

内容全面，涵盖Hadoop技术本身以及与其相关的Hive、HBase、Mahout、Pig、ZooKeeper、Avro、Chukwa等所有子项目

实战性强，为各个知识点精心设计了大量经典的案例，易于理解，可操作性强



陆嘉恒 著

Hadoop in Action

Hadoop 实战



机械工业出版社
China Machine Press

将网格计算、并行计算和虚拟化等技术融为一体的云计算技术已成为我们当下存储和处理海量数据的最佳选择之一。Hadoop的开源、高性能、高容错、跨平台等特点又使其成为架构云计算平台的首选。本书以实践为主，理论与实践相结合，全面阐述了整个Hadoop技术体系，适合读者系统地学习。强烈推荐！

—— Hadoop中文网

目前，国内的很多互联网企业都在使用或正准备使用Hadoop技术，这些企业都面临着一个共同的难题：Hadoop方面的人才难求。Hadoop方面的人才之所以难找，一方面是因为Hadoop在国内应用的时间不长，从业人员并不多；另一方面是因为Hadoop技术本身难以掌握，而且与涉及海量数据处理的实际生产环境密切相关。本书很好地把握住了当前Hadoop从业者的核心需求，不仅理论知识全面，更重要的是包含大量与实际生产环境相结合的案例，极具指导意义。

—— Hadoop用户社区

在全球范围内，已经有数量庞大的大中型互联网公司开始使用Hadoop，国外的Amazon、Facebook、Yahoo!，国内的腾讯、百度、淘宝、阿里巴巴等都是成功应用Hadoop的典范。然而，Hadoop技术本身却极为复杂，而且涉及众多其他的技术，学习门槛比较高。本书从初中级读者的需求出发，以实践为导向，全面而系统地讲解了Hadoop技术本身，以及与之相关的其他各种技术。对于想系统学习Hadoop和想增加实战经验的读者来说，本书不可多得！

—— 51CTO

客服热线：(010) 88378991, 88361066
购书热线：(010) 68326294, 88379649, 68995259
投稿热线：(010) 88379604
读者信箱：hzsj@hzbook.com



华章网站 <http://www.hzbook.com>

网上购书：www.china-pub.com

上架指导：计算机/程序设计

ISBN 978-7-111-35944-9



9 787111 359449

定价：69.00元



Hadoop in Action

Hadoop实战

陆嘉恒 著



机械工业出版社
China Machine Press



本书是一本系统且极具实践指导意义的 Hadoop 工具书和参考书。内容全面，对 Hadoop 整个技术体系进行了全面的讲解，不仅包括 HDFS 和 MapReduce 这两大核心内容，而且还包括 Hive、HBase、Mahout、Pig、ZooKeeper、Avro、Chukwa 等与 Hadoop 相关的子项目的内容。实战性强，为各个知识点精心设计了大量经典的小案例，易于理解，可操作性强。

全书一共 18 章：第 1 章全面介绍了 Hadoop 的概念、优势、项目结构、体系结构，以及它与分布式计算的关系；第 2 章详细讲解了 Hadoop 集群的安装和配置，以及常用的日志分析技巧；第 3 章分析了 Hadoop 在 Yahoo!、eBay、Facebook 和百度的应用案例，以及 Hadoop 平台海量数据的排序；第 4～7 章深入地讲解了 MapReduce 计算模型、MapReduce 应用的开发方法、MapReduce 的工作机制，同时还列出了多个 MapReduce 的应用案例，涉及单词计数、数据去重、排序、单表关联和多表关联等内容；第 8～11 章全面地阐述了 Hadoop 的 I/O 操作、HDFS 的原理与基本操作，以及 Hadoop 的各种管理操作，如集群的维护等；第 12～17 章详细而系统地讲解了 Hive、HBase、Mahout、Pig、ZooKeeper、Avro、Chukwa 等所有与 Hadoop 相关的子项目的原理及使用，以及这些子项目与 Hadoop 的整合使用；第 18 章以实例的方式讲解了常用 Hadoop 插件的使用和 Hadoop 插件的开发。

本书既适合没有 Hadoop 基础的初学者系统地学习，又适合有一定 Hadoop 基础但是缺乏实践经验的读者实践和参考。

封底无防伪标均为盗版

版权所有，侵权必究

本书法律顾问 北京市展达律师事务所

图书在版编目 (CIP) 数据

Hadoop 实战 / 陆嘉恒著. —北京：机械工业出版社，2011.9

(云计算技术系列丛书)

ISBN 978-7-111-35944-9

I. H… II. 陆… III. 数据处理—应用软件—网络编程 IV. TP274

中国版本图书馆 CIP 数据核字 (2011) 第 192844 号

机械工业出版社 (北京市西城区百万庄大街 22 号 邮政编码 100037)

责任编辑：杨绣国 陈佳媛

北京京师印务有限公司印刷

2011 年 10 月第 1 版第 1 次印刷

186mm×240mm·28.75 印张

标准书号：ISBN 978-7-111-35944-9

定价：69.00 元

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

客服热线：(010) 88378991；88361066

购书热线：(010) 68326294；88379649；68995259

投稿热线：(010) 88379604

读者信箱：hzsj@hzbook.com





目 录

前 言

第 1 章 Hadoop 简介 /1

- 1.1 什么是 Hadoop/2
 - 1.1.1 Hadoop 概述 /2
 - 1.1.2 Hadoop 的历史 /2
 - 1.1.3 Hadoop 的功能与作用 /2
 - 1.1.4 Hadoop 的优势 /3
 - 1.1.5 Hadoop 的应用现状和发展趋势 /3
- 1.2 Hadoop 项目及其结构 /3
- 1.3 Hadoop 的体系结构 /6
 - 1.3.1 HDFS 的体系结构 /6
 - 1.3.2 MapReduce 的体系结构 /7

资源分享网
PDG

- 1.4 Hadoop 与分布式开发 /7
- 1.5 Hadoop 计算模型——MapReduce/10
- 1.6 Hadoop 的数据管理 /10
 - 1.6.1 HDFS 的数据管理 /11
 - 1.6.2 HBase 的数据管理 /12
 - 1.6.3 Hive 的数据管理 /15
- 1.7 小结 /17

第2章 Hadoop 的安装与配置 /18

- 2.1 在 Linux 上安装与配置 Hadoop/19
 - 2.1.1 安装 JDK 1.6/19
 - 2.1.2 配置 SSH 免密码登录 /20
 - 2.1.3 安装并运行 Hadoop/21
- 2.2 在 Windows 上安装与配置 Hadoop/23
 - 2.2.1 安装 Cygwin/24
 - 2.2.2 配置环境变量 /24
 - 2.2.3 安装和启动 sshd 服务 /24
 - 2.2.4 配置 SSH 免密码登录 /24
- 2.3 安装和配置 Hadoop 集群 /25
 - 2.3.1 网络拓扑 /25
 - 2.3.2 定义集群拓扑 /25
 - 2.3.3 建立和安装 Cluster /26
- 2.4 日志分析及几个小技巧 /32
- 2.5 小结 /33

第3章 Hadoop 应用案例分析 /35

- 3.1 Hadoop 在 Yahoo! 的应用 /36
- 3.2 Hadoop 在 eBay 的应用 /38
- 3.3 Hadoop 在百度的应用 /40
- 3.4 Hadoop 在 Facebook 的应用 /43
- 3.5 Hadoop 平台上的海量数据排序 /46
- 3.6 小结 /53

第4章 MapReduce 计算模型 /54

- 4.1 为什么要用 MapReduce/55



VIII

- 4.2 MapReduce 计算模型 /56
 - 4.2.1 MapReduce Job/56
 - 4.2.2 Hadoop 中的 Hello World 程序 /56
 - 4.2.3 MapReduce 的数据流和控制流 /64
- 4.3 MapReduce 任务的优化 /65
- 4.4 Hadoop 流 /67
 - 4.4.1 Hadoop 流的工作原理 /68
 - 4.4.2 Hadoop 流的命令 /69
 - 4.4.3 实战案例：添加 Bash 程序和 Python 程序到 Hadoop 流中 /70
- 4.5 Hadoop Pipes/72
- 4.6 小结 /74

第 5 章 开发 MapReduce 应用程序 /75

- 5.1 系统参数的配置 /76
- 5.2 配置开发环境 /78
- 5.3 编写 MapReduce 程序 /79
 - 5.3.1 Map 处理 /79
 - 5.3.2 Reduce 处理 /80
- 5.4 本地测试 /81
- 5.5 运行 MapReduce 程序 /83
 - 5.5.1 打包 /84
 - 5.5.2 在本地模式下运行 /85
 - 5.5.3 在集群上运行 /86
- 5.6 网络用户界面 /87
 - 5.6.1 JobTracker 页面 /87
 - 5.6.2 工作页面 /88
 - 5.6.3 返回结果 /90
 - 5.6.4 任务页面 /93
 - 5.6.5 任务细节页面 /93
- 5.7 性能调优 /94
- 5.8 MapReduce 工作流 /96
 - 5.8.1 将问题分解成 MapReduce 工作 /97
 - 5.8.2 运行相互依赖的工作 /97
- 5.9 小结 /98



第6章 MapReduce 应用案例 /99

- 6.1 单词计数 /100
 - 6.1.1 实例描述 /100
 - 6.1.2 设计思路 /100
 - 6.1.3 程序代码 /101
 - 6.1.4 代码解读 /102
 - 6.1.5 程序执行 /103
 - 6.1.6 代码结果 /103
- 6.2 数据去重 /104
 - 6.2.1 实例描述 /104
 - 6.2.2 设计思路 /105
 - 6.2.3 程序代码 /105
- 6.3 排序 /106
 - 6.3.1 实例描述 /106
 - 6.3.2 设计思路 /107
 - 6.3.3 程序代码 /107
- 6.4 单表关联 /109
 - 6.4.1 实例描述 /109
 - 6.4.2 设计思路 /110
 - 6.4.3 程序代码 /110
- 6.5 多表关联 /113
 - 6.5.1 实例描述 /113
 - 6.5.2 设计思路 /114
 - 6.5.3 程序代码 /114
- 6.6 小结 /116

第7章 MapReduce 工作机制 /117

- 7.1 MapReduce 作业的执行流程 /118
 - 7.1.1 MapReduce 任务的执行总流程 /118
 - 7.1.2 提交作业 /119
 - 7.1.3 初始化作业 /121
 - 7.1.4 分配任务 /123
 - 7.1.5 执行任务 /125
 - 7.1.6 更新任务执行进度和状态 /126
 - 7.1.7 完成作业 /127



- 7.2 错误处理机制 /127
 - 7.2.1 硬件故障 /127
 - 7.2.2 任务失败 /128
- 7.3 作业调度机制 /128
- 7.4 shuffle 和排序 /129
 - 7.4.1 map 端 /130
 - 7.4.2 reduce 端 /131
 - 7.4.3 shuffle 过程的优化 /132
- 7.5 任务执行 /133
 - 7.5.1 推测式执行 /133
 - 7.5.2 任务 JVM 重用 /134
 - 7.5.3 跳过坏记录 /134
 - 7.5.4 任务执行环境 /135
- 7.6 小结 /136

第 8 章 Hadoop I/O 操作 /137

- 8.1 I/O 操作中的数据检查 /138
- 8.2 数据的压缩 /142
 - 8.2.1 Hadoop 对压缩工具的选择 /142
 - 8.2.2 压缩分割和输入分割 /143
 - 8.2.3 在 MapReduce 程序中使用压缩 /143
- 8.3 数据的 I/O 中序列化操作 /144
 - 8.3.1 Writable 类 /144
 - 8.3.2 实现自己的 Hadoop 数据类型 /152
- 8.4 针对 MapReduce 的文件类 /153
 - 8.4.1 SequenceFile 类 /154
 - 8.4.2 MapFile 类 /159
- 8.5 小结 /161

第 9 章 HDFS 详解 /162

- 9.1 Hadoop 的文件系统 /163
- 9.2 HDFS 简介 /165
- 9.3 HDFS 体系结构 /166
 - 9.3.1 HDFS 的相关概念 /166
 - 9.3.2 HDFS 的体系结构 /167



- 9.4 HDFS 的基本操作 /169
 - 9.4.1 HDFS 的命令行操作 /169
 - 9.4.2 HDFS 的 Web 界面 /171
- 9.5 HDFS 常用 Java API 详解 /173
 - 9.5.1 使用 Hadoop URL 读取数据 /173
 - 9.5.2 使用 FileSystem API 读取数据 /174
 - 9.5.3 创建目录 /176
 - 9.5.4 写数据 /177
 - 9.5.5 删除数据 /178
 - 9.5.6 文件系统查询 /178
- 9.6 HDFS 中的读写数据流 /182
 - 9.6.1 文件的读取 /182
 - 9.6.2 文件的写入 /184
 - 9.6.3 一致性模型 /185
- 9.7 HDFS 命令详解 /186
 - 9.7.1 通过 distcp 进行并行复制 /186
 - 9.7.2 HDFS 的平衡 /187
 - 9.7.3 使用 Hadoop 归档文件 /188
 - 9.7.4 其他命令 /190
- 9.8 小结 /194

第 10 章 Hadoop 的管理 /195

- 10.1 HDFS 文件结构 /196
- 10.2 Hadoop 的状态监视和管理工具 /200
 - 10.2.1 审计日志 /200
 - 10.2.2 监控日志 /200
 - 10.2.3 Metrics/201
 - 10.2.4 Java 管理扩展 /203
 - 10.2.5 Ganglia/204
 - 10.2.6 Hadoop 管理命令 /206
- 10.3 Hadoop 集群的维护 /210
 - 10.3.1 安全模式 /210
 - 10.3.2 Hadoop 的备份 /211
 - 10.3.3 Hadoop 的节点管理 /212
 - 10.3.4 系统升级 /214



10.4 小结 /216

第 11 章 Hive 详解 /217

11.1 Hive 简介 /218

11.1.1 Hive 的数据存储 /218

11.1.2 Hive 的元数据存储 /220

11.2 Hive 的基本操作 /220

11.2.1 在集群上安装 Hive/220

11.2.2 配置 Hive/222

11.3 Hive QL 详解 /224

11.3.1 数据定义 (DDL) 操作 /224

11.3.2 数据操作 (DML) /231

11.3.3 SQL 操作 /233

11.3.4 Hive QL 的使用实例 /235

11.4 Hive 的网络 (WebUI) 接口 /237

11.5 Hive 的 JDBC 接口 /238

11.6 Hive 的优化 /241

11.7 小结 /243

第 12 章 HBase 详解 /244

12.1 HBase 简介 /245

12.2 HBase 的基本操作 /245

12.2.1 HBase 的安装 /245

12.2.2 运行 HBase /249

12.2.3 HBase Shell/250

12.2.4 HBase 配置 /254

12.3 HBase 体系结构 /255

12.4 HBase 数据模型 /259

12.4.1 数据模型 /259

12.4.2 概念视图 /260

12.4.3 物理视图 /260

12.5 HBase 与 RDBMS/261

12.6 HBase 与 HDFS/262

12.7 HBase 客户端 /262

12.8 Java API /263



12.9 HBase 编程实例之 MapReduce /270

12.10 模式设计 /273

12.10.1 学生表 /273

12.10.2 事件表 /274

12.11 小结 /275

第 13 章 Mahout 详解 /276

13.1 Mahout 简介 /277

13.2 Mahout 的安装和配置 /277

13.3 Mahout API 简介 /278

13.4 Mahout 中的聚类和分类 /280

13.4.1 什么是聚类和分类 /280

13.4.2 Mahout 中的数据表示 /281

13.4.3 将文本转化成向量 /282

13.4.4 Mahout 中的聚类、分类算法 /283

13.4.5 算法应用实例 /288

13.5 Mahout 应用：建立一个推荐引擎 /292

13.5.1 推荐引擎简介 /292

13.5.2 使用 Taste 构建一个简单的推荐引擎 /292

13.5.3 简单分布式系统下基于产品的推荐系统简介 /294

13.6 小结 /297

第 14 章 Pig 详解 /299

14.1 Pig 简介 /300

14.2 Pig 的安装和配置 /300

14.2.1 Pig 的安装条件 /300

14.2.2 Pig 的下载、安装和配置 /301

14.2.3 Pig 运行模式 /301

14.3 Pig Latin 语言 /304

14.3.1 Pig Latin 语言简介 /304

14.3.2 Pig Latin 的使用 /305

14.3.3 Pig Latin 的数据类型 /307

14.3.4 Pig Latin 关键字 /308

14.4 用户定义函数 /313

14.4.1 编写用户定义函数 /313



- 14.4.2 使用用户定义函数 /315
- 14.5 Pig 实例 /315
 - 14.5.1 Local 模式 /316
 - 14.5.2 MapReduce 模式 /318
- 14.6 Pig 进阶 /319
 - 14.6.1 数据实例 /319
 - 14.6.2 Pig 数据分析 /320
- 14.7 小结 /324

第 15 章 ZooKeeper 详解 /326

- 15.1 ZooKeeper 简介 /327
 - 15.1.1 ZooKeeper 的设计目标 /327
 - 15.1.2 数据模型和层次命名空间 /328
 - 15.1.3 ZooKeeper 中的节点和临时节点 /328
 - 15.1.4 ZooKeeper 的应用 /329
- 15.2 ZooKeeper 的安装和配置 /329
 - 15.2.1 在集群上安装 ZooKeeper /329
 - 15.2.2 配置 ZooKeeper /334
 - 15.2.3 运行 ZooKeeper /336
- 15.3 ZooKeeper 的简单操作 /339
 - 15.3.1 使用 ZooKeeper 命令的简单操作步骤 /339
 - 15.3.2 ZooKeeper API 的简单使用 /340
- 15.4 ZooKeeper 的特性 /343
 - 15.4.1 ZooKeeper 的数据模型 /343
 - 15.4.2 ZooKeeper 会话及状态 /345
 - 15.4.3 ZooKeeper Watches /346
 - 15.4.4 ZooKeeper ACL /346
 - 15.4.5 ZooKeeper 的一致性保证 /347
- 15.5 ZooKeeper 的 Leader 选举 /348
- 15.6 ZooKeeper 锁服务 /348
 - 15.6.1 ZooKeeper 中的锁机制 /349
 - 15.6.2 ZooKeeper 提供的一个写锁的实现 /350
- 15.7 使用 ZooKeeper 创建应用程序 /351
- 15.8 小结 /355



第 16 章 Avro 详解 /356

- 16.1 Avro 简介 /357
 - 16.1.1 模式声明 /358
 - 16.1.2 数据序列化 /362
 - 16.1.3 数据排列顺序 /364
 - 16.1.4 对象容器文件 /365
 - 16.1.5 协议声明 /367
 - 16.1.6 协议传输格式 /368
 - 16.1.7 模式解析 /370
- 16.2 Avro 的 C/C++ 实现 /371
- 16.3 Avro 的 Java 实现 /382
- 16.4 GenAvro (Avro IDL) 语言 /385
- 16.5 Avro SASL 概述 /390
- 16.6 小结 /392

第 17 章 Chukwa 详解 /393

- 17.1 Chukwa 简介 /394
- 17.2 Chukwa 架构 /395
 - 17.2.1 客户端 (Agent) 及其数据模型 /395
 - 17.2.2 收集器 (Collector) 和分离解析器 (Demux) /396
 - 17.2.3 HICC/398
- 17.3 Chukwa 的可靠性 /399
- 17.4 Chukwa 集群搭建 /400
 - 17.4.1 基本配置要求 /400
 - 17.4.2 安装 Chukwa/400
- 17.5 Chukwa 数据流的处理 /407
- 17.6 Chukwa 与其他监控系统比较 /408
- 17.7 小结 /409

第 18 章 Hadoop 的常用插件与开发 /411

- 18.1 Hadoop Studio 简介和使用 /412
 - 18.1.1 Hadoop Studio 的安装和配置 /412
 - 18.1.2 Hadoop Studio 的使用举例 /413
- 18.2 Hadoop Eclipse 简介和使用 /419



- 18.2.1 Hadoop Eclipse 安装和配置 /420
- 18.2.2 Hadoop Eclipse 的使用举例 /420
- 18.2.3 Hadoop Eclipse 插件开发 /421
- 18.3 Hadoop Streaming 简介和使用 /422
 - 18.3.1 Hadoop Streaming 的使用举例 /426
 - 18.3.2 使用 Hadoop Streaming 时常见的问题 /428
- 18.4 Hadoop Libhdfs 简介和使用 /430
 - 18.4.1 Hadoop Libhdfs 安装和配置 /430
 - 18.4.2 Hadoop Libhdfs API 简介 /430
 - 18.4.3 Hadoop Libhdfs 的使用举例 /431
- 18.5 小结 /432

附录 A 云计算在线检测平台 /434

- A.1 平台介绍 /435
- A.2 结构和功能 /435
 - A.2.1 前台用户接口的结构和功能 /435
 - A.2.2 后台程序运行的结构和功能 /437
- A.3 检测流程 /437
- A.4 使用 /438
 - A.4.1 功能使用 /438
 - A.4.2 返回结果介绍 /439
 - A.4.3 使用注意事项 /440
- A.5 小结 /441





本章内容

- ☐ 什么是 Hadoop
- ☐ Hadoop 项目及其结构
- ☐ Hadoop 的体系结构
- ☐ Hadoop 与分布式开发
- ☐ Hadoop 计算模型——MapReduce
- ☐ Hadoop 的数据管理
- ☐ 小结



1.1 什么是 Hadoop

1.1.1 Hadoop 概述

Hadoop 是 Apache 软件基金会旗下的一个开源分布式计算平台。以 Hadoop 分布式文件系统（HDFS, Hadoop Distributed Filesystem）和 MapReduce（Google MapReduce 的开源实现）为核心的 Hadoop 为用户提供了系统底层细节透明的分布式基础架构。HDFS 的高容错性、高伸缩性等优点允许用户将 Hadoop 部署在低廉的硬件上，形成分布式系统；MapReduce 分布式编程模型允许用户在不了解分布式系统底层细节的情况下开发并行应用程序。所以用户可以利用 Hadoop 轻松地组织计算机资源，从而搭建自己的分布式计算平台，并且可以充分利用集群的计算和存储能力，完成海量数据的处理。

1.1.2 Hadoop 的历史

Hadoop 的源头是 Apache Nutch，该项目开始于 2002 年，是 Apache Lucene 的子项目之一。2004 年，Google 在“操作系统设计与实现”（OSDI, Operating System Design and Implementation）会议上公开发表了题为“MapReduce: Simplified Data Processing on Large Clusters”（MapReduce: 简化大规模集群上的数据处理）的论文，之后受到启发的 Doug Cutting 等人开始尝试实现 MapReduce 计算框架，并将它与 NDFS（Nutch Distributed File System）结合，以支持 Nutch 引擎的主要算法。由于 NDFS 和 MapReduce 在 Nutch 引擎中有着良好的应用，所以它们于 2006 年 2 月被分离出来，成为了一套完整而独立的软件，起名为 Hadoop。到了 2008 年年初，Hadoop 已成为 Apache 的顶级项目，它被包括 Yahoo! 在内的很多互联网公司所采用。现在，Hadoop 已经发展成为包含 HDFS、MapReduce、Pig、ZooKeeper 等子项目的集合，用于分布式计算。

1.1.3 Hadoop 的功能与作用

我们为什么需要 Hadoop 呢？众所周知，现代社会的信息量增长速度极快，这些信息里又积累着大量的数据，其中包括个人数据和工业数据。预计到 2020 年，每年产生的数字信息将会有超过 1/3 的内容驻留在云平台中或借助云平台处理。我们需要对这些数据进行分析 and 处理，以获取更多有价值的信息。那么我们如何高效地存储和管理这些数据，如何分析这些数据呢？这时可以选用 Hadoop 系统，它在处理这类问题时，采用了分布式存储方式，提高了读写速度，并扩大了存储容量。采用 MapReduce 来整合分布式文件系统上的数据，可保证分析和处理数据的高效。与此同时，Hadoop 还采用存储冗余数据的方式保证了数据的安全性。

Hadoop 中 HDFS 的高容错特性，以及它是基于 Java 语言开发的，这使得 Hadoop 可以部署在低廉的计算机集群中，同时不限于某个操作系统。Hadoop 中 HDFS 的数据管理能力，MapReduce 处理任务时的高效率，以及它的开源特性，使其在同类的分布式系统中大放异

彩，并在众多行业和科研领域中被广泛采用。

1.1.4 Hadoop 的优势

Hadoop 是一个能够让用户轻松架构和使用的分布式计算平台。用户可以轻松地在 Hadoop 上开发和运行处理海量数据的应用程序。它主要有以下几个优点：

- ❑ 高可靠性。Hadoop 按位存储和处理数据的能力值得人们信赖。
- ❑ 高扩展性。Hadoop 是在可用的计算机集簇间分配数据并完成计算任务的，这些集簇可以方便地扩展到数以千计的节点中。
- ❑ 高效性。Hadoop 能够在节点之间动态地移动数据，并保证各个节点的动态平衡，因此其处理速度非常快。
- ❑ 高容错性。Hadoop 能够自动保存数据的多个副本，并且能够自动将失败的任务重新分配。

1.1.5 Hadoop 的应用现状和发展趋势

由于 Hadoop 优势突出，基于 Hadoop 的应用已经遍地开花，尤其是在互联网领域。Yahoo! 通过集群运行 Hadoop，以支持广告系统和 Web 搜索的研究；Facebook 借助集群运行 Hadoop，以支持其数据分析和机器学习；百度则使用 Hadoop 进行搜索日志的分析和网页数据的挖掘工作；淘宝的 Hadoop 系统用于存储并处理电子商务交易的相关数据；中国移动研究院基于 Hadoop 的“大云”（BigCloud）系统用于对数据进行分析并对外提供服务。

2008 年 2 月，Hadoop 最大贡献者的 Yahoo! 构建了当时规模最大的 Hadoop 应用，它们 2000 个节点上面执行了超过 1 万个 Hadoop 虚拟机器来处理超过 5PB 的网页内容，分析大约 1 兆个网络连接之间的网页索引资料。这些网页索引资料压缩后超过 300TB。Yahoo! 正是基于这些为用户提供了高质量的搜索服务。

Hadoop 目前已经取得了非常突出的成绩。随着互联网的发展，新的业务模式还将不断涌现，Hadoop 的应用也会从互联网领域向电信、电子商务、银行、生物制药等领域拓展。相信在未来，Hadoop 将会在更多的领域中扮演幕后英雄，为我们提供更加快捷优质的服务。

1.2 Hadoop 项目及其结构

现在 Hadoop 已经发展成为包含多个子项目的集合。虽然其核心内容是 MapReduce 和 Hadoop 分布式文件系统（HDFS），但 Hadoop 下的 Common、Avro、Chukwa、Hive、HBase 等子项目也是不可或缺的。它们提供了互补性服务或在核心层上提供了更高层的服务。图 1-1 展现了 Hadoop 的项目结构图。

下面将对 Hadoop 的各个子项目进行更详细的介绍。

1) Core/Common: 从 Hadoop 0.20 版本开始, Hadoop Core 项目便更名为 Common。Common 是为 Hadoop 其他子项目提供支持的常用工具, 它主要包括 FileSystem、RPC 和串行化库, 它们为在廉价的硬件上搭建云计算环境提供基本的服务, 并且为运行在该平台上的软件开发提供了所需的 API。

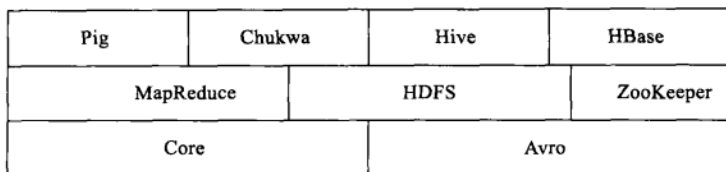


图 1-1 Hadoop 的项目结构图

2) Avro: Avro 是用于数据序列化的系统。它提供了丰富的数据结构类型、快速可压缩的二进制数据格式、存储持久性数据的文件集、远程调用 RPC 的功能和简单的动态语言集成功能。其中, 代码生成器既不需要读写文件数据, 也不需要实现或实现 RPC 协议, 它只是一个可选的对静态类型语言的实现。

Avro 系统依赖于模式 (Schema), Avro 数据的读和写是在模式之下完成的。这样就可以减少写入数据的开销, 提高序列化的速度并缩减其大小。同时, 也可以方便动态脚本语言的使用, 因为数据连同其模式都是自描述的。

在 RPC 中, Avro 系统的客户端和服务端通过握手协议进行模式的交换。因此当客户端和服务端拥有彼此全部的模式时, 不同模式下的相同命名字段、丢失字段和附加字段等信息的一致性就得到了很好的解决。

3) MapReduce: MapReduce 是一种编程模型, 用于大规模数据集 (大于 1TB) 的并行运算。“映射”(map)、“化简”(reduce) 等概念和它们的主要思想都是从函数式编程语言中借来的。它使得编程人员在不了解分布式并行编程的情况下也能方便地将自己的程序运行在分布式系统上。MapReduce 在执行时先指定一个 map (映射) 函数, 把输入键值对映射成一组新的键值对, 经过一定的处理后交给 reduce, reduce 对相同 key 下的所有 value 进行处理后再输出键值对作为最终的结果。

图 1-2 是 MapReduce 的任务处理流程图, 它展示了 MapReduce 程序将输入划分到不同的 map 上, 再将 map 的结果合并到 reduce, 然后进行处理的输出过程。详细介绍请参考本

