

Informatica PowerCenter 权威指南

杜绍森 著



中国工信出版集团



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
<http://www.phei.com.cn>

Informatica PowerCenter

权威指南

杜绍森 著 /

电子工业出版社
Publishing House of Electronics Industry
北京•BEIJING

内 容 简 介

在大数据时代，掌控数据首先需要掌握数据的处理能力。俗话说：“工欲善其事，必先利其器。”
Informatica PowerCenter 作为业界广泛使用的数据处理工具之一，被全球多数大型机构、组织认可并采用。

本书全面地介绍了 Informatica PowerCenter 的主要功能及高级特性。

本书分为 3 个部分：第一部分为基础篇，包括第 1~4 章，系统介绍了 PowerCenter 的基础组件和常用功能，并在其中穿插了大量实践案例；第二部分为高级篇，包括第 5~8 章，系统介绍了 PowerCenter 并行、集群、性能调优和字符集管理等高级内容；第三部分为扩展篇，包括第 9 章，简要介绍了 CDC 的相关知识，PowerCenter 与 SAP、MPP、Hadoop 集成，以及非结构化和半结构化数据处理能力。

本书适合 PowerCenter 的入门者及有一定 PowerCenter 使用经验的用户参考，也可作为各数据仓库、大数据专业培训机构的培训教材。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有，侵权必究。

图书在版编目（CIP）数据

Informatica PowerCenter 权威指南 / 杜绍森著. —北京：电子工业出版社，2015.9

ISBN 978-7-121-27045-1

I. ①I… II. ①杜… III. ①企业管理—数据管理 IV. ①F270.7

中国版本图书馆 CIP 数据核字（2015）第 203208 号

责任编辑：徐津平

特约编辑：赵树刚

印 刷：北京京师印务有限公司

装 订：北京京师印务有限公司

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本：787×980 1/16 印张：22.75 字数：475 千字

版 次：2015 年 9 月第 1 版

印 次：2015 年 9 月第 1 次印刷

印 数：3000 册 定价：69.00 元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，
联系及邮购电话：（010）88254888。

质量投诉请发邮件至 zlts@phei.com.cn，盗版侵权举报请发邮件至 dbqq@phei.com.cn。

服务热线：（010）88258888。

推 荐 序

犹豫了很久，以我现今的职位给作者写序，是否有些自吹自擂？但读完书稿，我决定了：一本好书，介绍一个好产品，既然与我受用，何不推荐给更多的人呢？

“IT”是信息（Information）和技术（Technology）的缩写，它的发展不过三十多年的时间。在 IT 发展前期的大部分时间里，其所有进步大部分体现在“T”上，例如计算性能、存储容量、网络拓展及打印效果等。直到近些年“大数据时代”的出现，人们才开始了对于信息数据，也就是“I”的关注。我相信，这个变化是 IT 发展的必然，是一个破茧成蝶的过程，并且这个关注也一定会延续很多年。

同样，正是因为“大数据”日益深入人心，企业的 IT 规划和发展越来越与“大数据”相关联，PowerCenter 才得以从一个 IT 人员得心应手的工具，蜕变为大数据应用的一个重要环节。记得在 2014 年国务院工业和信息化部颁布的大数据白皮书中，就明确地将“数据准备”定义为大数据发展的第一个环节。由此，作为在数据集成领域里历年排名第一的 PowerCenter，也就承担起了“帮助企业实现大数据应用的第一步”的重要使命。

本书前 6 章中规中矩，如同一本深入浅出的教科书，将具备一些基本 IT 知识的人士引进数据迁移的奇妙世界，加上作者风趣的调侃，学来丝毫不觉得枯燥单调。第 7 章开始是实战描述，实际上是一系列的应用经验分享，这些宝贵的经验之谈，可以让初学者在未来的实践中少走弯路，还可以将本书作为可以随时受教的参考书。更值得一提的是，不同于普通的产品手册，本书作者以其十几年的理论研究和教育培训，以及主导或参与诸多中外企业“数据集成项目”实施的经验，将 PowerCenter 的很多功能细节描述得淋漓尽致。本书对于有意进行 ETL 教学的教育培训机构，不失为一本经典的教材；而对于有意培养自己成为 ETL 应用高手的 IT 人士，则是一本有益而又有趣的读物。

