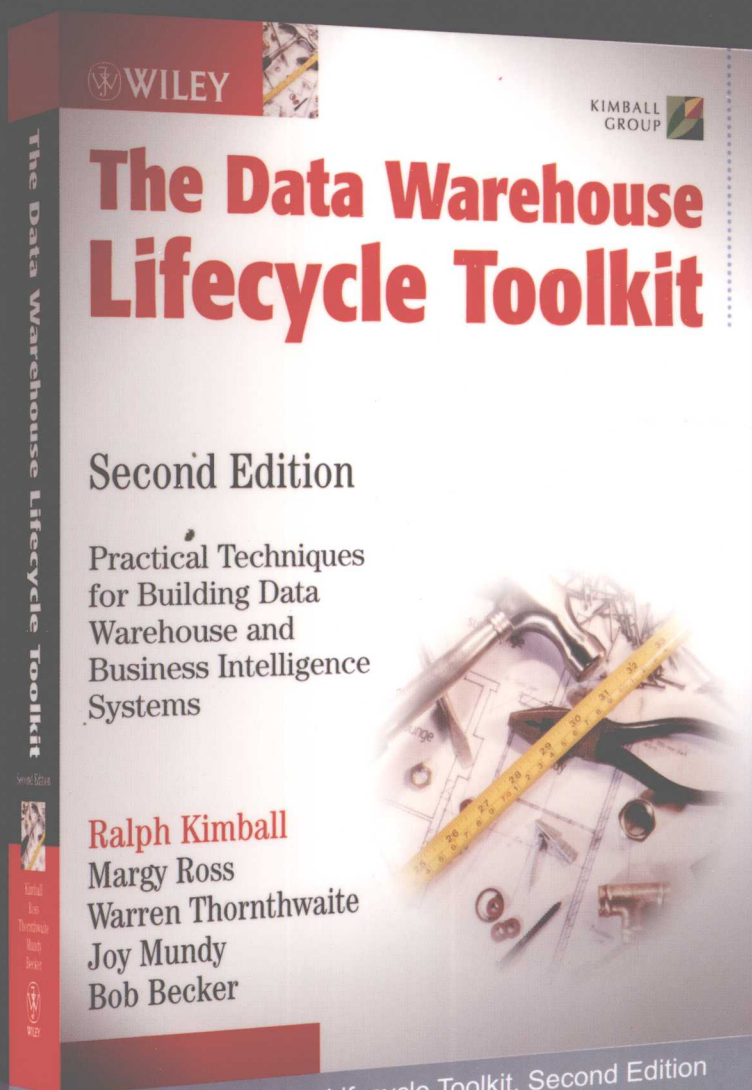


数据仓库生命周期 工具箱 (第二版)

(美) Ralph Kimball 等著 唐富年 孙媛媛 译



The Data Warehouse Lifecycle Toolkit, Second Edition



清华大学出版社

国外计算机科学经典教材

数据仓库生命周期工具箱

(第二版)

(美) Ralph Kimball 等著
唐富年 孙媛媛 译

清华大学出版社

北 京

Ralph Kimball, Margy Ross, Warren Thornthwaite, Joy Mundy, Bob Becker

The Data Warehouse Lifecycle Toolkit, Second Edition

EISBN: 978-0-470-14977-5

Copyright©2008 by Wiley Publishing, Inc.

All Rights Reserved. This translation published under license.

本书中文简体字版由 Wiley Publishing, Inc. 授权清华大学出版社出版。未经出版者书面许可, 不得以任何方式复制或抄袭本书内容。

北京市版权局著作权合同登记号 图字: 01-2009-4344

本书封面贴有 Wiley 公司防伪标签, 无标签者不得销售。

版权所有, 侵权必究。侵权举报电话: 010-62782989 13701121933

图书在版编目(CIP)数据

数据仓库生命周期工具箱(第二版)/(美)金博尔(Kimball, R.)等著;唐富年,孙媛媛译.

—北京:清华大学出版社,2009.9

书名原文: The Data Warehouse Lifecycle Toolkit, Second Edition

ISBN 978-7-302-20374-2

I. 数… II. ①金… ②唐… ③孙… III. 数据库系统 IV. TP311.13

中国版本图书馆 CIP 数据核字(2009)第 100738 号

责任编辑:王军 李楷平

装帧设计:孔祥丰

责任校对:成凤进

责任印制:王秀菊

出版发行:清华大学出版社

地 址:北京清华大学学研大厦 A 座

<http://www.tup.com.cn>

邮 编:100084

社 总 机:010-62770175

邮 购:010-62786544

投稿与读者服务:010-62776969, c-service@tup.tsinghua.edu.cn

质 量 反 馈:010-62772015, zhiliang@tup.tsinghua.edu.cn

印 刷 者:清华大学印刷厂

装 订 者:三河市溧源装订厂

经 销:全国新华书店

开 本:185×260 印 张:31 字 数:754 千字

版 次:2009 年 9 月第 1 版 印 次:2009 年 9 月第 1 次印刷

印 数:1~4000

定 价:68.00 元

本书如存在文字不清、漏印、缺页、倒页、脱页等印装质量问题,请与清华大学出版社出版部联系调换。联系电话:(010)62770177 转 3103 产品编号:030075-01