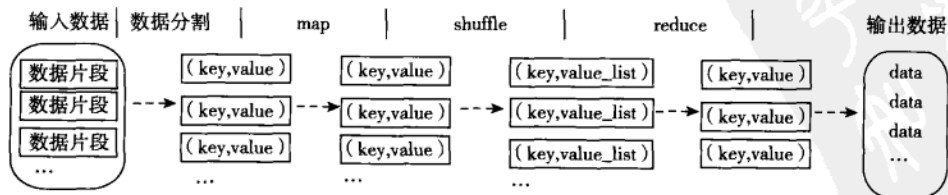


图 1-2 MapReduce 的任务处理流程图

章的 1.3 节。

4) HDFS：是一个分布式文件系统。由于 HDFS 具有高容错性 (fault-tolerant) 的特点，所以可以设计部署在低廉 (low-cost) 的硬件上。它可以通过提供高吞吐率 (high throughput) 来访问应用程序的数据，适合那些有着超大数据集的应用程序。HDFS 放宽了可移植操作系统接口 (POSIX, Portable Operating System Interface) 的要求，这样就可以实现以流的形式访问文件系统中的数据。HDFS 原本是开源的 Apache 项目 Nutch 的基础结构，最后它成为了 Hadoop 的基础架构之一。

以下是 HDFS 的设计目标：

- ❑ 检测和快速恢复硬件故障。硬件故障是常见的问题，整个 HDFS 系统由数百台或数千台存储着数据文件的服务器组成，而如此多的服务器意味着高故障率，因此，故障的检测和自动快速恢复是 HDFS 的一个核心目标。
- ❑ 流式的数据访问。HDFS 使应用程序能流式地访问它们的数据集。HDFS 被设计成适合进行批量处理，而不是用户交互式的处理。所以它重视数据吞吐量，而不是数据访问的反应速度。
- ❑ 简化一致性模型。大部分的 HDFS 程序操作文件时需要一次写入，多次读取。一个文件一旦经过创建、写入、关闭之后就不需要修改了，从而简化了数据一致性问题和高吞吐量的数据访问问题。
- ❑ 通信协议。所有的通信协议都在 TCP/IP 协议之上。一个客户端和明确配置了端口的名字节点 (NameNode) 建立连接之后，它和名称节点 (NameNode) 的协议便是客户端协议 (Client Protocol)。数据节点 (DataNode) 和名字节点 (NameNode) 之间则用数据节点协议 (DataNode Protocol)。

关于 HDFS 的具体介绍请参考本章的 1.3 节。

5) Chukwa：Chukwa 是开源的数据收集系统，用于监控和分析大型分布式系统的数据。Chukwa 是在 Hadoop 的 HDFS 和 MapReduce 框架之上搭建的，它同时继承了 Hadoop 的可扩展性和健壮性。Chukwa 通过 HDFS 来存储数据，并依赖于 MapReduce 任务处理数据。Chukwa 中也附带了灵活且强大的工具，用于显示、监视和分析数据结果，以便更好地利用所收集的数据。

6) Hive：Hive 最早是由 Facebook 设计的，是一个建立在 Hadoop 基础之上的数据仓库，它提供了一些用于数据整理、特殊查询和分析存储在 Hadoop 文件中的数据集的工具。Hive 提供的是一种结构化数据的机制，它支持类似于传统 RDBMS 中的 SQL 语言来帮助那些熟悉 SQL 的用户查询 Hadoop 中的数据，该查询语言称为 Hive QL。与此同时，那些传统的 MapReduce 编程人员也可以在 Mapper 或 Reducer 中通过 Hive QL 查询数据。Hive 编译器会把 Hive QL 编译成一组 MapReduce 任务，从而方便 MapReduce 编程人员进行 Hadoop 应用的开发。

7) HBase：HBase 是一个分布式的、面向列的开源数据库，该技术来源于 Google 的论文“Bigtable：一个结构化数据的分布式存储系统”。如同 Bigtable 利用了 Google 文件系统

(Google File System) 提供的分布式数据存储方式一样, HBase 在 Hadoop 之上提供了类似于 Bigtable 的能力。HBase 是 Hadoop 项目的子项目。HBase 不同于一般的关系数据库, 其一, HBase 是一个适合于存储非结构化数据的数据库; 其二, HBase 是基于列而不是基于行的模式。HBase 和 Bigtable 使用相同的数据模型。用户将数据存储在一个表里, 一个数据行拥有一个可选择的键和任意数量的列。由于 HBase 表示疏松的, 用户可以给行定义各种不同的列。HBase 主要用于需要随机访问、实时读写的大数据 (Big Data)。具体介绍请参考本书第 12 章 “HBase 详解”。

8) Pig: Pig 是一个对大型数据集进行分析和评估的平台。Pig 最突出的优势是它的结构能够经受住高度并行化的检验, 这个特性让它能够处理大型的数据集。目前, Pig 的底层由一个编译器组成, 它在运行的时候会产生一些 MapReduce 程序序列, Pig 的语言层由一种叫做 Pig Latin 的正文型语言组成。有关 Pig 的具体内容请参考本书第 14 章 “Pig 详解”。

上面讨论的 8 个子项目在本书中都有相应的章节进行详细的介绍。

1.3 Hadoop 的体系结构

如前面的内容所说, HDFS 和 MapReduce 是 Hadoop 的两大核心。而整个 Hadoop 的体系结构主要是通过 HDFS 来实现对分布式存储的底层支持的, 并且它会通过 MapReduce 来实现对分布式并行任务处理的程序支持。

1.3.1 HDFS 的体系结构

我们首先介绍 HDFS 的体系结构, HDFS 采用了主从 (Master/Slave) 结构模型, 一个 HDFS 集群是由一个 NameNode 和若干个 DataNode 组成的。其中 NameNode 作为主服务器, 管理文件系统的命名空间和客户端对文件的访问操作; 集群中的 DataNode 管理存储的数据。HDFS 允许用户以文件的形式存储数据。从内部来看, 文件被分成若干个数据块, 而且这若干个数据块存放在一组 DataNode 上。NameNode 执行文件系统的命名空间操作, 比如打开、关闭、重命名文件或目录等, 它也负责数据块到具体 DataNode 的映射。DataNode 负责处理文件系统客户端的文件读写请求, 并在 NameNode 的统一调度下进行数据块的创建、删除和复制工作。图 1-3 给出了 HDFS 的体系结构。

NameNode 和 DataNode 都被设计成可以在普通商用计算机上运行。这些计算机通常运行的是 GNU/Linux 操作系统。HDFS 采用 Java 语言开发, 因此任何支持 Java 的机器都可以部署 NameNode 和 DataNode。一个典型的部署场景是集群中的一台机器运行一个 NameNode 实例, 其他机器分别运行一个 DataNode 实例。当然, 并不排除一台机器运行多个 DataNode 实例的情况。集群中单一的 NameNode 的设计则大大简化了系统的架构。NameNode 是所有 HDFS 元数据的管理者, 用户数据永远不会经过 NameNode。

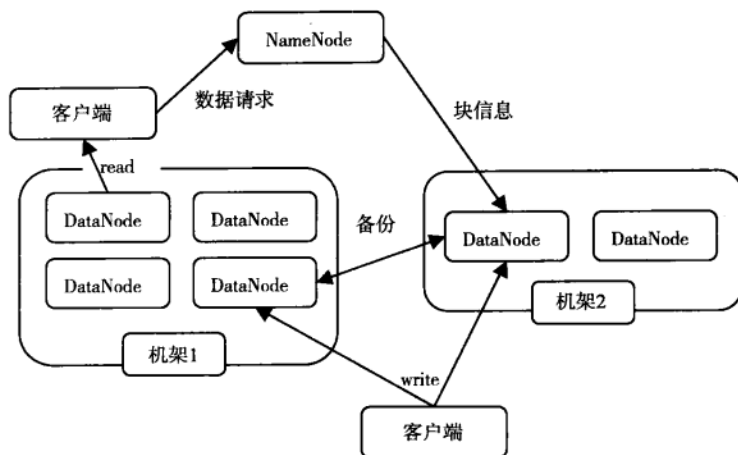


图 1-3 HDFS 的体系结构图

1.3.2 MapReduce 的体系结构

接下来介绍 MapReduce 的体系结构，MapReduce 是一种并行编程模式，这种模式使得软件开发人员可以轻松地编写出分布式并行程序。在 Hadoop 的体系结构中，MapReduce 是一个简单易用的软件框架，基于它可以将任务分发到由上千台商用机器组成的集群上，并以一种高容错的方式并行处理大量的数据集，实现 Hadoop 的并行任务处理功能。MapReduce 框架是由一个单独运行在主节点上的 JobTracker 和运行在每个集群从节点上的 TaskTracker 共同组成的。主节点负责调度构成一个作业的所有任务，这些任务分布在不同的从节点上。主节点监控它们的执行情况，并且重新执行之前失败的任务；从节点仅负责由主节点指派的任务。当一个 Job 被提交时，JobTracker 接收到提交作业和配置信息之后，就会将配置信息等分发给从节点，同时调度任务并监控 TaskTracker 的执行。

从上面的介绍可以看出，HDFS 和 MapReduce 共同组成了 Hadoop 分布式系统体系结构的核心。HDFS 在集群上实现了分布式文件系统，MapReduce 在集群上实现了分布式计算和任务处理。HDFS 在 MapReduce 任务处理过程中提供了文件操作和存储等支持，MapReduce 在 HDFS 的基础上实现了任务的分发、跟踪、执行等工作，并收集结果，二者相互作用，完成了 Hadoop 分布式集群的主要任务。

1.4 Hadoop 与分布式开发

我们通常说的分布式系统其实是分布式软件系统，即支持分布式处理的软件系统，它是在通信网络互联的多处理机体系结构上执行任务的，包括分布式操作系统、分布式程序设计语言及其编译（解释）系统、分布式文件系统和分布式数据库系统等。Hadoop 是分布式

软件系统中文件系统这一层的软件，它实现了分布式文件系统和部分分布式数据库的功能。Hadoop 中的分布式文件系统 HDFS 能够实现数据在计算机集群组成的云上高效的存储和管理，Hadoop 中的并行编程框架 MapReduce 能够让用户编写的 Hadoop 并行应用程序运行更加简化。下面简单介绍一下基于 Hadoop 进行分布式并发编程的相关知识，详细的介绍请参看后面有关 MapReduce 编程的章节。

Hadoop 上的并行应用程序开发是基于 MapReduce 编程框架的。MapReduce 编程模型的原理是：利用一个输入的 key/value 对集合来产生一个输出的 key/value 对集合。MapReduce 库的用户用两个函数来表达这个计算：Map 和 Reduce。

用户自定义的 map 函数接收一个输入的 key/value 对，然后产生一个中间 key/value 对的集合。MapReduce 把所有具有相同 key 值的 value 集合在一起，然后传递给 reduce 函数。

用户自定义的 reduce 函数接收 key 和相关的 value 集合。reduce 函数合并这些 value 值，形成一个较小的 value 集合。一般来说，每次 reduce 函数调用只产生 0 或 1 个输出的 value 值。通常我们通过一个迭代器把中间的 value 值提供给 reduce 函数，这样就可以处理无法全部放入内存中的大量的 value 值集合了。

图 1-4 是 MapReduce 的数据流图，这个过程简而言之就是将大数据集分解为成百上千个小数据集，每个（或若干个）数据集分别由集群中的一个节点（一般就是一台普通的计算机）进行处理并生成中间结果，然后这些中间结果又由大量的节点合并，形成最终结果。图 1-4 也指出了 MapReduce 框架下并行程序中的三个主要函数：map、reduce、main。在这个结构中，需要用户完成的工作仅仅是根据任务编写 map 和 reduce 两个函数。

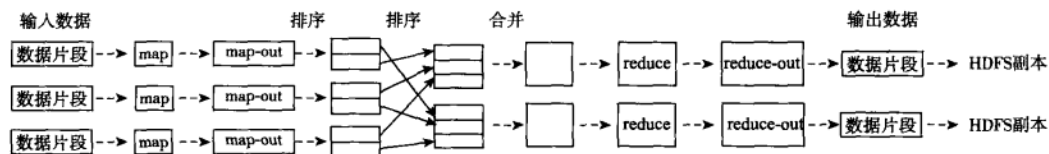


图 1-4 MapReduce 的数据流图

MapReduce 计算模型非常适合在大量计算机组成的大规模集群上并行运行。图 1-4 中的每一个 map 任务和每一个 reduce 任务均可以同时运行于一个单独的计算节点上，可想而知，其运算效率是很高的，那么这样的并行计算是如何做到的呢？下面将简单介绍一下其原理。

1. 数据分布存储

Hadoop 分布式文件系统（HDFS）由一个名称节点（NameNode）和 N 个数据节点（DataNode）组成，每个节点均是一台普通的计算机。在使用方式上 HDFS 与我们熟悉的单机文件系统非常类似，它可以创建目录，创建、复制和删除文件，以及查看文件的内容等。但 HDFS 底层把文件切割成了 Block，然后这些 Block 分散地存储于不同的 DataNode 上，每个 Block 还可以复制数份数据存储于不同的 DataNode 上，达到容错容灾的目的。NameNode 则是整个 HDFS 的核心，它通过维护一些数据结构来记录每一个文件被切割成了多少个

Block、这些 Block 可以从哪些 DataNode 中获得，以及各个 DataNode 的状态等重要信息。

2. 分布式并行计算

Hadoop 中有一个作为主控的 JobTracker，用于调度和管理其他的 TaskTracker，JobTracker 可以运行于集群中的任意一台计算机上。TaskTracker 则负责执行任务，它必须运行于 DataNode 上，也就是说 DataNode 既是数据存储节点，也是计算节点。JobTracker 将 map 任务和 reduce 任务分发给空闲的 TaskTracker，让这些任务并行运行，并负责监控任务的运行情况。如果某一个 TaskTracker 出了故障，JobTracker 会将其负责的任务转交给另一个空闲的 TaskTracker 重新运行。

3. 本地计算

数据存储在某一计算机上，就由哪台计算机进行这部分数据的计算，这样可以减少数据在网络上的传输，降低对网络带宽的需求。在 Hadoop 这类基于集群的分布式并行系统中，计算节点可以很方便地扩充，它所能提供的计算能力近乎无限，但是由于数据需要在不同的计算机之间流动，故网络带宽变成了瓶颈，“本地计算”是一种最有效的节约网络带宽的手段，业界把这形容为“移动计算比移动数据更经济”。

4. 任务粒度

把原始大数据集切割成小数据集时，通常让小数据集小于或等于 HDFS 中一个 Block 的大小（默认是 64MB），这样能够保证一个小数据集是位于一台计算机上的，便于本地计算。有 M 个小数据集待处理，就启动 M 个 map 任务，注意这 M 个 map 任务分布于 N 台计算机上，它们会并行运行，reduce 任务的数量 R 则可由用户指定。

5. 数据分割（Partition）

把 map 任务输出的中间结果按 key 的范围划分成 R 份（R 是预先定义的 reduce 任务的个数），划分时通常使用 hash 函数（如： $\text{hash}(\text{key}) \bmod R$ ），这样可以保证某一范围内的 key 一定是由一个 reduce 任务来处理的，可以简化 Reduce 的过程。

6. 数据合并（Combine）

在数据分割之前，还可以先对中间结果进行数据合并（Combine），即将中间结果中有相同 key 的 <key, value> 对合并成一对。Combine 的过程与 reduce 的过程类似，很多情况下可以直接使用 reduce 函数，但 Combine 是作为 map 任务的一部分，在执行完 map 函数后紧接着执行的。Combine 能够减少中间结果中 <key, value> 对的数目，从而降低网络流量。

7. Reduce

Map 任务的中间结果在做完 Combine 和 Partition 之后，以文件形式存于本地磁盘上。中间结果文件的位置会通知主控 JobTracker，JobTracker 再通知 reduce 任务到哪一个 DataNode 上去取中间结果。注意，所有的 map 任务产生的中间结果均按其 key 值用同一个 hash 函数划分成了 R 份，R 个 reduce 任务各自负责一段 key 区间。每个 reduce 需要向许多个 map 任务节点取得落在其负责的 key 区间内的中间结果，然后执行 reduce 函数，形成一

个最终的结果文件。

8. 任务管道

有 R 个 reduce 任务，就会有 R 个最终结果，很多情况下这 R 个最终结果并不需要合并成一个最终结果，因为这 R 个最终结果又可以作为另一个计算任务的输入，开始另一个并行计算任务，这也就形成了任务管道。

这里简要介绍了在并行编程方面 Hadoop 中 MapReduce 编程模型的原理、流程、程序结构和并行计算的实现，MapReduce 程序的详细流程、编程接口、程序实例等请参见后面章节。

1.5 Hadoop 计算模型——MapReduce

MapReduce 是 Google 公司的核心计算模型，它将运行于大规模集群上的复杂的并行计算过程高度地抽象为了两个函数：map 和 reduce。Hadoop 是 Doug Cutting 受到 Google 发表的关于 MapReduce 的论文的启发而开发出来的。Hadoop 中的 MapReduce 是一个使用简易的软件框架，基于它写出来的应用程序能够运行在由上千台商用机器组成的大型集群上，并以一种可靠容错的方式并行处理上 T 级别的数据集，实现了 Hadoop 在集群上的数据和任务的并行计算与处理。

一个 MapReduce 作业 (job) 通常会把输入的数据集切分为若干个独立的数据块，由 map 任务 (task) 以完全并行的方式处理它们。框架会先对 map 的输出进行排序，然后把结果输入给 reduce 任务。通常作业的输入和输出都会被存储在文件系统中。整个框架负责任务的调度和监控，以及重新执行已经失败的任务。

通常，MapReduce 框架和分布式文件系统是运行在一组相同的节点上的，也就是说，计算节点和存储节点在一起。这种配置允许框架在那些已经存好数据的节点上高效地调度任务，这可以使整个集群的网络带宽被非常高效的利用。

MapReduce 框架由一个单独的 master JobTracker 和集群节点上的 slave TaskTracker 共同组成。master 负责调度构成一个作业的所有任务，这些任务分布在不同的 slave 上。master 监控它们的执行情况，并重新执行已经失败的任务，而 slave 仅负责执行由 master 指派的任务。

在 Hadoop 上运行的作业需要指明程序的输入/输出位置 (路径)，并通过实现合适的接口或抽象类来提供 map 和 reduce 函数。同时还需要指定作业的其他参数，构成作业配置 (job configuration)。在 Hadoop 的 jobclient 提交作业 (jar 包/可执行程序等) 和配置信息给 JobTracker 之后，JobTracker 会负责分发这些软件和配置信息给 slave 及调度任务，并监控它们的执行，同时提供状态和诊断信息给 jobclient。

1.6 Hadoop 的数据管理

前面重点介绍了 Hadoop 及其体系结构和计算模型 MapReduce，现在开始介绍 Hadoop

的数据管理，主要包括 Hadoop 的分布式文件系统 HDFS、分布式数据库 HBase 和数据仓库工具 Hive 的数据管理。

1.6.1 HDFS 的数据管理

HDFS 是分布式计算的存储基石，Hadoop 分布式文件系统和其他分布式文件系统有很多类似的特质：

- ❑ 对于整个集群有单一的命名空间；
- ❑ 具有数据一致性。适合一次写入多次读取的模型，客户端在文件没有被成功创建之前是无法看到文件存在的；
- ❑ 文件会被分割成多个文件块，每个文件块被分配存储到数据节点上，而且会根据配置由复制文件块来保证数据的安全性。

从前面的介绍和图 1-3 可以看出，HDFS 通过三个重要的角色来进行文件系统的管理：NameNode、DataNode 和 Client。NameNode 可以看做是分布式文件系统管理者，主要负责管理文件系统的命名空间、集群配置信息和存储块的复制等。NameNode 会将文件系统的 Metadata 存储在内存中，这些信息主要包括文件信息、每一个文件对应的文件块的信息和每一个文件块在 DataNode 中的信息等。DataNode 是文件存储的基本单元，它将文件块（Block）存储在本地文件系统中，保存了所有 Block 的 Metadata，同时周期性地将所有存在的 Block 信息发送给 NameNode。Client 就是需要获取分布式文件系统文件的应用程序。以下通过三个具体的操作来说明 HDFS 对数据的管理。

（1）文件写入

- 1) Client 向 NameNode 发起文件写入的请求。
- 2) NameNode 根据文件大小和文件块的配置情况，返回给 Client 它所管理的 DataNode 的信息。
- 3) Client 将文件划分为多个 Block，根据 DataNode 的地址信息，按顺序将其写入每一个 DataNode 块中。

（2）文件读取

- 1) Client 向 NameNode 发起读取文件的请求。
- 2) NameNode 返回文件存储的 DataNode 信息。
- 3) Client 读取文件信息。

（3）文件块（Block）复制

- 1) NameNode 发现部分文件的 Block 不符合最小复制数这一要求或部分 DataNode 失效。
- 2) 通知 DataNode 相互复制 Block。
- 3) DataNode 开始直接相互复制。

HDFS 作为分布式文件系统在数据管理方面还有几个值得借鉴的功能：

- ❑ 文件块（Block）的放置：一个 Block 会有三份备份，一份放在 NameNode 指定的 DataNode 上，另一份放在与指定的 DataNode 不在同一台机器上的 DataNode 上，最

后一份放在与指定的 DataNode 在同一 Rack 上的 DataNode 上。备份的目的是为了数据安全，采用这种配置方式主要是考虑同一 Rack 失败的情况，以及不同 Rack 之间的数据拷贝会带来的性能问题。

- ❑ 心跳检测：用心跳检测 DataNode 的健康状况，如果发现问题就采取数据备份的方式来保证数据的安全性。
- ❑ 数据复制（场景为 DataNode 失败、需要平衡 DataNode 的存储利用率和平衡 DataNode 数据交互压力等情况）：使用 Hadoop 时可以用 HDFS 的 balancer 命令配置 Threshold 来平衡每一个 DataNode 的磁盘利用率。假设设置了 Threshold 为 10%，那么执行 balancer 命令的时候，首先会统计所有 DataNode 的磁盘利用率的平均值，然后判断如果某一个 DataNode 的磁盘利用率超过这个均值，那么将会把这个 DataNode 的 block 转移到磁盘利用率低的 DataNode 上，这对于新节点的加入来说十分有用。
- ❑ 数据校验：采用 CRC32 做数据校验。在写入文件 Block 的时候，除了写入数据外还会写入校验信息，在读取的时候则需要校验后再读入。
- ❑ 单个 NameNode：如果失败，任务处理信息将会记录在本地文件系统和远端的文件系统中。
- ❑ 数据管道性的写入：当客户端要写入文件到 DataNode 上时，客户端首先会读取一个 Block，然后写到第一个 DataNode 上，接着由第一个 DataNode 将其传递到备份的 DataNode 上，直到所有需要写入这个 Block 的 DataNode 都成功写入后，客户端才会开始写下一个 Block。
- ❑ 安全模式：分布式文件系统启动的时候会有安全模式（系统运行期间也可以通过命令进入安全模式），当分布式文件系统处于安全模式时，文件系统中的内容不允许修改也不允许删除，直到安全模式结束。安全模式主要是为了在系统启动的时候检查各个 DataNode 上的数据块的有效性，同时根据策略进行必要的复制或删除部分数据块。在实际操作过程中，若在系统启动时修改和删除文件会出现安全模式不允许修改的错误提示，只需要等待一会儿即可。

1.6.2 HBase 的数据管理

HBase 是一个类似 Bigtable 的分布式数据库，它的大部分特性和 Bigtable 一样，是一个稀疏的、长期存储的（存在硬盘上）、多维度的排序映射表。这张表的索引是行关键字、列关键字和时间戳。每个值是一个不解释的字符数组，数据都是字符串，没有类型。用户在表格中存储数据，每一行都有一个可排序的主键和任意多的列。由于是稀疏存储的，所以同一张表里面的每一行数据都可以有截然不同的列。列名字的格式是“<family>:<label>”，它是由字符串组成的，每一张表有一个 family 集合，这个集合是固定不变的，相当于表的结构，只能通过改变表结构来改变表的 family 集合，但是 label 值相对于每一行来说都是可以改变的。

HBase 把同一个 family 中的数据存储在同一个目录下，而 HBase 的写操作是锁行的，

每一行都是一个原子元素，都可以加锁。所有数据库的更新都有一个时间戳标记，每次更新都会生成一个新的版本，而 HBase 会保留一定数量的版本，这个值是可以设定的。客户端可以选择获取距离某个时间点最近的版本，或者一次获取所有版本。

以上从微观上介绍了 HBase 的一些数据管理措施。那么 HBase 作为分布式数据库在整体上从集群出发又是如何管理数据的呢？

HBase 在分布式集群上主要依靠由 HRegion、HMaster、HClient 组成的体系结构从整体上管理数据。

HBase 体系结构的三大重要组成部分是：

- ❑ HBaseMaster：HBase 主服务器，与 Bigtable 的主服务器类似。
 - ❑ HRegionServer：HBase 域服务器，与 Bigtable 的 Tablet 服务器类似。
 - ❑ HBaseClient：HBase 客户端是由 `org.apache.hadoop.HBase.client.HTable` 定义的。
- 下面将对这三个组件进行详细的介绍。

(1) HBaseMaster

一个 HBase 只部署一台主服务器，它通过领导选举算法（Leader Election Algorithm）确保只有唯一的主服务器是活跃的，ZooKeeper 保存主服务器的服务器地址信息。如果主服务器瘫痪，可以通过领导选举算法从备用服务器中选择新的主服务器。

主服务器初始化集群。当主服务器第一次启动时，会试图从 HDFS 获取根或根域目录，若获取失败则创建根或根域目录，以及第一个元域目录。下次启动时，主服务器就可以获取集群和集群中所有域的信息了。

主服务器负责域的分配工作。首先，主服务器会分配根域，并存储指向根域所在域服务器地址的指针。其次，主服务器遍历根域查询元域，并分配元域到域服务器中。每个元域中包含了所有的用户域，用户域中存储了多个用户表。如果所有的元域分配完毕，主服务器将会分配用户域到相应的域服务器，以保持域服务器间的负载平衡。

主服务器时刻监视着域服务器的运行状态。一旦主服务器检测到某一域服务器不可达时，它将分离出域服务器上的每个域的预写日志（write-ahead log）文件。之后，主服务器会将域重新分配到其他域服务器上，并运行之。如果主服务器发现一个域服务器超负荷运行，则会取消或关闭该域服务器的一些域，并将这些域分配到其他低负载的域服务器上。

主服务器还负责表格的管理。例如，调整表格的在线/离线状态和改变表格的模式（增加或删除列族）等。此外，客户端还可以请求本地域直接从域服务器上读取数据。

在 Bigtable 中，当主服务器和域服务器的连接断开时，域服务器仍然可以继续服务，因为 Bigtable 提供了一种额外的锁管理机制，这种机制中的锁管理器（Chubby）保证了域服务器服务的可用性。而在 HBase 中，由于没有提供锁管理机制，当主服务器崩溃时，整个集群系统都要重新启动，因为主服务器是所有域服务器的管理中心。

下面介绍元表和根表的概念：

元表（Meta Table）包含了所有用户域的基本信息，域信息包括起始关键字、结束关键字、域是否在线、域所在的域服务器地址等。元表会随着用户域的增长而增长。

根表 (Root Table) 被定义为存储单一域的信息, 并指向元表中的所有域。与元表一样, 根表也包含每个元域的信息和元域所在的域服务器地址。

根表和元表中的每行大约为 1KB。域默认大小为 256MB, 根域可以映射 2.6×10^5 个元域。同样, 元域可以映射相应数量的用户域。因此, 根域可以映射 6.9×10^{10} 个用户域, 大约可以存储 1.9×10^{19} 字节的数据。

(2) HRegionServer

HBase 域服务器主要有服务于主服务器分配的域、处理客户端的读写请求、缓冲区回写、压缩和分割域等功能。

每个域只能由一台域服务器来服务。当它开始服务于某域时, 它会从 HDFS 文件系统中读取该域的日志和所有存储文件。同时它还会管理操作 HDFS 文件的持久性存储工作。

客户端通过与主服务器通信获取域和域所在域服务器的列表信息后, 就可以直接向域服务器发送域读写请求了。域服务器收到写请求时, 首先将写请求信息写入一个预写日志文件中, 该文件取名为 HLog。同一个域的所有写请求都被记录在同一个 HLog 文件中。一旦写请求被记录在 HLog 中之后, 它将被缓存在存储缓存区 (MemCache) 中。每个 HStore 对应一个存储缓存区。对于读请求, 域服务器先要检测请求数据在存储缓存区中是否被命中, 如果没有命中, 域服务器再去查找相关的映射文件。

当存储缓存区的大小达到一定阈值后, 需要将存储缓存区中的数据回写到磁盘上, 形成映射文件, 并在 HLog 日志文件中标记。因此当再次执行时, 可以跳跃到最后一次回写之前的操作上。回写也可能因域服务器存储器压力而被触发。

当映射文件的数量达到一定阈值时, 域服务器会将最近常写入的映射文件进行轻度的合并压缩。此外, 域服务器还会周期性地对所有的映射文件进行压缩, 使其成为单一的映射文件。之所以周期性地压缩所有的映射文件, 是因为最早的映射文件通常都比较大, 而最近的映射文件则要小很多, 压缩要消耗很多的时间, 具体消耗的时间主要取决于读取、合并和写出最大映射文件所需要的 I/O 操作次数。压缩和处理读写请求是同时进行的。在一个新的映射文件移入之前, 读写操作将被挂起, 直到映射文件被加入 HStore 的活跃映射文件列表中, 且已合并的旧映射文件被删除后, 才会释放读写操作。

当 HStore 中映射文件的大小达到一定的阈值时 (目前默认的阈值为 256MB), 域服务器就要对域进行分割了。域被均分为两个子域, 分割操作执行速度很快, 因为子域是直接从父域中读取数据的。之后, 父域处于离线状态。域服务器在元域中记录新的子域, 并通知主服务器可以将子域分配给其他域服务器。如果域分割消息在网络传输中丢失, 主服务器可以在周期性扫描元域中未被分配的域信息时发现分割操作。一旦父域被关闭, 所有对父域的读写操作将被挂起。客户端则会探测域的分割信息, 当新的子域在线时, 客户端再发出读写请求。当子域触发压缩操作时, 父域的数据将复制到子域中。父域将会在两个子域都完成压缩操作时被回收。