曾经有不少朋友问我：当成了 ETL 的行家里手以后，下一个职业目标会有哪些发展方

向？所以，我想借此序的一角，分享一些我的认知，供大家参考。

第一，云数据集成和管理。根据 IDC 的预测，2017 年全球 SaaS 和云软件模式将占软件开支的 1/6。越来越多的云应用系统承诺并交付更简单、更快捷和更智能的业务营运方法，所以，掌握云数据集成，会让你在不可阻挡的云服务趋势下游刃有余。

第二，下一代数据洞察。不同于第一代商业智能（BI）对展示形态和分析过去的重视，大数据时代的数据洞察，更加关注数据的质量而不是数据的展现形式，更重视预测未来的行为模式而非过去的行为分析。所以，要想成为大数据分析专家，你必须懂得数据质量管理和服务前瞻性的分析。当然，保障分析结果正确的前提是确保数据的统一性、完整性，并找到数据的关联性。

第三，数据治理。大数据时代，越来越多的企业将数据纳入其固定资产；在金融和医疗行业，数据相关的合规性成为政府监管的重要指标；为了应对客户需求和市场业务模式的变化，许多企业开始考虑应用整合和迁移……这些巨大的变化，不断催生出数据治理的高手，他们必须在行业规范、企业应用系统、数据的关联性和安全性方面具备独特的技能。因此，了解行业特性、行业应用，使之与数据集成相结合，便成为你进行数据治理的更高境界。

近年来，关于大数据的定义一直在调整，而大数据应用的目标却始终没变，那就是：发现数据价值，帮助企业降低成本并实现业务创新。在过去短短的两三年里，中国作为自然的大数据国家，已经在大数据的理论研究和实际应用方面取得了巨大的进展。大数据的应用会推动各行各业诞生越来越多的数据科学家，那是行业知识和数据治理兼备的卓越人才。IT 的发展已经实现了由“计算机科学”向“数据科学”的转换，近年来，“数据科学”又开始向行业应用进行大规模迁移。所以，数据科学家既是数据价值的挖掘者，更是行业产品和流程的创新者，他们的价值不是向企业的高管提供分析报表，他们本身就是企业的高管，他们在用数据作为依据，实现企业面向客户、市场、产品和流程方面的创新。

千里之行，始于足下。与各位读者共勉。

Informatica 大中国区总经理 王晨杰

自序

初识 Informatica，大概是在十四五年前的一个偶然的机会，公司接到一个叫作决策支持系统（DSS）的项目。尽管当时作为工程师和客户一起整理了项目的需求，完成了需求的确认和签字，但我现在几乎无法记起任何关于需求的内容，对项目实施过程的某些环节却仍然记忆犹新。项目开始时，公司安排了两位工程师参与 ETL（Extract Transformation Load），一位是我，另一位与我现在仍是同事，我们当时使用了一个叫作 PowerMart 的工具，版本是 5.1。这就是我和 Informatica 的第一次亲密接触。当时 Informatica 的总代理也是曾经大名鼎鼎的 Sybase，据说我们的这个项目是 Informatica 进入中国后的第二个项目。

从此，我开始了自己漫长的 Informatica 之路。当时我所就职的公司敏锐地察觉到数据仓库/商业智能是未来的趋势之一，开始着手准备发展数据仓库方面的业务。当年有个著名的第三方调研机构，叫 IDC。我所就职的公司通过查阅 IDC 报告，发现 Informatica 是当时 ETL 市场份额最高的公司，于是果断决定采用这个工具作为自己的数据仓库的开发平台。当年公司的果断、决心，让我至今想起，仍非常钦佩。在 IT 人才严重短缺的年代，虽然年纪很轻、经验不足，但我还是作为经营分析项目的项目经理、技术经理等开始了自己的数据管理生涯。

此后一段时间，我在不经意间进入悲催的计费岁月，加班、加班、通宵、通宵……每个项目都以几年来计算，历经两个完整的移动计费项目，在此期间认识了很多好朋友和师长。这是我与 Informatica 断绝联系的几年，也是在技术方面拓展能力的几年。

我与 Informatica 的缘分并没有结束。有一天，原来的同事告诉我 Informatica 在招售前工程师，我就毅然决定去应聘，满足自己转向咨询领域的一点梦想，后来发现售前和咨询还是有所区别的，这是后话。因此，8 年前我加入了 Informatica 中国，开始成为一名专职的售前工程师。当年的 Informatica 只有这一个产品，人不是很多。我仍清晰地记着当时的

