

大数据的架构技术 与应用实践的探究

李佐军 编著

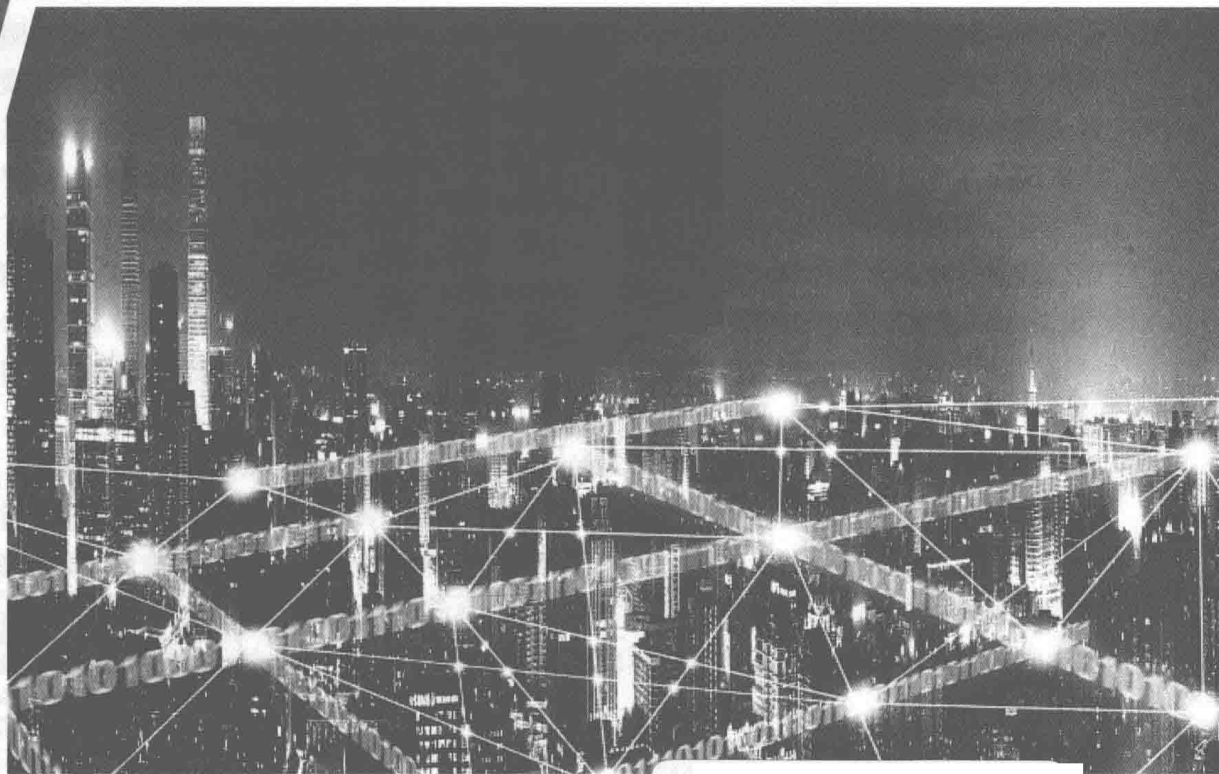


NORTHEAST NORMAL UNIVERSITY PRESS
WWW.NBNUP.COM

东北师范大学出版社

大数据的架构技术 与应用实践的探究

李佐军 编著



NORTHEAST NORMAL UNIVERSITY PRESS

WWW.NNUP.COM

东北师范大学出版社

图书在版编目(CIP)数据

大数据的架构技术与应用实践的探究 / 李佐军编著.
-- 长春: 东北师范大学出版社, 2019.4
ISBN 978-7-5681-5633-2

I. ①大… II. ①李… III. ①数据处理—研究 IV.
①TP274

中国版本图书馆CIP数据核字(2019)第060931号

责任编辑: 卢永康
 责任校对: 田 旺
 策划编辑: 王春彦
 封面设计: 优盛文化
 责任印制: 张允豪

东北师范大学出版社出版发行
长春市净月经济开发区金宝街118号(邮政编码: 130117)
销售热线: 0431-84568036
传真: 0431-84568036
网址: <http://www.nenup.com>
电子函件: sdcbcs@mail.jl.cn
定州启航印刷有限公司印装
2019年4月第1版 2019年4月第1次印刷
幅画尺寸: 185mm×260mm 印张: 17.25 字数: 302千
定价: 79.00元

随着 2013 年“大数据元年”的开启，各行各业都已经将大数据视为推动企业发展、推进行业进步、加快产业升级、促进民生繁荣、巩固社会安全甚至提升国家竞争力的核心武器。大数据从数据挖掘、商业智能发展而来，是信息技术发展的必然产物，是中国经济新常态下创新驱动的发动机和产业转型的助推器，带动了技术研发体系创新、管理方式改革、商业模式创新和产业价值链体系重构，推动了跨领域、跨行业的数据融合和协同创新。

大数据时代已经来临，它将在众多领域掀起变革的巨浪。但我们要冷静地看到，大数据的核心是为客户挖掘数据中蕴藏的价值，而不是软硬件的堆砌。因此，针对不同领域的大数据应用模式、商业模式研究将是大数据产业健康发展的关键。我们相信，在国家的统筹规划与支持下，通过各地方政府因地制宜地制定大数据产业发展策略，国内外 IT 龙头企业以及众多创新企业的积极参与，大数据产业未来的发展前景十分广阔。然而，大数据的实践应用不能一蹴而就，必须遵循科学的方法，循序渐进。无论是从业务的角度还是从技术的角度，将大数据应用讲清楚都不容易，尤其是要使非本领域的专家能对大数据有一个全面的了解更非易事。