出版说明

近年来，我国的高等教育特别是计算机学科教育，进行了一系列大的调整和改革，亟需一批门类齐全、具有国际先进水平的计算机经典教材，以适应我国当前计算机科学的教學需要。通过使用国外优秀的计算机科学经典教材，可以了解并吸收国际先进的教学思想和教学方法，使我国的计算机科学教育能够跟上国际计算机教育发展的步伐，从而培养出更多具有国际水准的计算机专业人才，增强我国计算机产业的核心竞争力。为此，我们从国外多家知名的出版机构 Pearson、McGraw-Hill、John Wiley & Sons、Springer、Cengage Learning 等精选、引进了这套“国外计算机科学经典教材”。

作为世界级的图书出版机构，Pearson、McGraw-Hill、John Wiley & Sons、Springer、Cengage Learning 通过与世界级的计算机教育大师携手，每年都为全球的计算机高等教育奉献大量的优秀教材。清华大学出版社和这些世界知名的出版机构长期保持着紧密友好的合作关系，这次引进的“国外计算机科学经典教材”便全是出自上述这些出版机构。同时，为了组织该套教材的出版，我们在国内聘请了一批知名的专家和教授，成立了专门的教材编审委员会。

教材编审委员会的运作从教材的选题阶段即开始启动，各位委员根据国内外高等院校计算机科学及相关专业的现有课程体系，并结合各个专业的培养方向，从上述这些出版机构出版的计算机系列教材中精心挑选针对性强的题材，以保证该套教材的优秀性和领先性，避免出现“低质重复引进”或“高质消化不良”的现象。

为了保证出版质量，我们为这套教材配备了一批经验丰富的编辑、排版、校对人员，制定了更加严格的出版流程。本套教材的译者，全部由对应专业的高校教师或拥有相关经验的 IT 专家担任。每本教材的责编在翻译伊始，就定期不间断地与该书的译者进行交流与反馈。为了尽可能地保留与发扬教材原著的精华，在经过翻译、排版和传统的三审三校之后，我们还请编审委员或相关的专家教授对文稿进行审读，以最大程度地弥补和修正在前面一系列加工过程中对教材造成的误差和瑕疵。

由于时间紧迫和受全体制作人员自身能力所限，该套教材在出版过程中很可能还存在一些遗憾，欢迎广大师生来电来信批评指正。同时，也欢迎读者朋友积极向我们推荐各类优秀的国外计算机教材，共同为我国高等院校计算机教育事业贡献力量。

清华大学出版社

国外计算机科学经典教材

编审委员会

主任委员：

孙家广 清华大学教授

副主任委员：

周立柱 清华大学教授

委员（按姓氏笔画排序）：

王成山	天津大学教授
王 珊	中国人民大学教授
冯少荣	厦门大学教授
冯全源	西南交通大学教授
刘乐善	华中科技大学教授
刘腾红	中南财经政法大学教授
吉根林	南京师范大学教授
孙吉贵	吉林大学教授
阮秋琦	北京交通大学教授
何 晨	上海交通大学教授
吴百锋	复旦大学教授
李 彤	云南大学教授
沈钧毅	西安交通大学教授
邵志清	华东理工大学教授
陈 纯	浙江大学教授
陈 钟	北京大学教授
陈道蕃	南京大学教授
周伯生	北京航空航天大学教授
孟祥旭	山东大学教授
姚淑珍	北京航空航天大学教授
徐佩霞	中国科学技术大学教授
徐晓飞	哈尔滨工业大学教授
秦小麟	南京航空航天大学教授
钱培德	苏州大学教授
曹元大	北京理工大学教授
龚声蓉	苏州大学教授
谢希仁	中国人民解放军理工大学教授

关于作者

本书几位作者的职业生涯都有着非常相似的轨迹。每个作者在数据仓库和商业智能(DW/BI)领域从事咨询和培训的时间都超过了15年。大多数作者都曾经在Metaphor计算机系统公司(20世纪80年代一家具有开创性的提供决策支持系统产品的厂商)中共事过。所有的作者都是Kimball集团的成员,都在Kimball大学讲授课程。他们经常向*Intelligent Enterprise*杂志和其他的业界刊物供稿,而且大多数作者以前都参与过工具箱系列丛书的撰写。

Ralph Kimball 创建了 Kimball 集团。从 20 世纪 80 年代中期开始,他就是 DW/BI 行业维度方法的领导者,并且培养了 10 000 多位 IT 专业人士。Ralph 在斯坦福大学(Stanford University)获得过电子工程专业博士学位。

Margy Ross 是 Kimball 集团的总裁。她从 1982 年就开始专门从事 DW/BI 方面的研究,特别强调业务需求分析和维度建模。Margy 在西北大学(Northwestern University)获得过工业管理学学士学位。

Warren Thornthwaite 从 1980 年开始从事 DW/BI 方面的研究。他刚开始负责管理 Metaphor 公司的咨询机构,后来为斯坦福大学和 WebTV 工作。Warren 在密西根大学(University of Michigan)获得过传媒研究学士学位,并且从宾夕法尼亚大学沃顿商学院(University of Pennsylvania's Wharton School)获得了 MBA 学位。

