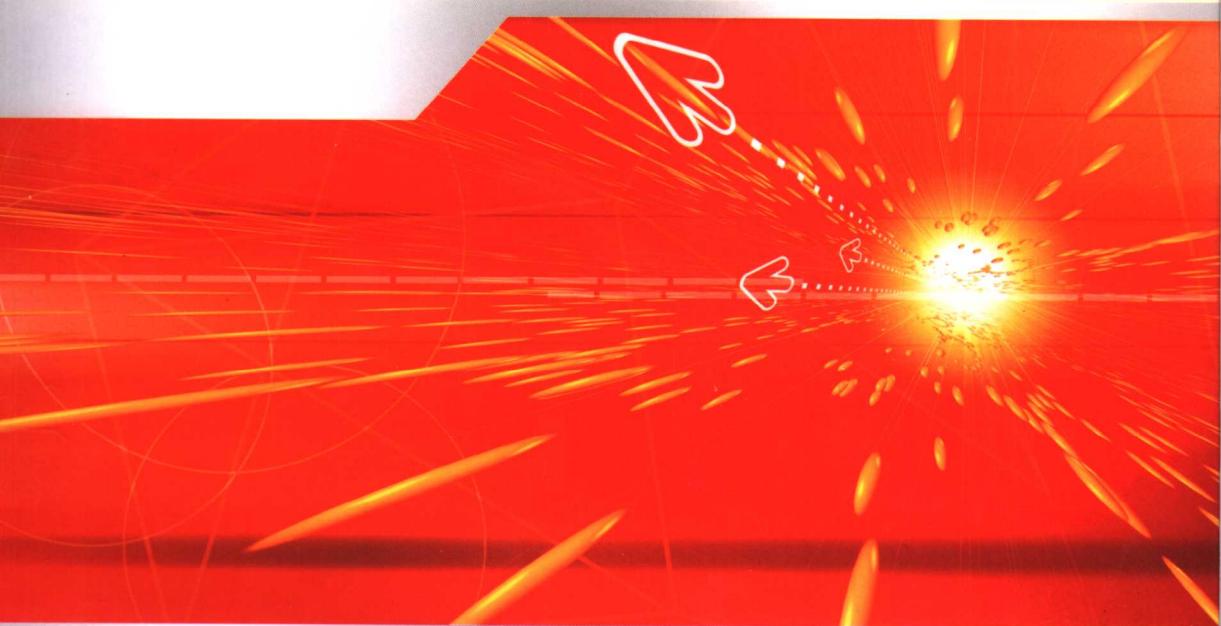


李晓波等 著

科学数据共享 关键技术

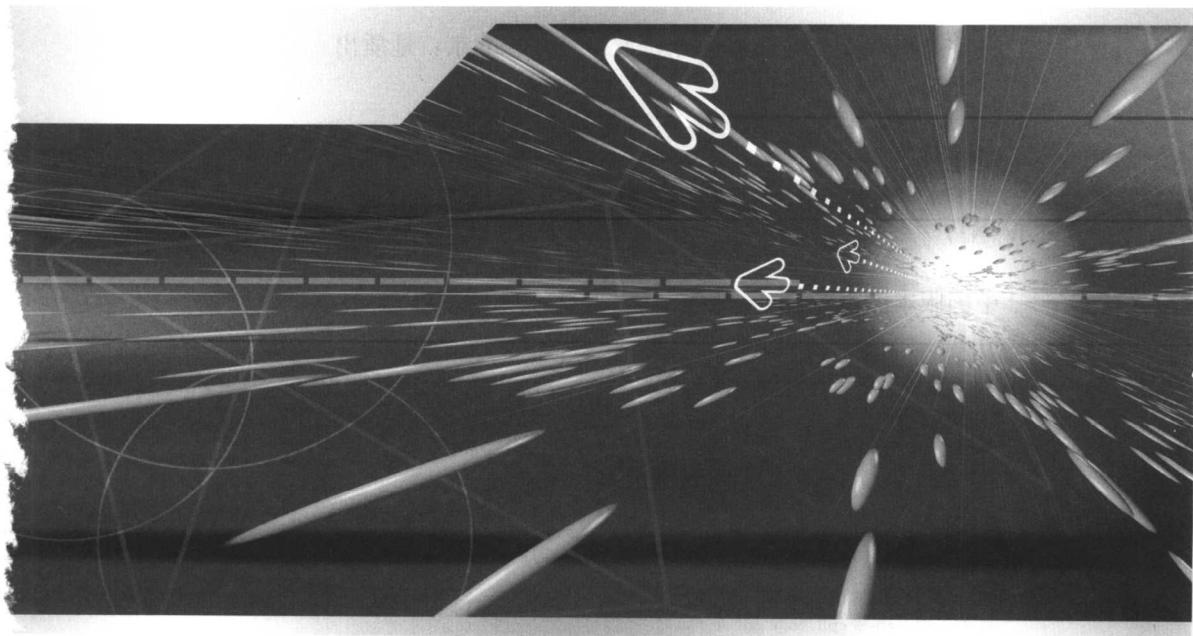


KEXUE SHUJU GONGXIANG
GUANJIAN JISHU

地 质 出 版 社

李晓波等 著

科学数据共享 关键技术



地 质 出 版 社

· 北 京 ·

内 容 提 要

本书系统总结了科学数据共享关键技术研究成果，以科学数据共享理论为指导，分析了科学数据共享工程技术需求和相关技术国内外应用发展情况，提出了科学数据共享若干关键技术和技术平台框架，论述了科学数据汇交、海量数据管理的技术模型与实现方法，阐明了科学数据网上定位查询、浏览分发等共享服务的功能和技术路线，明确了系统安全保障策略与技术方案，探讨了网格技术在科学数据共享中的应用前景。

本书可供科学数据共享试点单位及信息技术等相关专业科研人员、教师、学生等参阅。

图书在版编目（CIP）数据

科学数据共享关键技术 / 李晓波等著. —北京：地质出版社，2007. 11

ISBN 978 - 7 - 116 - 05496 - 7

I. 科… II. 李… III. 自然科学－数据管理－研究
IV. N37

中国版本图书馆 CIP 数据核字（2007）第 152343 号

责任编辑：柳 青

责任校对：黄苏晔

出版发行：地质出版社

社址邮编：北京海淀区学院路 31 号，100083

咨询电话：(010) 82324508 (邮购部)；(010) 82324573 (编辑室)

网 址：<http://www.gph.com.cn>

电子邮箱：zbs@gph.com.cn

传 真：(010) 82310759

印 刷：北京北林印刷厂

开 本：787mm×960mm 1/16

印 张：17.25

字 数：431 千字

印 数：1—1600 册

版 次：2007 年 11 月北京第 1 版·第 1 次印刷

定 价：45.00 元

书 号：ISBN 978 - 7 - 116 - 05496 - 7