为了帮助读者建立起对大数据应用全面、系统的认识，而不只是知道一些零散的技术或服务术语，本书站在系统论的高度对大数据的架构技术与实践应用做了全面论述。首先，介绍大数据的概念、发展现状与趋势，让读者对大数据有一个初步认知。其次，在分析大数据架构的基础上，展开介绍大数据的相关技术，涵盖技术支撑、数据获取技术、机器学习和数据挖掘技术、交互式分析技术、批处理技术、安全与隐私等，并结合大数据技术的各种应用实例，力图为读者呈现一个大数据的全景画卷。

在本书的编写过程中，参考和借鉴了一些学者的研究成果，在此对这些学者表示衷心的感谢。另外，由于时间及作者水平有限，本书难免存在疏漏与不妥之处，在本书出版之际，真诚地欢迎各位读者对本书提出宝贵的意见和建议。

第一章	大数据的初步认知	001
	第一节 何谓大数据	001
	第二节 大数据的发展现状	003
	第三节 大数据的发展趋势	018
第二章	大数据的架构分析	023
	第一节 大数据架构概述	023
	第二节 大数据架构设计	025
	第三节 大数据架构的参考模型	037
第三章	大数据的技术支撑	043
	第一节 云计算与大数据	043
	第二节 云资源的管理与调度	051
	第三节 开源云管理平台——OpenStack	059
	第四节 虚拟化技术的发展	063
	第五节 云存储系统的技术与分类	069
第四章	大数据的数据获取	096
	第一节 数据分类与数据获取组件	096
	第二节 探针在数据获取中的原理作用	098
	第三节 网页采集与日志收集	105
	第四节 数据分发中间件的作用分析	113
第五章	机器学习和数据挖掘技术	120
	第一节 机器学习与数据挖掘的关系	120
	第二节 机器学习的方式与类型	121
	第三节 机器学习与数据挖掘的应用	124
	第四节 深度学习的实践与发展	126

第六章	交互式分析	143
第一节	交互式分析的含义与特点	143
第二节	SQL on Hadoop	144
第三节	MPP DP 技术及其实际应用	148
第四节	大数据仓库的技术分析	167
第七章	批处理技术	171
第一节	批处理技术的含义与特点	171
第二节	MapReduce 模式解析	172
第三节	Spark 架构详解	176
第四节	BSP 框架模型及其实现	198
第八章	大数据的安全与隐私	208
第一节	大数据时代面临的安全挑战	208
第二节	安全问题解决的框架、技术与体系	214
第三节	大数据隐私保护的政策法规与技术介绍	222
第九章	电子商务与社会化网络中大数据分析	226
第一节	推荐系统分类与算法介绍	226
第二节	计算广告中大数据的应用	229
第三节	社交网络中大数据的应用	233
第十章	大数据架构的实现与应用实例	238
第一节	大数据架构的实现途径	238
第二节	基于 MATLAB 的图像处理与应用	242
第三节	基于 IT 治理的区域医疗卫生大数据架构的实现	245
第四节	微博用户兴趣建模系统架构设计	251
第五节	铁路客运旅游大数据平台分析	254
参考文献		264
后记		268



第一章

大数据的初步认知

第一节 何谓大数据

全球知名咨询公司麦肯锡在研究报告中指出，数据已经渗透到每个行业和业务职能领域，逐渐成为重要的生产因素，而人们对海量数据的运用将预示着新一波生产率增长和消费者盈余浪潮的到来。近年来，大数据概念的提出为中国数据分析行业的发展提供了无限的空间，越来越多的人认识到数据的价值。

那么，什么是大数据呢？大数据是规模非常巨大和复杂的数据集，传统数据库管理工具处理起来面临很多问题，如获取、存储、检索、共享、分析和可视化，数据量达到 PB、EB 和 ZB 的级别。大数据的特点可用 4 个 V 来概括：

- (1) 数据量 (Volume) 是持续快速增加的；
- (2) 高速度 (Velocity) 的数据 I/O；
- (3) 多样化 (Variety) 的数据类型和来源；
- (4) 数据价值 (Value) 大。

一、大数据的应用

大约从 2009 年开始，“大数据”才成为互联网信息技术行业的流行词语。美国互联网数据中心指出，互联网上的数据每年将增长 50%，每两年翻一番，而目前世界上 90% 以上的数据是最近几年才产生的。此外，数据并非单纯指人们在互联网上发布的信息，全世界的工业设备、汽车、电表上有着无数的数码传感器，随时测量和传递着有关位置、运动、震动、温度、湿度乃至空气中化学物质的变化，也产生了海量的数据信息。

数据充斥所带来的影响远远超出了企业界。贾斯汀·格里莫将数学与政治科学联系起来，他研究的内容涉及对博客文章、国会演讲和新闻稿进行计算机自动化分析等，希望借此洞察政治观点是如何传播的。在科学和体育、广告和公共卫生等其他领域，也有类似的情况——朝着数据驱动型的发现和决策的方向发生转变。