(3) HBaseClient

HBase 客户端负责查找用户域所在的域服务器地址。HBase 客户端会与 HBase 主机交换

消息以查找根域的位置，这是两者之间唯一的交流。

定位根域后，客户端连接根域所在的域服务器，并扫描根域获取元域信息，元域包含所需用户域的域服务器地址。客户端再连接元域所在的域服务器，扫描元域来获取所需用户域所在的域服务器地址。定位用户域后，客户端连接用户域所在的域服务器并发出读写请求。用户域的地址将在客户端中被缓存，后续的请求无须重复上述过程。

不管是由于主服务器为了负载均衡而重新分配域还是域服务器崩溃，客户端都会重新扫描元表来定位新的用户域地址。如果元域被重新分配，客户端将扫描根域来定位新的元域地址。如果根域也被重新分配，客户端将会连接主机定位新的根域地址，并通过重复上述过程来定位用户域地址。

综上所述，在 HBase 的体系结构中，HBase 主要由主服务器、域服务器和客户端三部分组成。主服务器作为 HBase 的中心，管理着整个集群中的所有域，监控每个域服务器的运行情况等；域服务器接收来自服务器的分配域，处理客户端的域读写请求并回写映射文件等；客户端主要用来查找用户域所在的域服务器地址信息。

1.6.3 Hive 的数据管理

Hive 是建立在 Hadoop 上的数据仓库基础构架。它提供了一系列的工具，用来进行数据提取、转化、加载，这是一种可以存储、查询和分析存储在 Hadoop 中的大规模数据的机制。Hive 定义了简单的类 SQL 查询语言，称为 QL，它允许熟悉 SQL 的用户查询数据。作为一个数据仓库，Hive 的数据管理按照使用层次可以从元数据存储、数据存储和数据交换三个方面来介绍。

(1) 元数据存储

Hive 将元数据存储于 RDBMS 中，有三种模式可以连接到数据库：

- ❑ Single User Mode：此模式连接到一个 In-memory 的数据库 Derby，一般用于 Unit Test。
- ❑ Multi User Mode：通过网络连接到一个数据库中，这是最常用的模式。
- ❑ Remote Server Mode：用于非 Java 客户端访问元数据库，在服务器端启动一个 MetaStoreServer，客户端则利用 Thrift 协议通过 MetaStoreServer 来访问元数据库。

(2) 数据存储

首先，Hive 没有专门的数据存储格式，也没有为数据建立索引，用户可以非常自由地组织 Hive 中的表，只需要在创建表的时候告诉 Hive 数据中的列分隔符和行分隔符，它就可以解析数据了。

其次，Hive 中所有的数据都存储在 HDFS 中，Hive 中包含 4 种数据模型：Table、External Table、Partition、Bucket。

Hive 中的 Table 和数据库中的 Table 在概念上是类似的，每一个 Table 在 Hive 中都有一个相应的目录来存储数据。例如，一个表 pvs，它在 HDFS 中的路径为：/wh/pvs，其中，wh 是在 hive-site.xml 中由 \${hive.metastore.warehouse.dir} 指定的数据仓库的目录，所有的

Table 数据（不包括 External Table）都保存在这个目录中。

Partition 对应于数据库中 Partition 列的密集索引，但是 Hive 中 Partition 的组织方式与数据库中的很不相同。在 Hive 中，表中的一个 Partition 对应于表下的一个目录，所有的 Partition 数据都存储在对应的目录中。例如：pvs 表中包含 ds 和 city 两个 Partition，则对应于 ds = 20090801, city = US 的 HDFS 子目录为：/wh/pvs/ds=20090801/city=US；对应于 ds = 20090801, city = CA 的 HDFS 子目录为：/wh/pvs/ds=20090801/city=CA。

Buckets 对指定列计算 hash，根据 hash 值切分数据，目的是为了便于并行，每一个 Buckets 对应一个文件。将 user 列分散至 32 个 Bucket 上，首先对 user 列的值计算 hash，比如，对应 hash 值为 0 的 HDFS 目录为：/wh/pvs/ds=20090801/city=US/part-00000；对应 hash 值为 20 的 HDFS 目录为：/wh/pvs/ds=20090801/city=US/part-00020。

External Table 指向已经在 HDFS 中存在的数 据，可以创建 Partition。它和 Table 在元数据的组织结构上是相同的，而在实际数据的存储上则有较大的差异。

在 Table 的创建过程和数据加载过程（这两个过程可以在同一个语句中完成）中，实际数据会被移动到数据仓库目录中。之后对数据的访问将会直接在数据仓库的目录中完成。删除表时，表中的数据和元数据将会被同时删除。

External Table 只有一个过程，因为加载数据和创建表是同时完成的。实际数据是存储在 Location 后面指定的 HDFS 路径中的，它并不会移动到数据仓库目录中。

（3）数据交换

数据交换主要分为以下几个部分（如图 1-5 所示）：

- ❑ 用户接口：包括客户端、Web 界面和数据库接口。
- ❑ 元数据存储：通常是存储在关系数据库中的，如 MySQL、Derby 等。
- ❑ 解释器、编译器、优化器、执行器。
- ❑ Hadoop：用 HDFS 进行存储，利用 MapReduce 进行计算。

用户接口主要有三个：客户端、数据库接口和 Web 界面，其中最常用的是客户端。Client 是 Hive 的客户端，当启动 Client 模式时，用户会想要连接 Hive Server，这时需要指出 Hive Server 所在的节点，并且在该节点启动 Hive Server。Web 界面是通过浏览器访问 Hive 的。

Hive 将元数据存储存储在数据库中，如 MySQL、Derby 中。Hive 中的元数据包括表的名字、

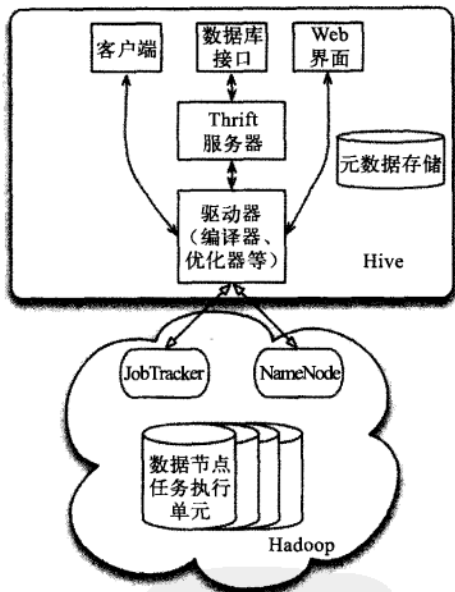


图 1-5 Hive 的数据交换图

表的列和分区及其属性、表的属性（是否为外部表等）、表数据所在的目录等。

解释器、编译器、优化器完成 HQL 查询语句从词法分析、语法分析、编译、优化到查询计划的生成。生成的查询计划存储在 HDFS 中，并在随后由 MapReduce 调用执行。

Hive 的数据存储在 HDFS 中，大部分的查询由 MapReduce 完成（包含 * 的查询不会生成 MapReduce 任务，比如 `select * from tbl`）。

以上从 Hadoop 的分布式文件系统 HDFS、分布式数据库 HBase 和数据仓库工具 Hive 入手介绍了 Hadoop 的数据管理，它们都通过自己的数据定义、体系结构实现了数据从宏观到微观的立体化管理，完成了 Hadoop 平台上大规模的数据存储和任务处理。

1.7 小结

本章首先介绍了 Hadoop 分布式计算平台：它是由 Apache 软件基金会开发的一个开源分布式计算平台。以 Hadoop 分布式文件系统（HDFS）和 MapReduce 为核心的 Hadoop 为用户提供了系统底层细节透明的分布式基础架构。由于 Hadoop 拥有可计量、成本低、高效、可信等突出特点，基于 Hadoop 的应用已经遍地开花，尤其是在互联网领域。

本章接下来介绍了 Hadoop 项目及其结构，现在 Hadoop 已经发展成为一个包含多个子项目的集合，被用于分布式计算，虽然 Hadoop 的核心是 Hadoop 分布式文件系统和 MapReduce，但 Hadoop 下的 Common、Avro、Chukwa、Hive、HBase 等子项目提供了互补性服务或在核心层之上提供了更高层的服务。紧接着，简要介绍了以 HDFS 和 MapReduce 为核心的 Hadoop 体系结构。

本章最后介绍了 Hadoop 的数据管理，主要包括 Hadoop 的分布式文件系统 HDFS、分布式数据库 HBase 和数据仓库工具 Hive 的数据管理。它们都有自己完整的数据定义和体系结构，以及实现数据从宏观到微观的立体管理数据方法，这都为 Hadoop 平台的数据存储和任务处理打下了基础。本章中的许多内容在本书后面的章节中都会详细展开介绍。





第 2 章

Hadoop 的安装与配置

本章内容

- ☐ 在 Linux 上安装与配置 Hadoop
- ☐ 在 Windows 上安装与配置 Hadoop
- ☐ 安装和配置 Hadoop 集群
- ☐ 日志分析及几个小技巧
- ☐ 小结

资源分享

PDG

Hadoop 的安装非常简单，大家可以在官网下载到最近的几个版本，网址为 <http://apache.etoak.com/hadoop/core/>。

Hadoop 最早是为了在 Linux 平台上使用而开发的，但是 Hadoop 在 UNIX、Windows 和 Mac OS X 系统上也运行良好。不过，在 Windows 上运行 Hadoop 稍显复杂，首先必须安装 Cygwin 以模拟 Linux 环境，然后才能安装 Hadoop。

在 Unix 上安装 Hadoop 的过程与在 Linux 上安装基本相同，因此下面不会对其进行详细介绍。

2.1 在 Linux 上安装与配置 Hadoop

在 Linux 上安装 Hadoop 之前，需要先安装两个程序：

- ❑ JDK 1.6 或更高版本；
- ❑ SSH（安全外壳协议），推荐安装 OpenSSH。

下面简述一下安装这两个程序的原因：

- ❑ Hadoop 是用 Java 开发的，Hadoop 的编译及 MapReduce 的运行都需要使用 JDK。
- ❑ Hadoop 需要通过 SSH 来启动 slave 列表中各台主机的守护进程，因此 SSH 也是必须安装的，即使是安装伪分布式版本（因为 Hadoop 并没有区分集群式和伪分布式）。对于伪分布式，Hadoop 会采用与集群相同的处理方式，即依次序启动文件 `conf/slaves` 中记载的主机上的进程，只不过伪分布式中 slave 为 `localhost`（即为自身），所以对于伪分布式 Hadoop，SSH 一样是必须的。

2.1.1 安装 JDK 1.6

安装 JDK 的过程很简单，下面以 Ubuntu 为例。

（1）下载和安装 JDK

确保可以连接到互联网，输入命令：

```
sudo apt-get install sun-java6-jdk
```

输入密码，确认，然后就可以安装 JDK 了。

这里先解释一下 `sudo` 与 `apt` 这两个命令，`sudo` 这个命令允许普通用户执行某些或全部需要 root 权限命令，它提供了详尽的日志，可以记录下每个用户使用这个命令做了些什么操作；同时 `sudo` 也提供了灵活的管理方式，可以限制用户使用命令。`sudo` 的配置文件为 `/etc/sudoers`。

`apt` 的全称为 the Advanced Packaging Tool，是 Debian 计划的一部分，是 Ubuntu 的软件包管理软件，通过 `apt` 安装软件无须考虑软件的依赖关系，可以直接安装所需要的软件，`apt` 会自动下载有依赖关系的包，并按顺序安装，在 Ubuntu 中安装有 `apt` 的一个图形化界面程序 `synaptic`（中文译名为“新立得”），大家如果有兴趣也可以使用这个程序来安装所需要的软

件。(如果大家想了解更多,可以查看一下关于 Debian 计划的资料。)

(2) 配置环境变量

输入命令:

```
sudo gedit /etc/profile
```

输入密码,打开 profile 文件。

在文件的最下面输入如下内容:

```
#set Java Environment
export JAVA_HOME=(你的JDK安装位置,一般为/usr/lib/jvm/java-6-sun)
export CLASSPATH=".:$JAVA_HOME/lib:$CLASSPATH"
export PATH="$JAVA_HOME/bin:$PATH"
```

这一步的意义是配置环境变量,使你的系统可以找到 JDK。

(3) 验证 JDK 是否安装成功

输入命令:

```
java -version
```

查看信息:

```
java version "1.6.0_14"
Java(TM) SE Runtime Environment (build 1.6.0_14-b08)
Java HotSpot(TM) Server VM (build 14.0-b16, mixed mode)
```

2.1.2 配置 SSH 免密码登录

同样以 Ubuntu 为例,假设用户名为 u。

1) 确认已经连接上互联网,输入命令

```
sudo apt-get install ssh
```

2) 配置为可以无密码登录本机。

首先查看在 u 用户下是否存在 .ssh 文件夹(注意 ssh 前面有“.”,这是一个隐藏文件夹),输入命令:

```
ls -a /home/u
```

一般来说,安装 SSH 时会自动在当前用户下创建这个隐藏文件夹,如果没有,可以手动创建一个。

接下来,输入命令:

```
ssh-keygen -t dsa -P '' -f ~/.ssh/id_dsa
```

解释一下,ssh-keygen 代表生成密钥;-t(注意区分大小写)表示指定生成的密钥类型;dsa 是 dsa 密钥认证的意思,即密钥类型;-P 用于提供密码;-f 指定生成的密钥文件。(关于密钥密码的相关知识这里就不详细介绍了,里面会涉及 SSH 的一些知识,如果读者有兴趣,可以自行查阅资料。)

在 Ubuntu 中, ~ 代表当前用户文件夹, 这里即 /home/u。

这个命令会在 .ssh 文件夹下创建两个文件 id_dsa 及 id_dsa.pub, 这是 SSH 的一对私钥和公钥, 类似于钥匙及锁, 把 id_dsa.pub (公钥) 追加到授权的 key 里面去。

输入命令:

```
cat ~/.ssh/id_dsa.pub >> ~/.ssh/authorized_keys
```

这段话的意思是把公钥加到用于认证的公钥文件中, 这里的 authorized_keys 是用于认证的公钥文件。

至此无密码登录本机已设置完毕。

3) 验证 SSH 是否已安装成功, 以及是否可以无密码登录本机。

输入命令:

```
ssh -version
```

显示结果:

```
OpenSSH_5.1p1 Debian-6ubuntu2, OpenSSL 0.9.8g 19 Oct 2007
Bad escape character 'rsion'.
```

显示 SSH 已经安装成功了。

输入命令:

```
ssh localhost
```

会有如下显示:

```
The authenticity of host 'localhost (::1)' can't be established.
RSA key fingerprint is 8b:c3:51:a5:2a:31:b7:74:06:9d:62:04:4f:84:f8:77.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'localhost' (RSA) to the list of known hosts.
Linux master 2.6.31-14-generic #48-Ubuntu SMP Fri Oct 16 14:04:26 UTC 2009 i686
```

```
To access official Ubuntu documentation, please visit:
http://help.ubuntu.com/
```

```
Last login: Mon Oct 18 17:12:40 2010 from master
admin@Hadoop:~$
```

这说明已经安装成功, 第一次登录时会询问你是否继续链接, 输入 yes 即可进入。

实际上, 在 Hadoop 的安装过程中, 是否无密码登录是无关紧要的, 但是如果不配置无密码登录, 每次启动 Hadoop, 都需要输入密码以登录到每台机器的 DataNode 上, 考虑到一般的 Hadoop 集群动辄数百台或上千台机器, 因此一般来说都会配置 SSH 的无密码登录。

2.1.3 安装并运行 Hadoop

介绍 Hadoop 的安装之前, 先介绍一下 Hadoop 对各个节点的角色定义。

Hadoop 分别从三个角度将主机划分为两种角色。第一，划分为 master 和 slave，即主人与奴隶；第二，从 HDFS 的角度，将主机划分为 NameNode 和 DataNode（在分布式文件系统中，目录的管理很重要，管理目录的就相当于主人，而 NameNode 就是目录管理者）；第三，从 MapReduce 的角度，将主机划分为 JobTracker 和 TaskTracker（一个 job 经常被划分为多个 task，从这个角度不难理解它们之间的关系）。

Hadoop 有官方发行版与 cloudera 版，其中 cloudera 版是 Hadoop 的商用版本，这里先介绍 Hadoop 官方发行版的安装方法。

Hadoop 有三种运行方式：单节点方式、单机伪分布方式与集群方式。乍看之下，前两种方式并不能体现云计算的优势，在实际应用中并没有什么意义，但是在程序的测试与调试过程中，它们还是很有意义的。

你可以通过以下地址获得 Hadoop 的官方发行版：

<http://www.apache.org/dyn/closer.cgi/Hadoop/core/>

下载 Hadoop-0.20.2.tar.gz 并将其解压，这里会解压到用户目录下，一般为：`/home/[你的用户名]/`。

□ 单节点方式配置：

安装单节点的 Hadoop 无须配置，在这种方式下，Hadoop 被认为是一个单独的 Java 进程，这种方式经常用来调试。

□ 伪分布式配置：

你可以把伪分布式的 Hadoop 看做是只有一个节点的集群，在这个集群中，这个节点既是 master，也是 slave；既是 NameNode 也是 DataNode；既是 JobTracker，也是 TaskTracker。

伪分布式的配置过程也很简单，只需要修改几个文件，如下所示。

进入 conf 文件夹，修改配置文件：

```
Hadoop-env.sh:
export JAVA_HOME="你的 JDK 安装地址"
```

指定 JDK 的安装位置：

```
conf/core-site.xml:
<configuration>
  <property>
    <name>fs.default.name</name>
    <value>hdfs://localhost:9000</value>
  </property>
</configuration>
```

这是 Hadoop 核心的配置文件，这里配置的是 HDFS 的地址和端口号。

```
conf/hdfs-site.xml:
<configuration>
  <property>
    <name>dfs.replication</name>
```




```

        <value>1</value>
    </property>
</configuration>

```

这是 Hadoop 中 HDFS 的配置，配置的备份方式默认为 3，在单机版的 Hadoop 中，需要将其改为 1。

```

conf/mapred-site.xml:
<configuration>
    <property>
        <name>mapred.job.tracker</name>
        <value>localhost:9001</value>
    </property>
</configuration>

```

这是 Hadoop 中 MapReduce 的配置文件，配置的是 JobTracker 的地址和端口。

需要注意的是，如果安装的是 0.20 之前的版本，那么只有一个配置文件，即为 Hadoop-site.xml。

接下来，在启动 Hadoop 前，需格式化 Hadoop 的文件系统 HDFS（这点与 Windows 是一样的，重新分区后的卷总是需要格式化的）。进入 Hadoop 文件夹，输入下面的命令：

```
bin/Hadoop NameNode -format
```

格式化文件系统，接下来启动 Hadoop。

输入命令：

```
bin/start-all.sh (全部启动)
```

最后，验证 Hadoop 是否安装成功。

打开浏览器，分别输入网址：

```

http://localhost:50030 (MapReduce 的 Web 页面)
http://localhost:50070 (HDFS 的 Web 页面)

```

如果都能查看，说明 Hadoop 已经安装成功。

对于 Hadoop 来说，安装 MapReduce 及 HDFS 都是必须的，但是如果有必要，你依然可以只启动 HDFS (start-dfs.sh) 或 MapReduce (start-mapred.sh)。

关于完全分布式的 Hadoop 会在 2.3 节详述。

2.2 在 Windows 上安装与配置 Hadoop

相对于 Linux，Windows 版本的 JDK 安装过程更容易，你可以在 http://www.java.com/zh_CN/download/manual.jsp 下载到最新版本的 JDK。这里再次申明，Hadoop 的编译及 MapReduce 程序的运行，很多地方都需要使用 JDK 的相关工具，因此只安装 JRE 是不够的。

安装过程十分简单，运行即可，程序会自动配置环境变量（在之前的版中还没有这项功能，新版本的 JDK 中已经可以自动配置环境变量了）。

2.2.1 安装 Cygwin

Cygwin 是在 Windows 平台下模拟 Unix 环境的一个工具，只有通过它才可以在 Windows 环境下安装 Hadoop。可以通过这个链接下载 Cygwin：

<http://www.cygwin.cn/setup.exe>

双击运行安装程序，选择 `install from internet`。

根据网络状况，选择合适的源下载程序。

进入 `select packages` 界面，然后进入 `Net`，勾选 `openssl` 及 `openssh`（如图 2-1 所示）。

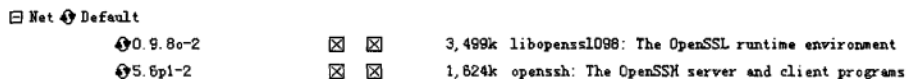


图 2-1 勾选 openssl 及 openssh

如果打算在 Eclipse 上编译 Hadoop，还必须安装“Base Category”下的“sed”（如图 2-2 所示）。



图 2-2 勾选 sed

另外建议安装“Editors Category”下的“vim”，以便在 Cygwin 上直接修改配置文件。

2.2.2 配置环境变量

依次点击我的电脑→属性→高级系统设置→环境变量，修改环境变量里的 path 设置，在其后添加 Cygwin 的 bin 目录和 Cygwin 的 usr\bin 目录。

2.2.3 安装和启动 sshd 服务

点击桌面上的 Cygwin 图标，启动 Cygwin，执行 `ssh-host-config` 命令，当要求输入 Yes/No 时，选择输入 No。当看到“Have fun”时，表示 sshd 服务安装成功。

在桌面上的“我的电脑”图标上右击，点击“管理”菜单，启动 CYGWIN sshd 服务。

2.2.4 配置 SSH 免密码登录

执行 `ssh-keygen` 命令生成密钥文件。按如下命令生成 `authorized_keys` 文件：

```
cd ~/.ssh/
cp id_rsa.pub authorized_keys
```

完成上述操作后, 执行 `exit` 命令先退出 Cygwin 窗口, 如果不执行这一步操作, 下面的操作可能会遇到错误。

接下来, 重新运行 Cygwin, 执行 `ssh localhost` 命令, 在第一次执行时会有提示, 然后输入 `yes`, 直接回车即可。

另外, 在 Windows 上安装 Hadoop 的过程与 Linux 一样, 这里就不再赘述了。

2.3 安装和配置 Hadoop 集群

2.3.1 网络拓扑

通常来说, 一个 Hadoop 的集群体系结构由两层网络拓扑组成, 如图 2-1 所示。结合实际的应用来看, 每个机架中会有 30 ~ 40 台机器, 这些机器共享一个 1GB 带宽的网络交换机。在所有的机架之上还有一个核心交换机或路由器, 通常来说其网络交换能力为 1GB 或更高。可以很明显地看出, 同一个机架中机器节点之间的带宽资源肯定要比不同机架中机器节点间丰富。这也是 Hadoop 随后设计数据读写分发策略要考虑的一个重要因素。

2.3.2 定义集群拓扑

在实际应用中, 为了使 Hadoop 集群获得更高的性能, 读者需要配置集群使 Hadoop 能够感知其所在的网络拓扑结构。当然如果集群中机器数量很少, 而且它们存在于一个机架中, 那么就不用做太多额外的工作, 而当集群中存在多个机架时, 就要使 Hadoop 清晰地知道每台机器所在的机架。随后, 在处理 MapReduce 任务时, Hadoop 会优先选择在机架内部进行数据传输, 而不是在机架间, 这样就可以更充分地使用网络带宽资源。同时, HDFS 可以更加智能地部署数据副本, 并在性能和可靠性间寻找到最优的平衡。

在 Hadoop 中, 网络的拓扑结构、机器节点及机架的网络位置定位都是通过树结构来描述的。通过它来确定节点间的距离, 这个距离是 Hadoop 做决策判断时的参考因素。NameNode 也是通过这个距离来决定应该把数据副本放到哪里的。当一个 map 任务到达时, 它会被分配到一个 TaskTracker 上运行, JobTracker 节点则会使用网络位置来确定 map 任务执行的机器节点。

在图 2-3 中, 笔者使用树结构来描述网络拓扑结构, 主要包括两个网络位置: 交换机 / 机架 1 和交换机 / 机架 2。因为在图中的集群只有一个最高级别的交换机, 所以此网络拓扑可简化描述为 / 机架 1 和 / 机架 2。

在配置 Hadoop 时, Hadoop 会确定节点地址和其网络位置的映射, 此映射在代码中通过 Java 接口 `DNSToSwitchMapping` 实现, 代码如下:

```
public interface DNSToSwitchMapping {
    public List<String> resolve(List<String> names);
}
```

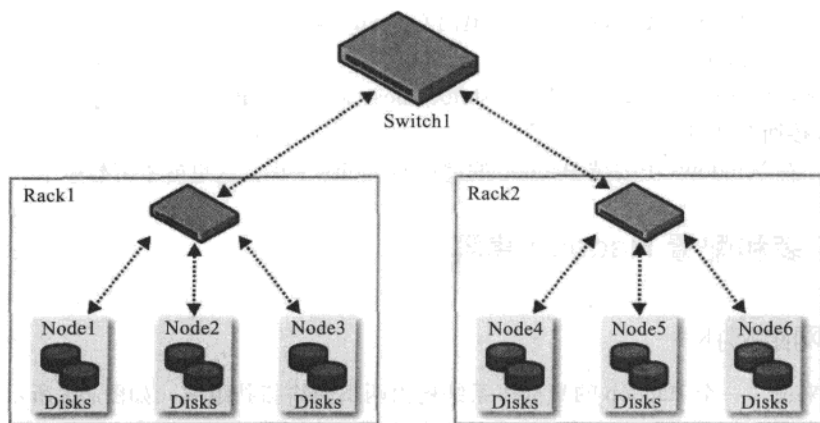


图 2-1 Hadoop 的网络拓扑结构

其中参数 `names` 是 IP 地址的一个 List 数据，这个函数返回的值为对应网络位置的字符串列表。在 `opology.node.switch.mapping.impl` 的配置参数中定义了一个 `DNSToSwitchMapping` 接口的实现，`NameNode` 通过它确定完成任务的机器节点所在的网络位置。

在图 2-1 的实例中，可以将节点 1、节点 2、节点 3 映射到 / 机架 1 中，节点 4、节点 5、节点 6，映射到 / 机架 2 中。事实上在实际应用中，管理员可能不需要手动做额外的工作去配置这些映射关系，系统有一个默认的接口实现 `ScriptBasedMapping`。它可以运行用户自定义的一个脚本区完成映射，如果用户没有定义，它则会将所有的机器节点映射到一个单独的网络位置中默认的机架上。如果用户定义了映射，则这个脚本的位置由 `topology.script.file.name` 的属性控制。脚本必须获取一批主机的 IP 地址作为参数进行映射，同时生成一个标准的网络位置给输出。

2.3.3 建立和安装 Cluster

想要建立 Hadoop 集群，首先要做的就是选择并购买机器，在机器到手之后，就要进行网络部署并安装软件了。安装和配置 Hadoop 有很多方法，这一部分内容在前文已经详细讲解（见 2.1 节和 2.2 节），同时还告诉了读者在实际部署时应该考虑的情况。

为了简化我们在每个机器节点上安装和维护相同软件的过程，通常会采用自动安装法，比如 Red Hat Linux 下的 Kickstart 或 Debian 的全程自动化安装。这些工具先会记录你的安装过程，以及你对选择项的选择，然后根据记录来自动安装软件。同时它会在每个进程的结尾提供一个钩子执行脚本，在对那些不包含在标准安装中的最终系统进行调整和自定义时这是非常有用的。

下面我们将具体介绍一下如何部署和配置 Hadoop。Hadoop 为了应对不同的使用需求（不管是开发、实际应用还是研究），有着不同的运行方式，包括单节点、单机伪分布、集群等。前面已经详细介绍了 Hadoop 在 Windows 和 Linux 平台下的安装与配置。下面将对

Hadoop 的分布式配置做具体的介绍。

1. Hadoop 集群的配置

在配置伪分布式的过程中，大家也许会觉得 Hadoop 的配置很简单，但那只是最基本的配置。

Hadoop 的配置文件分为两类：

❑ 只读类型的默认文件：src/core/core-default.xml、src/hdfs/hdfs-default.xml、src/mapred/mapred-default.xml、conf/mapred-queues.xml。

❑ 定位 (site-specific) 设置：conf/core-site.xml、conf/hdfs-site.xml、conf/mapred-site.xml、conf/mapred-queues.xml。

除此之外，也可以通过设置 conf/Hadoop-env.sh 来为 Hadoop 的守护进程设置环境变量（在 bin/ 文件夹内）。

Hadoop 是通过 org.apache.Hadoop.conf.Configuration 来读取配置文件的，在 Hadoop 的设置中，Hadoop 的配置是通过资源 (resource) 定位的，每个资源由一系列 name/value 对以 XML 文件的形式构成，它以一个字符串命名或以 Hadoop 定义的 Path 类命名（这个类是用于定义文件系统中的文件或文件夹的）。如果是以字符串命名的，Hadoop 会通过 classpath 调用此文件。如果是以 Path 类命名的，那么 Hadoop 会直接在本地文件系统中搜索文件。

资源设定有两个特点，如下所示。

❑ Hadoop 允许定义最终参数 (final parameters)，如果任意资源声明了 final 这个值，那么之后加载的任何资源都不能改变这个值，定义最终资源的格式是这样的：

```
<property>
  <name>dfs.client.buffer.dir</name>
  <value>/tmp/Hadoop/dfs/client</value>
  <final>true</final>    // 注意这个值
</property>
```

❑ Hadoop 允许参数传递，如下所示，当 tempdir 被调用时，basedir 会作为值被调用。

```
<property>
  <name>basedir</name>
  <value>/user/${user.name}</value>
</property>

<property>
  <name>tempdir</name>
  <value>${basedir}/tmp</value>
</property>
```

刚才提到，读者可以通过设置 conf/Hadoop-env.sh 为 Hadoop 的守护进程设置环境变量。一般来说，大家至少需要在这里设置在主机上安装 JDK 的位置 (JAVA_HOME)，以使 Hadoop 找到 JDK。大家也可以在这里通过 HADOOP_*_OPTS 对不同的守护进程分别进行设置，如表 2-1 所示。