版本为 PowerCenter 8.1.1。现在的 Informatica 已经与早期差别非常大了，但是很多人还是习惯把 Informatica 的数据集成产品 PowerCenter 叫作 Informatica。

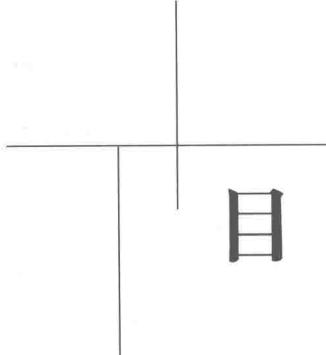
为什么要写一本关于 PowerCenter 的书呢？其实我内心一直有这样的冲动。PowerCenter 是一个非常好的产品，在国内也有近千家用户，有大量的开发者和管理者，随着大数据的推广，还有大量的后继者会陆续开始 PowerCenter 之旅。一本中文材料会帮助所有的用户更加快速、更加全面地了解 PowerCenter，充分利用自己在 PowerCenter 方面的投资。这个冲动持续了很久，包括期间陆续说服数位同事参与，但是大家都有繁忙的本职工作，一次一次被耽搁。直到 2015 年春节前，内心的冲动促使我开始动笔了。前两天坐在我对面的北区销售总监说，我写这本书像写回忆录，想想确实有道理，我的自序也是回忆录的样子，希望大家能够谅解。

如何写好这本书？这也是我非常纠结的一个问题。什么样的深度？适合什么样的人群？是否有读者愿意来读？如何帮助读者了解 PowerCenter？和同事们讨论过很多次，似乎还是没有下定决心。当我下笔的时候，尤其是写到 50 多页的时候，我觉得我已经坚定了这本书的方向：不求全面，但求让读者快速地掌握 PowerCenter；不求精深，但求将最常用的功能展现给读者；不求华丽的词藻，但求读者能读懂。

PowerCenter 是什么？它是 ETL 工具。什么是 ETL？大数据及数据仓库 70% 左右的工作都在做 ETL，在 Gartner 报告中它被划为 Data Integration 产品。Informatica 也曾定义自己是一家 Data Integration 公司。我是这样解释 ETL 工具的：它是 Data Integration 产品在数据仓库、大数据项目中的一个应用场景，它同时还有其他的应用场景，比如数据交换、数据安全，这些也是 PowerCenter 在后期的扩展。

希望本书能够成为分享我这些年掌握的 PowerCenter 相关知识的一个载体，成为初学者的入门教材，成为有经验者的开发人员的一本参考书。

杜绍森
2015 年 8 月



目 录

第 1 章 PowerCenter Hello World 世界	1
1.1 Informatica Hello World.....	1
1.2 PowerCenter 架构和客户端简介	3
1.2.1 PowerCenter 架构	3
1.2.2 PowerCenter 客户端	5
1.3 PowerCenter Hello World	7
第 2 章 PowerCenter 基础组件	27
2.1 Source.....	27
2.1.1 数据库源	28
2.1.2 文本文件源	30
2.2 Target.....	33
2.2.1 数据库目标	33
2.2.2 文本文件目标	34
2.3 Expression 表达式	35
Expression 中的变量端口 (Variable Port)	40

2.4	Filter	41
2.5	Source Qualifier	43
2.5.1	Source Qualifier 的作用	43
2.5.2	数据库数据源的 Source Qualifier	44
2.5.3	Source Qualifier 自定义 SQL	47
2.5.4	Source Qualifier 复杂关联	48
2.6	Sorter	49
2.7	Joiner	51
2.7.1	关联类型	52
2.7.2	Sorted Joiner	54
2.7.3	Joiner 的独特作用	55
2.7.4	自关联 (Self-Join)	56
2.8	Lookup	57
2.8.1	Lookup Caching enabled	59
2.8.2	非连接的 Lookup	61
2.8.3	Lookup SQL Override	63
2.8.4	共享 Lookup Cache	65
2.8.5	Dynamic Lookup	65
2.8.6	Lookup、Source Qualifier 和 Joiner 的对比	69
2.9	Stored Procedure	70
2.9.1	Connected Stored Procedure	70
2.9.2	Unconnected Stored Procedure	72
2.9.3	Pre- or Post-Session Stored Procedure	74
2.10	Union	76
2.11	Transaction Control	78
2.11.1	Transaction Control 有效性问题	79
2.11.2	Transaction Control 组件	80

