 475001)

摘要: 网格思想起源于 20 世纪 60 年代, 近十年来, 基于计算机技术及其相关学科的迅猛发展, 网格技术的研究进入了一个超速发展时期. 相应地, 地球信息科学也因计算机和遥感技术而产生了革命性的进展, 每天获得的数以 Tb、Pb 的地质数据得不到有效利用的问题日益困扰广大地质工作者. 将网格技术引入地质数据存储和共享系统将有助于解决这一难题.

关键词: 数据共享; 网格计算; 数据网格; 地球信息科学.

中图分类号: TP393.0 **文献标识码:** A

文章编号: 1000-2383(2002)05-0539-05

作者简介: 孙九林(1937—), 男, 中国工程院院士, 研究员, 主要从事信息科学技术在农业与资源环境中的应用及地球信息科学与信息化农业的研究.

网格、网格计算、数据网格是互联网上的新名词, 也是众多科技人员关注的新动向, 更是地质数据工作者关注的焦点. 网格、网格计算、数据网格理论与技术的发展给各行各业带来了新的机遇和挑战. 随着计算机软硬件技术的飞速发展, 地质观测和监控数据每日以 Tb、Pb 的速度增加, 如何充分发挥这些数据的能量, 使之能更好地为国民经济建设服务, 已成为地球科学数据建设与管理的关键. 网格技术的逐步成熟, 为地球科学数据的保值增值提供了新的机遇.

1996 年中国科学院地学部发布了《关于进一步做好我国地球科学、资源环境科学研究基础资料与数据共享的建议》, 建议数据资料共享工作由国家科委负责宏观指导, 制定政策、法规, 统一规划资料数据中心和全国信息网络的建设^[1]. 随着科学数据共享机制障碍的逐步解除, 数据共享与数据积累的技术矛盾日益突出. 基于 TCP/IP 协议的万维网并不能很好地解决人们在数据共享方面所面临的问题, 主要是由于分布式数据环境中 TCP/IP 协议的点对点(peer-to-peer)传输优点变成了缺点, 使万维网上出现了大量的信息/数据孤岛.

1 网格与网格计算

20 世纪 60 年代 J. C. R. Licklider 在人与机器共生(man-machine symbiosis)的论文中, 最先提出了网格概念. 但只是近十年才开始大规模关于网格的研究. 与网格研究相关的概念有: 元计算(metacomputing)、对等计算(peer-to-peer computing)、分布计算(distributed computing)、计算网格(computing grid)、信息网格(information grid)、知识网格(knowledge grid)、数据网格(data grid)、访问网格(access grid)、万维网服务(web service)等^[2-3].

网格实际上是继传统因特网、万维网之后的第三代因特网应用. Internet 经历了 3 个发展阶段: 第一代 Internet 是将计算机连起来, 能收发 e-mail; 第二代 Internet 是将网页连起来, 构成万维网(Web); 第三代 Internet 应该是把所有资源全面联通, 包括计算资源、存储资源、通信资源、软件资源、信息资源、知识资源等, 最终实现网络虚拟环境上的资源共享和协同工作, 消除信息孤岛和资源孤岛, 形成网格(grid)或信息网格(information grid). 下一代 Internet 可看成一台分布式计算机, 通过网格管理整个网络上的各种资源. 通过网格使用网络上的各种资源、通过网格使网络上的各种资源得到创新^[4].

网格计算(grid computing)比较通俗的解释是

收稿日期: 2002-06-26

基金项目: 科技部基础性工作“中国地球科学数据中心完善与服务”专项资金(No. 2001DEA30027).