在公共卫生、经济发展和经济预测等领域中，大数据的预见能力正在被开发，而且已经崭露头角。研究者发现，曾有一次，“流感症状”和“流感治疗”等词语在谷歌上的搜索查询量

增加，而在几个星期以后，到某个地区医院急诊室就诊的流感病人数量就有所增加。

二、大数据的战略意义

大数据技术的战略意义不在于掌握庞大的数据信息，而在于对这些含有意义的数据进行专业化处理。换言之，如果把大数据比作一种产业，那么这种产业实现盈利的关键在于提高对数据的“加工能力”；通过“加工”实现数据的“增值”。中国物联网校企联盟认为，物联网的发展离不开大数据，依靠大数据可以提供足够有利的资源。

随着云时代的来临，大数据（Big data）也吸引了越来越多的关注。《著云台》的分析师团队认为，大数据通常用来形容一个公司创造的大量非结构化和半结构化数据，这些数据在下载关系数据库用于分析时会花费过多时间和金钱。大数据分析常和云计算联系到一起，因为实时的大型数据集分析需要像 MapReduce 一样的框架来向数十、数百甚至数千台计算机分配工作。

大数据分析相比于传统的数据仓库应用，具有数据量大、查询分析复杂等特点。

三、大数据的作用

第一，对大数据的处理分析正成为新一代信息技术融合应用的结点。移动互联网、物联网、社交网络、数字家庭、电子商务等是新一代信息技术的应用形态，这些应用不断产生大数据。云计算为这些海量、多样化的大数据提供存储和运算平台。通过对不同来源数据的管理、处理、分析与优化，将结果反馈到上述应用中，将创造出巨大的经济和社会价值。大数据具有催生社会变革的能量。但释放这种能量，需要严谨治理、富有洞见的数据分析和激发管理创新环境。

第二，大数据是信息产业持续高速增长的新引擎。面向大数据市场的新技术、新产品、新服务、新业态会不断涌现。在硬件与集成设备领域，大数据将对芯片、存储产业产生重要影响，还将催生一体化数据存储处理服务器、内存计算等市场。在软件与服务领域，大数据将引发数据快速处理分析、数据挖掘技术和软件产品的发展。

第三，大数据利用将成为提高核心竞争力的关键因素。各行各业的决策正在从“业务驱动”转向“数据驱动”。对大数据的分析，可以使零售商实时掌握市场动态并迅速做出应对；可以为商家制定更加精准有效的营销策略提供决策支持；可以帮助企业为消费者提供更加及时和个性化的服务。在医疗领域，大数据可提高诊断准确性和药物有效性；在公共事业领域，可发挥促进经济发展、维护社会稳定等方面的重要作用。

四、大数据与传统数据库

传统数据库 / 数据仓库是 GB/TB 级、高质量、较干净、强结构化、Top-down、重交易、确定解。大数据是 PB 级以上、有噪声、有冗余、非结构化、Bottom-up、重交互、满意解。大数据出现后，NoSQL 模式变得非常流行。大数据引发了一些问题，如对数据库高并发读写要求、对海量数据的高效存储和访问需求、对数据库高可扩展性和高可用性的需求，使传统 SQL 主要性能没有用武之地。互联网巨头对 NoSQL 数据模式应用非常广泛，如谷歌的 Big Table、Facebook 的 Cassandra、Oracle 的 NoSQL 及亚马逊的 Dynamo 等。从大数据处理角度来看，MapReduce

成为事实的标准。大数据的存储和处理,已有了成熟解决方案,这对于在系统软件中占较大比重的操作系统来说没有太大变化。一些重要的命题还没有解决,例如,操作系统对新兴计算资源的直接抽象的调度(GPU、APU),分布式文件系统下的统一数据视图、全数据中心范围内能耗管理、大数据下的安全性等,还不成熟,有待研发。

五、大数据与 Web

大多数研究大数据的商业公司都有明确的商业目的,即更好地支撑 Web 服务,如谷歌搜索引擎服务、Facebook SNS 网站、新浪微博网站等。在大数据驱动下的 Web 服务特征:更加流畅的网页交互体验,更加快速的社会资讯获取,更加便捷的日常工作和生活,更加深入的人、机、物融合。

回顾一下 Web 的发展,Web1.0 时代 Web 内容主要由网站服务商提供,Web2.0 时代用户大量参与 Web 内容的贡献,如博客和微博。到了 Web3.0 时代,特征就是人、机、物共同参与 Web 内容贡献,使 Web 形成对真实世界的全面映射。

大数据来源于人、机、物,同时服务于人、机、物。大数据时代系统软件,特别是操作系统有待进一步发展。人、机、物融合大数据将推动 Web 进入崭新 Web3.0 时代。

第二节 大数据的发展现状

一、大数据对社会的影响

数据用户行为反映真实需求。一切行为皆有前兆。未来的不确定性是人类产生工具类的根源之一。简单来说,从各种各样的数据中快速获取有价值的信息的能力,即大数据技术。大数据时代,软件价值体现在它所带来的数据规模、流量与活性;公司价值在于其拥有数据的规模、活性以及收集、运用数据的能力,这些决定了公司的核心竞争力。从国家层面看,国家数据主权体现在对数据的占有和控制。数字主权将是继边防、海防、空防之后,另一个大国博弈的空间。