（如对本书有建议或意见，敬请致电本社；如本书有印装问题，本社负责调换）

前 言

信息与数据资源是知识经济时代重要的生产要素、无形资产和社会财富。新中国成立以来，我国各部門在生产、科研活动中，积累了海量科学数据，它们是国家科技创新和经济社会发展的重要战略资源。为加强我国科学数据資源的开发利用，充分发挥其基础性、公益性和战略性作用，科技部决定组织实施科学数据共享工程。

科学数据共享工程是在国家统筹规划下，应用现代信息技术，整合集成各部門、各单位的公益性、基础性科学数据資源，实现其公开与共用，并充分利用国际科学数据資源；通过制定和完善共享政策、法规和管理体系，把各部門、各单位乃至个人由政府投入所获取与积累的科学数据資源，纳入国家科学数据共享管理的统一框架；通过科学数据中心和科学数据网的建设及共享技术的研究开发与应用，形成跨部門、跨地区、跨学科、多层次、分布式的国家科学数据管理与共享服务体系，大力提高科学数据的管理与共享服务水平，为国家科技整体水平的提高和经济社会发展，提供可靠的科学数据資源的支持。工程的实施，对推进国家创新体系建设、促进经济结构调整和经济增长方式的转变具有十分重要的意义。

实施科学数据共享工程的关键要素是科学数据資源、共享政策法规和共享技术支撑。其中，共享技术是实现科学数据广泛共享的最有效途径。为深入研究和探索科学数据共享的技术途径，推动各领域科学数据中心（网）的建设，科技部立项开展了科学数据共享关键技术研究，针对科学数据共享若干关键技术问题，研究提出一些解决方案。本书系统总结了科学数据共享关键技术研究成果，以科学数据共享理论为指导，分析了科学数据共享工程技术需求和相关技术国内外应用发展情况，提出了科学数据共享若干关键技术和技术平台框架，论述了科学数据汇交、海量数据管理的技术模型与实现方法，阐明了科学数据网上定位查询、浏览分发等共享服务的功能和技术路线，明确了系统安全保障策略与技术方案，探讨了网格技术在科学数据共享中的应用前景。目前，根据科学数据共享关键技术研究成果开发的部分原型系统，已应用于科学数据共享工程的试点项目建设，并将逐步加以推广。