Joy Mundy 从 1992 年开始从事 DW/BI 方面的研究,在斯坦福、Web TV 和微软公司的 SQL Server 产品开发机构中都工作过一段时间。Joy 在塔夫斯大学(Tufts University)获得过经济学学士学位,并且在斯坦福大学获得了工程经济系统硕士学位。

Bob Becker 从 1989 年就开始帮助多个行业的客户解决 DW/BI 方面的难题,并提供相应的解决方案,还为卫生保健机构做了大量的工作。Bob 从明尼苏达州商学院(Minnesota's School of Business)获得了市场营销学商学士学位。

译者序

本书是数据仓库和商业智能领域的又一部经典著作，作者 Kimball 等人在数据仓库领域享有很高的声誉。本书的作者都长期工作在数据仓库/商业智能系统开发的第一线，他们将自己多年的经验和感悟写入到了本书的字里行间。全书讲述了整个 Kimball 生命周期过程各个环节的具体工作，从业务需求的视角，引导读者全面认识数据仓库/商业智能系统的开发。

本书并不是一本为刚刚涉足数据仓库领域的新手准备的入门教材，我们建议您在阅读本书之前最好能有一些数据仓库或数据库领域的基本知识，这对把握本书的内容有很多好处。在书中，作者使用了一种通俗而明晰的手法进行内容的讲述，因此译者在翻译时也尽量使用比较通俗的语言。本书的写作风格匠心独运，在内容安排上也体现出了 Kimball 生命周期中的“迭代”思想。全书在章节上有总有分，一层层展开，一步步深入。一些读者刚刚拿到这本书时也许难以习惯它的内容布局，但是正如作者所说，只要坚持读完全书，最终必然会理解数据仓库/商业智能系统开发工作的总体流程和关键环节。

本书在每一章开始都说明了这一章的内容适合于哪些读者，在每一章的结尾都进行了详细而具有针对性的总结。读者可以通读全书，也可以根据自身的需要选读其中的部分章节。目前市场上有关数据仓库的书籍非常多，而且各自的侧重点都有所不同，恐怕难以评论孰优孰劣。但是本书的出现，也许会让您有耳目一新的感觉，因为本书的出发点是站在业务的角度，而并没有单纯从技术的角度来安排各个章节的内容，这就使书中的内容更贴近实际，也更具有可操作性。

在本书的翻译过程中，我们参考了很多数据仓库和商业智能方面的书籍，这对我们准确把握本书的内容起到了重要作用。遗憾的是，在数据仓库和商业智能领域至今仍然有大量的术语没有统一的标准译法，在翻译过程中我们对这些术语尽量都使用了比较常用的译法。在此，非常感谢本书第一版的译者们，他们的工作具有很高的借鉴价值。

本书由唐富年、孙媛媛翻译。同时，也非常感谢国防科技大学信息系统与管理学院的张文煜博士，他为本书的翻译提出了一些很宝贵的意见。肖国尊负责翻译质量和进度的控制，以及翻译思想的指导，在此予以衷心感谢。鉴于译者水平有限，译文中难免存在错漏之处，还望广大读者谅解并不吝指正。敬请读者将反馈意见发至邮箱 be-flying@sohu.com。

前 言

在《数据仓库生命周期工具箱》第一版出版之后的九年里，数据仓库产业已经发生了显著的变化。现在，数据仓库产业已经变得十分成熟，并且得到了商业界的接受和认可。在这九年里，硬件和软件的发展取得了令人难以置信的成就。我们已经开始谈论“TB”字节而不再是“GB”字节。但是，数据仓库的任务基本上没有什么改变。

许多人所在的机构里有数千位的数据仓库用户，从业务决策者到一般的数据仓库用户，再到市场营销和财务用户中的骨干成员。事实上，操作方面的迫切需求是数据仓库研究中的最热点问题，而且每个人都坚持认为他们需要“实时”数据。在数据仓库变得越来越重要、越来越直观的同时，用户反复提出保密性、安全性和合规性方面的需求。业务用户正在逐步意识到高质量数据的价值，这 and 传统制造业注重质量管理是一样的道理。最后，可能也是最重要的，我们为自己所从事的这个行业起了一个新的名称，这个名称反映了我们的真正目的。它就是：商业智能(Business Intelligence)。为了强调这一点，在本书的大多数地方，我们都把您要创建的整个系统叫做 DW/BI 系统。

商业智能的这种转变将主动权移交到了业务用户手中，而不再是由 IT 人员掌握主动权。但同时这个转变将全部注意力都集中到了数据仓库的使命上：它是商业智能必需的平台。数据仓库需要做繁重的工作，它从源系统中取得数据，对数据进行清洗，并将数据组织起来使普通的业务用户能够看懂它。当然，我们力争实现世界级的商业智能，但是世界级的商业智能需要您拥有一个世界级的数据仓库。反过来说，一个数据仓库没有商业智能将会遭遇彻底的失败。