(一) 泛互联网化

泛互联网化是收集用户数据的唯一低成本方式,能够带来数据规模和数据活性。泛互联网化带来软件使用的三个变化:跨平台、碎片化和门户化。

1. 跨平台

应用软件深度整合网络浏览器功能,桌面、移动终端(手机、平板电脑)拥有相同的体验和协同的功能。

2. 门户化

用户无须启用其他软件即可完成绝大多数的工作和沟通需求,对于个性化的用户需求,可以直接调用第三方应用或者插件完成。

3. 碎片化

把原来大型臃肿的软件拆分成多个独立的功能组件,用户可以按需下载使用。

这三个变化的核心意义在于收集用户行为资源，提高客户黏性；降低软件总体拥有成本，改变商业模式。

（二）行业垂直整合

开源软件加剧了基础软件的同质化趋势，而软、硬件一体化的趋势，进一步弱化了产业链上游的发言权。大数据产业结构发展趋势有两个维度：第一维度是大数据产业链，围绕数据的采集、整理、分析和反馈。第二维度是垂直的行业，像媒体、零售、金融服务、医疗和电信。

从这两个维度来看，大数据有三类商业模式：第一，大数据价值链环节，专注价值链的高附加值环节。第二，垂直产业大数据整合，利用大数据提高垂直产业效率。第三，大数据赋能者，提供大数据基础设施、技术和工具。

（三）数据成为资产

未来企业的竞争，将是拥有数据规模和活性的竞争，将是对数据解释和运用的竞争。围绕数据，可以演绎出六种新的商业模式：租售数据模式、租售信息模式、数据媒体模式、数据使用模式、数据空间运营模式、大数据技术提供商。

1. 租售数据模式

简单来说，即租售广泛收集、精心过滤、时效性强的数据。

2. 租售信息模式

一般聚焦某个行业，广泛收集相关数据，深度整合萃取信息，以庞大的数据中心加上专用传播渠道，也可成为一方霸主。此处，信息指的是经过加工处理，承载一定行业特征的数据集合。

3. 数据媒体模式

全球广告市场空间为 5 000 亿美元，具备培育千亿级公司的土壤和成长空间。这类公司的核心资源是获得实时、海量、有效的数据，立身之本是大数据分析技术，盈利源于精准营销。

4. 数据使用模式

如果没有大量的数据，缺乏有效的数据分析技术，此类公司的业务其实难以开展。通过在线分析小微企业的交易数据、财务数据，甚至可以计算出应提供多少贷款，多长时间可以收回等关键问题，把坏账风险降到最低。

5. 数据空间运营模式

从历史上看，传统的 IDC 即这种模式，互联网巨头都在提供此类服务。海外的 Dropbox、国内微盘都是此类公司的代表。这类公司的想象空间在于可以成长为数据聚合平台，盈利模式将趋于多元化。

6. 大数据技术提供商

从数据量上来看，非结构化数据是结构化数据的 5 倍以上，任何一个种类的非结构化数据处理都可以重现现有结构化数据的辉煌。语音数据处理领域、视频数据处理领域、语义识别领域、图像数据处理领域都可能出现大型的、高速成长的公司。

（四）云平台数据更加完善

企业越来越希望能将自己的各类应用程序及基础设施转移到云平台上。就像其他 IT 系统那样，大数据的分析工具和数据库也将走向云计算。

云计算能为大数据带来哪些变化呢?

首先,云计算为大数据提供了可以弹性扩展、相对便宜的存储空间和计算资源,使中小企业可以像亚马逊一样通过云计算来完成大数据分析。其次,云计算 IT 资源庞大、分布较为广泛,是异构系统较多的企业及时准确处理数据的有力方式,甚至是唯一的方式。

当然,大数据要走向云计算,还有赖于数据通信带宽的提高和云资源池的建设,需要确保原始数据能迁移到云环境以及资源池可以按需弹性扩展。

二、大数据的挑战与研究现状

大数据分析相比于传统的数据仓库应用,具有数据量大、查询分析复杂等特点。为了设计适合大数据分析的数据仓库架构,下面列举了大数据分析平台需要具备的几个重要特性,对当前的主流实现平台——并行数据库、MapReduce 及基于两者的混合架构进行了分析归纳,指出了各自的优势及不足,同时对各个方向的研究现状及大数据分析方面进行了介绍。

(一) 概述

最近几年,数据仓库又成为数据管理研究的热点领域,主要原因是当前数据仓库系统面临的需求在数据源、需提供的数据服务和所处的硬件环境等方面发生了根本性的变化,并且这些变化是我们必须面对的。

1. 三个变化

(1) 数据量。由 TB 级升到 PB 级,并仍在持续爆炸式增长。2011 年经调查显示,最大的数据仓库中的数据量,每两年增加 3 倍(年均增长率为 173%),其增长速度远超摩尔定律增长速度。照此增长速度计算,最近几年最大数据仓库中的数据量将逼近 100 PB。

(2) 分析需求。由常规分析转向深度分析(Deep Analytics)。数据分析日益成为企业利润必不可少的支撑点。根据 TDWI(中国商业智能网)对大数据分析的报告,企业已经不满足对现有数据的分析和监测,而是期望能对未来趋势有更多的分析和预测,以增强企业竞争力(图 1-1)。这些分析操作包括移动平均线分析、数据关联关系分析、回归分析、市场分析等复杂统计分析,我们称之为深度分析。

(3) 硬件平台。由高端服务器转向由中低端硬件构成的大规模机群平台。由于数据量的迅速增加,并行数据库的规模不得不随之增大,从而导致其成本的急剧上升。出于成本的考虑,越来越多的企业将应用由高端服务器转向由中低端硬件构成的大规模机群平台。