构建稳定、先进的科学数据共享技术支撑体系，是一项迫切的任务和长期的

工作。本书提出的有关科学数据共享的若干关键技术问题和解决方案尚为粗浅，需要在实践和技术的发展中逐步加以修正和完善。

本书第一章、第二章由李晓波、汪志强、张垚垚编写，第三章由宦茂盛、丁龙玺编写，第四章由徐枫、宦茂盛、石雯雯编写，第五章由徐宝龙、丁龙玺、李上珠编写，第六章由徐茂智、赵彦慧、沈浔浔编写，第七章由刘海涛、管海兵编写；全书由李晓波审阅统稿。本书编写过程中，得到了孙九林、滕绵震、黄鼎成、傅小锋、李俊、李集明、戴爱德、孔瑞等专家的指导，在此致以诚挚谢意！

2006 年 12 月 20 日

目 次

前 言

1 数据共享相关技术应用发展概况	(1)
1.1 相关信息技术发展概况	(1)
1.1.1 信息技术发展概况	(1)
1.1.2 数据共享服务相关技术	(2)
1.2 国内外应用发展情况	(10)
1.2.1 国外应用发展情况	(10)
1.2.2 国内应用发展情况	(16)
2 科学数据共享技术需求分析	(21)
2.1 科学数据资源分布与特点	(21)
2.1.1 科学数据资源分布	(21)
2.1.2 科学数据资源特点	(23)
2.2 科学数据共享工程建设目标和框架	(25)
2.3 功能需求	(26)
2.4 技术平台框架	(28)
2.5 关键技术需求	(30)
2.5.1 科学数据汇交	(30)
2.5.2 科学数据资源综合管理	(32)
2.5.3 网上科学数据共享服务	(36)
2.5.4 科学数据共享安全	(38)
3 科学数据汇交技术	(41)
3.1 概述	(41)
3.2 技术架构	(41)
3.2.1 应用层	(42)
3.2.2 支撑平台层	(42)
3.2.3 数据层	(42)
3.3 成果登记	(42)
3.3.1 成果登记业务流程	(42)
3.3.2 成果登记主要技术	(45)

3.3.3 科技成果登记系统示例.....	(50)
3.4 网上数据上传.....	(61)
3.4.1 数据汇交流程.....	(61)
3.4.2 汇交数据的形式.....	(61)
3.4.3 数据汇交关键技术.....	(63)
3.4.4 网上数据上传系统基本功能设计.....	(65)
3.5 元数据汇交.....	(69)
3.5.1 元数据汇交流程.....	(69)
3.5.2 汇交元数据的内容和形式.....	(70)
3.5.3 核心元数据汇交技术.....	(72)
3.5.4 科学数据共享核心元数据汇交系统实现.....	(73)
4 科学数据综合管理技术.....	(76)
4.1 概述.....	(76)
4.2 技术架构.....	(77)
4.2.1 系统组成.....	(78)
4.2.2 系统逻辑架构.....	(79)
4.3 集中式数据管理.....	(81)
4.3.1 集中式数据管理概述.....	(81)
4.3.2 数据抽取、转换和加载.....	(83)
4.3.3 集中式数据集目录管理.....	(99)
4.3.4 集中式数据字典管理	(103)
4.4 分布式数据管理	(109)
4.4.1 分布式数据管理概述	(109)
4.4.2 分布式数据目录管理	(110)
4.4.3 分布式数据字典管理	(112)
4.5 数据存储策略	(115)
4.5.1 存储策略	(116)
4.5.2 存储介质及方式	(118)
5 科学数据共享服务技术	(121)
5.1 体系结构	(121)
5.1.1 目录服务	(122)
5.1.2 数据服务	(122)
5.1.3 延伸服务	(123)
5.2 目录服务	(124)

5.2.1 目录服务体系	(124)
5.2.2 目录服务功能模型	(126)
5.2.3 目录服务关键技术	(128)
5.3 数据服务	(135)
5.3.1 数据查询	(136)
5.3.2 数据浏览	(139)
5.3.3 数据下载	(142)
5.3.4 数据服务关键技术	(145)
5.4 延伸服务	(151)
5.4.1 数据挖掘	(151)
5.4.2 数据可视化	(156)
5.4.3 数据专题服务	(159)
5.4.4 延伸服务关键技术	(161)
6 科学数据共享安全技术	(164)
6.1 安全体系结构	(164)
6.1.1 共享服务安全结构	(166)
6.1.2 数据安全	(170)
6.2 身份认证	(174)
6.2.1 身份认证模型综述	(174)
6.2.2 单点认证模型综述	(178)
6.2.3 相关理论和 SSO 实现	(180)
6.2.4 数据共享的单点登录模型	(185)
6.3 访问控制	(190)
6.3.1 自主访问控制	(191)
6.3.2 强制访问控制	(192)
6.3.3 基于角色的访问控制	(193)
6.3.4 其他访问控制模型及比较	(198)
6.3.5 共享系统的访问控制	(200)
6.4 安全审计及其他措施	(201)
6.4.1 安全审计	(201)
6.4.2 计算机取证	(206)
6.4.3 其他安全措施	(210)
6.5 数据安全	(214)
6.5.1 数据共享安全模型	(214)