计算供应网,即由众多的用户在大范围的网络上共享处理性能、文件以及应用软件。“grid computing”源自“power grid”(电力供应网)的专业术语。“power grid”的原意是电力供应商根据用户的需要供应电力,消费者只需支付自己使用的那部分电费。鉴于此,网格研究者们希望整个 Internet 可以像电力网一样供人们使用,尤其在电脑的计算能力(computing)方面,网格计算的目标是根据用户的需求通过网络提供必要的处理量,使用户只需支付相应的使用费。同时,电力网格模式是计算机网络发展的目标(用户所需要的数据可能是分布式异构数据库),即用户只专注于其所感兴趣的数据而不必关心其存储位置及存储方式,通过网格计算使用户透明方便地存取网上数据资源。网格计算最基础的项目是以 Globus 协议(以美国能源部及 NASA 等政府的研发机构为主推动 grid computing 的主要项目,他们开发的网络协议被称为“Globus”)为基础推动的。Globus 协议作为自由软件已经在因特网上公开^[5]。

2 网格研究进展

2.1 国际研究进展

目前,网格的研究主要在美国和欧洲。英国政府已投资 1 亿英镑,用来研制“UK national grid”(英国国家网格)。美国政府用于网格技术基础研究的经费则已达 5 亿美元。美国军方正规划实施一个宏大的网格计划,叫做“全球信息网格”(global information grid),预计在 2020 年完成。

随着网格研究在学术界的加速,信息产业界的大公司也相继公布了与网格目标一致的研究开发计划。惠普、IBM、微软、Sun 等公司最近取得共识,支持 XML、SOAP、UDDI 等万维网标准,从而更有利于开发新一代的网络应用,即万维网服务。其目的是将因特网上的资源和信息汇聚在一起,组合成企业和消费者所需要的服务。惠普推出了 eSpeak 万维网服务平台;IBM 用它的 Web Sphere 平台和一系列中间件实现万维网服务;微软的路线是通过其 .net 计划和 C# 语言实现万维网服务;Sun 则通过 open network environment (Sun ONE)计划和 Java 平台来实现它。另外,IBM 最近宣布,将投资 40 亿美元启动一个全公司的“网格计算创新计划”;Sun 则在 2000 年 9 月公布了其网格引擎软件。

作的模式,其中全球网格论坛(global grid forum)是目前主要的合作组织。目前比较有影响的研究计划有 Globus、Legion、Information Power Grid、Euro-Grid、Distributed Terascale Facility 等^[5]。

2.2 国内研究进展

我国对网格技术的研究起步较晚,相关工作开始于 1998 年。但由于网格技术是一项刚起步的研究,因此我们在网格研究的关键技术方面与国外差距不大,基本处于相同的起跑线上。目前,我国的网格技术研究主要集中于中科院计算所、国防科大、江南计算所、清华大学等几家在高性能计算方面有较强实力的研究单位。这些单位在高性能计算研究方面有很好的技术积累和很强的科研能力。其中,中科院计算所在高性能计算领域的主要成果是曙光 3000 超级服务器;其他单位的主要成果有银河巨型机、同方探索机群系统等。

中国科学院计算技术研究所对网格技术的研究已较为深入。其研究目标更多定位于“服务网络”,如果地学数据工作者与之进行很好的相互沟通,正好可以解决目前地学数据共享、存储所面临的困难。中科院计算所的网格研究工作统称为“Vega”(织女星网格),它不像某些公司那样以开发某种专业产品为目的,而是把网格做成像因特网和 Web 一样的开放性产品,使用并参与国际标准制定。我们要做的创新工作,在符合我国国情的同时,又必须融入国际标准^[9]。

从 1999 年底到 2001 年初,中科院计算所联合国内十几家科研单位,共同承担了“863”重点项目——“国家高性能计算环境”(national high performance computing environment,简称 NHPCE)的研发任务。该项目的目标是建立一个计算资源广域分布、支持异构特性的计算网格示范系统,它把我国的 8 个高性能计算中心通过 Internet 连接起来,进行统一的资源管理、信息管理和用户管理,并在此基础上开发了多个需要高性能计算能力的网格应用系统,取得了一系列研究成果。

2.3 数据网格应用研究

数据网格(data grid)是欧共体国家正在研制的下一代互联网应用,该项目以欧洲 6 个著名的研究组织为主来进行“超网”(supernet)的开发,其中包括欧洲空间管理局(ESA)、法国国家数据科学研究中心(CNRS)、荷兰国家核物理和高能物理研究所(NIKHEF)、英国著名的粒子物理与天文研究委员