2. 两个问题

图 1-2 为一个典型的数据仓库架构。

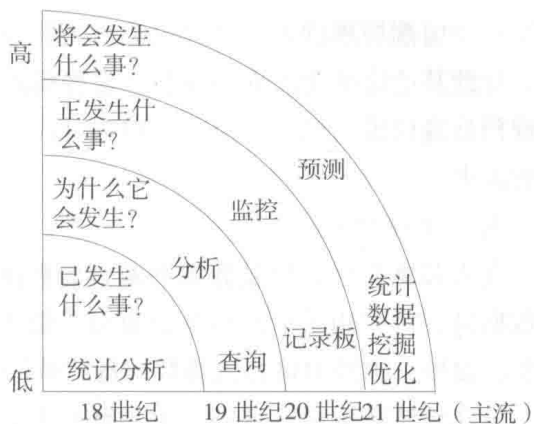


图 1-1 分析的趋势图

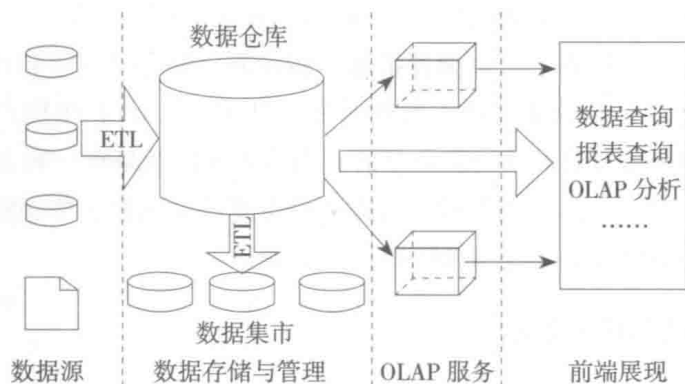


图 1-2 典型的数据仓库架构

由图 1-2 可以看出，传统的数据仓库将整个实现划分为 4 个层次，数据源中的数据先通过 ETL 工具被抽取到数据仓库中进行集中存储和管理，再按照星形模型或雪花模型组织数据，然后由 OLAP 工具从数据仓库中读取数据，生成数据立方体（MOLAP）或者直接访问数据仓库进行数据分析（ROLAP）。在大数据时代，此种计算模式存在以下两个问题。

（1）数据移动代价过高。在数据源层和分析层间引入一个存储管理层，可以提升数据质量并针对查询进行优化，同时付出了较大的数据迁移代价和执行时的连接代价。数据先通过复杂且耗时的 ETL 过程存储到数据仓库中，再在 OLAP 服务器中转化为星形模型或者雪花模型；执行分析时，又通过连接方式将数据从数据库中取出。这些代价在 TB 级时也许可以接受，但面对大数据，其执行时间至少会增长几个数量级。更为重要的是，对于大量的即时分析，这种数据移动的计算模式是不可取的。

（2）不能快速适应变化。传统的数据仓库假设主题是较少变化的，其应对变化的方式是对数据源到前端展现的整个流程中的每个部分进行修改，然后再重新加载数据，甚至重新计算数据，导致其适应变化的周期较长。这种模式比较适合对数据质量和查询性能要求较高，而不太计较预处理代价的场合。但在大数据时代，分析处在变化的业务环境中，这种模式将难以适应新的需求。

3. 一个鸿沟

在大数据时代，巨量数据与系统的数据处理能力间将会产生一个鸿沟：一边是至少 PB 级的数据量，另一边是面向传统数据分析能力设计的数据仓库和各种 BI 工具。如果这些系统工具发展缓慢，该鸿沟将会随着数据量的持续爆炸式增长而逐步拉大。

虽然传统数据仓库可以采用舍弃不重要数据或者建立数据集市的方式来缓解此问题，但毕竟只是权宜之策，并非系统级解决方案，而且舍弃的数据在未来可能会重新使用，以发掘出更大的价值。

（二）期望特性

数据仓库系统需具备几个重要特性，如表 1-1 所示。

表 1-1 数据仓库系统需要具备的特性

特 性	说 明
高度可扩展性	横向大规模可扩展，大规模并行处理
高性能	快速响应复杂查询与分析
高度容错性	对硬件平台一致性要求不高，适应能力强
支持异构环境	业务需求变化时，能快速反应
较低的分析延迟	既方便查询，又能处理复杂分析
较低成本	较高的性价比
向下兼容性	支持传统的商务智能工具

1. 高度可扩展性

一个明显的事实是，数据库不能依靠一台或少数几台机器的升级（scale-up 纵向扩展）满足数据量的爆炸式增长，而是希望能方便地做到横向可扩展（scale-out）来实现此目标。

普遍认为无共享结构 shared-nothing，每个节点都拥有私有内存和磁盘，并且通过高速网络与其他节点互连，具备较好的扩展性。分析型操作往往涉及大规模的并行扫描、多维聚集及星形连接操作，这些操作也比较适合在无共享结构的网络环境下运行。Teradata 即采用此结构，Oracle 在其新产品 Exadata 中也采用了此结构。

2. 高性能

数据量的增长并没有降低对数据库性能的要求，反而有所提高。软件系统性能的提升可以降低企业对硬件的投入成本，节省计算资源，提高系统吞吐量。巨量数据的效率优化，并行是必由之路。1 PB 数据在 50 MB/s 速度下串行扫描一次，需要 230 天；而在 6 000 块磁盘上，并行扫描 1 PB 数据只需要 1 小时。

3. 高度容错性

大数据的容错性要求在查询执行过程中，一个参与节点失效时，不需要重做整个查询，而机群节点数的增加会带来节点失效概率的增加。在大规模机群环境下，节点的失效将不再是稀有事件（根据谷歌报告，平均每个 MapReduce 数据处理任务即有 1.2 个工作节点失效）。因此，在大规模机群环境下，系统不能依赖硬件来保证容错性，要更多地考虑软件容错。