6.5.2 数据传输安全	(218)
6.5.3 数据存储安全	(223)
6.5.4 版权保护	(225)
7 科学数据共享网格初探	(229)
7.1 背景知识	(230)
7.1.1 网格的历史与定义	(231)
7.1.2 网格架构与标准	(233)
7.1.3 网格的应用领域	(234)
7.1.4 网格标准与技术	(235)
7.1.5 网格支撑工具 Globus	(237)
7.2 科学数据共享网格	(239)
7.2.1 体系结构	(239)
7.2.2 软件框架	(240)
7.2.3 系统架构	(241)
7.3 元数据管理	(243)
7.4 信息服务	(245)
7.4.1 信息服务系统基本功能	(246)
7.4.2 监控和发现服务	(248)
7.5 资源管理	(251)
7.5.1 资源代理	(251)
7.5.2 存储资源代理	(252)
7.6 数据管理	(255)
7.6.1 全局二级存储服务 GASS	(255)
7.6.2 数据传输	(256)
7.6.3 数据访问	(258)
7.6.4 数据复制	(259)
7.7 安全模型	(260)
7.7.1 网格安全性基础结构 GSI	(261)
7.7.2 OGSA 安全体系结构	(262)
7.7.3 Web 服务安全性模型	(264)
7.7.4 访问控制机制	(265)
参考文献	(267)

数据共享相关技术应用发展概况

数据资源的高速增长和知识经济的发展，对分布式海量数据资源管理、整合与共享产生了巨大需求。当前数据共享在技术上主要采用两种方法：一种将存储在硬盘上的数据，用硬盘或光盘拷贝进行交换；另一种则是依赖于计算机网络进行共享。随着计算机、网络、数据库等现代信息技术，特别是基于网络的信息管理与服务技术的发展，使人们在广域范围内随时随地获取自己感兴趣的信息成为可能。充分发挥计算机与网络的功能，在共同遵守数据共享的原则下建立全球性和区域性的数据网络，深入开发和应用数据共享相关技术，形成强大的数据共享平台，实现数据、信息乃至知识的广泛交流与共享，是国内外数据资源管理与共享应用的发展趋势。

1.1 相关信息技术发展概况

1.1.1 信息技术发展概况

信息技术是当前发展最快的技术领域，是推进国家经济发展的核心技术之一。计算机技术、网络技术、数据库技术、通信技术等现代信息技术不断创新、融合，为数据技术的发展和数据共享应用提供了坚实的基础和广阔的发展空间。

在计算机技术方面，CPU 芯片技术不断提升，运算速度已超过 3 GHz。65 纳米制作工艺、64 位计算、双核技术等的最新发展，推动 CPU 性能持续扩展和提高。处于计算领域最高端的高性能计算机仍是 21 世纪 IT 领域争夺的制高点。高性能计算机运算速度不断突破，标量型产品成为主流，高性能计算机向网络化、体系结构主流化、开放和标准化、应用的多样化等方面发展。高性能计算的网络化发展使网格计算成为新的研究热点，并发展成一门重要的新兴技术。

网络技术的发展深刻改变着人类生产和生活方式。光纤通讯技术、无线网络技术、超宽带技术竞相发展，互联网应用迅速增长，并进入新的发展阶段，下一代网络（Next Generation Network，NGN）、语义网（the Semantic Web）等新技术已经浮出水面。所谓下一代网络泛指以 IP 为中心，同时可以支持语音、数据和多媒体业务的互联网、移动通信网络、固定电话通信网络的融合网络。NGN 在

网络容量和资源方面将具有高带宽和大容量、足够的地址资源等重要特征。IPv6技术已成为下一代网络的地址问题解决方案，它将IP地址的长度由32位扩展到128位，能够满足互联网飞速发展的需求。所谓“语义网”，是按照能够表达网页内容的“词语”链接起来的全球信息网，换言之，是用机器很容易理解和处理的方式链接起来的全球数据库。语义网是对万维网本质的变革，它能够根据语义进行判断，是一种能理解人类语言的智能网络，它使人们利用互联网巨量信息资源变得更为轻松和有效，有力地提升全球范围的信息与知识共享水平，为科技创新提供无尽的资源。