表 2-1 Hadoop 的守护进程配置表

守护进程 (Daemon)	配置选项 (Configure Options)
NameNode	HADOOP_NameNode_OPTS
DataNode	HADOOP_DataNode_OPTS
SecondaryNameNode	HADOOP_SECONDARYNameNode_OPTS
JobTracker	HADOOP_JOBTRACKER_OPTS
TaskTracker	HADOOP_TASKTRACKER_OPTS

例如，如果想设置 NameNode 使用 parallelGC，那么可以这样写：

```
export HADOOP_NameNode_OPTS="-XX:+UseParallelGC ${HADOOP_NameNode_OPTS}"
```

在这里也可以进行其他设置，比如设置 Java 的运行环境 (HADOOP_OPTS)，设置日志文件的存放位置 (HADOOP_LOG_DIR)，或者 SSH 的配置 (HADOOP_SSH_OPTS) 等。

关于 conf/core-site.xml、conf/hdfs-site.xml、conf/mapred-site.xml 的配置可查看表 2-2、表 2-3 和表 2-4。

表 2-2 conf/core-site.xml 的配置

参数 (Parameter)	值 (Value)
fs.default.name	NameNode 的 IP 地址及端口

表 2-3 conf/hdfs-site.xml 的配置

参数 (Parameter)	值 (Value)
dfs.name.dir	NameNode 存储名字空间及汇报日志的位置
dfs.data.dir	DataNode 存储数据块的位置

表 2-4 conf/mapred-site.xml 的配置

参数 (Parameter)	值 (Value)
mapreduce.jobtracker.address	JobTracker 的 IP 地址及端口
mapreduce.jobtracker.system.dir	MapReduce 在 HDFS 上存储文件的位置，例如 /Hadoop/mapred/system/
mapreduce.cluster.local.dir	MapReduce 的缓存数据存储在文件系统上的位置
mapred.tasktracker.{map/reduce}.tasks.maximum	每台 TaskTracker 所能运行的 Map 或 Reduce 的 task 最大数量
dfs.hosts/dfs.hosts.exclude	允许或禁止的 DataNode 列表
mapreduce.jobtracker.hosts.filename/ mapreduce.jobtracker.hosts.exclude.filename	允许或禁止的 TaskTrackers 列表
mapreduce.cluster.job-authorization-enabled	布尔类型，标志着 job 存取控制列表是否支持对 job 的观察和修改

配置并不复杂，一般而言，除了规定端口、IP 地址、文件的存储位置外，其他配置都不是必须修改的，可以根据需要决定是采用默认配置还是自己修改。还有一点需要注意的是，

以上配置都被默认为最终参数（final parameters），这些参数都不可以在程序中再次修改。

接下来可以看一下 conf/mapred-queues.xml 的配置列表（如表 2-5 所示）。

表 2-5 conf/mapred-queues.xml 的配置

标签或属性（Tag/Attribute）	值（Value）	是否可刷新
queues	配置文件的根元素	无意义
aclsEnabled	布尔类型 <queues> 标签的属性，标志着存取控制列表是否支持控制 job 的提交及所有 queues 的管理	是
queue	<queues> 的子元素，定义系统中的 queue	无意义
name	<queue> 的子元素，代表名字	否
state	<queue> 的子元素，代表 queue 的状态	是
acl-submit-job	<queue> 的子元素，定义一个能提交 job 到该 queue 的用户或组的名单列表	是
acl-administer-job	<queue> 的子元素，定义一个能更改 job 的优先级或能杀死已提交到该 queue 的 job 用户或组的名单列表	是
properties	<queues> 的子元素，定义优先调度规则	无意义
property	<properties> 的子元素	无意义
key	<property> 的子元素	调度程序指定
value	<property> 的属性	调度程序指定

相信大家能猜出上表中的 mapred-queues.xml 文件是用来做什么的，这个文件就是用来设置 MapReduce 系统的队列顺序的。queues 是 JobTracker 中的一个抽象概念，可以在一定程度上管理 job，因此它为管理员提供了一种管理 job 的方式。这种控制是常见且有效的，例如通过这种管理可以把不同的用户划分为不同的组，或分别赋予他们不同的级别，并且会优先执行高级别用户提交的 job。

按照这个思路，很容易想到三种原则：

- ☐ 同一类用户提交的 job 统一提交到同一个 queue 中；
- ☐ 运行时间较长的 job 可以提交到同一个 queue 中；
- ☐ 把很快就能运行完成的 job 划分到一个 queue 中，并且限制好 queue 中 job 的数量上限。

queues 的有效性很依赖在 JobTracker 中通过 mapreduce.jobtracker.taskscheduler 设置的调度规则（scheduler），一些调度算法可能只需要一个 queue，不过有些调度算法可能会很复杂，需要设置很多 queue。

queues 的大部分设置的更改都不需要重新启动 MapReduce 系统就可以生效，不过也有一些需要重启系统的，具体可见表 2-5。

conf/mapred-queues.xml 的文件配置与其他文件略有不同，配置格式如下所示：

```
<queues aclsEnabled="$aclsEnabled">
  <queue>
    <name>$queue-name</name>
    <state>$state</state>
  </queue>
```

```

<name>$child-queue1</name>
<properties>
  <property key="$key" value="$value"/>
  ...
</properties>
<queue>
  <name>$grand-child-queue1</name>
  ...
</queue>
</queue>
<queue>
  <name>$child-queue2</name>
  ...
</queue>
...
...
...
<queue>
  <name>$leaf-queue</name>
  <acl-submit-job>$acls</acl-submit-job>
  <acl-administer-jobs>$acls</acl-administer-jobs>
  <properties>
    <property key="$key" value="$value"/>
    ...
  </properties>
</queue>
</queue>
</queues>

```

以上这些就是 Hadoop 配置的主要内容，其他还有一些诸如内存配置方面的信息，如有兴趣可以参阅官方的配置文档。

2. 一个具体的配置

为了方便阐述，这里只搭建一个有三台主机的小集群。

相信读者还没有忘记 Hadoop 对主机的三种定位方式，分别为 master 和 slave，JobTracker 和 TaskTracker，NameNode 和 DataNode。为了方便，在分配 IP 地址时顺便规定一下角色。

下面是为这三台机器分配的 IP 地址及相应的角色：

```

10.37.128.2-master,NameNode,jobtracker-master (主机名)
10.37.128.3-slave,DataNode,tasktracker-slave1 (主机名)
10.37.128.4-slave,DataNode,tasktracker-slave2 (主机名)

```

首先在三台主机上创建相同的用户（这是 Hadoop 的基本要求）：

- 1) 在三台主机上安装 JDK 1.6，并设置环境变量。
- 2) 在这三台主机上安装 OpenSSH，并配置 SSH 可以无密码登录。

安装方式不再赘述，建立 ~/.ssh 文件夹，如已存在，则无须创建。“~”代表 Ubuntu 的

当前用户文件夹。

生成密钥并配置 SSH 无密码登录本机，输入命令：

```
ssh-keygen -t dsa -P '' -f ~/.ssh/id_dsa
cat ~/.ssh/id_dsa.pub >> ~/.ssh/authorized_keys
```

将文件拷贝到两台 slave 主机相同的文件夹内，输入命令：

```
scp authorized_keys slave1:~/.ssh/
scp authorized_keys slave2:~/.ssh/
```

查看是否可以从 master 主机无密码登录 slave，输入命令：

```
ssh slave1
ssh slave2
```

3) 在三台主机上分别设置 /etc/hosts 及 /etc/hostname。

hosts 这个文件用于定义主机名与 IP 地址之间的对应关系。

/etc/hosts:

```
127.0.0.1 localhost
10.37.128.2 master
10.37.128.3 slave1
10.37.128.4 slave2
```

hostname 这个文件用于定义 Ubuntu 的主机名。

/etc/hostname:

你的主机名（如 master, slave1 等）

4) 配置三台主机的 Hadoop 文件，内容如下：

conf/Hadoop-env.sh:

```
export JAVA_HOME="你的Java安装地址"
```

conf/core-site.xml:

```
<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>

<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
  <name>fs.default.name</name>
  <value>hdfs://master:9000</value>
</property>
<property>
  <name>Hadoop.tmp.dir</name>
  <value>你希望Hadoop存储数据块的位置</value> // 此文件夹需手动创建
</property>
```

```
</configuration>
```

conf/hdfs-site.xml:

```
<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>

<!-- Put site-specific property overrides in this file. -->

<configuration>
  <property>
    <name>dfs.replication</name>
    <value>2</value>
  </property>
</configuration>
```

conf/mapred-site.xml:

```
<?xml version="1.0"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>

<!-- Put site-specific property overrides in this file. -->

<configuration>
  <property>
    <name>mapred.job.tracker</name>
    <value>master:9001</value>
  </property>
</configuration>
```

conf/masters:

```
master
```

conf/slaves:

```
slave1
slave2
```

5) 启动 Hadoop。

```
bin/Hadoop NameNode -format
bin/start-all.sh
```

你可以通过命令:

```
Hadoop dfsadmin -report
```

查看集群状态, 或者通过 <http://master:50070> 及 <http://master:50030> 查看集群状态。

2.4 日志分析及几个小技巧

如果大家在安装的时候遇到问题, 或者按步骤安装完后却不能运行 Hadoop, 那么建议

仔细查看日志信息，Hadoop 记录了详尽的日志信息，日志文件保存在 logs 文件夹内。

无论是启动，还是以后会经常用到的 MapReduce 中的每一个 job，以及 HDFS 等相关信息，Hadoop 均存有日志文件以供分析。

例如：

NameNode 和 DataNode 的 namespaceID 不一致，这个错误是很多人在安装时会遇到的，日志信息为：

```
java.io.IOException: Incompatible namespaceIDs in /root/tmp/dfs/data: NameNode
namespaceID = 1307672299; DataNode namespaceID = 389959598
```

若 HDFS 一直没有启动，读者可以查询日志，并通过日志进行分析，以上提示信息显示了 NameNode 和 DataNode 的 namespaceID 不一致。

这个问题一般是由于两次或两次以上的格式化 NameNode 造成的，有两种方法可以解决，第一种方法是删除 DataNode 的所有资料；第二种方法是修改每个 DataNode 的 namespaceID（位于 /dfs/data/current/VERSION 文件中）或修改 NameNode 的 namespaceID（位于 /dfs/name/current/VERSION 文件中），使其一致。

下面这两种方法在实际应用中也可能用到。

1) 重启坏掉的 DataNode 或 JobTracker。当 Hadoop 集群的某单个节点出现问题时，一般不必重启整个系统，只须重启这个节点，它会自动连入整个集群。

在坏死的节点上输入如下命令即可：

```
bin/Hadoop-daemon.sh start DataNode
bin/Hadoop-daemon.sh start jobtracker
```

2) 动态加入 DataNode 或 TaskTracker。这个命令允许用户动态将某个节点加入集群中。

```
bin/Hadoop-daemon.sh --config ./conf start DataNode
bin/Hadoop-daemon.sh --config ./conf start tasktracker
```

2.5 小结

本章主要讲解了 Hadoop 的安装和配置过程。Hadoop 的安装过程并不复杂，基本配置也简单明了，其中有以下几个关键点：

- ❑ Hadoop 主要是用 Java 语言编写的，它无法使用 Linux 预装的 OpenJDK，因此在安装 Hadoop 前要先安装 Oracle 公司的 JDK（版本要在 1.6 以上）；
- ❑ 作为分布式系统，Hadoop 需要通过 SSH 的方式启动处于 slave 上的程序，因此必须安装和配置 SSH。

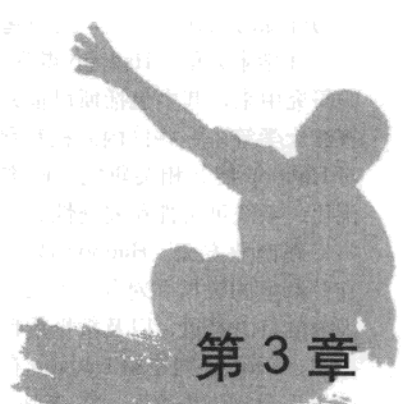
因此，在安装 Hadoop 前需要安装 JDK 和 SSH。

在 Windows 系统上安装 Hadoop 与在 Linux 系统上安装有一点不同，就是在 Windows 系统上需要通过 Cygwin 模拟 Linux 环境，而 SSH 的安装也需要在安装 Cygwin 时进行选择，

请不要忘了这一点。

集群配置只须记住 `conf/Hadoop-env.sh`、`conf/core-site.xml`、`conf/hdfs-site.xml`、`conf/mapred-site.xml`、`conf/mapred-queues.xml` 这 5 个文件的作用即可，另外 Hadoop 有些配置是可以在程序中修改的，这部分内容不是本章的重点，因此没有详细说明。





第 3 章

Hadoop 应用案例分析

本章内容

- ☐ Hadoop 在 Yahoo! 的应用
- ☐ Hadoop 在 eBay 的应用
- ☐ Hadoop 在百度的应用
- ☐ Hadoop 在 Facebook 的应用
- ☐ Hadoop 平台上的海量数据排序
- ☐ 小结

资源库
PDG

随着企业的数据量的迅速增长, 存储和处理大规模数据已成为企业的迫切需求。Hadoop 作为开源的云计算平台, 已引起了学术界和企业的普遍兴趣。

在学术方面, Hadoop 得到了各科研院所的广泛关注, 多所著名大学加入到 Hadoop 集群的研究中来, 其中包括斯坦福大学、加州大学伯克利分校、康奈尔大学、卡耐基·梅隆大学、普渡大学等。一些国内高校和科研院所如中科院计算所、清华大学、中国人民大学等也开始对 Hadoop 展开相关研究, 研究内容涉及 Hadoop 的数据存储、资源管理、作业调度、性能优化、系统可用性和安全性等多个方面。

在商业方面, Hadoop 技术已经在互联网领域得到了广泛的应用。互联网公司往往需要存储海量的数据并对其进行处理, 而这正是 Hadoop 的强项。如 Facebook 使用 Hadoop 存储内部的日志拷贝, 以及数据挖掘和日志统计; Yahoo! 利用 Hadoop 支持广告系统并处理网页搜索; Twitter 则使用 Hadoop 存储微博数据、日志文件和其他中间数据等。在国内, Hadoop 同样也得到了许多公司的青睐, 如百度主要将 Hadoop 应用于日志分析和网页数据库的数据挖掘; 阿里巴巴则将 Hadoop 用于商业数据的排序和搜索引擎的优化等。

下面我们将选取具有代表性的 Hadoop 应用案例进行分析, 让读者了解 Hadoop 在企业界的应用情况。

3.1 Hadoop 在 Yahoo! 的应用

关于 Hadoop 技术的研究和应用, Yahoo! 始终处于领先地位, 它将 Hadoop 应用于自己的各种产品中, 包括数据分析、内容优化、反垃圾邮件系统、广告和优化选择、大数据处理和 ETL 等; 同样, 在用户兴趣预测、搜索排名、广告定位等方面得到了充分的应用。

在 Yahoo! 主页个性化方面, 实时服务系统通过 Apache 从数据库中读取 user 到 interest 的映射, 并且每隔 5 分钟生产环境中的 Hadoop 集群就会基于最新数据重新排列内容, 每隔 7 分钟则在页面上更新内容。

在邮箱方面, Yahoo! 利用 Hadoop 集群根据垃圾邮件模式为邮件计分, 并且每隔几个小时就在集群上改进反垃圾邮件模型, 集群系统每天还可以推动 50 亿次的邮件投递。

目前 Hadoop 最大的生产应用是 Yahoo! 的 Search Webmap 应用, 它运行在超过 10 000 台机器的 Linux 系统集群里, Yahoo! 的网页搜索查询使用的就是它产生的数据。Webmap 的构建步骤如下: 首先进行网页的爬取, 同时产生包含所有已知网页和互联网站的数据库, 以及一个关于所有页面及站点的海量数据组; 然后将这些数据传输给 Yahoo! 搜索中心执行排序算法。在整个过程中, 索引中页面间的链接数量将会达到 1TB, 经过压缩的数据产出量会达到 300TB, 运行一个 MapReduce 任务就需使用超过 10 000 的内核, 而在生产环境中使用数据的存储量超过 5PB。

Yahoo! 在 Hadoop 中同时使用了 Hive 和 Pig, 在许多人看来, Hive 和 Pig 大体上相似而且 Pig Latin 与 SQL 也十分相似。那么 Yahoo! 为什么要同时使用这些技术呢? 主要是因为 Yahoo! 的研究人员在查看了它们的工作负载并分析了应用案例后认为不同的情况下需要

使用不同的工具。

先了解一下大规模数据的使用和处理背景。大规模的数据处理经常分为三个不同的任务：数据收集、数据准备和数据表示，这里并不打算介绍数据收集阶段，因为 Pig 和 Hive 主要用于数据准备和数据表示阶段。

数据准备阶段通常被认为是提取、转换和加载（Extract Transform Load, ETL）数据的阶段，或者认为这个阶段是数据工厂。这里的数据工厂只是一个类比，在现实生活中的工厂接收原材料后会生产出客户所需的产品，而数据工厂与之相似，它在接收原始数据后，可以输出供客户使用的数据集。这个阶段需要装载和清洗原始数据，并让它遵守特定的数据模型，还要尽可能地让它与其他数据源结合等。这一阶段的客户一般都是程序员、数据专家或研究者。

数据表示阶段一般指的都是数据仓库，数据仓库存储了客户所需要的产品，客户会根据需要选取合适的产品。这一阶段的客户可能是系统的数据工程师、分析师或决策者。

根据每个阶段负载和用户情况的不同，Yahoo！在不同的阶段使用不同的工具。结合了诸如 Oozie 等工作流系统的 Pig 特别适合于数据工厂，而 Hive 则适合于数据仓库。下面将分别介绍数据工厂和数据仓库。

Yahoo！的数据工厂存在三种不同的工作用途：流水线、迭代处理和科学研究。

经典的数据流水线包括数据反馈、清洗和转换。一个常见例子是 Yahoo！的网络服务器日志，这些日志需要进行清洗以去除不必要的信息，数据转换则是要找到点击之后所转到的页面。Pig 是分析大规模数据集的平台，它建立在 Hadoop 之上并提供了良好的编程环境、优化条件和可扩展的性能。Pig Latin 是关系型数据流语言，并且是 Pig 核心的一部分，基于以下的原因，Pig Latin 相比于 SQL 而言，更适合构建数据流。首先，Pig Latin 是面向过程的，并且 Pig Latin 允许流水线开发者自定义流水线中检查点的位置；其次，Pig Latin 允许开发者直接选择特定的操作实现方式而不是依赖于优化器；最后，Pig Latin 支持流水线的分支，并且 Pig Latin 允许流水线开发者在数据流水线的任何地方插入自己的代码。Pig 和诸如 Oozie 等工作流工具一起使用来创建流水线，一天可以运行数以万计的 Pig 作业。

迭代处理也是需要 Pig 的，在这种情况下通常需要维护一个大规模的数据集。数据集上的典型处理包括加入一小片数据后就会改变大规模数据集的状态。如考虑这样一个数据集，它存储了 Yahoo！新闻中现有的所有新闻。我们可以把它想象成一幅巨大的图，每个新闻就是一个节点，新闻节点若有边相连则说明这些新闻指的是同一个事件。每隔几分钟就会有新的新闻加入进来，这些工具需要将这些新闻节点加到图中，并找到相似的新闻节点用边连接起来，还要删除被新节点覆盖的旧节点。这和标准流水线不同的是它不断有小变化，这就需要增长处理模型在合理的时间范围内处理这些数据了。例如，所有的新节点加入图中后，又有一批新的新闻节点到达，在整个图上重新执行连接操作是不现实的，这也许会花费数个小时。相反，在新增加的节点上执行连接操作并使用全连接（full join）的结果是可行的，而且这个过程只需要花费几分钟时间。标准的数据库操作可以使用 Pig Latin 通过上述方式实现，这时 Pig 就会得到很好的应用。

Yahoo! 有许多的科研人员, 他们需要用网格工具处理千万亿大小的数据, 还有许多研究人员希望快速地写出脚本来测试自己的理论或获得更深的理解。但是在数据工厂中, 数据不是以一种友好的、标准的方式呈现的, 这时 Pig 就可以大显身手了, 因为它可以处理未知模式的数据, 还有半结构化和非结构化的数据。Pig 与 streaming 相结合使得研究者在小规模数据集上测试的 Perl 和 Python 脚本可以很方便地在大规模数据集上运行。

在数据仓库处理阶段, 有两个主要的应用: 商业智能分析和特定查询 (Ad-hoc query)。在第一种情况下, 用户将数据连接到商业智能 (BI) 工具 (如 MicroStrategy) 上来产生报告或深入的分析。在第二种情况下, 用户执行数据分析师或决策者的特定查询。这两种情况下, 关系模型和 SQL 都很好用。事实上, 数据仓库已经成为 SQL 使用的核心, 它支持多种查询并具有分析师所需的工具, Hive 作为 Hadoop 的子项目为其提供了 SQL 接口和关系模型, 现在 Hive 团队正开始将 Hive 与 BI 工具通过接口 (如 ODBC) 结合起来使用。

Pig 在 Yahoo! 得到了广泛应用, 这使得数据工厂的数据被移植到 Hadoop 上运行成为可能。随着 Hive 的深入使用, Yahoo! 打算将数据仓库移植到 Hadoop 上。在同一系统上部署数据工厂和数据仓库将会降低数据加载到仓库的时间, 这也使得共享工厂和仓库之间的数据、管理工具、硬件等成为可能。Yahoo! 在 Hadoop 上同时使用多种工具使 Hadoop 能够执行更多的数据处理。

3.2 Hadoop 在 eBay 的应用

在 eBay 上存储着上亿种商品的信息, 而且每天有数百万种的新商品增加, 因此需要用云系统来存储和处理 PB 级别的数据, 而 Hadoop 则是个很好的选择。

Hadoop 是建立在商业硬件上的容错、可扩展、分布式的云计算框架, eBay 利用 Hadoop 建立了一个大规模的集群系统——Athena, 它被分为五层 (如图 3-1 所示), 下面从最底层向上开始介绍:

1) Hadoop 核心层, 包括 Hadoop 运行时环境、一些通用设施和 HDFS, 其中文件系统为读写大块数据而做了一些优化, 如将块的大小由 128MB 改为 256MB。

2) MapReduce 层, 为开发和执行任务提供 API 和控件。

3) 数据获取层, 现在数据获取层的主要框架是 HBase、Pig 和 Hive:

□ HBase 是根据 Google BigTable 开发的按列存储的多维空间数据库, 通过维护数据的划分和范围提供有序的数据, 其数据储存在 HDFS 上。

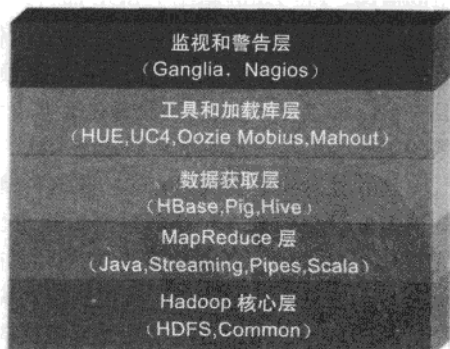


图 3-1 Athena 的层次

□ Pig (Latin) 是提供加载、筛选、转换、提取、聚集、连接、分组等操作的面向过程的语言, 开发者使用 Pig 建立数据管道和数据工厂。

□ Hive 是用于建立数据仓库的使用 SQL 语法的声明性语言。对于开发者、产品经理和分析师来说, SQL 接口使得 Hive 成为很好的选择。

4) 工具和加载库层, UC4 是 eBay 从多个数据源自动加载数据的企业级调度程序。加载库有: 统计库 (R)、机器学习库 (Mahout)、数学相关库 (Hama) 和 eBay 自己开发的用于解析网络日志的库 (Mobius)。

5) 监视和警告层, Ganglia 是分布式集群的监视系统, Nagios 则用来警告一些关键事件如服务器不可达、硬盘已满等。

eBay 的企业服务器运行着 64 位的 RedHat Linux :

□ NameNode 负责管理 HDFS 的主服务器;

□ JobTracker 负责任务的协调;

□ HBaseMaster 负责存储 HBase 存储的根信息, 并且方便与数据块或存取区域进行协调;

□ ZooKeeper 是保证 HBase 一致性的分布式锁协调器。

用于存储和计算的节点是 1U 大小的运行 Cent OS 的机器, 每台机器拥有 2 个四核处理器和 2TB 大小的存储空间, 每 38 ~ 42 个节点单元为一个 rack, 这组建成了高密度网络。有关网络方面, 顶层 rack 交换机到节点的带宽为 1Gbps, rack 交换机到核心交换机的带宽为 40Gbps。

这个集群是 eBay 内多个团队共同使用的, 包括产品和一次性任务。这里使用 Hadoop 公平调度器 (Fair Scheduler) 来管理分配、定义团队的任务池、分配权限、限制每个用户和组的并行任务、设置优先权期限和延迟调度。

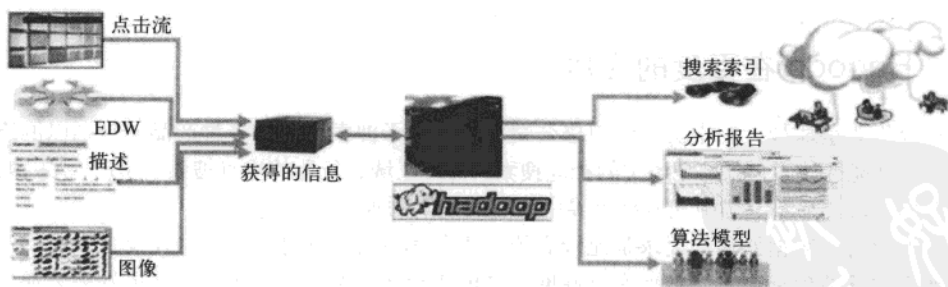


图 3-2 数据流

数据流的具体处理过程如图 3-2 所示, 系统每天需要处理 8TB 至 10TB 的新数据, 而 Hadoop 主要用于:

□ 基于机器学习的排序, 使用 Hadoop 计算需要考虑多个因素 (如价格、列表格式、卖家记录、相关性) 的排序函数, 并需要添加新因素来验证假设的扩展功能, 以增强

eBay 物品搜索的相关性。

- ❑ **对物品描述数据的挖掘**，在完全无人监管的方式下使用数据挖掘和机器学习技术将物品描述清单转化为与物品相关的键/值对，以扩大分类的覆盖范围。

eBay 的研究人员在系统构建和使用过程中遇到的挑战及一些初步计划有以下几个方面：

- ❑ **可扩展性**，当前主系统的 NameNode 拥有扩展的功能，随着集群的文件系统不断增长，需要存储大量的元数据，所以内存占有量也在不断增长。若是 1PB 的存储量则需要将近 1GB 的内存量，可能的解决方案是使用等级结构的命名空间划分，或者使用 HBase 和 ZooKeeper 联合对元数据进行管理。
- ❑ **有效性**，NameNode 的有效性对产品的工作负载很重要，开源社区提出了一些备用选择，如使用检查点和备份节点、从 Secondary NameNode 中转移到 Avatar 节点、日志元数据复制技术等。eBay 研究人员根据这些方法建立了自己的产品集群。
- ❑ **数据挖掘**，在存储非结构化数据的系统上建立支持数据管理、数据挖掘和模式管理的系统。新的计划提议将 Hive 的元数据和 Owl 添加到新系统中，并称为 Howl。eBay 研究人员努力将这个系统联系到分析平台上去，这样用户可以很容易地在不同的数据系统中挖掘数据。
- ❑ **数据移动**，eBay 研究人员考虑发布数据转移工具，这个工具可以支持在不同的子系统如数据仓库和 HDFS 之间进行数据的复制。
- ❑ **策略**，通过配额实现较好的归档、备份等策略（Hadoop 现有版本的配额需要改进）。eBay 的研究人员基于工作负载和集群的特点对不同的集群确定配额。
- ❑ **标准**，eBay 研究人员开发健壮的工具来为数据来源、消耗情况、预算情况、使用情况等进行度量。

同时 eBay 正在改变收集、转换、使用数据的方式，以提供更好的商业智能服务。

3.3 Hadoop 在百度的应用

百度作为全球最大的中文搜索引擎公司，提供基于搜索引擎的各种产品，包括以网络搜索为主的功能性搜索；以贴吧为主的社区搜索；针对区域、行业的垂直搜索、MP3 音乐搜索，以及百科等，几乎覆盖了中文网络世界中所有的搜索需求。

百度对海量数据处理的要求是比较高的，要在线下对数据进行分析，还要在规定的时间内处理完并反馈到平台上。百度在互联网领域的平台需求如图 3-3 所示，这里就需要通过性能较好的云平台进行处理了，Hadoop 就是很好的选择。在百度，Hadoop 主要应用于以下几个方面：

- ❑ 日志的存储和统计；
- ❑ 网页数据的分析和挖掘；
- ❑ 商业分析，如用户的行为和广告关注度等；
- ❑ 在线数据的反馈，及时得到在线广告的点击情况；

□ 用户网页的聚类，分析用户的推荐度及用户之间的关联度。

MapReduce 主要是一种思想，不能解决所有领域内与计算有关的问题，百度的研究人员认为比较好的模型应该如图 3-4 所示，HDFS 实现共享存储，一些计算使用 MapReduce 解决，一些计算使用 MPI 解决，而还有一些计算需要通过两者来共同处理。因为 MapReduce 适合处理数据很大且适合划分的数据，所以在处理这类数据时就可以用 MapReduce 做一些过滤，得到基本的向量矩阵，然后通过 MPI 进一步处理后返回结果，只有整合技术才能更好地解决问题。