本书是 DW/BI 系统的设计人员、管理人员和所有者在实际工作中的指南。我们尽可能使本书的内容非常具体和实用，以便将这本书与其他 DW/BI 书籍区分开来。被本书的内容搞得眼花缭乱并没有关系，我们希望您一直坚持读下去，最终必然会达到预定目标。这本书描述了一个条理清晰的 DW/BI 系统框架，这个框架从确定整个企业 DW/BI 系统的初始范围开始，经过详细的开发和部署步骤，一直到最后计划下一个阶段的工作。

全世界安装有好几个功能各异的数据仓库。很多 DW/BI 系统的所有者都是完全按照生命周期的思想来进行开发的。或许从生命周期思想中可以得到的最大收获就是：每个 DW/BI 系统都是不断发展变化的，它永远都不是静止的，也决不会停止转变。新的商业需

求会不断涌现,新的管理者和执行官也会对 DW/BI 系统提出一些不可预知的要求,还会有新的可用的数据源加入到系统中。至少, DW/BI 系统需要根据所处机构的变化而同步地变化。稳定的机构也会要求 DW/BI 系统取得适度的演化,不断变化的动态机构则会使 DW/BI 系统的任务变得更具有挑战性。

考虑到 DW/BI 系统具有不断演化的特性,我们需要灵活可变的、适应性强的设计方法,还必须同时扮演 DBA 和 MBA 的角色。我们需要将来自单个业务过程的小块数据连成大块数据,从而形成企业级的数据仓库。同时,还要求对 DW/BI 系统所做的改变始终是适度的。一个适度的改变不会使以前的数据或先前的应用程序无效。

本书结构

本书有两大基本主题。第一个主题就是 Kimball 生命周期方法。您可能会问:“是什么使得 Kimball 生命周期方法与其他方法不同?”最简单的答案是我们从业务用户的角度开始,找出他们完成工作需要什么,并以此来建立 DW/BI 系统。有了这些需求之后,我们逐步向下使报表、应用程序、数据库和软件系统地进行工作,最后再深入到底层设施的物理层。这与技术驱动的方法形成鲜明的对比——其顺序正好相反。在 20 世纪 90 年代初期,一些 IT 工作室并不知道怎样使用我们这种面向业务和用户的方法。但是随着 2008 年本书的出版,“商业智能”这个名称本身就说明了这一切,即应该由用户和业务来驱动数据仓库。

第二个主题是“总线架构”。本书中介绍了如何进行单个业务过程的连续迭代,使读者最终能够创建一个企业 DW/BI 系统。在本书中,您会看到我们将维度模型作为一种向业务用户表现数据的可靠方式。推荐这种方法只有一个原因:它确实能够满足业务用户的愿望,简单并且具有高效的查询性能。我们由衷地感谢您能够选择本书所讲述的维度建模方法。最后,您可以使用任何您认为比较合适的方式将数据呈现给用户。这不应当由我们来决定,而是应当由用户来决定。

本书中涵盖了上述观点,提供了能够辅助用户完成工作的各种具有实际价值的技能和实用的工具。通过这种方式,希望能将我们自 1982 年以来在建立 DW/BI 系统的过程中积累的想法和经验都讲授给您。

本书读者对象

本书的主要读者应该是那些真正需要在实现“作为商业智能应用平台的数据仓库”的过程中负责创建和管理工作的设计人员或者管理人员。因为“作为商业智能应用平台的数据仓库”这句话十分冗长而拗口,因此在提到整个系统时我们都使用“DW/BI”这一名称,它说明您需要负责从初始源系统获得数据一直到将数据显示到业务用户屏幕的整个过程。

尽管本书包含一些介绍性的内容,但是我们认为这本书对于已经对数据仓库技术有一定接触的 IT 专业人员来说将非常有用。在 2002 年出版的另一本相关书籍,是由 Ralph Kimball 和 Margy Ross 编写的 *The Data Warehouse Toolkit [Second Edition]*, 该书更加深入而具体地讲述了维度建模。

通过设计和交付一个真正的数据仓库，您可能已经积累了一些经验并且形成了自己的观点，这就是最好的知识背景！开发一个实际的数据仓库所积累的经验是任何其他方式都无法替代的。在将自己的“杰作”提交给一群要求苛刻的业务用户时，我们这些作者都曾有过一些令人感到羞愧的经历。令人难以接受的是，大多数用户的实际工作与技术毫不相干，他们甚至可能特别不喜欢技术。但是如果我们的技术易于使用，并且能够为用户提供明显的使用价值，那么业务用户还是会使用我们的技术的。

本书要求有一定的专业知识。其中，有关设计技巧和体系结构方面的论述毋庸置疑地会引入一些您未曾遇到过的专业术语。我们已经对本书进行了精心梳理，以确保大多数技术方面的主题都是读者应该能够理解的。我们尽量使本书不会因为内容本身的原因而使您陷入细节上的困扰。在本书后面的 DW/BI 的术语表中，将简要地解释我们在书中所使用的术语。

尽管我们希望读者能够完整地阅读本书来理解 Kimball 生命周期的全过程，但是我们在每一章的开始也会强调该章所主要针对的读者，这样您就可以更好地判断哪些内容需要精读，哪些内容可以跳过。希望您的经验和看法使您搭建起自己的框架，这样我们的观点就可以被串联在一起了。阅读完第一章以后，您将会看到在建立一个 DW/BI 系统时必须按照三条并行的路线推进：技术、数据和商业智能应用程序。我们还在每一章开头的“*You Are Here*”图中都指明了这三条路线。尽管这三条路线之间明显会互相影响，但是它们的开发应该以并行方式和异步方式进行。