数据库技术是计算机科学技术中发展最快、应用最广泛的重要分支之一。随着网络应用日益深入和数据对象的不断拓展，数据库技术与其他计算机技术相结合，迅速向分布式、并行、智能和多维结构方向发展。与分布处理技术相结合的分布式数据库，是分布在计算机网络上多个逻辑相关的数据库集合，具有数据的物理分布性、数据的逻辑整体性、数据的分布透明性、场地自治和协调、数据的冗余及冗余透明性等特点。与并行处理技术相结合的并行数据库，利用多处理器平台的能力，通过多种并行性，提供优化的响应时间与事务吞吐量，如数据库集群技术显著提高了数据处理速度、数据可用性和数据集可扩展性。数据库技术与人工智能技术相结合，使数据库在反映能力上具有主动性、快速性和智能化的特点。基于多维数据模型的后关系型数据库（Post-Relational Database，PRDB），将多维处理技术和面向对象技术集成在一起，提供了事务处理应用开发所需的高性能和灵活性，同时支持应用和数据的复杂性，并拥有比关系型技术更强的扩展性、更快的编程能力以及更便捷的使用特性。此外，多媒体数据库、空间数据库、数据仓库等数据库技术，提高了对复杂数据对象的处理、管理和集成的能力。

计算机、网络和数据库等现代信息技术的融合发展，为海量分布式数据资源管理与共享相关技术的发展提供了有力支撑。

1.1.2 数据共享服务相关技术

国际科学数据委员会第18届世界大会提出了科学数据领域若干前沿问题，包括数据综合与数据互操作问题、数据处理技术与数据显现工具问题、数据保护与档案管理问题等，这些问题也是科学数据共享面临的主要问题。在数据资源和数据需求不断增长的强劲推动下，元数据、Web服务、数据ETL、数据可视化、空间数据共享、网格等一系列相关技术发展迅速，令人注目。

1.1.2.1 元数据技术

元数据（Metadata）是实现数据共享，特别是实现网络数据共享的重要技

术。元数据是关于数据的数据，描述数据对象的各种属性及相关关系等内容，主要用于计算机系统之间共享数据时对数据进行说明，实现不同系统对同一数据的处理的一致性。元数据由元数据元素、结构和描述句法组成，类型包括描述性元数据、管理性元数据、结构性元数据或技术性元数据等。它在许多领域有具体的规定和应用。

通常，元数据是对信息资源或数据结构化的描述，它描述信息资源或数据本身的特征和属性，规定数字化信息的组织，具有定位、发现、证明、评估、选择等功能。根据数据拥有者或管理者发布的元数据，人们能够对数据的适用性进行详细、深入的了解，包括数据的内容、用途、格式、质量、处理方法和获取方法等各方面细节，从而确定它们是否适合自己的使用需要，并可按照元数据中提供的获取方法最终取得数据。

元数据应用领域相当广泛，在数据资源管理和共享应用中正发挥着越来越大的作用。在数字图书馆应用中，通过如 DC (Dublin Core) 元数据等元数据系统为分布的、种类繁多的数字化信息提供了整合工具与纽带，实现对海量图书资源的有效检索服务。在数据仓库应用中，元数据是描述数据及其环境的数据，包括了描述性元素、管理性元素和技术性元素，既有对业务数据的描述信息以帮助用户理解和使用数据，也有关于数据项存储方法等支持系统对数据管理和处理的信息。数据仓库系统中的元数据机制主要支持以下系统管理功能：描述哪些数据在数据仓库中；定义要进入数据仓库中的数据和从数据仓库中产生的数据；记录根据业务事件发生而随之进行的数据抽取工作时间安排；记录并检测系统数据一致性的要求和执行情况；衡量数据质量。在地理空间信息共享应用中，空间元数据用于描述地理数据集的内容、质量、表示方式、空间参考、管理方式以及数据集的其他特征，帮助人们理解和共享有关地理空间信息。以元数据为核心的目录查询，是利用元数据技术提供信息服务的一种标准模式。建立基于元数据库的分布式数据资源目录检索系统，通过元数据标准的核心元素将信息以动态分类的形式展现给用户，用户通过浏览门户网站提供的元数据摘要（核心元数据）可以快速确定自己所需的信息范围，然后要求门户网站在该范围内进一步搜索，从而实现对信息的发现、定位、访问与获取。