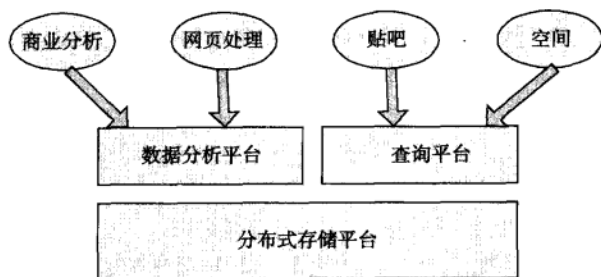


图 3-3 互联网领域的平台需求

百度现在拥有 3 个 Hadoop 集群，总规模在 700 台机器左右，其中有 100 多台新机器和 600 多台要淘汰的机器（它们的计算能力相当于 200 多台新机器），不过其规模还在不断的增加中。现在每天运行的 MapReduce 任务在 3000 个左右，处理数据约 120TB/天。

百度为了更好地用 Hadoop 进行数据处理，在以下几个方面做了改进和调整：

（1）调整 MapReduce 策略

- 限制作业处于运行状态的任务数；
- 调整预测执行策略，控制预测执行量，一些任务不需要预测执行；
- 根据节点内存状况进行调度；
- 平衡中间结果输出，通过压缩处理减少 I/O 负担。

（2）改进 HDFS 的效率和功能

- 权限控制，在 PB 级数据量的集群上数据应该是共享的，这样分析起来比较容易，但是需要对权限进行限制；
- 让分区与节点独立，这样，一个分区坏掉后节点上的其他分区还可以正常使用；
- 修改 DFSClient 选取块副本位置的策略，增加功能使 DFSClient 选取块时跳过出错的 DataNode；

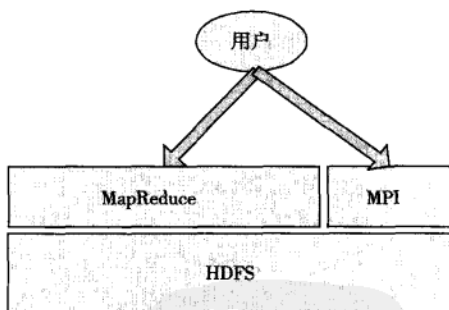


图 3-4 计算模型

❑ 解决 VFS (Virtual File System) 的 POSIX (Portable Operating System Interface of Unix) 兼容性问题。

(3) 修改 Speculative 的执行策略

- ❑ 采用速率倒数替代速率, 防止数据分布不均时经常不能启动预测执行情况的发生;
- ❑ 增加任务时必须达到某个百分比后才能启动预测执行的限制, 解决 reduce 运行等待 map 数据的时间问题;
- ❑ 只有一个 map 或 reduce 时, 可以直接启动预测执行。

(4) 对资源使用进行控制

- ❑ 对应用物理内存进行控制。如果内存使用过多会导致操作系统跳过一些任务, 百度通过修改 Linux 内核对进程使用的物理内存进行独立的限制, 超过阈值可以终止进程。
- ❑ 分组调度计算资源, 实现存储共享、计算独立, 在 Hadoop 中运行的进程是不可抢占的。
- ❑ 在大块文件系统中, X86 平台下一个页的大小是 4KB。如果页较小, 管理的数据就会很多, 会增加数据操作的代价并影响计算效率, 因此需要增加页的大小。

百度在使用 Hadoop 时也遇到了一些问题, 主要有:

- ❑ **MapReduce 的效率问题**: 比如, 如何在 shuffle 效率方面减少 I/O 次数以提高并行效率; 如何在排序效率方面设置排序为可配置的, 因为排序过程会浪费很多的计算资源, 而一些情况下是不需要排序的。
- ❑ **HDFS 的效率和可靠性问题**: 如何提高随机访问效率, 以及数据写入的实时性问题, 如果 Hadoop 每写一条日志就在 HDFS 上存储一次, 效率会很低。
- ❑ **内存使用的问题**: reducer 端的 shuffle 会频繁地使用内存, 这里采用类似 Linux 的 buddy system 来解决, 保证 Hadoop 用最小的开销达到最高的利用率; 当 Java 进程内容使用内存较多时, 可以调整垃圾回收 (GC) 策略; 有时存在大量的内存复制现象, 这会消耗大量 CPU 资源, 同时还会导致内存使用峰值极高, 这时需要减少内存的复制。
- ❑ **作业调度的问题**: 如何限制任务的 map 和 reduce 计算单元的数量, 以确保重要计算可以有足够的计算单元; 如何对 TaskTracker 进行分组控制, 以限制作业执行的机器, 同时还可以在用户提交任务时确定执行的分组并对分组进行认证。
- ❑ **性能提升的问题**: UserLogs cleanup 在每次 task 结束的时候都要查看一下日志, 以决定是否清除, 这会占用一定的任务资源, 可以通过将清理线程从子 Java 进程移到 TaskTracker 来解决; 子 Java 进程会对文本行进行切割而 map 和 reduce 进程则会重新切割, 这将造成重复处理, 这时需要关掉 Java 进程的切割功能; 在排序的时候也可以实现并行排序来提升性能; 实现对数据的异步读写也可以提升性能。
- ❑ **健壮性的问题**: 需要对 mapper 和 reducer 程序的内存消耗进行限制, 这就要修改 Linux 内核, 增加其限制进程的物理内存的功能; 也可以通过多个 map 程序共享一块内存, 以一定的代价减少对物理内存的使用; 还可以将 DataNode 和 TaskTracker 的 UGI 配置为普通用户并设置账号密码; 或者让 DataNode 和 TaskTracker 分账号启动, 确保 HDFS 数据的安全性, 防止 Tracker 操作 DataNode 中的内容; 在不能保证用户

的每个程序都很健壮的情况下，有时需要将进程终止掉，但要保证父进程终止后子进程也被终止。

- **Streaming 局限性的问题**：比如，只能处理文本数据，mapper 和 reducer 按照文本行的协议通信，无法对二进制的数据进行简单处理。为了解决这个问题，百度人员新写了一个类 Bistreaming (Binary Streaming)，这里的子 Java 进程 mapper 和 reducer 按照 (KeyLen, Key, ValLen, Value) 的方式通信，用户可以按照这个协议编写程序。
- **用户认证的问题**：这个问题的解决办法是让用户名、密码、所属组都在 NameNode 和 Job Tracker 上集中维护，用户连接时需要提供用户名和密码，从而保证数据的安全性。百度下一步的工作重点可能主要会涉及以下内容：
 - 内存方面，降低 NameNode 的内存使用并研究 JVM 的内存管理；
 - 调度方面，改进任务可以被抢占的情况，同时开发出自己的基于 Capacity 的作业调度器，让等待作业队列具有优先级且队列中的作业可以设置 Capacity，并可以支持 TaskTracker 分组；
 - 压缩算法，选择较好的方法提高压缩比、减少存储容量，同时选取高效率的算法以进行 shuffle 数据的压缩和解压；
 - 对 mapper 程序和 reducer 程序使用的资源进行控制，防止过度消耗资源导致机器死机。以前是通过修改 Linux 内核来进行控制的，现在考虑通过在 Linux 中引入 cgroup 来对 mapper 和 reducer 使用的资源进行控制；
 - 将 DataNode 的并发数据读写方式由多线程改为 select 方式，以支持大规模并发读写和 Hypertable 的应用。

百度同时也在使用 Hypertable，它是以 Google 发布的 BigTable 为基础的开源分布式数据存储系统，百度将它作为分析用户行为的平台，同时在元数据集中化、内存占用优化、集群安全停机、故障自动恢复等方面做了一些改进。

3.4 Hadoop 在 Facebook 的应用

Facebook 作为全球知名的社交网站，拥有超过 3 亿的活跃用户，其中约有 3 千万用户至少每天更新一次自己的状态；用户每月总共上传 10 亿余张照片、1 千万个视频；以及每周共享 10 亿条内容，包括日志、链接、新闻、微博等。因此 Facebook 需要存储和处理的数据量是非常巨大的，每天新增加 4TB 压缩后的数据，扫描 135TB 大小的数据，在集群上执行 Hive 任务超过 7500 次，每小时需要进行 8 万次计算，所以高性能的云平台对 Facebook 来说是非常重要的，而 Facebook 主要将 Hadoop 平台用于日志处理、推荐系统和数据仓库等方面。

Facebook 将数据存储在使用 Hadoop/Hive 搭建的数据仓库上，这个数据仓库拥有 4800 个内核，具有 5.5PB 的存储量，每个节点可存储 12TB 大小的数据，同时，它还具有两层网络拓扑，如图 3-5 所示。Facebook 中的 MapReduce 集群是动态变化的，它基于负载情况和

集群节点之间的配置信息可动态移动。

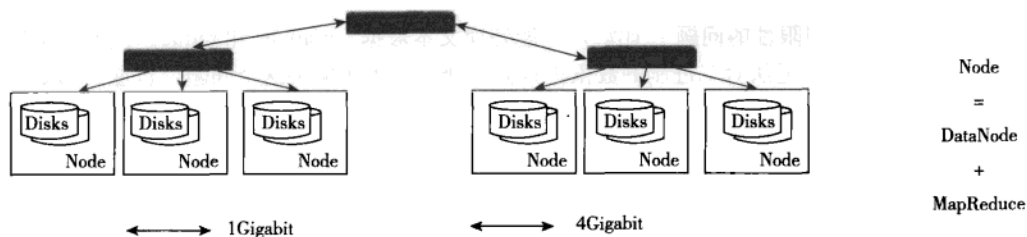


图 3-5 集群的网络拓扑

图 3-6 为 Facebook 的数据仓库架构，在这个架构中，网络服务器和内部服务生成日志数据，这里 Facebook 使用开源日志收集系统，它可以将数以百计的日志数据集存储在 NFS 服务器上，但大部分日志数据会复制到同一个中心的 HDFS 实例中，而 HDFS 存储的数据都会放到利用 Hive 构建的数据仓库中。Hive 提供了类 SQL 的语言来与 MapReduce 结合，创建并发布多种摘要和报告，以及在它们的基础上进行历史分析。Hive 上基于浏览器的接口允许用户执行 Hive 查询。Oracle 和 MySQL 数据库用来发布这些摘要，这些数据容量相对较小，但查询频率较高并需要实时响应。一些旧的数据需要及时归档，并存储在较便宜的存储设备上，如图 3-7 所示。

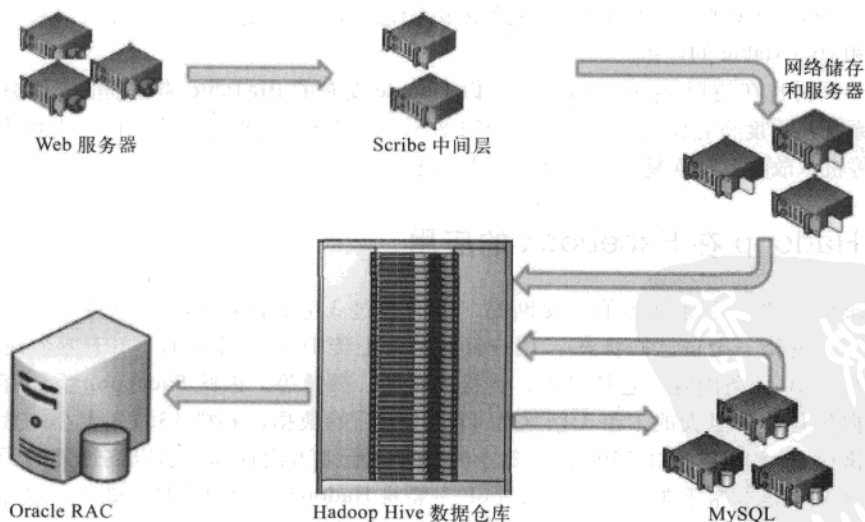


图 3-6 Facebook 数据仓库架构

下面介绍 Facebook 在 AvatarNode 和调度策略方面所做的一些工作。AvatarNode 主要用于 HDFS 的恢复和启动，若 HDFS 崩溃，原有技术恢复首先需要花 10 ~ 15 分钟来读取

12GB 的文件镜像并写回, 还要用 20 ~ 30 分钟处理来自 2000 个 DataNode 的数据块报告, 最后用 40 ~ 60 分钟来恢复崩溃的 NameNode 和部署软件。表 3-1 说明了 BackupNode 和 AvatarNode 的区别, AvatarNode 作为普通的 NameNode 启动, 处理所有来自 DataNode 的消息。AvatarDataNode 与 DataNode 相似, 支持多线程和针对多个主节点的多队列, 但无法区分原始和备份。人工恢复使用 AvatarShell 命令行工具, AvatarShell 执行恢复操作并更新 ZooKeeper 的 zNode, 恢复过程对用户来说是透明的。分布式 Avatar 文件系统实现在现有文件系统的上层。

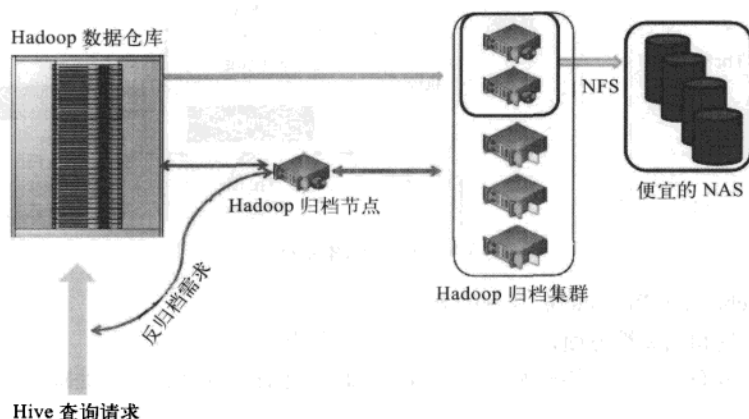


图 3-7 数据归档

表 3-1 BackupNode 和 AvatarNode 的区别

BackupNode (冷备份)	AvatarNode (热备份)
Namespace 状态与原始的同步 没有数据块和 DataNode 信息 在可用之前仍需 20 ~ 30 分钟	Namespace 状态与原始相比有几个事务的延迟 拥有全部的数据块和 DataNode 信息 6500 万个文件恢复不超过一分钟

基于位置的调度策略在实际应用中存在着一些问题：如需要高内存的任务可能会被分配给拥有低内存的 TaskTracker；CPU 资源有时未被充分利用；为不同硬件的 TaskTracker 进行配置也比较困难等。Facebook 采用基于资源的调度策略，即公平享有调度方法，实时监测系统并收集 CPU 和内存的使用情况，调度器会分析实时的内存消耗情况，然后在任务之间公平分配任务的内存使用量。它通过读取 /proc/ 目录解析进程树，并收集进程树上所有的 CPU 和内存的使用信息，然后通过 TaskCounters 在心跳 (heartbeat) 时发送信息。

Facebook 的数据仓库使用 Hive，其构架如图 3-8 所示，有关 Hive 查询语言的相关知识可查阅第 11 章的内容。这里 HDFS 支持三种文件格式：文本文件 (TextFile)，方便其他应用程序读写；顺序文件 (SequenceFile)，只有 Hadoop 能够读取并支持分块压缩；RCFile，使用顺序文件基于块的存储方式，每个块按列存储，这样有较好的压缩率和查询性能。Facebook

未来会在 Hive 上进行改进，以支持索引、视图、子查询等新功能。

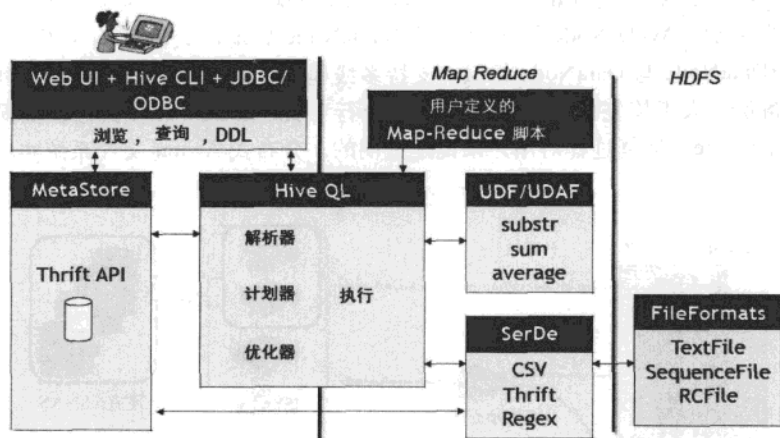


图 3-8 Hive 的体系结构

现在 Facebook 使用 Hadoop 遇到的挑战有：

- ☐ 服务质量和隔离性方面，较大的任务会影响集群性能；
- ☐ 安全性方面，如果软件漏洞导致 NameNode 事务日志崩溃该如何处理；
- ☐ 数据归档方面，如何选择归档数据，以及数据如何归档；
- ☐ 性能提升方面，如何有效地解决瓶颈等。

3.5 Hadoop 平台上的海量数据排序

Yahoo! 研究人员使用 Hadoop 完成了 Jim Gray 基准排序，此排序包含许多相关的基准，每个基准都有自己的规则。所有的排序基准都是通过测量不同记录的排序时间来制定的，每个记录为 100 字节，其中前面的 10 字节是键，剩余的部分是数值。MinuteSort 是比较在一分钟内所排序的数据量大小，GraySort 是比较在对大规模数据（至少 100TB）进行排序时的排序速率（TBs/minute）。基准规则具体如下：

- ☐ 输入数据必须与数据生成器生成的数据完全匹配；
- ☐ 任务开始的时候，输入数据不能在操作系统的文件缓存中。在 Linux 环境下，排序程序之间需要使用内存来交换其他内容；
- ☐ 输入和输出数据都是没有经过压缩的；
- ☐ 输出不能对输入进行重写；
- ☐ 输出文件必须存放到磁盘上；
- ☐ 必须计算输入和输出数据的每个键 / 值对的 CRC32，共 128 位校验和，当然，输入和输出必须对应相等；

- 输出如果分成多个输出文件，那么必须是完全有序的，也就是将这些输出文件连接以后必须是完全有序的输出；
- 开始和分布程序到集群上也要记入计算时间内；
- 任何抽样也要记入计算时间内。

Yahoo! 的研究人员使用 Hadoop 排列 1TB 数据用时 62 秒，排列 1PB 数据用时 16.25 个小时，具体如表 3-2 所示，它获得了 Daytona 类 GraySort 和 MinuteSort 级别的冠军。

表 3-2 数据规模与排序时间

数据大小 (Bytes)	节点数	副本数	时 间
500 000 000 000	1406	1	59 秒
1 000 000 000 000	1460	1	62 秒
100 000 000 000 000	3452	2	173 分钟
1 000 000 000 000 000	3658	2	975 分钟

下面的内容是根据基准排序的官方网站 (<http://sortbenchmark.org/>) 上有关使用 Hadoop 排序的相关内容整理而成。

Yahoo! 的研究人员编写了三个 Hadoop 应用程序来进行 TB 级数据的排序：

- TeraGen 是产生数据的 map/reduce 程序；
- TeraSort 进行数据取样，并使用 map/reduce 对数据进行排序；
- TeraValidate 是用来验证输出数据是否有序的 map/reduce 程序。

TeraGen 用来产生数据，它将数据按行排列并且根据执行任务的数目为每个 map 分配任务，每个 map 任务产生所分配行数范围内的数据。最后，TeraGen 使用 1800 个任务产生总共 100 亿行的数据存储在 HDFS 上，每个存储块的大小为 512MB。

TeraSort 是标准的 map/reduce 排序程序，但这里使用的是不同的分配方法。程序中使用 $N-1$ 个已排好序的抽样键值来为 reduce 任务分配排序数据的行数范围。例如，键值 key 在范围 $\text{sample}[i-1] \leq \text{key} < \text{sample}[i]$ 的数据会分配给第 i 个 reduce 任务。这样就保证了第 i 个 reduce 任务输出的数据都比第 $i+1$ 个 reduce 任务输出的数据小。为了加速分配过程，分配器在抽样键值上建立两层的索引结构树。TeraSort 在任务提交之前在输入数据中进行抽样，并将产生的抽样数据写入 HDFS 中。这里规定了输入输出格式，使得三个应用程序可以正确地读取并写入数据。reduce 任务的副本数默认是 3，这里设置为 1，因为输出数据不需要复制到多个节点上。这里配置的任务为 1800 个 map 任务和 1800 个 reduce 任务，并为任务的栈设置了充足的空间，防止产生的中间数据溢出到磁盘上。抽样器使用 100 000 个键值来决定 reduce 任务的边界，如图 3-9 所示，分布并不是很完美。

TeraValidate 保证输出数据是全部排好序的，它为输出目录的每个文件分配一个 map 任务（如图 3-10 所示），map 任务检查每个值是否大于等于前一个值，同时输出最大值和最小值给 reduce 任务，reduce 任务检查第 i 个文件的最小值是否大于第 $i-1$ 文件的最大值，如果不是则产生错误报告。

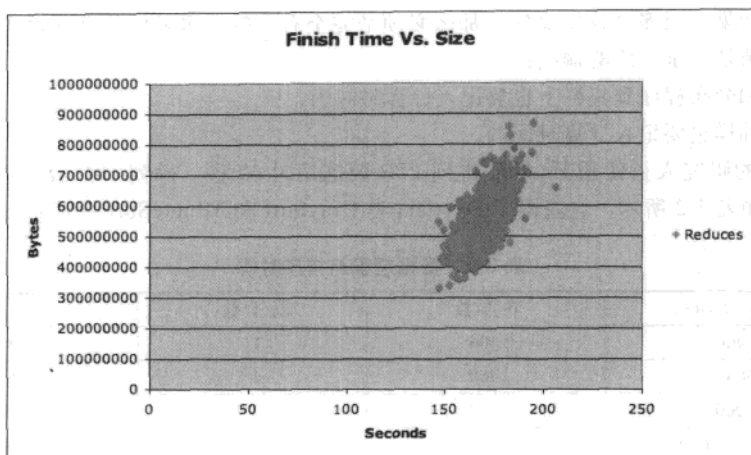


图 3-9 reduce 任务的输出大小和完成时间分布图

以上应用程序运行在雅虎搭建的集群上，其集群配置为：

- ❑ 910 个节点；
- ❑ 每个节点拥有 4 个英特尔双核 2.0GHz 至强处理器；
- ❑ 每个节点拥有 4 个 SATA 硬盘；
- ❑ 每个节点有 8GB 的内存；
- ❑ 每个节点有 1GB 的以太网带宽；
- ❑ 40 个节点一个 rack ；
- ❑ 每个 rack 到核心有 8GB 的以太网带宽；
- ❑ 操作系统为 Red Hat Enterprise Linux Server Release 5.1(kernel 2.6.18) ；
- ❑ JDK 为 Sun Java JDK 1.6.0_05-b13。

整个排序过程在 209 秒（3.48 分钟）内完成，尽管拥有 910 个节点，但是网络核心是与其他 2000 个节点的集群共享的，所以运行时间会因为其他集群的活动而有所变化。

使用 Hadoop 进行 GraySort 基准排序时，Yahoo! 的研究人员将上面的 map/reduce 应用程序稍加修改以适应新的规则，整个程序分为 4 个部分，分别为：

- ❑ TeraGen 是产生数据的 map/reduce 程序；
- ❑ TeraSort 进行数据取样，并使用 map/reduce 对数据进行排序；
- ❑ TeraSum 是 map/reduce 程序，用来计算每个键 / 值对的 CRC32，共 128 位校验和；
- ❑ TeraValidate 是用来验证输出数据是否有序的 map/reduce 程序，并且计算校验和的总和。

TeraGen 和 TeraSort 与上面介绍的一样，TeraValidate 除了增加了计算输出目录校验和总和的任务以外，其他都一样，这里不再赘述。

TeraSum 计算每个键 / 值对的 CRC32 的校验和，每个 map 任务计算输入的校验和并输出，然后一个 reduce 任务将每个 map 生成的校验和相加。这个程序用来计算输入目录下每

个键/值对校验和的和, 还用来检查排序输出后的正确性。

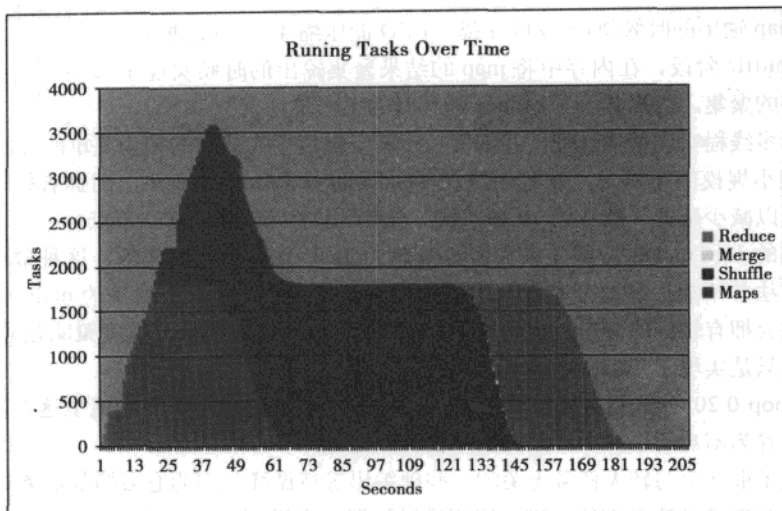


图 3-10 每个阶段的任务数

这次基准测试运行在 Yahoo! 的 Hammer 集群上, 集群的具体细节如下:

- 将近 3800 个节点 (在这样大规模的集群中, 一些节点会坏掉);
- 每个节点两个双核 2.5GHz 的 Xeons 处理器;
- 每个节点 4 个 SATA 硬盘;
- 每个节点 8GB 内存 (在 PB 级排序前会升级到 16GB);
- 每个节点 1GB 的以太网带宽;
- 每个 rack 拥有 40 个节点;
- 每个节点到核心有 8GB 的以太网带宽;
- 操作系统为 Red Hat Enterprise Linux Server Release 5.1 (kernel 2.6.18);
- JDK 为 Sun Java JDK (1.6.0 05-b13 and 1.6.0 13-b03) (32 and 64 bit)。

对于较大规模的排序, 这里 NameNode 和 JobTracker 使用的是 64 位的 JVM。排序测试所用的 Hadoop 平台也做了一些变化, 主要有:

- 重新实现了 Hadoop shuffle 阶段的 reducer 部分, 在重新设计后提高了 shuffle 的性能, 解除了瓶颈, 而且代码也更容易维护和理解了;
- 新的 shuffle 过程从一个节点获取多个 map 的结果, 而不是之前的一次只取一个结果。这样防止了多余的连接和传输开销;
- 允许配置 shuffle 连接的超时时间, 在小规模排序时则可以将其减小, 因为一些情况下 shuffle 会在超时时间到期后停止, 这会增加任务的延迟时间;
- 设置 TCP 为无延迟并增加 TaskTracker 和 TaskTracker 之间 ping 的频率, 以减少发现

问题的延迟时间：

- ❑ 增加一些代码，用来检测从 shuffle 传输数据的正确性，防止引起 reduce 任务的失败。
- ❑ 在 map 输出的时候使用 LZO 压缩，LZO 能压缩 45% 的数据量；
- ❑ 在 shuffle 阶段，在内存中将 map 的结果聚集输出的时候实现了 reduce 需要的内存到内存的聚集，这样减少了 reduce 运行时的工作量；
- ❑ 使用多线程实现抽样过程，并编写一个基于键值平均分布的较为简单的分配器；
- ❑ 在较小规模的集群上，配置系统在 TaskTracker 和 JobTracker 之间拥有较快的心跳频率，以减少延迟（默认为 10 秒/1000 节点，配置为 2 秒/1000 节点）；
- ❑ 默认的 JobTracker 按照先来先服务策略为 TaskTracker 分配任务，这种贪心的任务分配方法并不能很好地分布数据。从全局的角度来看，如果一次性为 map 分配好任务，系统会拥有较好的分布，但是为所有的 Hadoop 程序实现全局调度策略是非常困难的，这里只是实现了 TeraSort 的全局调度策略；
- ❑ Hadoop 0.20 增加了安装和清除任务的功能，但是在排序基准测试里这并不需要，可以设置为不启动来减少开始和结束任务的延迟；
- ❑ 删除了框架中与较大任务无关的一些硬编码等待循环，因为它会增加任务延迟时间；
- ❑ 允许为任务设置日志的级别，这样通过配置日志级别可以从 INFO 到 WARN 减少日志的内容，减少日志的内容对系统的性能有较大的提高，但是增加了调试和分析的困难；
- ❑ 优化任务分配代码，但还未完成。目前，对输入文件使用 RPC 请求到 NameNode 上会花费大量的时间。

Hadoop 与上面的测试相比有了很大的改进，可以在更短的时间内执行更多的任务。值得注意的是，在大集群和分布式应用程序中需要转移大量数据，这会导致执行时间有很大的变化。但是随着 Hadoop 的改进，它能够更好地处理硬件故障，这种时间变化也就微不足道了。不同规模的数据排序所需的时间如表 3-2 所示。

因为较小规模的数据需要更短的延迟和更快的网络，所以使用集群中的部分节点来进行计算。将较小规模计算的输出副本数设置为 1，因为整个过程较短且运行在较小的集群上，节点坏掉的可能性相对较小。而在较大规模的计算上，节点坏掉是难免的，于是将节点副本数设置为 2。HDFS 保证节点换掉后数据不会丢失，因为不同的副本放在不同的节点上。

Yahoo! 的研究人员统计了 JobTracker 上从任务提交状况获得的任务数随时间的变化，图 3-11、图 3-12、图 3-13、图 3-14 显示了每个时间点下的任务数。maps 只有一个阶段，而 reduces 拥有三个阶段：shuffle、merge 和 reduce。shuffle 是从 maps 中转移数据的，merge 在测试中并不需要；reduce 阶段进行最后的聚集并写到 HDFS 上。如果将这些图与图 3-6 进行比较，你会发现建立任务的速度变快了。图 3-6 中每次心跳建立一个任务，那么所有任务建立起来需要 40 秒，现在 Hadoop 每次心跳可以设置好一个 TaskTracker，可见减少任务建立的开销是非常重要的。

运行大规模数据时，数据传输的次数对任务性能的影响也是非常大的。在 PB 级数据排序中，每个 map 处理 15GB 的数据而不是默认的 128MB，每个 reduce 处理 50GB 的数据。