4. 支持异构环境

建设同构系统的大规模机群难度较大，原因在于计算机硬件更新较快，一次性购置大量同构的计算机是不可取的，而且会在未来添置异构计算资源。此外，不少企业已经积累了一些闲置的计算机资源。此种情况下，异构环境不同节点的性能是不一样的，可能出现“木桶效应”，即最慢节点的性能决定整体处理性能。因此，异构的机群需要特别关注负载均衡、任务调度等方面的设计。

5. 较低的分析延迟

分析延迟是分析前的数据准备时间。在大数据时代，分析所处的业务环境是变化的，因此要求系统能动态地适应业务分析需求。在分析需求发生变化时，减少数据准备时间，系统能尽快地做出反应，快速地进行数据分析。

6. 较低的成本

在满足需求的前提下数据仓库系统技术成本越低，其生命力就越强。值得指出的是，成本是一个综合指标，不仅是硬件或软件的代价，还应包括日常运维成本（网络费用、电费、建筑等）和管理人员成本等。据报告显示，数据中心的主要成本不是硬件的购置成本，而是日常运维成本，因此在设计系统时需要更多地关注此项内容。

7. 向下兼容性

数据仓库发展的30年，产生了大量面向客户业务的数据处理工具（如 Informatica、Data Stage 等）、分析软件（如 SPSS、R、Matlab 等）和前端展现工具（如水晶报表等）。这些软件是一笔宝贵的财富，已被分析人员所熟悉，是大数据时代中小规模数据分析的必要补充。因此，新的数据仓库需考虑同传统商务智能工具的兼容性。由于这些系统往往提供标准驱动程序，如 ODBC、JDBC 等，这项需求的实际要求是对 SQL 的支持。

总而言之，以较低的成本投入进行高效的数据分析是大数据分析的基本目标。

（三）并行数据库

并行数据库系统（Parallel Database System）是新一代高性能的数据库系统，是在 MPP 和集群并行计算环境的基础上建立的数据库系统。

并行数据库技术起源于20世纪70年代的数据库机（Database Machine）的研究，研究的内容主要集中在关系代数操作的并行化和实现关系操作的专用硬件设计上，希望通过硬件实现关系数据库操作的某些功能，该研究以失败告终。20世纪80年代后期，并行数据库技术的研究方向逐步转到了通用并行机方面，研究的重点是并行数据库的物理组织、操作算法、优化和调度策略。从20世纪90年代至今，随着处理器、存储、网络等相关基础技术的发展，并行数据库技术的研究上升到一个新的水平，研究的重点也转移到数据操作的时间并行性和空间并行性上。

并行数据库系统的目标是高性能（High Performance）和高可用性（High Availability），通过多个处理节点并行执行数据库任务，提高整个数据库系统的性能和可用性。

性能指标关注的是并行数据库系统的处理能力，具体的表现可以总结为数据库系统处理事务的响应时间。并行数据库系统的高性能可以从两个方面理解，即速度提升（Speed-up）和范围提升（Scale-up）。速度提升是指通过并行处理，可以使用更少的时间完成更多样的数据库事务。范围提升是指通过并行处理，在相同的处理时间内，可以完成更多的数据库事务。并行数据库系统基于多处理节点的物理结构，将数据库管理技术与并行处理技术有机结合，实现系统的高性能。

可用性指标关注的是并行数据库系统的健壮性，也就是当并行处理节点中的一个节点或多个节点部分失效或完全失效时，整个系统对外持续响应的能力。高可用性可以同时从硬件与软件两个方面提供保障。

在硬件方面,通过冗余的处理节点、存储设备、网络链路等硬件措施,可以保证当系统中某节点部分或完全失效时,其他的硬件设备可以接手处理,对外提供持续服务。

在软件方面,通过状态监控与跟踪、互相备份、日志等技术手段,可以保证当前系统中某节点部分或完全失效时,由它所进行的处理或由它所掌控的资源可以无损失或基本无损失地转移到其他节点,并由其他节点继续对外提供服务。

为了实现和保证高性能和高可用性,可扩充性也成为并行数据库系统的一个重要指标。可扩充性是指并行数据库系统通过增加处理节点或者硬件资源(处理器、内存等),使其可以平滑地或线性地扩展其整体处理能力的特性。

随着对并行计算技术研究的深入和 SMP、MPP 等处理机技术的发展,并行数据库的研究也进入了一个新的领域。集群已经成为并行数据库系统中最受关注的热点。目前,并行数据库领域主要还有下列问题需要进一步研究和解决。

(1) 并行体系结构及其应用,这是并行数据库系统的基础问题。为了达到并行处理的目的,参与并行处理的各个处理节点之间是否要共享资源、共享哪些资源、需要多大程度的共享,这些就需要研究并行处理的体系结构及有关实现技术。

(2) 并行数据库的物理设计,主要是在并行处理的环境下,对数据分布的算法的研究、数据库设计工具与管理工具的研究。

(3) 处理节点间通信机制的研究。为了实现并行数据库的高性能,并行处理节点要最大限度地协同处理数据库事务,因此节点间必不可少地存在通信问题。如何支持大量节点之间消息和数据的高效通信,也成为并行数据库系统中一个重要的研究课题。