国际上的网格研究主要采用开放源码、公开合

(PPARC)以及欧洲粒子物理研究所(CERN),其合作者来自捷克、芬兰、法国、德国、匈牙利、意大利、荷兰、西班牙、瑞典和英国^[7]。数据网格计划的目标是要开发出新一代速度更快、承受能力更大的数据网络。开发后的网络,将允许在欧洲范围内进行各种类型数据库的信息交流与参考,这种大范围跨组织间的合作被视为一种至关重要的举措,可以方便许多学科,如高能物理、化学、生物等领域的研究。

数据网格应用研究是一项规模宏伟的计划,由欧洲共同体投资 980 万欧元,历时 3 年时间(2001—2003 年),有超过 200 名科学家和研究人员参与研究,它的发展将使各行各业的技术人员和专家从中受益。该研究面临的首要问题是如何共享互联网上提供的大量分布式异构数据,虽然这并不是数据网格研究计划的初衷。要共享互联网数据有赖于即将出现的网格计算技术,网格计算技术可望建立具有可扩张的巨型计算环境以有效管理互联网上的文档、数据库、计算机、科学仪器和设备。

数据网格研究计划被分成 12 个子课题,有 4 个研究层面:试验床和基础结构层、应用研究层、数据网格中间件研究层和管理发布层。

3 数据网格与地学数据共享

3.1 中国地学数据共享现状

地球科学数据是研究地球形成演化、探讨人类生态环境及其变迁、减轻自然灾害、合理开发资源和促进社会可持续发展的重要科学数据,是宝贵的科学财富^[8]。在地学数据共享应用中,地学数据共享网格为“单一数据源”。自 20 世纪 70 年代以来,全球信息化的发展速度加快,科学数据积累迅速增加。据估算,人类社会最近 30 年所积累的科学数据总量已经超过了人类 5 000 年发展历史所积累的数据量总和^[9]。同时,由于科学数据在信息时代新经济中的特殊地位,使得数据资源可以直接或间接地给数据拥有者和数据使用者带来各种效益,因此,个人、团体、单位和行业从事数据积累和数据应用的行为迅速发展起来。随着数据资源的积累越来越多,社会各个层面对数据使用的呼声也越来越大。这就造成了一方面海量数据的拥有者苦于需为数据管理和存储付出巨额的支出,另一方面大量用户在到处找寻可以使用的数据,有时即使知道数据拥有者或数据的存储地址,但由于种种客观原因而无法得到或使用

数据。这其中固然有政策或科学数据本身一些特性的限制,但随着美国 20 世纪 80 年代国家科学数据共享框架方案的出台,各国在科学数据使用上的主观限制因素正日益减少,而更多面临的是物理或技术层面上的障碍,人们面对的是一座硕大无比的科学数据垃圾山。

3.2 数据网格在地学数据共享中的应用

网格是构筑在互联网上的一组新兴技术,它将高速互联网、高性能计算机、大型数据库、传感器、远程设备等融为一体,为科技人员和普通百姓提供更多的资源、功能和交互性。互联网主要为人们提供电子邮件、网页浏览等通信功能,而网格的功能则更多、更强,它能让人们透明地使用计算、存储等其他资源。

地学数据共享中数据网格的研究重点是如何消除信息孤岛和知识孤岛,实现信息资源和知识资源的智能共享。要解决的数据共享不是一般的文件交换与信息浏览,而是要把所有个人与单位连接成一个虚拟的社会组织(virtual organization),实现在动态变化环境中具有灵活控制的协作式信息资源共享。数据服务网格与 Web 最大的区别是一体化,即用户看到的不是数不清的门类繁多的网站,而是单一的入口和单一系统映像。比如一个用户需要某一方面的地学数据,他不必知道有哪些数据供应商或数据生产者,他只需通过数据网格提供的元数据库进行最简单的查询,即可找到他所需要的地学数据。同时他不需要知道数据处于何处以及数据的存储方式,只要查询到的数据符合研究要求,经过网格计算,他即可从数据网格中轻松获取所需要的数据和数据格式。

3.3 地学数据共享网格的逻辑模型

实现数据服务网格应用的关键在于网格管理软件。数据服务网格的服务包括文件消息、计算、数据内容、事务处理和知识服务等,因此数据服务网格可大致分为计算网格、数据网格与知识网格。网格管理软件在操作系统之上,可以看成是一种中间件。