如果按照 1.5GB/map 进行处理, 需要 40 个小时才能完成。因此, 为了增加吞吐量, 增加每个块的大小是非常重要的。

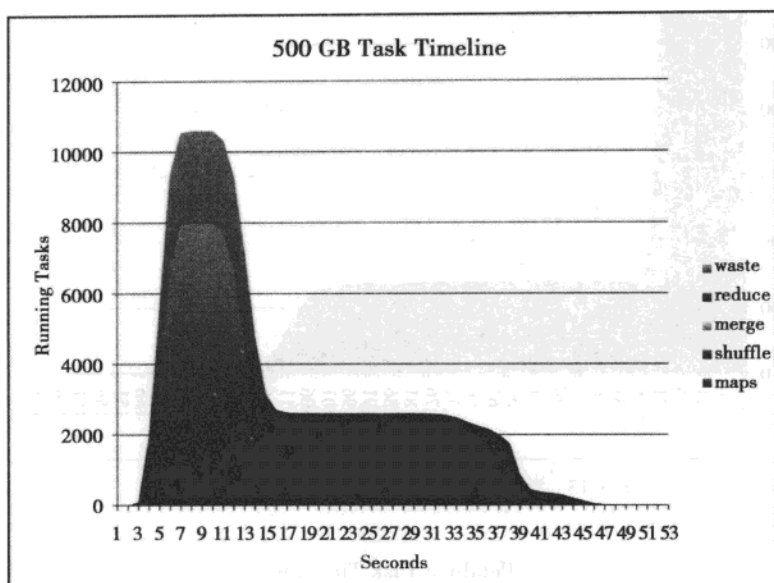


图 3-11 数据量为 500GB 时任务数随时间的变化

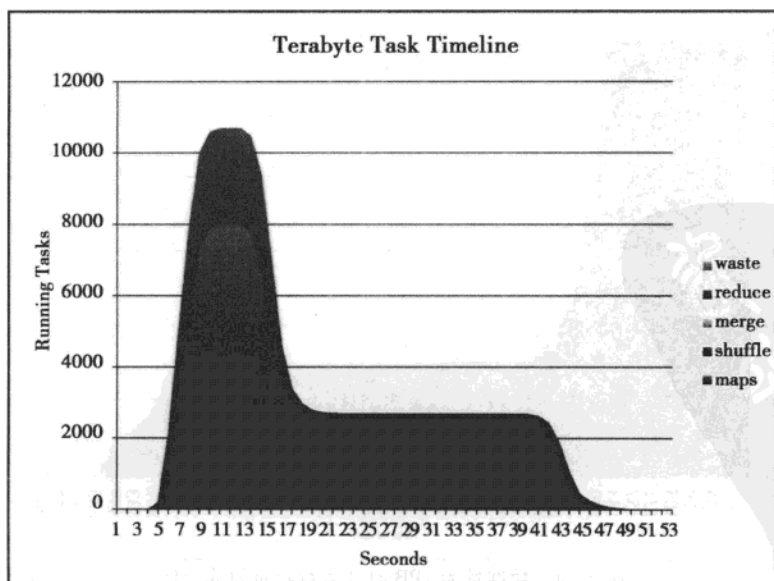


图 3-12 数据量为 1TB 时任务数随时间的变化

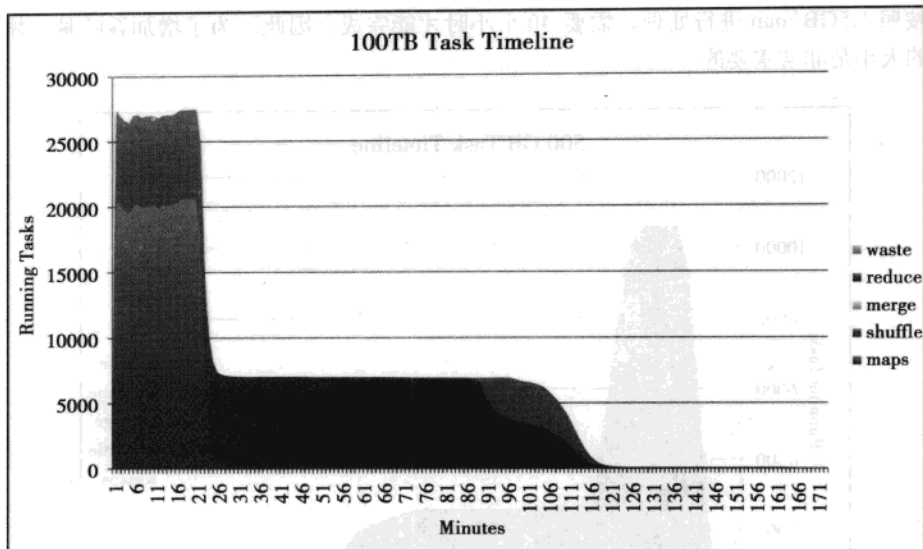


图 3-13 数据量为 100TB 时任务数随时间的变化

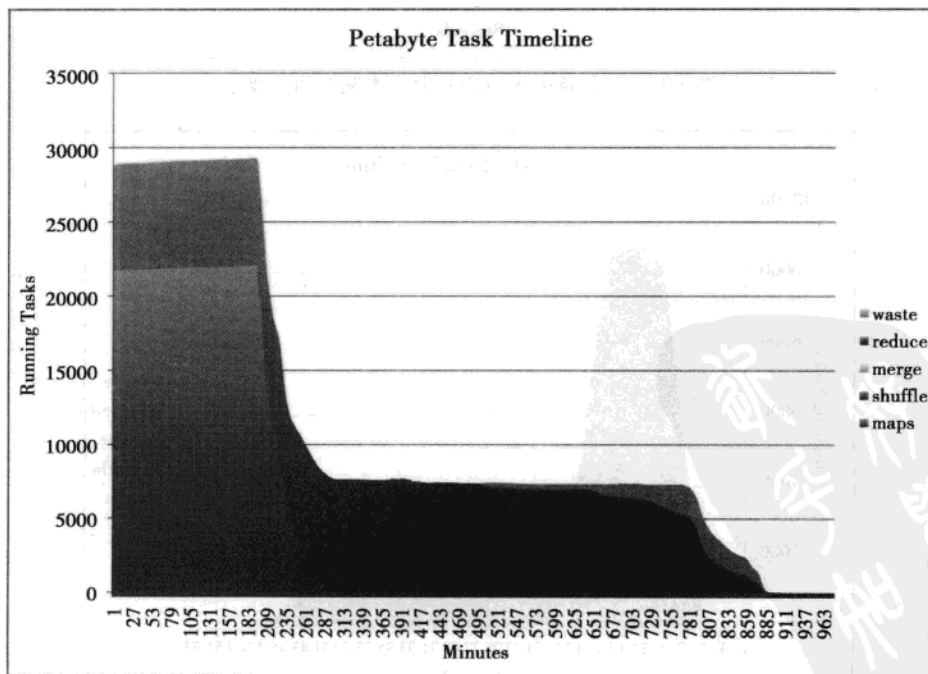


图 3-14 数据量为 1PB 时任务数随时间的变化

3.6 小结

本章主要介绍了 Hadoop 的具体使用案例,我们选取了 Yahoo!、百度、Facebook、eBay 和海量数据排序为例进行说明,主要介绍了商业公司如何使用 Hadoop 来增强自己的服务,以及它们在使用 Hadoop 中遇到的各种问题和改进的方法。Hadoop 是开源的系统,任何公司可以根据自己的业务需要对 Hadoop 进行修改或改进,同时也为 Hadoop 的改进贡献了自己的力量。

随着 Hadoop 的不断改进,其强大的分布式功能被越来越多的人熟知,使用 Hadoop 的公司队伍也在不断壮大中,具体可以登录 <http://wiki.apache.org/hadoop/PoweredBy> 查看。

参考文献

另外,本章关于 Hadoop 在 Yahoo! 的应用内容是根据 Hadoop 云计算大会上 Yahoo! 研究人员的报告整理而成的,Pig 和 Hive 应用相关内容来自 Yahoo! 研究人员的博客^[1],大家如果想要了解 Hadoop 在 Yahoo! 应用中的更多细节和进展,请关注 Yahoo! Hadoop 团队的博客(developer.yahoo.com/blogs/hadoop)。

Hadoop 在 eBay 的应用内容是根据 eBay 研究人员的技术博客^[2]整理而成的,其中参考了 eBay 分析平台开发部 Anil Madan 介绍的 Hadoop 在 eBay 的使用情况,大家如果想要了解 Hadoop 在 eBay 应用的更多信息,可以关注 eBay 研究人员的技术博客(www.ebaytechblog.com)。

百度使用 Hadoop 平台的情况则是根据近几届 Hadoop 中国云计算大会上百度研究人员的报告整理而成,大家如果想了解更详细的信息或 Hadoop 中国云计算大会的相关信息可登录 Hadoop in China 网站:<http://www.hadooper.cn>。

Facebook 使用 Hadoop 的情况是根据 Facebook 相关技术人员在各个云计算大会上所做的报告整理而成的。


Hadoop 平台上的海量数据排序的相关内容是根据 Hadoop 基准排序测试的报告^[3]、^[4]撰写而成,如果大家想要了解排序基准测试的更多细节或排名情况,可以登录 <http://sortbenchmark.org> 查看。

[1] Alan Gates, Pig and Hive at Yahoo!, http://developer.yahoo.com/blog/hadoop/posts/2010/08/pig_and_hive_at_yahoo/

[2] Anil Madan, Hadoop-The power of the Elephant, <http://www.ebaytechblog.com/2010/10/29/hadoop-the-power-of-the-elephant/>

[3] Owen O'Mallay, TeraByte Sort on Apache Hadoop, <http://sortbenchmark.org/YahooHadoop.pdf>

[4] Owen O'Mallay & Arun C.Murthy, Winning a 60 Second Dash with a Yellow Elephant, <http://sortbenchmark.org/Yahoo2009.pdf>



第 4 章

MapReduce 计算模型

本章内容

- ☐ 为什么要用 MapReduce
- ☐ MapReduce 计算模型
- ☐ MapReduce 任务的优化
- ☐ Hadoop 流
- ☐ Hadoop Pipes
- ☐ 小结

资源库

PDG

MapReduce 被广泛地应用于日志分析、海量数据的排序、在海量数据中查找特定模式等场景中。Hadoop 根据 Google 的论文实现了 MapReduce 这个编程框架，并将源代码完全贡献了出来。本章将要向大家介绍 MapReduce 这个流行的编程框架。

4.1 为什么要用 MapReduce

MapReduce 的流行是有理由的，它非常简单、易于实现且扩展性强。可以通过它轻易地编写出同时多台主机上运行的程序，可以使用 Ruby、Python、PHP 和 C++ 等非 Java 类语言编写 map 或 reduce 程序，还可以在任何安装 Hadoop 的集群中运行同样的程序，不论这个集群有多少台主机。MapReduce 适合于处理大量的数据集，因为它会同时被多台主机一起处理，这样通常会有较快的速度。

下面来看一个例子。

引文分析是评价论文好坏的一个非常重要的方面，本例只对其中最简单的一部分，即论文的被引用次数进行了统计。假设大家有很多篇论文（百万级），且每篇论文的引文形式如下所示：

References

David M. Blei, Andrew Y. Ng, and Michael I. Jordan.
2003. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3:993-1022.

Samuel Brody and Noemie Elhadad. 2010. An unsupervised aspect-sentiment model for online reviews. In *NAACL '10*.

Jaime Carbonell and Jade Goldstein. 1998. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In *SIGIR '98*, pages 335-336.

Dennis Chong and James N. Druckman. 2010. Identifying frames in political news. In Erik P. Bucy and R. Lance Holbert, editors, *Sourcebook for Political Communication Research: Methods, Measures, and Analytical Techniques*. Routledge.

Cindy Chung and James W. Pennebaker. 2007. The psychological function of function words. *Social Communication: Frontiers of Social Psychology*, pages 343-359.

Gunes Erkan and Dragomir R. Radev. 2004. Lexrank: graph-based lexical centrality as salience in text summarization. *J. Artif. Int. Res.*, 22(1):457-479.

Stephan Greene and Philip Resnik. 2009. More than words: syntactic packaging and implicit sentiment. In *NAACL '09*, pages 503-511.

Aria Haghighi and Lucy Vanderwende. 2009. Exploring content models for multi-document summarization. In *NAACL '09*, pages 362-370.

Sanda Harabagiu, Andrew Hickl, and Finley Lacatusu. 2006. Negation, contrast and contradiction in text processing.

在单机上运行时，想要完成这个任务，需要先切分出所有论文的名字并存入一个 hash 表中，然后遍历所有论文，查看引文信息，一一计数。因为文章数量很多，需要进行很多次内外存交换，这无疑会延长程序的执行时间。但在 MapReduce 中，这是一个 WordCount 就能解决的问题。

4.2 MapReduce 计算模型

要了解 MapReduce，首先需要了解 MapReduce 的载体是什么。在 Hadoop 中，用于执行 MapReduce 任务的机器角色有两个：一个是 JobTracker；另一个是 TaskTracker。JobTracker 是用于调度工作的，TaskTracker 是用于执行工作的。一个 Hadoop 集群中只有一台 JobTracker。

4.2.1 MapReduce Job

在 Hadoop 中，每个 MapReduce 任务都被初始化为一个 Job。每个 Job 又可以分为两个阶段：map 阶段和 reduce 阶段。这两个阶段分别用两个函数来表示，即 map 函数和 reduce 函数。map 函数接收一个 $\langle \text{key}, \text{value} \rangle$ 形式的输入，然后同样产生一个 $\langle \text{key}, \text{value} \rangle$ 形式的中间输出，Hadoop 会负责将所有具有相同中间 key 值的 value 集合到一起传递给 reduce 函数，reduce 函数接收一个如 $\langle \text{key}, (\text{list of values}) \rangle$ 形式的输入，然后对这个 value 集合进行处理，每个 reduce 产生 0 或 1 个输出，reduce 的输出也是 $\langle \text{key}, \text{value} \rangle$ 形式的。

为了方便理解，分别将 3 个 $\langle \text{key}, \text{value} \rangle$ 对标记为 $\langle k1, v1 \rangle$ 、 $\langle k2, v2 \rangle$ 、 $\langle k3, v3 \rangle$ ，那么上面的所述的过程就可以用图 4-1 来表示了。

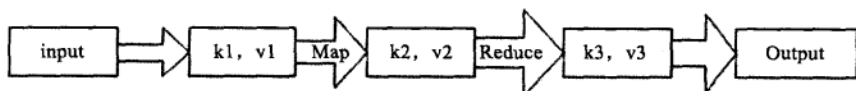


图 4-1 MapReduce 程序数据变化的基本模型

4.2.2 Hadoop 中的 Hello World 程序

上面的过程是 MapReduce 的核心，所有的 MapReduce 程序都具有如上的结构。下面我再举一个例子详述 MapReduce 的执行过程。

大家在初次接触编程语言时，看到的第一个示例程序可能都是“Hello World”。在 Hadoop 中也有一个类似于 Hello World 的程序，这就是 WordCount。本小节会结合这个程序具体讲解与 MapReduce 程序有关的所有类，这个程序的内容如下：

```
import java.io.IOException;
import java.util.*;

import org.apache.hadoop.fs.Path;
```

```

import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapred.*;
import org.apache.hadoop.util.*;

public class WordCount {

    public static class Map extends MapReduceBase implements Mapper<LongWritable,
        Text, Text, IntWritable> {
        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        public void map(LongWritable key, Text value, OutputCollector<Text,
            IntWritable> output, Reporter reporter) throws IOException {
            String line = value.toString();
            StringTokenizer tokenizer = new StringTokenizer(line);
            while (tokenizer.hasMoreTokens()) {
                word.set(tokenizer.nextToken());
                output.collect(word, one);
            }
        }
    }

    public static class Reduce extends MapReduceBase implements Reducer<Text,
        IntWritable, Text, IntWritable> {
        public void reduce(Text key, Iterator<IntWritable> values,
            OutputCollector<Text, IntWritable> output, Reporter reporter) throws
            IOException {
            int sum = 0;
            while (values.hasNext()) {
                sum += values.next().get();
            }
            output.collect(key, new IntWritable(sum));
        }
    }

    public static void main(String[] args) throws Exception {
        JobConf conf = new JobConf(WordCount.class);
        conf.setJobName("wordcount");

        conf.setOutputKeyClass(Text.class);
        conf.setOutputValueClass(IntWritable.class);

        conf.setMapperClass(Map.class);
        conf.setReducerClass(Reduce.class);

        conf.setInputFormat(TextInputFormat.class);
        conf.setOutputFormat(TextOutputFormat.class);

        FileInputFormat.setInputPaths(conf, new Path(args[0]));
    }
}

```

```

        FileOutputFormat.setOutputPath(conf, new Path(args[1]));

        JobClient.runJob(conf);
    }
}

```

同时，为了叙述方便，设定两个输入文件，如下：

```

echo "Hello World Bye World" > file01
echo "Hello Hadoop Goodbye Hadoop" > file02

```

看到这个程序，可能很多读者会对众多的预定义类感到很迷惑，其实这些类非常简单明了。首先，WordCount 程序的代码虽多，但是执行过程却很简单，在本例中，它首先将输入文件读进来，然后交由 map 程序处理，map 程序将输入读入后切出其中的单词，并标记它的数目为 1，形成 <word,1> 的形式，然后交由 reduce 处理，reduce 将相同 key 值（也就是 word）的 value 值收集起来，形成 <word, list of 1> 的形式，之后再将这些 1 值加起来，即为单词的个数，最后将这个 <key, value> 对以 TextOutputFormat 的形式输出到 HDFS 中。

针对这个数据流动过程，我们挑出了如下几句代码来表述它的执行过程：

```

JobConf conf = new JobConf(MyMapre.class);
conf.setJobName("wordcount");

conf.setInputFormat(TextInputFormat.class);
conf.setOutputFormat(TextOutputFormat.class);

conf.setMapperClass(Map.class);
conf.setReducerClass(Reduce.class);

FileInputFormat.setInputPath(conf, new Path(args[0]));
FileOutputFormat.setOutputPath(conf, new Path(args[1]));

```

首先讲解一下 Job 的初始化过程。main 函数调用 Jobconf 类来对 MapReduce Job 进行初始化，然后调用 setJobName() 方法命名这个 Job。对 Job 进行合理的命名有助于更快地找到 Job，以便在 JobTracker 和 TaskTracker 的页面中对其进行监视。接着就调用 setInputPath() 和 setOutputPath() 设置输入输出路径。下面会结合 WordCount 程序重点讲解 InputFormat()、OutputFormat()、map()、reduce() 这 4 种方法。

1. InputFormat() 和 InputSplit

InputSplit 是 Hadoop 定义的用来传送给每个单独的 map 的数据，InputSplit 存储的并非数据本身，而是一个分片长度和一个记录数据位置的数组。生成 InputSplit 的方法可以通过 InputFormat() 来设置。当数据传送给 map 时，map 会将输入分片传送到 InputFormat 上，InputFormat 则调用 getRecordReader() 方法生成 RecordReader，RecordReader 再通过 creatKey()、creatValue() 方法创建可供 map 处理的 <key, value> 对，即 <k1, v1>。简而言之，

InputFormat() 方法是用来生成可供 map 处理的 <key, value> 对的。

Hadoop 预定义了多种方法将不同类型的输入数据转化为 map 能够处理的 <key, value> 对，它们都继承自 InputFormat，分别是：

- ❑ BaileyBorweinPlouffe.BbpInputFormat
- ❑ ComposableInputFormat
- ❑ CompositeInputFormat
- ❑ DBInputFormat
- ❑ DistSum.Machine.AbstractInputFormat
- ❑ FileInputFormat

其中，FileInputFormat 又有多个子类，分别为：

- ❑ CombineFileInputFormat
- ❑ KeyValueTextInputFormat
- ❑ NLineInputFormat
- ❑ SequenceFileInputFormat
- ❑ TeraInputFormat
- ❑ TextInputFormat

其中，TextInputFormat 是 Hadoop 默认的输入方法，在 TextInputFormat 中，每个文件（或其一部分）都会单独地作为 map 的输入，而这是继承自 FileInputFormat 的。之后，每行数据都会生成一条记录，每条记录则表示成 <key, value> 形式：

- ❑ key 值是每个数据的记录在数据分片中的字节偏移量，数据类型是 LongWritable；
- ❑ value 值是每行的内容，数据类型是 Text。

也就是说，输入数据会以如下的形式被传入 map 中：

```
file01:
0 hello world bye world
file02
0 hello hadoop bye hadoop
```

因为 file01 和 file02 都会被单独输入到一个 map 中，因此它们的 key 值都是 0。

2. OutputFormat

每一种输入格式都有一种输出格式与其对应。同样，默认的输出格式是 TextOutputFormat，这种输出方式与输入类似，会将每条记录以一行的形式存入文本文件。不过，它的键和值可以是任意形式的，因为程序内部会调用 toString() 方法将键和值转换为 String 类型再输出。最后的输出形式如下所示：

```
Bye 2
Hadoop 2
Hello 2
World 2
```

3. map 和 reduce

map 方法和 reduce 方法是本章的重点，从前面的内容知道，map 函数接收经过 InputFormat 处理所产生的 <k1, v1>，然后输出 <k2, v2>。WordCount 的 map 函数如下：

```
public class MyMapre {
    public static class Map extends MapReduceBase implements Mapper<LongWritable,
        Text, Text, IntWritable> {
        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        public void map(LongWritable key, Text value, OutputCollector<Text,
            IntWritable> output, Reporter reporter) throws IOException {
            String line = value.toString();
            StringTokenizer tokenizer = new StringTokenizer(line);
            while (tokenizer.hasMoreTokens()) {
                word.set(tokenizer.nextToken());
                output.collect(word, one);
            }
        }
    }
}
```

map 函数继承自 MapReduceBase，并且它实现了 Mapper 接口，此接口是一个泛型类型，它有 4 种形式的参数，分别用来指定 map 的输入 key 值类型、输入 value 值类型、输出 key 值类型和输出 value 值类型。在本例中，因为使用的是 TextInputFormat，它的输出 key 值是 LongWritable 类型，输出 value 值是 Text 类型，所以 map 的输入类型即为 <LongWritable, Text>。如前面的内容所述，在本例中需要输出 <word, 1> 这样的形式，因此输出的 key 值类型是 Text，输出的 value 值类型是 IntWritable。

实现此接口类还需要实现 map 方法，map 方法会具体负责对输入进行操作，在本例中，map 方法对输入的行以空格为单位进行切分，然后使用 OutputCollect 收集输出的 <word, 1>，即 <k2, v2>。

下面来看 reduce：

```
public static class Reduce extends MapReduceBase implements Reducer<Text,
    IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text,
        IntWritable> output, Reporter reporter) throws IOException {
        int sum = 0;
        while (values.hasNext()) {
            sum += values.next().get();
        }
        output.collect(key, new IntWritable(sum));
    }
}
```

与 map 类似，reduce 函数也是继承自 MapReduceBase 的，需要实现 Reducer 接口。

reduce 函数以 map 的输出作为输入，因此 reduce 的输入类型是 `<Text, IntWritable>`。而 reduce 的输出是单词和它的数目，因此，它的输出类型是 `<Text, IntWritable>`。reduce 函数也要实现 reduce 方法，在此方法中，reduce 函数将输入的 key 值作为输出的 key 值，然后将获得的多个 value 值加起来，作为输出的 value 值。

4. 运行 MapReduce 程序

读者可以在 Eclipse 里运行 MapReduce 程序，也可以在命令行中运行 MapReduce 程序，但是在实际应用中，还是推荐到命令行中运行程序。按照第 2 章所介绍的，首先安装 Hadoop，然后编译、打包编写的 Java 程序，如下所示（以 hadoop-0.20.2 为例，安装路径是 `~/hadoop`）：

```
mkdir FirstJar
javac -classpath ~/hadoop/hadoop-0.20.2-core.jar -d FirstJar
WordCount.java
jar -cvf wordcount.jar -C FirstJar/ .
```

首先建立 FirstJar，然后编译文件生成 .class，存放到文件夹 FirstJar 中，并将 FirstJar 中的文件打包生成 wordcount.jar 文件。

接着上传输入文件（输入文件是 file01、file02，存放在 `~/input`）：

```
~/hadoop/bin/hadoop dfs -mkdir input
~/hadoop/bin/hadoop dfs -put ~/input/file0* input
```

在此上传过程中，先建立文件夹 input，然后上传文件 file01、file02 到 input 中。

最后运行生成的 jar 文件，为了叙述方便，先将生成的 jar 文件放入 Hadoop 的安装文件夹中（`HADOOP_HOME`），然后运行如下命令：

```
~/hadoop/bin/hadoop jar wordcount.jar WordCount input output
11/01/21 20:02:38 WARN mapred.JobClient: Use GenericOptionsParser for parsing the
arguments. Applications should implement Tool for the same.
11/01/21 20:02:38 INFO mapred.FileInputFormat: Total input paths to process : 2
11/01/21 20:02:38 INFO mapred.JobClient: Running job: job_201101111819_0002
11/01/21 20:02:39 INFO mapred.JobClient: map 0% reduce 0%
11/01/21 20:02:49 INFO mapred.JobClient: map 100% reduce 0%
11/01/21 20:03:01 INFO mapred.JobClient: map 100% reduce 100%
11/01/21 20:03:03 INFO mapred.JobClient: Job complete: job_201101111819_0002
11/01/21 20:03:03 INFO mapred.JobClient: Counters: 18
11/01/21 20:03:03 INFO mapred.JobClient: Job Counters
11/01/21 20:03:03 INFO mapred.JobClient: Launched reduce tasks=1
11/01/21 20:03:03 INFO mapred.JobClient: Launched map tasks=2
11/01/21 20:03:03 INFO mapred.JobClient: Data-local map tasks=2
11/01/21 20:03:03 INFO mapred.JobClient: FileSystemCounters
11/01/21 20:03:03 INFO mapred.JobClient: FILE_BYTES_READ=100
11/01/21 20:03:03 INFO mapred.JobClient: HDFS_BYTES_READ=46
11/01/21 20:03:03 INFO mapred.JobClient: FILE_BYTES_WRITTEN=270
11/01/21 20:03:03 INFO mapred.JobClient: HDFS_BYTES_WRITTEN=31
11/01/21 20:03:03 INFO mapred.JobClient: Map-Reduce Framework
11/01/21 20:03:04 INFO mapred.JobClient: Reduce input groups=4
```

```

11/01/21 20:03:04 INFO mapred.JobClient: Combine output records=0
11/01/21 20:03:04 INFO mapred.JobClient: Map input records=2
11/01/21 20:03:04 INFO mapred.JobClient: Reduce shuffle bytes=106
11/01/21 20:03:04 INFO mapred.JobClient: Reduce output records=4
11/01/21 20:03:04 INFO mapred.JobClient: Spilled Records=16
11/01/21 20:03:04 INFO mapred.JobClient: Map output bytes=78
11/01/21 20:03:04 INFO mapred.JobClient: Map input bytes=46
11/01/21 20:03:04 INFO mapred.JobClient: Combine input records=0
11/01/21 20:03:04 INFO mapred.JobClient: Map output records=8
11/01/21 20:03:04 INFO mapred.JobClient: Reduce input records=8

```

Hadoop 命令（注意不是 Hadoop 本身）会启动一个 JVM 来运行这个 MapReduce 程序，并自动获得 Hadoop 的配置，同时把类的路径（及其依赖关系）加入到 Hadoop 的库中。以上就是 Hadoop Job 的运行记录，从这里面可以看到，这个 Job 被赋予了一个 ID 号：job_201101111819_0002，而且得知输入文件有两个（Total input paths to process : 2），同时还可以了解 map 的输入输出记录（record 数及字节数），以及 reduce 的输入输出记录。比如说，在本例中，map 的 task 数量是 2 个，reduce 的 task 数量是 1 个。map 的输入 record 数是 2 个，输出 record 数是 8 个等信息。

可以通过命令查看输出文件为：

```

bye 2
hadoop 2
hello 2
world 2

```

5. 新的 API

从 0.20.2 开始，Hadoop 提供了一个新的 API，新的 API 是在 org.apache.hadoop.mapreduce 中的，旧版的 API 则在 org.apache.hadoop.mapred 中。新的 API 不兼容旧的 API，WordCount 程序用新的 API 重写如下：

```

import java.io.IOException;
import java.util.*;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.*;
import org.apache.hadoop.mapreduce.lib.output.*;
import org.apache.hadoop.util.*;

public class WordCount extends Configured implements Tool {
    public static class Map extends Mapper<LongWritable, Text, Text, IntWritable> {
        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();
        public void map(LongWritable key, Text value, Context context) throws

```



```

        IOException, InterruptedException {
            String line = value.toString();
            StringTokenizer tokenizer = new StringTokenizer(line);
            while (tokenizer.hasMoreTokens()) {
                word.set(tokenizer.nextToken());
                context.write(word, one);
            }
        }
    }

    public static class Reduce extends Reducer<Text, IntWritable, Text,
        IntWritable> {
        public void reduce(Text key, Iterable<IntWritable> values, Context context)
            throws IOException, InterruptedException {
            int sum = 0;
            for (IntWritable val : values) {
                sum += val.get();
            }
            context.write(key, new IntWritable(sum));
        }
    }

    public int run(String [] args) throws Exception {
        Job job = new Job(getConf());
        job.setJarByClass(WordCount.class);
        job.setJobName("wordcount");

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        job.setMapperClass(Map.class);
        job.setReducerClass(Reduce.class);

        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);

        FileInputFormat.setInputPaths(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        boolean success = job.waitForCompletion(true);
        return success ? 0 : 1;
    }

    public static void main(String[] args) throws Exception {
        int ret = ToolRunner.run(new WordCount(), args);
        System.exit(ret);
    }
}