(4) 并行操作算法。为提高并行处理的效率,需要在数据分布算法研究的基础上,深入研究链接、聚集、统计、排序等具体的数据操作在多节点上的并行操作算法。

(5) 并行操作的优化和同步。为获得高性能,如何将一个数据库处理事务合理地分解成相对独立的并行操作步骤,如何将这些步骤以最优的方式在多个处理节点间进行分配,如何在多个处理节点的同一步骤和不同步骤之间进行消息和数据的同步,这些问题都值得深入研究。

(6) 并行数据库中数据的加载和再组织技术。为了保证高性能和高可用性,并行数据库系统中的处理节点可能需要进行扩充(或者调整),这就需要考虑如何对原有数据进行卸载、加载,以及如何合理地在各个节点重新组织数据。

(四) MapReduce

MapReduce 的编程模型不同于以前学过的大多数编程模型,它是一种用于大规模数据集(大于 1 TB)的并行运算的编程模型。其概念 Map(映射)和 Reduce(化简),及它们的主要思想,都是从函数式编程语言里借来的和从矢量编程语言里借来的特性。它极大地方便了编程人员在不会分布式并行编程的情况下,将自己的程序在分布式系统上运行。当前的软件实现是指定一个 Map(映射)函数,用来把一组键值对映射成一组新的键值对;指定并发的 Reduce(化简)函数,用来保证所有映射的键值对中的每一个共享相同的键组。

MapReduce 采用“分而治之”的思想,把对大规模数据集的操作,分发给一个主节点管理下的各分节点共同完成,接着通过整合各分节点的中间结果,得到最终的结果。简单来说,

MapReduce 就是“任务的分散与结果的汇总”。

MapReduce 处理过程被 MapReduce 高度地抽象为两个函数：Map 和 Reduce。Map 负责把任务分解成多个任务，Reduce 负责把分解后多任务处理的结果汇总起来。至于在并行编程中的各种复杂问题，如分布式存储、工作调度、负载均衡、容错处理、网络通信等，均由 MapReduce 框架负责处理，可以不用程序员烦心。值得注意的是，用 MapReduce 来处理的数据集（或任务）必须具备这样的特点：待处理的数据集可以分解成许多小的数据集，且每一小数据集都可以完全并行地进行处理。

图 1-3 给出了使用 MapReduce 处理数据集的过程。

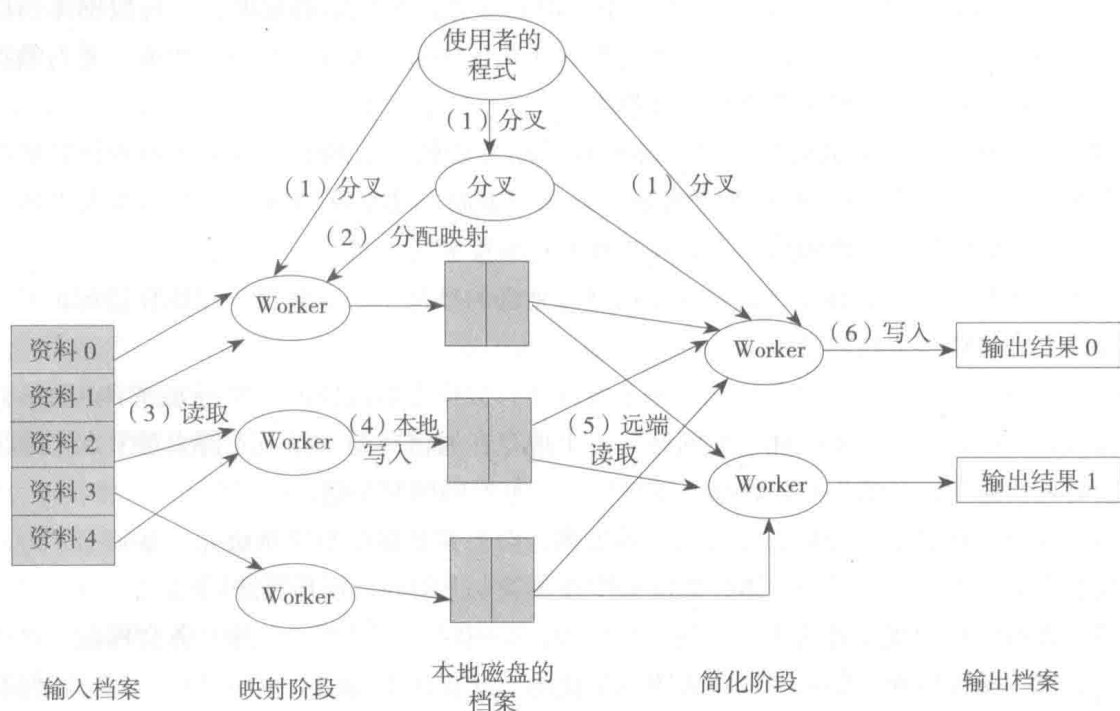


图 1-3 MapReduce 处理数据集的过程

由图 1-3 可以看出，该计算模型的核心部分是 Map 和 Reduce 函数。这两个函数的具体功能由用户设计实现，只要能够按照用户自定义的规则，将输入的 $\langle \text{key}, \text{value} \rangle$ 对转换成另一个或一批 $\langle \text{key}, \text{value} \rangle$ 对输出即可。