建立包括多种元素内容、结构完整的元数据，对数据组织、信息集成和知识共享有十分重要的意义。

元数据格式类型复杂，用途殊异，在用不同元数据格式描述的资源体系之间进行检索、资源描述和资源利用，就存在元数据互操作问题。元数据的互操作性要求在由不同的组织制定与管理且技术规范不尽相同的元数据环境下，能够做到对用户保持一致性的服务，也就是说对一个应用或用户来说，能够保证一个统一

的数据界面，保证一致性与对用户的透明。元数据互操作通常有两种方式：一是通过在不同系统间创建元数据映射，进行元数据格式转换，以实现跨系统的检索；二是采用核心元数据集，提供跨学科和格式的语义互操作性，用于各种应用软件、数据格式或者主题领域的字段和描述，它是适合于任何 Web 资源、标准的元数据，与现存其他元数据兼容。通过建立一个标准的资源描述框架 RDF，用以描述所有元数据格式，那么只要一个系统能够解析这个标准描述框架，就能解读相应的元数据格式。RDF 是一个能对结构化元数据进行编码，交换及再利用的体系框架，采用 XML 作为交换和处理元数据的通用语法结构体系，通过对通常意义上的语义、语法和结构的支持，提供各种不同元数据体系之间的互操作性。

包括元数据标准及有关元数据采集、编辑、管理、发布、互操作等在内的元数据技术，是当前对广域分布的、异构的信息资源进行整合与共享服务的有效手段。

1.1.2.2 Web 服务技术

Web 服务（Web Service）技术是互联网时代出现的最引人注目的关键技术之一，它对企业、用户甚至共同文化将产生巨大影响。

Web 服务是一种新的 Web 应用程序类型，是一种自包含、自解释、模块化的应用，能够被发布、定位，并且可从 Web 上的任何位置进行调用。一个 Web 服务就是一个可以被 URI 识别的软件应用，它的接口和绑定可以被 XML 描述与发现，并且可以通过基于 Internet 的协议直接支持与其他基于 XML 消息的软件应用的交互。

Web 服务的主要目标是，通过使用统一标准，能够统一封装数据、消息、行为等，在无需考虑具体应用环境下让不同系统跨越平台，彼此兼容，进行无缝通信和数据共享。Web 服务实现的功能可以是响应一个简单的请求，也可以是完成一个复杂的商务流程。理论上讲，一旦对 Web 服务进行了部署，其他 Web 服务应用程序就可以发现并调用它部署的服务。

简单来说，Web 服务就是可远程调用的应用程序组件，其本质目的是提供一个与操作系统无关、与程序设计语言无关、与机器类型无关、与运行环境无关的平台，实现 Internet 中应用程序的共享。

以往的分布式应用程序逻辑可通过使用分布式组件对象模型（DCOM）、公用对象请求代理程序体系结构（CORBA）和远程方法调用（RMI）等中间件平台，将服务置于远程系统中。但这些系统要求服务客户端与系统提供的服务本身之间必须进行紧密耦合，即要求一个同类基本结构，故无法扩展到互联网上。这些系统往往十分脆弱，如果一端的执行机制发生变化，如服务器端应用程序的接

口发生更改，那么另一端客户端便会崩溃。而 Web 服务彼此是松散耦合的，连接中的任何一方均可更改执行机制，却不影响应用程序的正常运行。

Web 服务是基于 HTTP、XML、SOAP（简单对象访问协议）、WSDL（Web 服务描述语言）和 UDDI（统一描述、发现和集成协议）等一系列标准协议的。基于这些标准协议，Web 服务实现了服务的跨平台、跨语言的共享。Web 服务是软件重用和应用集成的一种非常有力的形式。

Web 服务的平台架构主要由服务请求者、服务注册中心和服务提供者三部分组成。其实现过程包括服务发布与注册、服务查询与发现、服务绑定与调用。服务请求者指查询、调用服务的客户端程序；服务提供者即服务的所有者和部署服务的平台；服务注册中心是指用来存储服务信息的信息库，服务提供者在这里发布、注册服务，而服务请求者在这里查询、绑定服务，最终实现调用服务提供者的服务。

Web 服务有两大核心优势，即分布性和互操作性，这是数据及其服务共享的关键问题。在 Web 服务架构下，服务提供者和服务请求者都可以是分布式的，一个服务请求者可以远程调用多个服务提供者的服务，服务提供者也可以同时为多个服务请求者提供服务，这为服务共享提供了一个最佳的方式。Web 服务的信息表达基于标准通用的 XML 语言，在 XML 语言的基础上，使用 WSDL 和 UDDI 实现服务注册与发现，使用 SOAP 进行服务调用，这使不同的服务之间实现了互操作性。