```

从这个程序，可以看到新旧 API 的几个区别：

- 在新的 API 中, Mapper 与 Reducer 已经不是接口而是抽象类。而且 map 函数与 reduce 函数也已经不再实现 Mapper 和 Reducer 接口, 而是继承 Mapper 和 Reducer 抽象类。这样做更容易扩展, 因为添加方法到抽象类中更容易。
- 新的 API 中更广泛地使用了 context 对象, 并使用 MapContext 进行 MapReduce 间的通信, MapContext 同时充当 OutputCollector 和 Reporter 的角色。
- Job 的配置统一由 Configuration 来完成, 而不必额外地使用 JobConf 对守护进程进行配置。
- 由 Job 类来负责 Job 的控制, 而不是 JobClient, JobClient 在新的 API 中已经被删除。

此外, 新的 API 同时支持“推”和“拉”式的迭代方式, 在以往的操作中, `<key, value>` 对是被推入到 map 中的, 但是在新的 API 中, 允许程序将数据拉入 map 中, reduce 也一样, 这样做更加方便程序分批处理数据。

4.2.3 MapReduce 的数据流和控制流

前面已经提到了 MapReduce 的数据流和控制流的关系, 本节将结合 WordCount 实例具体解释它们的含义。图 4-2 是上例中 WordCount 程序的执行流程:

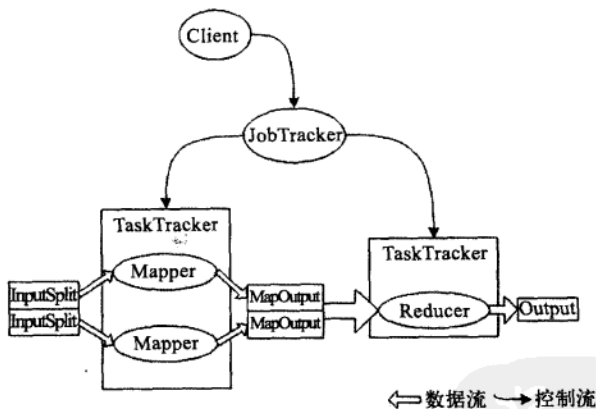


图 4-2 MapReduce 工作的简易图

由前面的内容可知, 负责控制及调度 MapReduce 的 job 的是 JobTracker, 负责运行 MapReduce 的 job 是 TaskTracker。当然, MapReduce 在运行时是分成 map task 和 reduce task 来处理的, 而不是完整的 job。简单的控制流大概是这样的: JobTracker 调度任务给 TaskTracker, TaskTracker 执行任务时, 会返回进度报告。JobTracker 则会记录进度的进行状况, 如果某个 TaskTracker 上的任务执行失败, 那么 JobTracker 会把这个任务分配给另一台 TaskTracker, 直到任务执行完成。

这里更详细地解释一下数据流。上例中有两个 map 任务及一个 reduce 任务。数据首先按照 TextInputFormat 形式被处理成两个 InputSplit, 然后输入到两个 map 中, map 程序会读

取 InputSplit 指定的位置的数据，然后按照设定的方式处理此数据，最后写入本地磁盘中。注意，这里并不是写到 HDFS 上，这应该很好理解，因为 map 的输出在 job 完成后即可删除了，因此不需要存储到 HDFS 上，虽然存储到 HDFS 上会更安全。但是，由于网络传输会降低 MapReduce 任务的执行效率，因此 map 的输出文件是写在本地磁盘上的。如果 map 程序在没来得及将数据传送给 reduce 时就崩溃了（程序出错或机器崩溃），那么 JobTracker 只需要另选一台机器重新执行这个 task 就可以了。

Reduce 会读取 map 的输出数据，合并 value，然后将它们输出到 HDFS 上。Reduce 的输出会占用很多的网络带宽，不过这与上传数据一样，是不可避免的。如果大家还是不能很好地理解数据流，下面有一个更具体的图（WordCount 执行时的数据流），如图 4-3 所示：

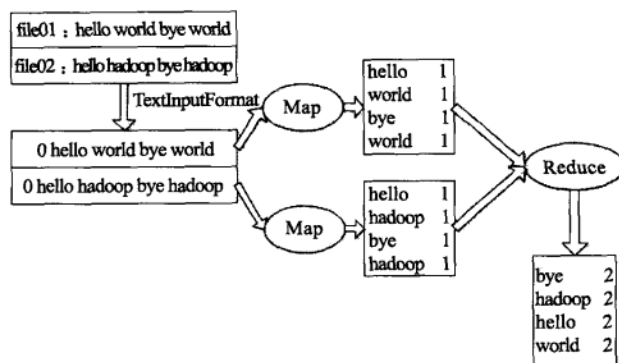


图 4-3 WordCount 数据流程图

相信看到图 4-3，大家应该能对 MapReduce 的执行过程有更深刻的了解了。

除此之外，还有两种情况需要注意：

1) MapReduce 在执行过程中往往不止一个 reduce task，reduce task 的数量是可以通过程序指定的，当存在多个 reduce task 时，每个 reduce 会收集一个或多个 key 值。需要注意的是，当出现多个 reduce task 时，每个 reduce task 都会生成一个输出文件。

2) 另外，没有 reduce 任务的时候，系统会直接将 map 的输出结果作为最终结果，同时 map task 的数量可以看做是 reduce task 的数量，即有多少个 map task 就有多少个输出文件。

4.3 MapReduce 任务的优化

相信每个程序员在编程时都会问自己两个问题：“我如何完成这个任务”，以及“如何能让程序运行得更快”。同样，MapReduce 计算模型的多次优化也是为了更好地解答这两个问题。

MapReduce 计算模型的优化涉及了方方面面的内容，但是主要集中在两个方面：一是计算性能方面的优化；二是 I/O 操作方面的优化。下面我们将重点讲解以下几个方面的内容。

1. 任务调度

任务调度是 Hadoop 中非常重要的一环，它的优化则又涉及两个方面的内容。计算方面：Hadoop 总会优先将任务分配给空闲的机器，使所有的任务能公平地分享系统资源；I/O 方面：Hadoop 会尽量将 map 任务分配给 InputSplit 所在的机器，以减少网络 I/O 的消耗。

2. 数据预处理与 InputSplit 的大小

MapReduce 任务擅长处理少量的大数据，而在处理大量的小数据时，MapReduce 的性能就要逊色很多。因此在提交 MapReduce 任务前可以先对数据进行一次预处理，将数据合并以提高 MapReduce 任务的执行效率，这招往往很有效。如果这还不够，可以参考 map 任务的运行时间，当一个 map 任务只需要运行几秒钟就可以结束时，就需要考虑是否应该给它分配更多的数据了。通常而言，一个 map 任务的运行时间在一分钟左右比较合适，可以通过设置 map 的输入数据大小来调节 map 的运行时间。在 FileInputFormat 中（除了 CombineFileInputFormat），Hadoop 会在处理每个 block 后将其作为一个 InputSplit，因此合理地设置 block 块大小是很重要的调节方式。除此以外，也可以通过合理地设置 map 任务的数量来调节 map 任务的数据输入。

3. map 和 reduce 任务的数量

合理地设置 map 任务与 reduce 任务的数量对提高 MapReduce 任务的效率是非常重要的，默认的设置往往不能很好地体现出 MapReduce 任务的需求，不过，设置它们的数量也要有一定的实践经验。

首先，定义两个概念，就是 map 任务槽和 reduce 任务槽。map/reduce 任务槽就是这个集群能够同时运行的 map/reduce 任务的最大数量。比如，在一个具有 1200 台机器的集群中，设置每台机器最多可以同时运行 10 个 map 任务，5 个 reduce 任务，那么这个集群的 map 任务槽就是 12000，reduce 任务槽是 6000。

可以通过任务槽对任务调度进行设置。设置 MapReduce 任务的 map 数量时主要参考的还是 map 的运行时间，不过设置 reduce 任务的数量时就只需要参考任务槽的设置了。一般来说，reduce 任务的数量应该是 reduce 任务槽的 0.95 倍或是 1.75 倍，这是基于不同的考虑而决定的。当 reduce 任务的数量是任务槽的 0.95 倍时，如果一个 reduce 任务失败，Hadoop 可以很快地找到一台空闲的机器重新执行这个任务。当 reduce 任务的数量是任务槽的 1.75 倍时，执行速度快的机器可以获得更多的 reduce 任务，因此可以使负载更加均衡，以提高任务的处理速度。

4. combine 函数

combine 函数是用于在本地合并数据的函数，在有些情况下，map 函数产生的中间数据会有很多重复的数据，比如在一个简单的 WordCount 程序中，因为词频是接近于一个 zipf 分布的，每个 map 任务可能会产生成千上万个 <the, 1> 记录，若将这些记录一一传送给 reduce 任务是会很耗时的，所以，MapReduce 框架运行用户写的一个 combine 函数，用于本地合并，这会大大减少网络 I/O 操作的消耗。那么，此时就可以利用 combine 函数先计算出在这个

block 中单词 the 的个数。合理地设计 combine 函数会有效地减少网络传输的数据量, 提高 MapReduce 的效率。

在 MapReduce 程序中使用 combine 很简单, 只需在程序中添加如下内容:

```
job.setCombinerClass(combine.class);
```

在 WordCount 程序中, 可以指定 reduce 函数为 combine 函数, 如下所示:

```
job.setReducerClass(Reduce.class);
```

5. 压缩

编写 MapReduce 程序时, 可以选择对 map 的输出和最终的输出结果进行压缩 (同时也可以选择压缩方式)。在一些情况下, map 的中间输出可能会很大, 对其进行压缩可以有效地减少网络上的数据传输量。对最终结果的压缩虽然会减少数据写 HDFS 的时间, 但是也会对该取产生一定的影响, 因此要根据实际情况来选择 (第 8 章中提供了一个小实验来验证压缩的效果)。

6. 自定义 comparator

在 Hadoop 中, 可以自定义数据类型以实现更复杂的目的, 比如, 当读者想实现 k-means 算法 (一个基础的聚类算法) 时可以定义 k 个整数的集合。自定义 Hadoop 数据类型时, 推荐自定义 comparator 来实现数据的二进制比较, 这样可以省去数据序列化和反序列化的时间, 提高程序的运行效率 (具体会在第 8 章中讲解)。

4.4 Hadoop 流

Hadoop 流提供了一个 API, 允许用户使用任何脚本语言编写 map 函数或 reduce 函数。Hadoop 流的关键是, 它使用 UNIX 标准流作为程序与 Hadoop 之间的接口。因此, 任何程序只要可以从标准输入流中读取数据, 并且可以写入数据到标准输出流, 那么就可以通过 Hadoop 流使用其他语言编写 MapReduce 程序的 map 函数或 reduce 函数。

举个最简单的例子 (Ubuntu, hadoop-0.20.2):

```
bin/hadoop jar contrib/streaming/hadoop-0.20.2-streaming.jar -input input -output output -mapper /bin/cat -reducer usr/bin/wc
```

从这个例子中可以看到, Hadoop 流引入的包是 hadoop-0.20.2-streaming.jar, 并且具有如下命令:

-input	指明输入文件路径
-output	指明输出文件路径
-mapper	指定 map 函数
-reducer	指定 reduce 函数

Hadoop 流的操作还有其他参数, 稍后会一一列出。

4.4.1 Hadoop 流的工作原理

先来看 Hadoop 流的工作原理。在上例中，map 和 reduce 都是 Linux 内的可执行文件，更重要的是，它们接收的都是标准输入（stdin），输出的都是标准输出（stdout）。如果大家熟悉 Linux，那么对它们一定不会陌生。

执行这个程序，如下所示：

程序的输入与 WordCount 程序是一样的，如下：

```
file01:
hello world bye world
file02
hello hadoop bye hadoop
```

输入命令：

```
bin/hadoop jar contrib/streaming/hadoop-0.20.2-streaming.jar -input input -output
output -mapper /bin/cat -reducer /usr/bin/wc
```

显示：

```
packageJobJar: [/root/tmp/hadoop-unjar7103575849190765740/] [] /tmp/
streamjob2314757737747407133.jar tmpDir=null
11/01/23 02:07:36 INFO mapred.FileInputFormat: Total input paths to process : 2
11/01/23 02:07:37 INFO streaming.StreamJob: getLocalDirs(): [/root/tmp/mapred/
local]
11/01/23 02:07:37 INFO streaming.StreamJob: Running job: job_201101111819_0020
11/01/23 02:07:37 INFO streaming.StreamJob: To kill this job, run:
11/01/23 02:07:37 INFO streaming.StreamJob: /root/hadoop/bin/hadoop job -Dmapred.
job.tracker=localhost:9001 -kill job_201101111819_0020
11/01/23 02:07:37 INFO streaming.StreamJob: Tracking URL: http://localhost:50030/
jobdetails.jsp?jobid=job_201101111819_0020
11/01/23 02:07:38 INFO streaming.StreamJob: map 0% reduce 0%
11/01/23 02:07:47 INFO streaming.StreamJob: map 100% reduce 0%
11/01/23 02:07:59 INFO streaming.StreamJob: map 100% reduce 100%
11/01/23 02:08:02 INFO streaming.StreamJob: Job complete: job_201101111819_0020
11/01/23 02:08:02 INFO streaming.StreamJob: Output: output
```

程序的输出是：

```
2      8      46
```

稍微解释一下这个结果，wc 命令用来统计文件中的行数、单词数与字节数，可以看到，这个结果是正确的。

Hadoop 流的工作原理并不复杂，其中 map 的工作原理如图 4-4 所示（reduce 与其相同）。

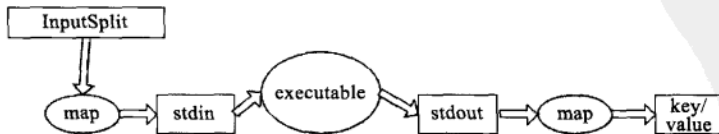


图 4-4 Hadoop 流 map 流程图

当一个可执行文件作为 Mapper 时，每一个 map 任务会以一个独立的进程启动这个可执行文件，然后在 map 任务运行时，会把输入切分成行提供给可执行文件，并作为它的标准输入（stdin）内容。当可执行文件运行出结果时，map 从标准输出（stdout）中收集数据，并将其转化为 <key, value> 对，作为 map 的输出。

reduce 与 map 相同，如果可执行文件作为 Reducer 时，reduce 任务会启动这个可执行文件，并且将 <key, value> 转化为行来作为这个可执行文件的标准输入（stdin）。然后 reduce 会收集这个可执行文件的标准输出（stdout）内容，并把每一行转化为 <key, value> 对，作为 reduce 的输出。

map 与 reduce 将输出转化为 <key,value> 对的默认方法是，将每行的第一个 tab 符号（制表符）之前的内容作为 key，之后的内容作为 value。如果没有 tab 符号，那么这一行的所有内容会作为 key，而 value 值为 null。当然这是可以更改的。

值得一提的是，可以使用 Java 类作为 map，而用一个可执行程序作为 reduce；或使用 Java 类作为 reduce，而用可执行程序作为 map。如下例所示：

```
/bin/hadoop jar contrib/streaming/hadoop-0.20.2-streaming.jar
-input myInputDirs -output myOutputDir -mapper
org.apache.hadoop.mapred.lib.IdentityMapper -reducer /bin/wc
```

4.4.2 Hadoop 流的命令

Hadoop 流提供自己的流命令选项及一个通用的命令选项，用于设置 Hadoop 流任务。首先介绍一下流命令。

1. Hadoop 流命令选项

Hadoop 流命令的具体内容如表 4-1 所示：

表 4-1 Hadoop 流命令

参 数	可选 / 必选	参 数	可选 / 必选
-input	必选	-cmdenv	可选
-output	必选	-inputreader	可选
-mapper	必选	-verbose	可选
-reducer	必选	-lazyoutput	可选
-file	可选	-numreduce tasks	可选
-inputformat	可选	-mapdebug	可选
-outputformat	可选	-reducedebg	可选
-partitioner	可选	-io	可选
-combiner	可选		

流命令中，必选的 4 个很好理解，分别用于指定输入输出文件的位置及 map 函数和 reduce 函数。在其他的可选命令中，下面只解释几个常用的。

❑ -file

-file 指令用于将文件加入到 Hadoop 的 Job 中。在上例中, cat 和 wc 都是 Linux 系统中的命令, 而在 Hadoop 流的使用中, 往往需要使用自己写的文件 (作为 map 或 reduce), 一般而言, 这些文件是 Hadoop 集群中的机器上没有的, 这时就需要使用 Hadoop 流中的 -file 命令将这个可执行文件加入到 Hadoop 的 Job 中了。

❑ -combiner

这个命令用来加入 combine 程序。

❑ -inputformat 和 -outputformat

这两个命令用来设置输入输出文件的处理方法, 这两个命令后面的参数必须是 Java 类。

2. Hadoop 流通用的命令选项

Hadoop 流的通用命令是用来配置 Hadoop 流的 Job 的。需要注意的是, 如果使用这部分配置, 则必须将其置于流命令配置之前, 否则命令会失败。这里简要列出命令列表 (如表 4-2 所示), 供大家参考。

表 4-2 Hadoop 流的 Job 设置命令

参数	可选 / 必选	参数	可选 / 必选
-conf	可选	-files	可选
-D	可选	-libjars	可选
-fs	可选	-archives	可选
-jt	可选		

4.4.3 实战案例: 添加 Bash 程序和 Python 程序到 Hadoop 流中

从上面的内容可以知道, Hadoop 流的 API 是一个扩展性非常强的框架, 因为它与程序相连的部分只有数据, 因此可以接收任何适用于 UNIX 标准输入输出的脚本语言。比如 Bash、PHP、Ruby、Python 等。

下面举出两个非常简单的例子来进一步说明它的特性。

1. 添加 Bash 程序到 Hadoop 流中

MapReduce 框架是一个非常适合在大规模的非结构化数据中查找数据的编程模型, grep 就是一个这种类型的例子。

在 Linux 中, grep 命令用来在一个或多个文件中查找某个字符模式 (这个字符模式可以代表字符串, 多用正则表达式表示)。

下面尝试在如下的数据中查找带有 Hadoop 字符串的行, 如下所示:

输入文件为:

file01:

hello world bye world


```
file02:
hello      hadoop bye hadoop
```

```
reduce 文件为:
reduce.sh:
grep hadoop
```

输入命令为:

```
bin/hadoop jar contrib/streaming/hadoop-0.20.2-streaming.jar -input input -output
output -mapper /bin/cat -reducer ~/Desktop/test/reducer.sh -file ~/Desktop/test/
reducer.sh
```

结果为:

```
hello      hadoop bye hadoop
```

显然, 这个结果是正确的。

2. 添加 Python 程序到 Hadoop 流中

对于 Python 来说, 情况稍微有些特殊, 因为 Python 是可以编译为 jar 包的, 如果将程序编译为 jar 包, 那么就可以采用运行 jar 包的方式来运行了。

不过, 同样也可以用流的方式运行 Python 程序。请看如下代码:

```
Reduce.py
#!/usr/bin/python

import sys;

def generateLongCountToken(id):
    return "LongValueSum:" + id + "\t" + "1"

def main(argv):
    line = sys.stdin.readline();
    try:
        while line:
            line = line[:-1];
            fields = line.split("\t");
            print generateLongCountToken(fields[0]);
            line = sys.stdin.readline();
        except "end of file":
            return None
    if __name__ == "__main__":
        main(sys.argv)
```

使用如下命令来运行:

```
bin/hadoop jar contrib/streaming/hadoop-0.20.2-streaming.jar -input input -output
pyoutput -mapper reduce.py -reducer aggregate -file reduce.py
```

注意其中的 aggregate, 它是 Hadoop 提供的一个包, 它提供一个 reduce 函数和一个 combine 函数。这个函数实现一些简单的类似求和、取最大值最小值等功能。

4.5 Hadoop Pipes

Hadoop Pipes 提供了一个在 Hadoop 上运行 C++ 程序的方法。与流不同的是，流使用的是标准输入输出作为可执行程序与 Hadoop 相关进程间通信的工具，而 Pipes 使用的是 Sockets。

先看示例程序：

```
wordcount.cpp
#include "hadoop/Pipes.hh"
#include "hadoop/TemplateFactory.hh"
#include "hadoop/StringUtils.hh"

const std::string WORDCOUNT = "WORDCOUNT";
const std::string INPUT_WORDS = "INPUT_WORDS";
const std::string OUTPUT_WORDS = "OUTPUT_WORDS";

class WordCountMap: public HadoopPipes::Mapper {
public:
    HadoopPipes::TaskContext::Counter* inputWords;

    WordCountMap(HadoopPipes::TaskContext& context) {
        inputWords = context.getCounter(WORDCOUNT, INPUT_WORDS);
    }

    void map(HadoopPipes::MapContext& context) {
        std::vector<std::string> words =
            HadoopUtils::splitString(context.getInputValue(), " ");
        for(unsigned int i=0; i < words.size(); ++i) {
            context.emit(words[i], "1");
        }
        context.incrementCounter(inputWords, words.size());
    }
};

class WordCountReduce: public HadoopPipes::Reducer {
public:
    HadoopPipes::TaskContext::Counter* outputWords;

    WordCountReduce(HadoopPipes::TaskContext& context) {
        outputWords = context.getCounter(WORDCOUNT, OUTPUT_WORDS);
    }

    void reduce(HadoopPipes::ReduceContext& context) {
        int sum = 0;
        while (context.nextValue()) {
            sum += HadoopUtils::toInt(context.getInputValue());
        }
        context.emit(context.getInputKey(), HadoopUtils::toString(sum));
    }
};
```

```

        context.incrementCounter(outputWords, 1);
    }
};

int main(int argc, char *argv[]) {
    return HadoopPipes::runTask(HadoopPipes::TemplateFactory<WordCountMap,
        WordCountReduce>());
}

```

这个程序连接的是一个 C++ 库，这个程序的结构类似于 Java 编写的程序。如新版 API 一样，这个程序使用 context 方法读入和收集 <key, value> 对。在使用时要重写 HadoopPipes 名字空间下的 mapper 函数和 reducer 函数，并用 context.emit() 方法输出 <key, value> 对。main 函数是应用程序的入口，它调用 HadoopPipes::runTask 方法，这个方法有一个 TemplateFactory 参数来创建 map 和 reduce 实例，也可以重载 factory 设置 combiner()、partitioner()、record reader、record writer。

接下来，编译这个程序。这个编译命令需要用到 g++，读者可以使用 apt 自动安装这个程序，如下所示：

```
apt-get install g++
```

然后建立文件 Makefile，如下所示：

```

HADOOP_INSTALL="你的hadoop安装文件夹"
PLATFORM=Linux-i386-32 (如果是AMD的CPU, 请使用Linux-amd64-64)

CC = g++
CPPFLAGS = -m32 -I$(HADOOP_INSTALL)/c++/$(PLATFORM)/include

```

```

wordcount: wordcount.cpp
    $(CC) $(CPPFLAGS) $< -Wall
    -L$(HADOOP_INSTALL)/c++/$(PLATFORM)/lib -lhadooppipes
    -lhadooputils -lpthread -g -O2 -o $@
注意在 $(CC) 前有一个 <tab> 符号，这个分隔符是很关键的。

```

这会在当前目录下建立一个 WordCount 可执行文件。

然后，上传可执行文件到 HDFS 上，这是为了 TaskTracker 能够获得这个可执行文件。这里上传到 bin 文件夹内。

```

~/hadoop/bin/hadoop fs -mkdir bin
~/hadoop/bin/hadoop dfs -put wordcount bin

```

然后，就可以运行这个 MapReduce 程序了，可以采用两种配置方式运行这个程序。一种方式是直接在命令中运行指定配置，如下所示：

```

~/hadoop/bin/hadoop pipes\
-D hadoop.pipes.java.recordreader=true\
-D hadoop.pipes.java.recordwriter=true\
-input input\

```

```
-output Coutput\  
-program bin/wordcount
```

另一种方式是预先配置写入配置文件中，如下所示：

```
<?xml version="1.0"?>  
<configuration>  
  <property>  
    // Set the binary path on DFS  
    <name>hadoop.pipes.executable</name>  
    <value>bin/wordcount</value>  
  </property>  
  <property>  
    <name>hadoop.pipes.java.recordreader</name>  
    <value>true</value>  
  </property>  
  <property>  
    <name>hadoop.pipes.java.recordwriter</name>  
    <value>true</value>  
  </property>  
</configuration>
```

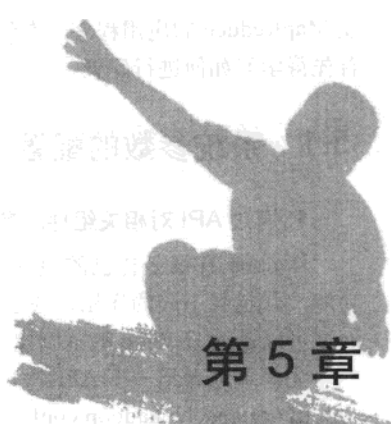
然后通过如下命令运行这个程序：

```
~/hadoop/bin/hadoop pipes -conf word.xml -input input -output output
```

将参数 `hadoop.pipes.executable` 和 `hadoop.pipes.java.recordreader` 设置为 `true` 表示使用 Hadoop 默认的输入输出方式（即 Java 的）。同样的，也可以设置一个 Java 语言编写的 mapper 函数、reducer 函数、combiner 函数和 partitioner 函数。实际上，在任何一个作业中，都可以混用 Java 类和 C++ 类。

4.6 小结

本章主要介绍了 MapReduce 的计算模型。其中的关键内容是一个流程和四个方法。一个流程指的是数据流程，输入数据到 $\langle k1, v1 \rangle$ 、 $\langle k1, v1 \rangle$ 到 $\langle k2, v2 \rangle$ 、 $\langle k2, v2 \rangle$ 到 $\langle k3, v3 \rangle$ 、 $\langle k3, v3 \rangle$ 到输出数据。四个方法就是这个数据转换过程中使用的方法（分别是 `InputFormat`、`map`、`reduce`、`OutputFormat`），以及其对应的转换过程。除此之外，还介绍了这个编程框架的几个优化方法，以及 Hadoop 流和 Hadoop Pipes，后者是在 Hadoop 中使用脚本文件及 C++ 编写 MapReduce 程序的方法。



第 5 章

开发 MapReduce 应用程序

本章内容

- ☐ 系统参数的配置
- ☐ 配置开发环境
- ☐ 编写 MapReduce 程序
- ☐ 本地测试
- ☐ 运行 MapReduce 程序
- ☐ 网络用户界面
- ☐ 性能调优
- ☐ MapReduce 工作流
- ☐ 小结

资源分享

PDG

在前面的章节中，已经介绍了 MapReduce 模型。在本章中，将介绍如何在 Hadoop 中开发 MapReduce 的应用程序。在编写 MapReduce 程序之前，需要安装和配置开发环境。因此，首先得学习如何进行配置。

5.1 系统参数的配置

1. 通过 API 对相关组件的参数进行配置

Hadoop 有很多自己的组件（例如 HBase、Chukwa 等），每一种组件都可以实现不同的功能，并起着不同的作用，通过多种组件的配合使用，Hadoop 就能够实现非常强大的功能。而这可以通过 Hadoop 的 API 来进行配置。

先简单地介绍一下 API^①，它被分成了以下几个部分（也就是几个不同的包）：

- org.apache.hadoop.conf：定义了系统参数的配置文件处理 API；
- org.apache.hadoop.fs：定义了抽象的文件系统 API；
- org.apache.hadoop.hdfs：Hadoop 分布式文件系统（HDFS）模块的实现；
- org.apache.hadoop.mapred：Hadoop 分布式计算系统（MapReduce）模块的实现，包括任务的分发调度等；
- org.apache.hadoop.ipc：用于网络服务端和客户端的工具，封装了网络异步 I/O 的基础模块；
- org.apache.hadoop.io：定义了通用的 I/O API，用于针对网络、数据库、文件等数据对象进行读写操作等。

在此我们需要用到的是 org.apache.hadoop.conf，用它来定义系统参数的配置。Configurations 类由源来设置，每个源包含以 XML 形式出现的一系列属性/值对。每个源以一个字符串或一个路径来命名。如果是字符串命名，则类路径检查该字符串代表的路径存在与否；如果是路径命名的，则直接通过本地文件系统进行检查，而用不着类路径。

下面举一个配置文件的例子，如代码清单 5-1 所示。

代码清单 5-1 configuration-default.xml

```
<? xml version="1.0"? >
<configuration>

  <property>
    <name>hadoop.tmp.dir</name>
    <value>/tmp/hadoop-${usr.name}</value>
    <description>A base for other temporary directories.</description>
  </property>

  <property>
```

① 可以参考 <http://hadoop.apache.org/common/docs/current/api>。

```

    <name>io.file.buffer.size</name>
    <value>4096</value>
    <description>the size of buffer for use in sequence file.</description>
</property>

<property>
    <name>height</name>
    <value>tall</value>
    <final>true</final>
</property>
</configuration>

```

这个文件里的信息可以通过以下的方式进行抽取：

```

Configuration conf = new Configuration();
Conf.addResource("configuration-default.xml");
assertThat(conf.get("hadoop.tmp.dir"), is("/tmp/hadoop-${usr.name}"));
assertThat(conf.get("io.file.buffer.size"), is("4096"));
assertThat(conf.get("height"), is("tall"));

```

2. 多个配置文件的整合

现在假设还有另外一个配置文件 configuration-site.xml, 它的具体代码细节如代码清单 5-2 所示。

代码清单 5-2 configuration-site.xml

```

<? xml version="1.0"? >
<configuration>

    <property>
        <name>io.file.buffer.size</name>
        <value>5000</value>
        <description>the size of buffer for use in sequence file.</description>
    </property>

    <property>
        <name>height</name>
        <value>short</value>
        <final>true</final>
    </property>
</configuration>

```

现在使用两个资源 configuration-default.xml 和 configuration-site.xml 来定义配置。将资源按顺序添加到 Configuration 中, 如下所示:

```

Configuration conf = new Configuration();
conf.addResource("configuration-default.xml");

```

```
conf.addResource("configuration-site.xml");
```

现在不同的资源当中有了相同属性，但是它们的取值却不一样，这时这些属性的取值应该如何确定呢？可以遵循这样一个原则：后添加进来的属性取值覆盖掉前面所添加资源中的属性取值。因此，此处的属性 `io.file.buffer.size` 取值应该是 5000 而不是先前的 4096，即：

```
assertThat(conf.get("io.file.buffer.size"),is("5000"));
```