依据国际网格理论技术研究的最新成果,地学数据网格系统可以分为 3 个基本层次:数据资源层、中间件层和应用层。数据网格资源层是构成网格系统的硬件基础,它包括各种计算资源,如个人电脑、工作站、超级计算机、贵重仪器、可视化设备、现有应用软件等,这些计算资源通过网络设备(路由器、集线器)连接起来。数据网格资源层仅仅实现了计算资源在物理上的连通,但从逻辑上看,这些资源仍然是

孤立的, 资源共享问题仍然没有得到解决. 因此, 必须在数据网格资源层的基础上通过数据网格中间件层来完成广域计算资源的有效共享. 数据网格中间件层是指一系列工具和协议软件, 其功能是屏蔽网格资源层中计算资源的分布、异构特性, 向数据网格应用层提供透明、一致的使用接口. 数据网格中间件层也称为数据网格操作系统 (grid operating system), 它同时需要提供用户编程接口和相应环境, 以支持网格应用的开发. 数据网格应用层是用户需求的具体体现. 在数据网格操作系统的支持下, 网格用户可以使用其提供的工具或环境开发各种应用系统. 能否在数据网格系统上开发应用系统以解决各种大型计算问题, 是衡量数据网格系统优劣的关键.

3.4 中国地学数据资源共享网格系统建设

在国外, 最著名的数据网格研究是美国的 Globus 项目. 该项目的主要研究目标有 2 个: (1) 网格技术的研究; (2) 相应软件的开发和标准的制定. 同时, Globus 项目还涉及到网格应用的开发及试验床的建立. 依据 Globus 项目提出的数据网格体系结构模型, 在此初步提出中国地学数据资源共享网格建设框架(图 1).

地学数据资源共享网格体系结构主要分为以下几个部分: (1) 数据网格结构 (grid fabric) 层. 它是一个本地控制的接口, 提供与资源相关的基本功能, 便于高层分布式网格服务的实现. 它提供共享获取的各种资源的入口, 它们是物理或逻辑实体, 包括计算资源、存储系统、目录、物理网络资源等; 这里的资源可以是一个逻辑实体, 例如一个分布式的文件系统、分布式集群计算机. 实现资源的共享, 需要使用 Internet 网络协议. 基于结构层, 高层协议可以如同操作本地资源一样操作其他主机的共享资源. 实现这层协议至少需要实现一个允许外界发现和查询资源结构和状态的机制和一个能控制服务质量的资源管

理机制; (2) 数据网格服务 (grid service) 层. 实现与数据资源无关和应用无关的功能, 网格服务的实现涉及到地域和机构的分布, 为高层协议提供了简单而且安全的通讯方式. 安全协议主要是为了解决安全问题的复杂性, 并提供一个有效的解决方案. 这个问题来源于不同系统安全策略的不同和用户数据对安全需求的不同. 从基础的用户登陆与权限认证到用户程序的权限赋予, 从本地资源的安全策略到异地主机基于账号的信任机制, 要解决这些不同的安全策略, 并提供可靠的、统一的接口, 是网格技术要解决的一个重要问题; (3) 数据网格应用工具 (grid application toolkit) 层. 提供更为专业化的服务和组件用于不同类型的网格数据应用; (4) 应用 (application) 层. 位于数据网格体系的最顶层, 是由用户开发的应用系统组成, 它为应用程序提供统一的接口和服务. 数据网格用户可以使用其他层次的接口和服务完成网格应用的开发. 应用程序集成应用层定义的语言框架, 通过应用层的协议, 对底层的资源进行访问, 而不再需要关心访问的复杂繁琐的实现机制.