2.12 Sequence.....	80
2.12.1 Sequence 的常规用法.....	80
2.12.2 共享 Sequence.....	82
2.12.3 可重用的 Sequence.....	83
2.13 Aggregator.....	84
2.13.1 条件聚合	85
2.13.2 使用 Aggregator 进行行列转换	86
2.14 Rank	88
2.15 Update strategy.....	90
2.15.1 Treat source rows as 属性的使用	91
2.15.2 Update strategy 使用	93
2.15.3 如何实现 Update else Insert	94
2.15.4 Update Stagety 案例: 缓慢变化维	98
2.16 SQL Transformation.....	104
2.16.1 Script Mode	104
2.16.2 Static Query Mode.....	106
2.16.3 Dynamic Query Mode	108
2.17 Java Transformation	109
2.17.1 Java Transformation 简介	109
2.17.2 Passive Java Transformation.....	114
2.17.3 Active Java Transformation.....	121
2.17.4 常见错误说明	123
2.18 Normalizer.....	124
2.19 Router	126
2.20 Custom Transformation.....	128
2.21 HTTP Transformation	129
2.22 XML 组件组	132

2.23 Transformation 中的一些概念	135
2.23.1 Connect 与 Unconnect	135
2.23.2 Active 与 Passive	136

第 3 章 Workflow 执行、监控 138

3.1 Session	139
3.1.1 Reusable Session	139
3.1.2 非 Reusable Session	141
3.2 最简单、最常用的 Workflow	143
3.2.1 并行执行	143
3.2.2 串行执行	144
3.2.3 调度	146
3.3 Worklet	147
3.4 Command	148
3.5 Control	150
3.6 发送 E-mail	151
3.6.1 配置发送 E-mail	151
3.6.2 在 Workflow 中使用 E-mail	151
3.7 Event Tasks	155
3.7.1 用户自定义事件使用	156
3.7.2 预定义事件使用	158
3.8 Timer	159
3.9 Decision	159
3.10 Assignment	160

第 4 章 常用功能汇集

163

4.1	Debugger	163
4.2	Mapplet/Reusable Transformation	165
4.2.1	Reusable Transformation.....	165
4.2.2	Mapplet.....	167
4.3	使用 Shortcut	169
4.3.1	Local Shortcut	170
4.3.2	Global Shortcut.....	171
4.4	Session 相关属性.....	173
4.4.1	Properties Tab 相关属性	173
4.4.2	Config Object Tab 相关属性.....	174
4.5	参数和变量	176
4.5.1	Mapping 参数.....	176
4.5.2	Mapping 变量.....	180
4.5.3	系统/Session 参数与变量	184
4.5.4	Workflow/Worklet 变量	189
4.5.5	Local 变量 (Local Variables)	191

第 5 章 PowerCenter 高级应用

193

5.1	任务分区 (Partition)	193
5.1.1	Database Partitioning.....	196
5.1.2	Hash Partitioning	201
5.1.3	Key Range Partitioning	204
5.1.4	Pass Through Partitioning.....	205
5.1.5	Round-Robin Partitioning.....	211
5.2	内存管理	214

5.2.1 DTM 内存	215
5.2.2 Transformation Cache	216
5.3 网格计算	219
5.3.1 Grid 架构	219
5.3.2 Grid 负载均衡	221
5.3.3 Grid 与任务分区 (Partition)	224
5.4 高可用性 (HA)	227
5.4.1 PowerCenter 自带的 HA 方案	228
5.4.2 依托第三方厂商的 HA 方案	229
5.4.3 两种 HA 方案对比	230
5.5 Web Service 应用	230
5.5.1 Web Service Hub	231
5.5.2 Web Service 调度/监控接口	232
5.5.3 Web Service Provider	234
5.5.4 Web Service Consumer	246
5.6 Pushdown Optimization	251
5.6.1 Pushdown 优化是什么	252
5.6.2 Pushdown 优化类型	252
5.7 版本控制及部署	256
5.7.1 Check In/Check Out	256
5.7.2 Team-Based 开发的一些有用功能	258
5.7.3 Label 与 Deployment Group	260
5.7.4 复制对象从开发 Repository 到生产 Repository	264