由于图书的内容必然是按照线性方式进行编排的，所以我们不得不线性地介绍 Kimball 生命周期中的所有步骤，就像这些步骤是以某种固定的次序发生一样。希望在读完本书以后，您能够想象出这些步骤在现实世界中是具有更现实、更复杂的关系的。读完本书以后，当您的项目进行到某个特定阶段时，请再返回到相应的章节并重新仔细地阅读其内容。这也就是为什么我们将其称作生命周期工具箱的原因。

这本书与第一版有什么不同

与第一版相比，第二版生命周期工具箱的内容有了明显的更新和重组。前三章可以帮助您理解整个 Kimball 生命周期过程，并确保您的工作已经满足继续向前推进所必需的条件。然后，我们努力使有关复杂架构的讨论更具有实际价值，并且更加紧密地将架构和 Kimball 生命周期中各项活动的次序联系起来。在第 4 章我们细致地讲述了 DW/BI 系统的完整架构，包括从原始数据的提取到最后将数据显示到业务用户屏幕上的整个过程。在第 5 章我们讲述了怎样为这个技术架构创建一个详细的计划和如何进行产品选择。然后，从第 6 章到第 12 章，我们沿着三条主线(数据库设计、ETL 系统和 BI 应用程序)系统地展开，先从概念上进行介绍，随后又从物理上进行了介绍。在最后的两章里，讲述了如何将这个精心设计的庞然大物部署到实际业务环境中，并且讨论了在第一轮实现之后怎样扩展 DW/BI 系统。

希望我们对数据仓库和商业智能的热情能够贯穿本书的始终。DW/BI 所面临的挑战令人着迷，而且也值得探索。毋庸置疑的是，数年之后当 DW/BI 厂商对他们的产品加以更新时，所有以前的东西都会被取代，相应的名称也会改变。但是我们的任务仍然不会改变，那就是：为业务用户提供数据和分析结果，使他们能够更好地进行业务决策。

目 录

第 1 章 Kimball 生命周期导论 1	
1.1 生命周期的历史..... 1	
1.2 生命周期里程碑..... 3	
1.2.1 项目/项目群规划..... 3	
1.2.2 项目/项目群管理..... 4	
1.2.3 业务需求定义..... 4	
1.2.4 技术路线..... 4	
1.2.5 数据路线..... 5	
1.2.6 商业智能应用路线..... 6	
1.2.7 部署..... 6	
1.2.8 维护..... 6	
1.2.9 增长..... 6	
1.3 使用生命周期图..... 7	
1.4 生命周期导航帮助..... 7	
1.5 生命周期相关术语简介..... 8	
1.5.1 数据仓库与商业智能..... 8	
1.5.2 ETL 系统..... 9	
1.5.3 业务过程维度模型..... 9	
1.5.4 商业智能应用程序..... 10	
1.6 小结..... 11	
第 2 章 项目/项目群的启动与管理 13	
2.1 确定项目..... 14	
2.1.1 评估 DW/BI 项目的准备 就绪情况..... 14	
2.1.2 弥补不足并确定下步工作..... 15	
2.1.3 确定初步范围和章程..... 18	
2.1.4 建立商业报告和合理性 证明..... 22	
2.2 项目规划..... 26	
2.2.1 确立项目标识..... 26	
2.2.2 项目人员配备..... 26	
2.2.3 制定项目计划..... 32	
2.2.4 制定沟通计划..... 35	
2.3 项目管理..... 37	
2.3.1 召开项目团队启动会议..... 38	
2.3.2 监控项目状态..... 39	
2.3.3 维护项目计划..... 40	
2.3.4 整理项目文档..... 40	
2.3.5 范围管理..... 40	
2.3.6 期望管理..... 42	
2.3.7 辨识项目陷入困境的征兆..... 42	
2.4 项目群管理..... 43	
2.4.1 确立管理职责和管理过程..... 43	
2.4.2 将数据管理员的地位 提升到企业层..... 44	
2.4.3 利用高效的方法和 架构最优方法..... 45	
2.4.4 进行定期评估..... 45	
2.4.5 沟通, 沟通, 沟通..... 46	