随着有关技术开发的不断深入，Web 服务的潜能将逐步充分地发挥出来，将来具有多种前沿的功能性，可跨越时区和语言的障碍。例如可以想象，位于欧洲或环太平洋地区的商家能够使用母语发起一次电话会议，呼叫他在中国的说汉语的生意合作伙伴，因为通过使用无线电通讯和广域网链接，人们可以利用 Web 服务体系执行必要的即时翻译。

1.1.2.3 数据 ETL 技术

数据 ETL（Extraction、Transformation、Loading），是数据抽取、数据转换和数据加载的简称，是一类用于从多个不同类型结构的数据源中抽取数据，进行清理转换并加载到目标数据库中，以实现多源异构数据集成的工具，是构成数据仓库、数据挖掘及商业智能等技术的基础。

数据抽取是指从不同的网络、应用、操作系统平台、数据库及数据格式中抽取数据的过程。数据抽取要在充分理解和认真规划数据源基础上，制定抽取规则，设计数据接口，实现数据高效读取、传输和转移。

数据转换是指对抽取数据的合并拆分、汇总计算、类型转换、清洗过滤、重

新格式化，包括关键数据的重新构建等等，实现数据从源数据结构到目标数据结构的转换。在数据转换过程中，难点是对存在有二义性、重复、不完整、违反业务规则等问题的所谓“脏”数据进行检测和清洗，以保证目标数据库及后续应用的数据质量。数据转换操作视具体情况可在数据抽取时进行，也可在数据加载时进行，或在数据抽取和数据加载时均有数据转换操作。

数据加载是将经过转换和清洗的数据加载到目标数据库中。

数据 ETL 是构建数据仓库的第一步，它在数据仓库体系结构中占有非常重要的地位。现在支持数据仓库应用的数据 ETL 产品有很多，如 DataStage、Informatica、Data Junction、OWB、Data Builder 等。数据 ETL 工具软件一般由元数据（规则）的创建、数据源到目标系统的映射、数据源接口界面、数据抽取、数据转换、数据分发和加载等功能模块组成。针对空间数据还有 Spatial ETL 工具，如 FME。

较之数据仓库应用，数据 ETL 在更广泛的数据集成应用中，面临更复杂的数据源，且数据抽取和数据加载两个操作往往分属不同部门，需要远程传输，彼此完全独立，高度自治。面向数据集成的 ETL 应用可能需要着重解决三个问题：一是将抽取的源数据转换为公共数据格式，然后再转化为目标数据格式；二是数据抽取的增量复制、动态更新；三是有效的过程管理与调度。

现有数据资源及应用具有显著的异构性和分布性特征，整合在系统、语法、结构和语义各层次上的异构数据是一个很大挑战，且由于滥用缩写词、惯用语、数据输入错误、重复记录、丢失值等原因造成数据质量问题，产生很多不“洁净”数据，整合利用困难，这是数据 ETL 产生和应用的重要背景。数据 ETL 按照统一的规则集成并提高数据的价值，把数据转换为信息、知识，它屏蔽了数据源复杂的业务逻辑，从而为各种基于数据仓库或集成数据的分析与应用提供了统一的数据接口，这为在多源数据中进行信息挖掘和开展有针对性的信息服务，提供了很好的技术手段。

在现实需求的推动下，数据 ETL 正在向通用化、高效化和智能化方向发展，它将具有更良好的可扩展性和对异构性及多样性的支持能力，数据源管理、ETL 规则定制、数据质量保证等将更加自动化和智能化。

数据 ETL 在海量分布式数据整合集成、知识发现和总体数据质量管理等应用中有广阔发展空间。

1.1.2.4 网格技术

网格技术是动态多机构虚拟组织中的资源共享和协同问题解决技术。网格(Grid)就是作为这样一种技术被提出的，它作为“黏合”中间件，来实现系统

的用户管理、资源信息管理、作业管理和安全认证管理等功能，保障计算系统的可靠运行。

关于网格的准确定义目前还没有，网格研究权威、美国 Globus 项目领袖 Ian Foster 认为，网格必须同时满足三个条件：① 协调非集中的控制资源——网格整合在不同控制域中的各种资源，协调在不同控制域中的各种使用者；② 使用标准、开放和通用的协议和接口；③ 取得非凡的服务质量——网格允许它的资源被协调使用，以得到多种服务质量，满足不同使用者需求，使得联合系统的功效比其各部分的功效总和要大得多。

网格可以分为计算网格、信息网格和知识网格等不同层次。

计算网格由通过高速网络连接在一起的超级计算机、大型服务器和专用设备组成，它的主要功能包括负载均衡、数据源地址解析、安全、复制和消息转发。计算网格是网格的系统层，它为应用层（信息网格和知识网格等）提供系统基础设施。