4 结语

将网格、网格技术、网格计算、数据网格思想引入到地学数据存储与管理中, 是目前最先进的互联网技术与思想同地球信息科学的紧密结合. 采用数据网格技术, 可以解决当前困扰地学数据工作者的一大难题, 即如何有效存储与发布所采集的地学数据, 同时又可以解决广大研究人员对数据的迫切需求. 进一步研究中, 急需解决的技术关键点在于:

(1) 如何保障安全性. 这是因为, 如果人们在诸如因特网这样的全球网络上实现共享计算处理资源, 那么, 罪犯就会有机可乘, 攻击系统漏洞. 加上这是一项十分庞大的计划, 因此, 如何确保安全性也就成为最困难的一个课题, 也是数据拥有者最担心的问题之一; (2) 由于数据资源的需求与供给都在动态变化, 而且分布在全球的各个角落, 完成用户要求的一项服务可能要调用北京的超级服务器、上海的数据库或安装在西安的某台计算机上的软件, 因此对服务器的响应时间、网络的带宽 (主要是带宽)、特别是网格管理软件的复杂性与灵活性及网络上各种设备的互操作性都有很高的要求; (3) 如何解决网格使用模式问题, 在现有的操作系统上, 计算机用户可以



图 1 地学资源共享网格体系结构模型

Fig. 1 Structure model of geo-resources sharing grid

使用各种软件工具来完成各种任务.而在网格环境下,用户可能需要通过新的方式来利用网格系统资源.因此,在网格操作系统上设计开发各种工具、应用软件是网格使用模式研究需要解决的关键问题.

参考文献:

- [1] 中科院地学部. 关于进一步做好我国地球科学、资源与环境科学研究基础资料与数据共享的建议[J]. 地理科学进展, 1996, 11(1): 122-123.
- The Geo-Science Department of Chinese Academy of Sciences. The basic datum and data sharing advice of improving Chinese geo-science, resource and environment science research [J]. Progress in Geography, 1996, 11(1): 122-123.
- [2] Catlett C. Global grid forum documents and recommendations: process and requirements [EB/OL]. Copyright by Global Grid Forum. <http://www.gridforum.org/Documents/GFD/GFD-C.1.doc>, 2002.
- [3] Foster I, Kesselman C, Nick J, et al. The physiology of the grid: an open grid service architecture for distributed systems integration [EB/OL]. <http://www.eu-datagrid.org>, 2002-01.
- [4] 刘世昕. 下一代因特网让人们使用资源像用电一样简单 [N]. 中国青年报, 2002-04-12.

- LIU S X. The next Internet makes people using resource as simple as using electricity [N]. China Youth Daily, 2002-04-12.
- [5] Foster I, Kesselman C. Globus: a metacomputing infrastructure toolkit [J]. Intl J Supercomputer Applications, 1997, 11(2): 115-128.
- [6] 徐扬, 王宏, 宋宇, 等. 网格计算网罗一切 [N]. 中国计算机报, 2002-04-15.
- XU Y, WANG H, SONG Y, et al. Grid computing collects everything [N]. Chinese Computer, 2002-04-15.
- [7] Mauro D, Gianfranco M, Roberto P. Project-Presentation [EB/OL]. <http://web.datagrid.cnr.it/introdocs/DataGrid-11-NOT-0103-1-1-Project-Presentation.pdf>, 2001-07-13.
- [8] 李军. 地球科学数据研究的初步探讨 [J]. 地理学报, 1996, (增刊1): 封底.
- LI J. The principium discuss of geo-science data research [J]. Acta Geographica Sinica, 1996, (Suppl 1): back cover.
- [9] 刘闯, 王正兴. 科学数据共享调研组的系列报告的一部分 [R]. 北京: 中国科技部, 2001.
- LIU C, WANG Z X. Part of investigation and research group series report: science data sharing [R]. Beijing: the Ministry of Science and Technology, 2001.

Geo-Data Sharing and Data-Grid

SUN Jiu-lin¹, LI Shuang^{2,3}

(1. Institute of Geographic Science and Natural Resources Research, Beijing 100101, China; 2. Department of Geography, Liaocheng University, Liaocheng 252059, China; 3. College of Environment and Planning, Henan University, Kaifeng 475001, China)

Abstract: The conception of grid started in 1960's. Due to the rapid development of computer technology-based academic disciplines for the past ten years, the research into the grid technology has entered into an excessively fast-speed period. In addition, the global information science has already experienced many revolutionary changes with the advancement of the computer and remote-sensing technology. Today, the insufficient use of the thousands of Tb or Pb geo-data acquired daily has troubled many earth scientists. However, this trouble can be easily shot by the introduction of grid technology to the geo-data storing and sharing system.

Key words: data sharing; grid computing; data grid; earth information science.