2.5	小结	46	3.8	项目层需求的调整	75
2.6	管理工作和降低风险	46	3.8.1	走近项目层	75
2.7	质量保证	46	3.8.2	为项目需求访谈做准备	76
2.8	关键角色	47	3.8.3	进行访谈	77
2.9	关键提交内容	47	3.8.4	深入调查数据	79
2.10	作量估计	47	3.8.5	审查访谈结果	79
2.11	站资源	48	3.8.6	准备和发布项目提交材料	80
2.12	任务列表	48	3.8.7	协商下一步工作并结束 本轮访谈	80
第3章	收集业务需求	51	3.9	应对富有挑战性的受访者	80
3.1	需求定义的各种方法	53	3.9.1	受过打击的用户	81
3.1.1	个别访谈 VS 集体 座谈会	53	3.9.2	超负荷的用户/替换用户	81
3.1.2	收集业务需求应避免 使用的方法	54	3.9.3	昏昏欲睡的用户	81
3.2	访谈准备	55	3.9.4	过分热心的用户	81
3.2.1	确定访谈小组	55	3.9.5	自以为无所不知的用户	82
3.2.2	研究业务机构	56	3.9.6	一窍不通的用户	82
3.2.3	选择受访者	57	3.9.7	用户的缺位	82
3.2.4	设计访谈问卷	58	3.10	小结	82
3.2.5	确定访谈时间表	60	3.11	管理工作和降低风险	83
3.2.6	通知受访者做好准备	61	3.12	保证质量	83
3.2.7	访谈中的基本规则综述	63	3.13	关键角色	83
3.3	进行访谈	65	3.14	关键提交内容	84
3.3.1	项目群层面的业务访谈	66	3.15	工作量估计	84
3.3.2	项目群层面上的 IT 访谈	67	3.16	网站资源	84
3.3.3	项目群合规性/安全性访谈	67	3.17	任务列表	85
3.4	总结访谈	67	第4章	技术架构介绍	87
3.4.1	确定项目群成功的标准	67	4.1	架构的价值	88
3.4.2	致谢并告辞	68	4.2	技术架构综述	89
3.5	审查访谈结果	69	4.2.1	从源系统到用户桌面 的流程	91
3.6	准备和发布项目群需求文档	70	4.2.2	常见架构特征	91
3.6.1	访谈书面说明	70	4.2.3	DW/BI 架构评估	94
3.6.2	项目群需求调查结果文档	71	4.3	后台架构	94
3.7	区分业务优先次序和 商定下步工作	73	4.3.1	ETL 一般性需求	95
3.7.1	以优先级的审查和 确定结束会议	73	4.3.2	创建与购买	95
3.7.2	结束本轮访谈	74	4.3.3	后台 ETL 流程	95
			4.3.4	源系统	97
			4.3.5	抽取	100

4.3.6	清洗和一致化	100	4.10	小结	137
4.3.7	提交	101	第 5 章	创建架构计划和选择产品	139
4.3.8	ETL 管理服务	101	5.1	创建架构	139
4.3.9	其他后台服务和趋势	102	5.1.1	架构开发过程	140
4.3.10	ETL 数据存储	102	5.1.2	设计应用程序架构计划	142
4.3.11	ETL 元数据	103	5.2	选择产品	149
4.3.12	后台总结	104	5.2.1	保留一个业务关注点	149
4.4	呈现服务器架构	105	5.2.2	主要 DW/BI 评估领域	149
4.4.1	信息方面的业务需求	105	5.2.3	评估供选方案并挑选产品	150
4.4.2	细节原子数据	106	5.2.4	后台和呈现服务器方 面的考虑事项	158
4.4.3	聚集	106	5.2.5	前台考虑事项	160
4.4.4	呈现服务器设计规定	108	5.2.6	管理元数据	161
4.4.5	调整呈现服务器架构	109	5.2.7	任命元数据管理员	162
4.4.6	机构考虑事项	109	5.2.8	创建元数据策略	162
4.4.7	呈现服务器元数据	110	5.3	保护系统安全	163
4.4.8	呈现服务器总结	110	5.3.1	保护硬件和操作系统 的安全	164
4.5	前台架构	111	5.3.2	保护开发环境的安全	164
4.5.1	BI 应用程序类型	112	5.3.3	保护网络安全	165
4.5.2	BI 管理服务	112	5.3.4	用户验证	167
4.5.3	BI 数据存储	118	5.3.5	数据保护	168
4.5.4	桌面工具架构方法	120	5.3.6	监视使用情况和保证 合规性	171
4.5.5	BI 元数据	120	5.3.7	备份和恢复计划	171
4.5.6	前台总结	121	5.3.8	创建底层设施图	172
4.6	底层设施	121	5.4	安装硬件和软件	174
4.6.1	底层设施驱动因素	122	5.5	小结	175
4.6.2	后台和呈现服务器底 层设施因素	122	5.6	管理工作和降低风险	175
4.6.3	并行处理硬件架构	124	5.7	质量保证	176
4.6.4	硬件性能推进器	127	5.8	关键角色	176
4.6.5	数据库平台因素	128	5.9	关键提交内容	177
4.6.6	前台底层设施要素	130	5.10	工作量估计	177
4.6.7	底层设施总结	132	5.10.1	创建架构计划	177
4.8	元数据	132	5.10.2	选择产品	177
4.8.1	元数据集成的价值	132	5.10.3	元数据	177
4.8.2	元数据集成的供选方案	133	5.10.4	安全性	178
4.8.3	元数据总结	134	5.11	网站资源	178
4.9	安全性	134			
4.9.1	安全方面的弱点	135			
4.9.2	安全性总结	137			