信息网格建立在计算网格的基础上，利用元数据和中间件来实现异构信息源的存取和数据共享。知识网格是使用知识发现技术从数据库和其他非结构化数据中提炼和生成知识。信息网格是基于互联网的为公众提供各种一体化信息服务的信息基础设施，实现在动态变化环境中有灵活控制的协作式信息资源共享。在 Internet 和 Web 上，数据和信息资源零散地分布在各个网络站点，而在信息网格中，资源被统一管理和使用，用户可以通过网格门户透明地使用整个网格资源，他们看到的是一个逻辑门户上的若干与自己相关的频道，而不同于在数不清的门类繁多的网站中搜索自己想要的信息。空间信息网格是信息网格的一种应用，是汇集和共享空间信息资源，具有按需服务能力的一种空间信息基础设施，它提供一体化的空间信息获取、处理与应用服务的基本技术框架，以及智能化的空间信息处理平台和应用环境，实现空间信息资源和知识资源的智能共享。

网格相关技术主要包括 Web Service、OGSA（开放网格服务结构）、WSRF（Web 服务资源框架）等。Web Service 是对分布在互联网上程序的定位与相互访问技术，提供永久的无状态服务，实现了程序间远程访问透明性。OGSA 是 Globus 项目设计出的一种网格体系结构，提供临时的有状态网格服务，实现了网格功能的多样化。Web Service 解决了发现和激发永久服务的问题，但是在网络中，大量的是临时服务，因此 OGSA 对 Web Service 进行了扩展，提出了网格服务（是广义的服务，包括各种计算资源、存储资源、网络、程序、数据库等等）的概念，使得它可以支持临时服务实例，并且能够动态创建和删除。WSRF（Web Services Resource Framework）提出了提供持久数据的方式，定义了使用 Web 服务来访问有状态资源的一系列规范。WSRF 在 Web Service 的永久无状态

服务基础上加入有状态的（本质上是临时的）资源，从 WSRF 观点来看的资源可以被理解为任何的具有扩展的生命周期的设备或者应用程序模块，而不只是一个简单的请求或者响应。这种设备或者应用程序模块是通过 Web services 来提供的。

目前，网格主要应用于分布式超级计算、分布式仪器系统、数据密集型计算、远程沉浸（网络环境下可视化与虚拟现实）和信息集成等领域。应用网格技术实现网络资源共享是信息网络发展的重要方向。国内外网格技术研究与开发利用发展很快，应用前景令人鼓舞。

1.1.2.5 信息安全技术

现代意义的信息安全技术发展总是跟计算机技术的发展紧密结合在一起的。信息安全最初指的是信息保密，随着网络时代的到来，信息安全的涵盖面越来越广，包括信息的保密性、完整性和可用性、可控性、不可否认性、可靠性等等。人们对信息安全的重视程度亦越来越高。1985 年，美国国防部颁布《可信计算机安全评价标准》(TCSEC)，用以指导计算机产品的制造和应用中的安全防护；欧洲一些国家在 20 世纪 90 年代颁布了相应的《信息技术安全评价标准》(ITSEC)；1999 年，我国发布了《计算机信息系统安全保护等级划分准则》(GB 17859—1999)；1999 年，国际标准化组织 ISO 提出了《信息技术安全评价通用准则》(CC)。这些标准的颁布为信息安全的研究和工程技术的发展提供了指导。

当前安全技术研究的着眼点主要在密码基础研究、基础设施安全、安全协议、安全服务等方面。密码学的研究热点在椭圆曲线、数字签名、Hash 函数等领域，这些领域的研究奠定信息安全的基础，为信息的完整性、机密性和可用性做理论保障。基础设施安全包含了网络、服务器、物理设备等软硬件设施，在我国特别需要这方面知识产权自主化的研究。安全协议是在密码学和网络研究的基础上，定义安全的交互协议，保障信息的动态安全。

在信息安全技术中，身份认证和访问控制对保证科学数据共享过程中的信息安全具有突出的重要意义。

身份认证的研究主要在身份认证的凭证和认证协议上。认证的凭证主要表现在用户的操作界面上，比如口令、证书、SmartCard 以及生物特征（指纹、虹膜、静脉等）等。这几方面始终都是身份认证的热点，比如如何将较弱的口令做成一个强大的不易被攻击的认证凭证，利用 SmartCard 和生物特征给用户登录提供方便等。目前的认证协议研究很多是跟认证凭证直接相关，热点多在利用公钥密码体制的不对称性，设计凭证交互和密钥交互协议，达到认证的方便，得到诸如双方互相认证，甚至匿名认证等结果。