在 Map 阶段，MapReduce 框架将任务的输入数据分隔成固定大小的片段，随后将每个片段进一步分解成一批键值对 $\langle K1, V1 \rangle$ 。Hadoop 为每一个片段创建一个 Map 任务（可简称为 Mapper）用于执行用户自定义的 Map 函数，并将对应片段中的 $\langle K1, V1 \rangle$ 对作为输入，得到计算的中间结果 $\langle K2, V2 \rangle$ 。接着，将中间结果按照 $K2$ 进行排序，并将 key 值相同的 value 放在一起形成一个新列表，形成 $\langle K2, \text{list}(V2) \rangle$ 元组。最后，根据 key 值的范围将这些组进行分组，对应不同的 Reduce 任务（可简称为 Reducer）。

在 Reduce 阶段，Reducer 把从不同 Mapper 接收来的数据整合在一起并进行排序，然后调用用户自定义的 Reduce 函数，对输入的 $\langle K2, \text{list}(V2) \rangle$ 对进行相应的处理，得到键值对 $\langle K3, V3 \rangle$ 。

并输出到 HDFS 上。既然 MapReduce 框架为每个片段创建一个 Mapper，那么谁来确定 Reducer 的数目呢？其答案是用户。Mapped-site.XML 配置文件中有一个表示 Reducer 数目的属性 mapped.reduce.tasks，该属性的默认值为 1，开发人员可通过 job.setNumReduceTasks() 方法重新设置该值。

MapReduce 架构结构组成部分主要有以下类型，下面给予介绍。

1. 组成部分

(1) JobClient。每一个 job 都会在用户端通过 JobClient 类将应用程序以及配置参数打包成 jar 文件存储在 HDFS，并把路径提交到 JobTracker，然后由 JobTracker 创建每一个 Task（即 MapTask 和 ReduceTask）并将它们分发到各个 TaskTracker 服务中去执行。

(2) JobTracker。JobTracker 是一个 master 服务，JobTracker 负责调度 job 的每一个子任务 task 运行于 TaskTracker 上，并监控它们，如果发现有失败的 task 就重新运行它。一般应该把 JobTracker 部署在单独的机器上。

(3) TaskTracker。TaskTracker 是运行于多个节点上的 slaver 服务。TaskTracker 则负责直接执行每一个 task。TaskTracker 都需要运行在 HDFS 的 DataNode 上。

MapReduce 架构结构中的各角色运行过程如图 1-4 所示。

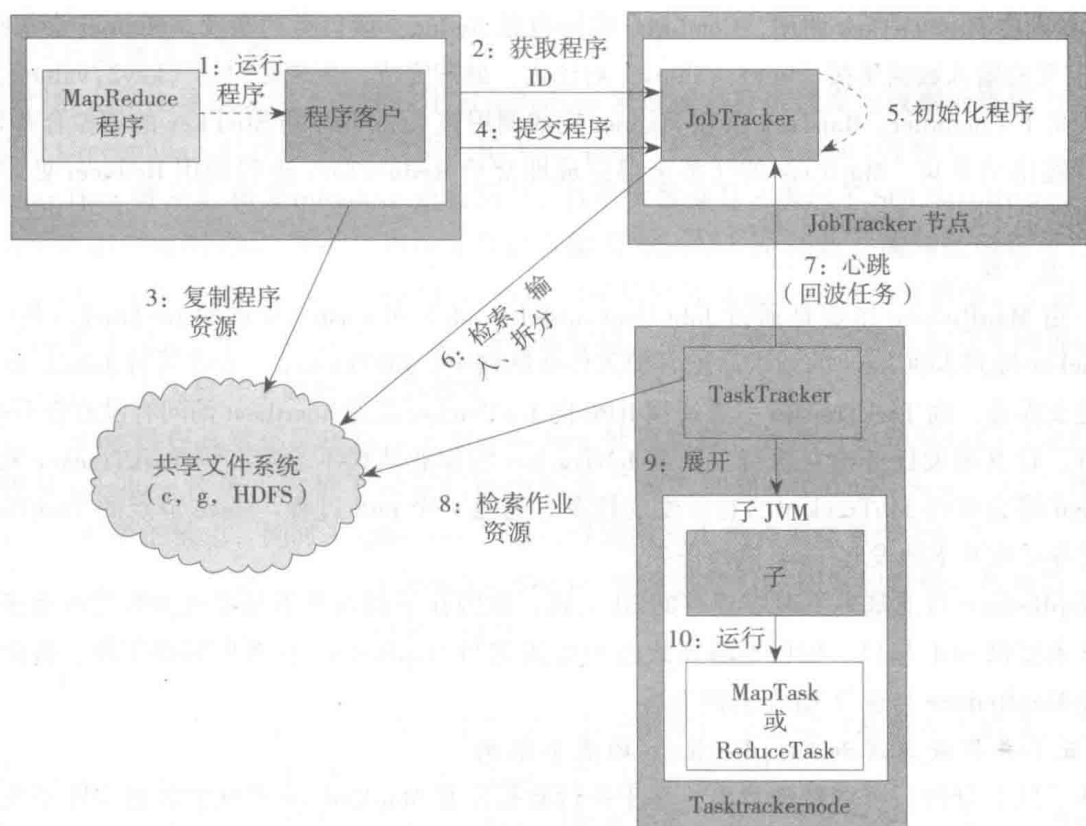


图 1-4 MapReduce 架构结构中的各角色运行过程图

2. 数据结构

(1) Mapper 和 Reducer。运行于 Hadoop 的 MapReduce 应用程序最基本的组成部分包括一个 Mapper 和一个 Reducer 类，以及一个创建 JobConf 的执行程序，在一些应用中还可以包括一