5.12 任务列表.....	178	第7章 维度模型设计.....	227
第6章 维度建模介绍.....	183	7.1 建模过程综述.....	227
6.1 使用维度建模的场合.....	184	7.2 组建团队.....	229
6.1.1 什么是维度建模.....	184	7.2.1 确定参加设计的人员.....	229
6.1.2 怎样进行规范化建模?.....	185	7.2.2 回顾需求.....	231
6.1.3 维度建模的好处.....	186	7.2.3 使用建模工具.....	231
6.2 维度建模入门.....	187	7.2.4 确立命名约定.....	233
6.2.1 事实表.....	187	7.2.5 为源数据调查和数据 探查做准备.....	234
6.2.2 维度表.....	189	7.2.6 获取场所和用品.....	234
6.2.3 四步维度设计过程.....	193	7.3 再论四步建模过程.....	234
6.3 企业数据仓库总线架构.....	194	7.3.1 第1步:选择业务过程.....	235
6.3.1 规划危机.....	194	7.3.2 第2步:声明粒度.....	235
6.3.2 总线架构.....	195	7.3.3 第3步:识别维度.....	236
6.3.3 价值链的意义.....	196	7.3.4 第4步:识别事实.....	237
6.3.4 通用矩阵的常见问题.....	197	7.4 设计维度模型.....	237
6.3.5 坚持使用一致性维度.....	198	7.4.1 建立高层维度模型.....	238
6.4 对维度的深入讨论.....	198	7.4.2 开发详细的维度模型.....	240
6.4.1 日期和时间.....	199	7.4.3 审查和验证模型.....	250
6.4.2 退化维.....	201	7.4.4 设计文档定稿.....	251
6.4.3 缓慢变化维.....	202	7.5 拥抱数据管理.....	252
6.4.4 角色扮演维.....	205	7.6 小结.....	253
6.4.5 杂项维.....	206	7.7 管理工作和降低风险.....	253
6.4.6 雪花型和支架.....	208	7.8 保证质量.....	254
6.4.7 处理层次结构.....	211	7.9 关键角色.....	254
6.4.8 使用桥接表的多值维.....	212	7.10 关键提交内容.....	254
6.5 更多关于事实的讨论.....	214	7.11 工作量估计.....	255
6.5.1 三个基本粒度.....	215	7.12 网站资源.....	255
6.5.2 不同粒度的事实及其分配.....	217	7.13 任务列表.....	255
6.5.3 多种货币和度量单位.....	219	第8章 物理数据库设计与性能规划... 257	
6.5.4 无事实的事实表.....	221	8.1 制定标准.....	258
6.5.5 合并事实表.....	221	8.1.1 遵守命名约定.....	259
6.6 有关维度建模的错觉和误区... 222		8.1.2 为空还是不为空.....	259
6.6.1 将关注点集中在部门报表 上导致的错误观点.....	222	8.1.3 设置登台表.....	259
6.6.2 提前汇总导致的错误观点... 223		8.1.4 制定文件位置标准.....	260
6.6.3 过于重视规范化导致 的错误观点.....	224	8.1.5 对用户访问的表使用 代用名或者视图.....	260
6.7 小结.....	225	8.1.6 主键.....	261

8.1.7 外键	262	8.13 保证质量	285
8.2 设计物理数据模型	263	8.14 关键角色	285
8.2.1 设计物理数据结构	263	8.15 关键提交内容	285
8.2.2 确定源到目标的映射	264	8.16 工作量估计	285
8.2.3 星型 VS 雪花型	265	8.17 网站资源	286
8.2.4 使用数据建模工具	266	8.18 任务列表	286
8.2.5 进行初步的规模估计	267		
8.3 创建开发数据库	268	第 9 章 抽取、转换和装载介绍	289
8.4 设计处理数据存储	269	9.1 归拢需求	290
8.5 设计初始索引方案	270	9.1.1 业务需求	290
8.5.1 索引和查询策略综述	270	9.1.2 合规性	290
8.5.2 为维度表建立索引	272	9.1.3 数据质量	291
8.5.3 为事实表建立索引	272	9.1.4 安全性	291
8.5.4 为装载数据(loads)		9.1.5 数据集成	291
建立索引	273	9.1.6 数据等待时间	292
8.5.5 为 OLAP 建立索引	273	9.1.7 存档和沿袭	292
8.5.6 在装载之后分析表和索引	273	9.1.8 用户提交界面	292
8.6 设计 OLAP 数据库	274	9.1.9 可用的技能	292
8.6.1 OLAP 数据粒度和		9.1.10 遗留许可证	293
深入钻取	274	9.2 ETL 系统的 34 个子系统	293
8.6.2 完善 OLAP 维度	274	9.3 抽取数据	293
8.6.3 定义 OLAP 计算	275	9.3.1 子系统 1——数据探查	294
8.7 建立测试数据库	276	9.3.2 子系统 2——变化数据	
8.8 设计聚集	276	捕捉系统	294
8.8.1 确定如何聚集	276	9.3.3 子系统 3——抽取系统	296
8.8.2 确定聚集的内容	277	9.4 数据的清洗和一致化	297
8.8.3 维护聚集	278	9.4.1 改进数据质量文化和过程	297
8.8.4 完成索引	279	9.4.2 子系统 4——数据清洗	
8.9 设计和构建数据库实例	279	系统	298
8.9.1 内存	280	9.4.3 子系统 5——错误事件	
8.9.2 块大小	280	模式	299
8.9.3 保存数据库构建脚本		9.4.4 子系统 6——审计维装	
和参数文件	280	配器	300
8.10 设计物理存储结构	281	9.4.5 子系统 7——重复数据	
8.10.1 计算表和索引的大小	281	删除系统	301
8.10.2 设计分区方案	281	9.4.6 子系统 8——一致化系统	302
8.10.3 设置存储	282	9.5 向呈现层交付数据	303
8.11 小结	284	9.5.1 子系统 9——缓慢变化	
8.12 管理工作和降低风险	284	维管理器	303