第 6 章 PowerCenter 实战汇总

266

6.1 PowerCenter 字符集	266
6.1.1 Oracle 数据库	267

6.1.2 DB2 字符集	268
6.1.3 AS/400 字符集	268
6.1.4 ODBC 字符集	269
6.1.5 文本文件字符集	270
6.1.6 Repository Service 字符集	271
6.1.7 Integration Service 字符集	272
6.1.8 Data Movement Mode	273
6.2 UNIX ODBC 配置	274
6.2.1 ODBC 常规配置	274
6.2.2 MySQL 社区版 ODBC 配置	276
6.3 使用 Mapping 动态分发文件	277
6.4 超越 EDW，商品自动价格跟踪	279
6.5 pmcmd 命令	283
6.6 pmrep 命令	284
6.7 infasetup 命令	284
6.8 Mapping Architect for Visio	286
6.9 MX View 语句	293
6.10 PowerCenter 与其他工具集成	294
第 7 章 性能调优	297
7.1 性能调优过程	298
7.2 发现瓶颈	299
7.2.1 定位目标写瓶颈及调优	301
7.2.2 定位源读瓶颈及调优	302
7.2.3 定位 Mapping/Session 瓶颈	303
7.2.4 定位系统瓶颈	305

7.3 Mapping 调优	305
7.3.1 Transformation 优化	305
7.3.2 列级别的优化	310
7.3.3 其他方面的优化	312
7.4 Session 调优	313
7.4.1 内存调优	313
7.4.2 PowerCenter 高级特性支持高性能	313
7.4.3 其他手段	314
7.5 SQL Override 调优	316

第 8 章 PowerCenter Troubleshooting 317

8.1 安装、启动过程的错误	317
8.2 开发过程的错误	319
8.3 Session 运行错误	320
8.4 源读或者目标写的错误	321

第 9 章 PowerCenter 扩展能力 322

9.1 PowerExchange CDC (变化数据捕捉)	322
9.1.1 PowerExchange CDC 的 3 种模式	323
9.1.2 开放数据库 CDC 基本原理	325
9.1.3 CDC 常见的一些讨论	326
9.1.4 CDC Real-Time for Oracle 安装配置 (实例)	327
9.1.5 CDC 定义注册组和添加捕获注册 (实例续)	331
9.1.6 CDC Mapping 开发及运行 (实例)	334

9.2 PowerCenter 与 SAP	336
9.2.1 R/3、mySAP、ECC	337
9.2.2 PowerCenter 与 BW	338
9.3 PowerCenter 与 MPP 数据库	339
9.4 PowerCenter 与 Hadoop	340
9.4.1 接口能力	341
9.4.2 PowerCenter on Hadoop	344
9.5 元数据管理与业务术语管理	345
9.5.1 元数据的血缘分析	346
9.5.2 元数据影响分析	346
9.5.3 业务数据管理	347
9.6 B2B Data Transformation	347

第1章

PowerCenter Hello World 世界

Hello World 程序是学习 IT 的入门术语，是初步了解一种技术的有效途径，因此本书也从 Hello World 开始。除了 PowerCenter Hello World，本书还提供了 Informatica Hello World 的个人版。

1.1 Informatica Hello World

Informatica 曾经是一个 ETL 工具的供应商，也是最好的 ETL 产品的供应商。在开始使用 ETL 工具，并加入到 Informatica 之后的很长一段时间内，我都认为 Informatica 是一家这样的公司。那时的 Informatica 只有一个产品，叫作 PowerCenter。那时，它的 Logo 下有一句话 “The Data Integration Company”。

此后，“幸福的日子”结束了，Informatica 陆续有了很多新产品，包括除 PowerCenter 外的 Data Quality、Data Archive、Master Data Management、Data Masking、Ultra Message、Data Integration Hub 等，这时，其 Logo 下的 “The Data Integration Company” 显著的 IT 蓝依然没有变化。这时的我也在思考什么是 Data Integration，并有了自己对 Data Integration 的诠释。我将 Infomratica 所谓的 Data Integration 诠释为三类数据集成，即下游集成、中游集成和上游集成。

下游集成指的是数据仓库，这是典型的数据集成项目，它有一个显著特点：从数据流