9.5.2	子系统 10——代理键生成器	306	9.6.8	子系统 29——沿袭和依赖分析器	320
9.5.3	子系统 11——层次管理器	306	9.6.9	子系统 30——问题自动调整系统	321
9.5.4	子系统 12——专用维度管理器	307	9.6.10	子系统 31——并行/管道系统	321
9.5.5	子系统 13——事实表构建器	308	9.6.11	子系统 32——安全系统	322
9.5.6	子系统 14——代理键管道	310	9.6.12	子系统 33——合规性管理器	322
9.5.7	子系统 15——多值维度桥接表构建器	311	9.6.13	子系统 34——元数据知识库管理器	323
9.5.8	子系统 16——延迟到达数据处理器	312	9.7	实时的意义	323
9.5.9	子系统 17——维度管理理系统	312	9.7.1	实时的分类	323
9.5.10	子系统 18——事实提供系统	313	9.7.2	实时的权衡	325
9.5.11	子系统 19——聚集构建器	313	9.7.3	呈现服务器上的实时分区	326
9.5.12	子系统 20——OLAP 多维数据集构建器	314	9.8	小结	327
9.5.13	子系统 21——数据传播管理器	314	第 10 章	设计和开发 ETL 系统	329
9.6	管理 ETL 环境	315	10.1	ETL 过程综述	329
9.6.1	子系统 22——作业调度器	315	10.2	启动	330
9.6.2	子系统 23——备份系统	316	10.3	设计 ETL 计划	331
9.6.3	子系统 24——恢复和重启系统	317	10.3.1	步骤 1——制订高层计划	331
9.6.4	子系统 25——版本控制系统	318	10.3.2	步骤 2——选择 ETL 工具	332
9.6.5	子系统 26——版本迁移系统	318	10.3.3	步骤 3——制定默认策略	333
9.6.6	子系统 27—— workflow 监视器	319	10.3.4	步骤 4——由目标表向下钻取	334
9.6.7	子系统 28——排序系统	320	10.3.5	设计 ETL 说明文档	336
			10.3.6	开发沙盒源系统	337
			10.4	设计一次性的历史装载处理	338
			10.4.1	步骤 5——使用历史数据填充维度表	339
			10.4.2	步骤 6——执行事实表历史装载	346
			10.5	设计增量 ETL 处理过程	352

10.5.1	步骤 7——维度表增量 处理	352	11.4	通过 BI 门户导航应用 程序	385
10.5.2	步骤 8——事实表增量 处理	355	11.4.1	考虑密度	387
10.5.3	步骤 9——聚集表和 OLAP 装载	359	11.4.2	基于业务过程的导航 结构	387
10.5.4	步骤 10——ETL 系统操作 和自动化	360	11.4.3	附加门户功能	388
10.6	小结	362	11.4.4	应用程序界面供选 方案	389
10.7	管理工作和降低风险	362	11.5	小结	389
10.8	保证质量	363	第 12 章	设计和开发商务智能应用 程序	391
10.9	关键角色	363	12.1	商业智能应用程序资源 规划	392
10.10	关键交付内容	363	12.1.1	BI 应用程序开发人员 的角色	392
10.11	工作量估计	363	12.1.2	谁来完成商业智能 工作	392
10.12	网站资源	364	12.1.3	生命周期时间安排	392
10.13	任务列表	364	12.2	BI 应用程序规范	393
第 11 章	商务智能应用程序介绍	367	12.2.1	创建应用程序标准和 模板	393
11.1	商业智能应用程序的重 要性	367	12.2.2	确定初始应用程序集	396
11.2	商业智能分析周期	369	12.2.3	制定详细的应用程序 规范	398
11.2.1	第一阶段: 监视活动	370	12.2.4	设计导航框架和门户	401
11.2.2	第二阶段: 识别异常	370	12.2.5	审查以及确认应用程序 和模型	402
11.2.3	第三阶段: 确定构成原因 的因素	370	12.2.6	与业务人士一同审查	403
11.2.4	第四阶段: 模型供选 方案	371	12.3	BI 应用程序开发	403
11.2.5	第五阶段: 采取行动并 跟踪结果	371	12.3.1	准备应用程序开发	403
11.2.6	分析周期的更多意义	371	12.3.2	构建应用程序	405
11.3	商业智能应用程序的类型	372	12.3.3	应用程序和数据的测试 和验证	411
11.3.1	直接访问查询和报表 工具	372	12.3.4	完成文档	412
11.3.2	标准报表	377	12.3.5	部署计划	412
11.3.3	分析性应用程序	378	12.4	BI 应用程序维护	412
11.3.4	仪表板和记分卡	379	12.5	小结	413
11.3.5	运营商业智能	381	12.6	管理工作并降低风险	413
11.3.6	数据挖掘	382			