但是，有一个特例，被标记为 `final` 的属性不能被后面定义的属性覆盖。`configuration-default.xml` 中的属性 `height` 被标记为 `final`，因此在 `configuration-site.xml` 中重写 `height` 并不会成功，它依然会从 `configuration-default.xml` 中取值，如下所示：

```
assertThat(conf.get("height"),is("tall"));
```

重写标记为 `final` 的属性通常情况下会报告配置错误，同时会有警告信息被记录下来以便诊断所用。管理员将守护进程地址文件之中的属性标记为 `final`，可防止用户在客户端配置文件中或在作业提交参数中改变其取值。

Hadoop 默认使用两个源进行配置，并按顺序加载 `core-default.xml` 和 `core-site.xml`。在实际应用中可能会添加其他的源，应按照它们添加的顺序进行加载。其中 `core-default.xml` 定义系统默认的属性，`core-site.xml` 定义在特定的地方重写。

5.2 配置开发环境

首先下载准备使用的 Hadoop 版本，然后将其解压到开发机器上面（详细过程见附录 A）。接下来，在集成开发环境中创建一个新的工程，然后将解压后的文件夹根目录下的 JAR 文件和 `lib` 目录之下的 JAR 文件加入到 `classpath` 中。之后就可以编译 Hadoop 程序了，并且可以在集成开发环境中以本地模式运行。

Hadoop 有三种不同的运行方式：本地模式、伪分布模式、完全分布模式。三种不同的运行方式各有各的好处与不足之处：本地模式安装与配置比较简单，运行在本地文件系统中，便于程序的调试，可及时查看程序运行的效果，但是当数据量比较大时，运行的速度会比较慢，并且没有体现出 Hadoop 分布式的优点；伪分布模式同样是在本地文件系统中运行，与本地模式的不同之处在于它运行的文件系统为 HDFS，好处是能够模仿完全分布模式，看到一些分布式处理的效果；完全分布模式则运行在多台机器的 HDFS 之上，完完全全地体现出了分布式的优点，但是调试程序会比较麻烦。

所以在实际运用中，可以结合这三种不同模式的优点，比如，编写和调试程序在本地模式和伪分布模式上进行，而实际处理大数据，则运行在完全分布模式下。而这就会涉及三种不同模式的配置与管理了，相关配置和管理会有相应的章节重点讲解，下面只分别简单介绍一下三种不同模式的配置。三种不同模式的配置文件分别取名叫 `hadoop-local.xml`（对应本地模式）、`hadoop-localhost.xml`（对应伪分布模式）、`hadoop-cluster.xml`（对应完全分布模式）。

hadoop-local.xml 中包含 Hadoop 默认的文件系统和 JobTracker, 如下所示:

```
<?xml version="1.0">
<configuration>

  <property>
    <name>fs.default.name</name>
    <value>file:///</value>
  </property>

  <property>
    <name>mapred.job.tracker</name>
    <value>local</value>
  </property>

</configuration>
```

hadoop-localhost.xml 指出 NameNode 和 JobTracker 都在本地文件系统上, 如下所示:

```
<?xml version="1.0">
<configuration>

  <property>
    <name>fs.default.name</name>
    <value>hdfs://localhost:9000</value>
  </property>

  <property>
    <name>mapred.job.tracker</name>
    <value>localhost:9001</value>
  </property>

  <property>
    <name>dfs.replication</name>
    <value>1</value>
  </property>

</configuration>
```

hadoop-cluster.xml 包含了集群中 NameNode 和 JobTracker 的详细地址, 这个应该根据实际的环境而定。

5.3 编写 MapReduce 程序

下面将通过一个计算学生平均成绩的例子, 来讲解开发 MapReduce 程序的流程。程序主要包括两部分的内容: Map 部分和 Reduce 部分, 分别实现了 map 和 reduce 的功能。

5.3.1 Map 处理

Map 处理的是一个纯文本文件, 文件中存放的数据是每一行表示一个学生的姓名和他相

应一科的成绩，如果有多门学科，则每个学生就存在多行数据。代码如下所示：

```
public static class Map
    extends Mapper<LongWritable, Text, Text, IntWritable> {
    public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {
        String line = value.toString(); // 将输入的纯文本文件的数据转化成 String
        System.out.println(line); // 为了便于程序的调试，输出读入的内容

        // 将输入的数据首先按行进行分割
        StringTokenizer tokenizerArticle = new StringTokenizer(line, "\n");
        // 分别对每一行进行处理
        while(tokenizerArticle.hasMoreTokens()){
            // 每行按空格划分
            StringTokenizer tokenizerLine = new StringTokenizer(tokenizerArticle.
                nextToken());
            String strName = tokenizerLine.nextToken(); // 学生姓名部分
            String strScore = tokenizerLine.nextToken(); // 成绩部分

            Text name = new Text(strName); // name of student
            int scoreInt = Integer.parseInt(strScore); // score of student
            context.write(name, new IntWritable(scoreInt)); // 输出姓名和成绩
        }
    }
}
```

通过数据集进行测试，结果显示完全可以将文件中的姓名和相应的成绩提取出来。需要解释的是：Mapper 处理的数据是由 InputFormat 分解过的数据集，其中 InputFormat 的作用是将数据集切割成小数据集 InputSplits，每一个 InputSplit 将由一个 Mapper 负责处理。此外，InputFormat 中还提供了一个 RecordReader 的实现，并将一个 InputSplit 解析成 <key,value> 对提供给了 map 函数。InputFormat 的默认值是 TextInputFormat，它针对文本文件，按行将文本切割成 InputSplits，并用 LineRecordReader 将 InputSplit 解析成 <key,value> 对，key 是行在文本中的位置，value 是文件中的一行。

本程序中的 InputFormat 使用的是默认值 TextInputFormat，因此结合上述程序的注释部分不难理解整个程序的处理流程和正确性。

5.3.2 Reduce 处理

Map 的结果会通过 partition 分发到 Reducer，Reducer 做完 Reduce 操作后，将通过 OutputFormat 输出，代码如下所示：

```
public static class Reduce
    extends Reducer<Text, IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterable<IntWritable> values,
        Context context) throws IOException, InterruptedException {

        int sum = 0;
```

```

int count=0;
Iterator<IntWritable> iterator = values.iterator();
while (iterator.hasNext()) {
    sum += iterator.next().get(); // 计算总分
    count++; // 统计总的科目数
}
int average = (int) sum/count; // 计算平均成绩
context.write(key, new IntWritable(average));
}
}

```

Mapper 最终处理的结果对 <key,value>, 会送到 Reducer 中进行合并, 合并的时候, 有相同 key 的键/值对则送到同一个 Reducer 上。Reducer 是所有用户定制 Reducer 类的基类, 它的输入是 key 和这个 key 对应的所有 value 的一个迭代器, 同时还有 Reducer 的上下文。Reduce 的结果由 Reducer.Context 的 write 方法输出到文件中。

5.4 本地测试

Score_Process 类继承于 Configured 的实现接口 Tool, 上述的 Map 和 Reduce 是 Score_Process 的内部类, 它们分别实现了 map 和 reduce 功能, 主函数存在于 Score_Process 中。下面创建一个 Score_Process 实例对程序进行测试。

Score_Process 的 run() 方法实现如下:

```

public int run(String [] args) throws Exception {
    Job job = new Job(getConf());
    job.setJarByClass(Score_Process.class);
    job.setJobName("Score_Process");

    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);

    job.setMapperClass(Map.class);
    job.setCombinerClass(Reduce.class);
    job.setReducerClass(Reduce.class);

    job.setInputFormatClass(TextInputFormat.class);
    job.setOutputFormatClass(TextOutputFormat.class);

    FileInputFormat.setInputPaths(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));

    boolean success = job.waitForCompletion(true);
    return success ? 0 : 1;
}

```

下面给出 main() 函数, 对程序进行测试:



```

public static void main(String[] args) throws Exception {
    int ret = ToolRunner.run(new Score_Process(), args);
    System.exit(ret);
}

```

程序在 Eclipse 中执行的时候，要在 run configuration 中设置好参数，输入的文件夹名为 input，输出的文件夹名为 output，请参考图 5-1。

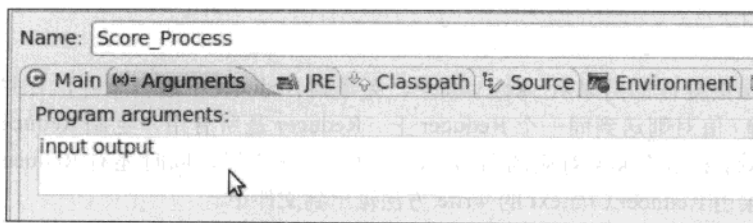


图 5-1 运行参数设置

其中输入 input 文件夹中的文件 input 的数据如图 5-2 所示。

陈洲立	67
陈东伟	90
李宁	87
杨森	86
刘东奇	78
谭果	94
盖盖	83
陈洲立	68
陈东伟	96
李宁	82
杨森	85
刘东奇	72
谭果	97
盖盖	82
陈洲立	46
陈东伟	48
李宁	67
杨森	33
刘东奇	28
谭果	78
盖盖	87

图 5-2 输入文件中的数据

运行过程的截图如图 5-3 所示。

运行的结果保存在 output 文件夹下，执行结果如图 5-4 所示。

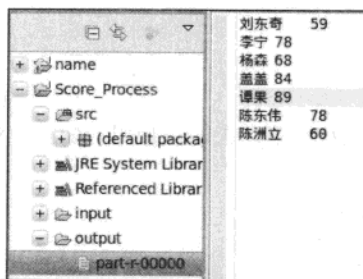


```

信息: Task 'attempt local_0001_r_000000_0' done.
2011-1-9 11:30:21 org.apache.hadoop.mapred.JobClient monitorAndPrintJob
信息: map 100% reduce 100%
2011-1-9 11:30:21 org.apache.hadoop.mapred.JobClient monitorAndPrintJob
信息: Job complete: job local_0001
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log
信息: Counters: 12
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log
信息: FileSystemCounters
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log
信息: FILE_BYTES_READ=27628
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log
信息: FILE_BYTES_WRITTEN=54859
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log
信息: Map-Reduce Framework
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log
信息: Reduce input groups=7
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log
信息: Combine output records=7
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log
信息: Map input records=23
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log
信息: Reduce shuffle bytes=0
2011-1-9 11:30:21 org.apache.hadoop.mapred.Counters log

```

图 5-3 部分运行过程



刘东奇	59
李宁	78
杨森	68
盖盖	84
谭栗	89
陈东伟	78
陈洲立	60

图 5-4 输出文件中的数据

5.5 运行 MapReduce 程序

想要测试人体的健康状况，先要知道各个人体组织的健康状况，然后再综合评价人体的健康状况。假设每个组织的健康指标是一个 0 到 100 之间的数字，得到综合身体健康状况的方法是计算所有组织健康指标的平均数。由于测试的人数众多，因此存储数据的格式为：姓名 + 得分 + #（代表一个人单个人体组织的健康状况），每个组织的健康状况分别用一个文件存储。现在一共有 1000 个组织参与了评估，即用 1000 个文件分别存储。

此例中对数据的处理与前面对学生成绩所进行的简单处理有一些区别，下面先将程序的主要部分列举出来。

Mapper 部分的代码如下：

```
public static class Map
```

```

    extends Mapper<LongWritable, Text, Text, IntWritable> {

        public void map(LongWritable key, Text value, Context context)
            throws IOException, InterruptedException {
            String line = value.toString();

            // 以“#”为分隔符, 将输入的文件分割成单个记录
            StringTokenizer tokenizerArticle = new StringTokenizer(line, "#");
            // 对每个记录进行处理
            while(tokenizerArticle.hasMoreTokens()){
                // 将每个记录分成姓名和分数两个部分
                StringTokenizer tokenizerLine = new StringTokenizer(tokenizerArticle.
                    nextToken());
                while(tokenizerLine.hasMoreTokens()){
                    String strName = tokenizerLine.nextToken();
                    if(tokenizerLine.hasMoreTokens()){
                        String strScore = tokenizerLine.nextToken();
                        Text name = new Text(strName); // 姓名
                        int scoreInt = Integer.parseInt(strScore); // 该组织的状况得分
                        context.write(name, new IntWritable(scoreInt));
                    } // if
                } // while
            } // map
        } // Mapper
    }

```

由于上述程序比较简单并且和单节点上的很相似, 配合注释就能够很好的理解, 因此就不再多讲解了。

下面是 Reducer 部分的代码:

```

public static class Reduce
    extends Reducer<Text, IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterable<IntWritable> values,
        Context context) throws IOException, InterruptedException {
        int sum = 0;
        int count=0;
        Iterator<IntWritable> iterator = values.iterator();
        while (iterator.hasNext()) {
            sum += iterator.next().get();
            count++;
        } //while
        int average = (int) sum/count;
        context.write(key, new IntWritable(average));
    } //reduce
} //Reducer

```

由于其主函数部分和上一部分的主函数完全一样, 在此处就不列举了。

5.5.1 打包

为了能够在命令行中运行, 首先需要对程序进行编译和打包, 下面就分别展示编译和打

包的过程。

编译代码如下所示：

```
javac -classpath /usr/local/hadoop/hadoop-0.20.2/hadoop-0.20.2-core.jar -d
ScoreProcessFinal_classes ScoreProcessFinal.java
```

编译后会生成以下几个文件，如图 5-5 所示。



图 5-5 编译生成的 class 文件

打包：

```
/usr/local/hadoop/hadoop-0.20.2/bin$ jar -cvf
/usr/local/hadoop/hadoop-0.20.2/bin/ScoreProcessFinal.jar -C ScoreProcessFinal_
classes/ .
```

标明清单 (manifest)

增加: ScoreProcessFinal\$Map.class(读入 = 1899) (写出 = 806) (压缩了 57%)

增加: ScoreProcessFinal\$Reduce.class(读入 = 1671) (写出 = 707) (压缩了 57%)

增加: ScoreProcessFinal.class(读入 = 2374) (写出 = 1183) (压缩了 50%)

打包后的结果如图 5-6 所示。



图 5-6 打包后的 jar 文件

5.5.2 在本地模式下运行

下面显示上面的程序在本地模式下执行的情况：

```
11/01/23 08:27:25 INFO input.FileInputFormat: Total input paths to process : 1000
11/01/23 08:27:27 INFO mapred.JobClient: Running job: job_201101230823_0001
11/01/23 08:27:28 INFO mapred.JobClient: map 0% reduce 0%
11/01/23 08:28:17 INFO mapred.JobClient: map 1% reduce 0%
11/01/23 08:28:47 INFO mapred.JobClient: map 2% reduce 0%
11/01/23 08:29:17 INFO mapred.JobClient: map 3% reduce 0%
11/01/23 08:30:23 INFO mapred.JobClient: map 5% reduce 0%
11/01/23 08:30:29 INFO mapred.JobClient: map 5% reduce 1%
11/01/23 08:30:59 INFO mapred.JobClient: map 6% reduce 1%
... ..
11/01/23 08:46:19 INFO mapred.JobClient: map 36% reduce 11%
11/01/23 08:46:28 INFO mapred.JobClient: map 36% reduce 12%
11/01/23 08:47:52 INFO mapred.JobClient: map 39% reduce 12%
11/01/23 08:48:02 INFO mapred.JobClient: map 39% reduce 13%
11/01/23 08:48:23 INFO mapred.JobClient: map 40% reduce 13%
```

```

11/01/23 08:48:53 INFO mapred.JobClient: map 41% reduce 13%
...
11/01/23 09:18:49 INFO mapred.JobClient: map 99% reduce 32%
11/01/23 09:18:55 INFO mapred.JobClient: map 99% reduce 33%
11/01/23 09:19:19 INFO mapred.JobClient: map 100% reduce 33%
11/01/23 09:19:40 INFO mapred.JobClient: map 100% reduce 71%
11/01/23 09:19:46 INFO mapred.JobClient: map 100% reduce 100%
11/01/23 09:19:48 INFO mapred.JobClient: Job complete: job_201101230823_0001
11/01/23 09:19:48 INFO mapred.JobClient: Counters: 17
11/01/23 09:19:48 INFO mapred.JobClient: Job Counters
11/01/23 09:19:48 INFO mapred.JobClient: Launched reduce tasks=1
11/01/23 09:19:48 INFO mapred.JobClient: Launched map tasks=1000
11/01/23 09:19:48 INFO mapred.JobClient: Data-local map tasks=1000
11/01/23 09:19:48 INFO mapred.JobClient: FileSystemCounters
11/01/23 09:19:48 INFO mapred.JobClient: FILE_BYTES_READ=36935228
11/01/23 09:19:48 INFO mapred.JobClient: HDFS_BYTES_READ=2376120337
11/01/23 09:19:48 INFO mapred.JobClient: FILE_BYTES_WRITTEN=857654430
11/01/23 09:19:48 INFO mapred.JobClient: HDFS_BYTES_WRITTEN=2134960
11/01/23 09:19:48 INFO mapred.JobClient: Map-Reduce Framework
11/01/23 09:19:48 INFO mapred.JobClient: Reduce input groups=213748
11/01/23 09:19:48 INFO mapred.JobClient: Combine output records=220587936
11/01/23 09:19:48 INFO mapred.JobClient: Map input records=1000
11/01/23 09:19:48 INFO mapred.JobClient: Reduce shuffle bytes=819866657
11/01/23 09:19:48 INFO mapred.JobClient: Reduce output records=213748
11/01/23 09:19:48 INFO mapred.JobClient: Spilled Records=227427872
11/01/23 09:19:48 INFO mapred.JobClient: Map output bytes=2394443000
11/01/23 09:19:48 INFO mapred.JobClient: Combine input records=431123764
11/01/23 09:19:48 INFO mapred.JobClient: Map output records=218872000
11/01/23 09:19:48 INFO mapred.JobClient: Reduce input records=8336172

```

上面的命令以 `inputOfScoreProcessFinal` 为输入, 同时以 `outputOfScoreProcessFinal` 为输出。到此, 已经将编译打包和在本地模式下运行的情况讲解完了。

5.5.3 在集群上运行

接下来会讲解程序如何在集群上运行。在笔者的实验环境中, 一共有 4 台机器, 其中一台同时担当 `JobTracker` 和 `NameNode` 的角色, 但不担当 `TaskTracker` 和 `DataNode` 的角色, 另外三台机器则同时担当 `TaskTracker` 和 `DataNode` 的角色。

首先, 将输入的文件复制到 HDFS 中, 用以下命令完成该功能:

```

u@master-chukwa:~/hadoop-0.20.2/bin$ hadoop dfs -copyFromLocal
/home/u/Desktop/inputOfScoreProcessFinal inputOfScoreProcessFinal

```

下面, 在命令行中运行程序:

```

u@master-chukwa:~/hadoop-0.20.2/bin$ hadoop jar /home/u/TG/ScoreProcessFinal.jar
ScoreProcessFinal inputOfScoreProcessFinal outputOfScoreProcessFinal

```

START


```

11/01/23 12:04:10 WARN mapred.JobClient: Use GenericOptionsParser for parsing the
arguments. Applications should implement Tool for the same.
11/01/23 12:04:10 INFO input.FileInputFormat: Total input paths to process : 1000
11/01/23 12:04:12 INFO mapred.JobClient: Running job: job_201101221603_0099
11/01/23 12:04:13 INFO mapred.JobClient: map 0% reduce 0%
11/01/23 12:04:16 INFO mapred.JobClient: Task Id :
attempt_201101221603_0099_m_001001_0, Status : SUCCEEDED

11/01/23 12:04:24 INFO mapred.JobClient: Task Id :
attempt_201101221603_0099_m_000000_0, Status : SUCCEEDED
11/01/23 12:04:24 INFO mapred.JobClient: Task Id :
... ..
11/01/23 12:09:54 INFO mapred.JobClient: map 25% reduce 2%
11/01/23 12:09:54 INFO mapred.JobClient: Task Id :
attempt_201101221603_0099_m_000244_0, Status : SUCCEEDED
11/01/23 12:09:56 INFO mapred.JobClient: Task Id :
attempt_201101221603_0099_m_000262_0, Status : SUCCEEDED

... ..
11/01/23 12:33:50 INFO mapred.JobClient: map 100% reduce 84%
11/01/23 12:33:58 INFO mapred.JobClient: map 100% reduce 85%
11/01/23 12:35:11 INFO mapred.JobClient: map 100% reduce 86%
11/01/23 12:36:23 INFO mapred.JobClient: map 100% reduce 87%
11/01/23 12:37:58 INFO mapred.JobClient: map 100% reduce 88%
11/01/23 12:39:24 INFO mapred.JobClient: Task Id :
attempt_201101221603_0099_r_000005_1, Status : SUCCEEDED
11/01/23 12:39:26 INFO mapred.JobClient: map 100% reduce 100%
... ..

```

上述命令运行 ScoreProcessFinal.jar 中的 ScoreProcessFinal 类，并且将 inputOfScoreProcessFinal 作为输入，outputOfScoreProcessFinal 作为输出。

5.6 网络用户界面

Hadoop 自带的网络用户界面在查看工作的信息时很方便（在 <http://jobtracker-host:50030/> 能找到用户界面）。当 Job 在运行的时候，它对跟踪 Job 工作进程很有用，同样在工作完成后查看工作统计和日志时也会很有用。

5.6.1 JobTracker 页面

下面显示一个 JobTracker 的页面，如图 5-7 所示。

页面的第一部分给出了 Hadoop 安装的详细信息，比如版本号、编译完成时间、JobTracker 当前的运行状态（下图显示的为“运行”状态）和开始时间。

接下来是集群的一个总结信息：集群容量及它的使用情况，以及集群上运行的 map 和 reduce 的数量、提交的工作总量、当前可用的 TaskTracker 节点数、集群的容量（用集群上可用的 map 和 reduce 任务槽的数量表示（“map 任务容量”和“reduce 任务容量”））和每个

节点平均可用槽的数量。

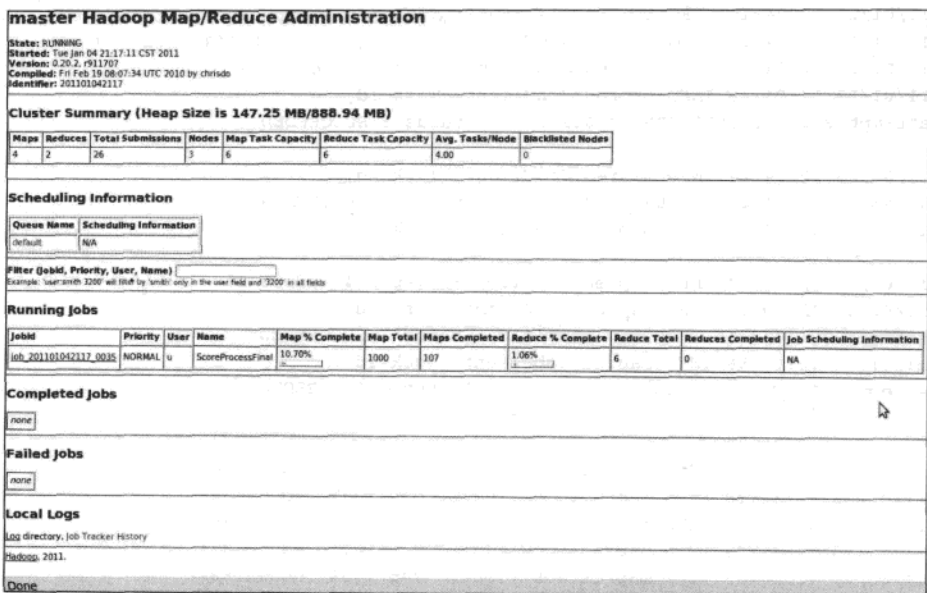


图 5-7 JobTracker 页面

总结的下面是一个正在运行的工作日程表，打开后能看到工作的序列。

再下面显示的是正在运行、完成、失败的工作，而且是通过表来体现的，在表中每一行显示的是一个工作并且显示了工作的 ID 号、所属者、名字和进程信息。

最后，在页面的最下面是 JobTracker 日志的链接和 JobTracker 的历史信息（JobTracker 运行的所有工作信息）。在将它们提交到历史页面之前，主要显示 100 个工作（可以通过 `mapred.job.name` 进行配置）。注意，历史记录是永久保存的，因此可以从 JobTracker 以前运行的工作中找到相关的记录。

5.6.2 工作页面

点击一个工作 ID 将看到一个工作页面，由于页面太大，下面将分别截图显示，首先显示的是页面最前面的内容（如图 5-8 所示）。

图 5-8 中，在页面的顶部是一个关于工作的一些总结性基本信息，比如工作所属者、名字、工作文件、工作已经执行了多长时间等。工作文件是工作的加强配置文件，包含在工作运行期间的所有有效属性和它们的取值。如果不确信某个属性的取值，可以点击进一步查看文件。

当工作运行时，可以在页面上监控它的进展情况，因为页面会周期性的更新。在总结信息的下面是一张表，它显示了 map 和 reduce 的进展情况。“任务栏”显示了该工作的 map 和

reduce 任务的总数 (map 和 reduce 各占一行)。其他列显示了这些任务的状态：“暂停”（等待执行）、“正在执行”、“完成”（运行成功）、“终止”（准确地说应该叫做“失败”），最后一列显示了失败或终止的任务所尝试的总数。从图中可以看出来一共有 1000 个 map 任务，因为一共有 1000 个输入文件，每个 map 任务会处理一个文件。同时，任何时候处于“正在执行”状态的 map 任务最多有 6 个，因为默认的设置是每个 TaskTracker 节点只有 2 个 map 任务，而我的实验环境中一共有 3 个 TaskTracker 节点，所以正常情况下处于“正在运行”状态的 map 任务数应该是 6，如果某一时刻处于“正在执行”状态的 map 任务数少于 6，则说明存在某些 map 任务尝试失败的情况。

Hadoop job_201101042117_0035 on master							
User: u							
Job Name: ScoreProcessFinal							
Job File: hdfs://master:9000/home/u/tmp/mapred/system/job_201101042117_0035/job.xml							
Job Setup: Successful							
Status: Running							
Started at: Fri Jan 14 20:01:20 CST 2011							
Running for: 17mins, 48sec							
Job Cleanup: Pending							
Kind	% Complete	Num Tasks	Pending	Running	Complete	Killed	Failed/Killed Task Attempts
map	80.09%	1000	193	6	801	0	0 / 0
reduce	8.82%	6	4	2	0	0	0 / 0

图 5-8 工作页面

图 5-9 显示的是工作页面最下面的内容。

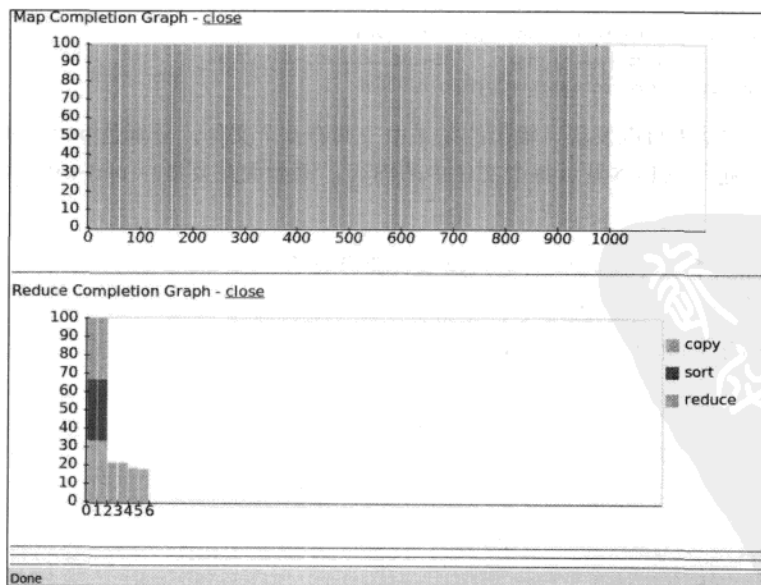


图 5-9 每个 map 和 reduce 任务的执行进度

图 5-9 是每个任务完成情况的一个图形化表示。Reduce 完成图分为 3 个阶段：复制（发生在 map 输出转交给 reduce 的 TaskTracker 时）、排序（发生在 reduce 输入合并的时候）和 reduce（发生在 reduce 函数起作用并产生最终输出的时候）。

5.6.3 返回结果

执行完毕后，可以通过以下几种方式得到结果。

(1) 用命令行直接显示

命令行如下所示：

```
u@master-chukwa:~/hadoop-0.20.2/bin$ hadoop dfs -ls outputOfScoreProcessFinal
```

显示的结果如下所示：

```
Found 7 items
drwxr-xr-x  - u supergroup      0 2011-01-23 12:04
/user/u/outputOfScoreProcessFinal/_logs
-rw-r--r--  2 u supergroup      352890 2011-01-23 12:24
/user/u/outputOfScoreProcessFinal/part-r-00000
-rw-r--r--  2 u supergroup      358387 2011-01-23 12:24
/user/u/outputOfScoreProcessFinal/part-r-00001
-rw-r--r--  2 u supergroup      355221 2011-01-23 12:30
/user/u/outputOfScoreProcessFinal/part-r-00002
-rw-r--r--  2 u supergroup      354239 2011-01-23 12:31
/user/u/outputOfScoreProcessFinal/part-r-00003
-rw-r--r--  2 u supergroup      359453 2011-01-23 12:31
/user/u/outputOfScoreProcessFinal/part-r-00004
-rw-r--r--  2 u supergroup      354770 2011-01-23 12:38
/user/u/outputOfScoreProcessFinal/part-r-00005
```

通过以上的结果可以发现，输出的结果中一共有 6 个文件，分别是 part-r-00000 到 part-r-00005。并且还可以具体显示每个文件中的内容，例如现在要显示 part-r-00000 中的内容，命令如下所示。

```
u@master-chukwa:~/hadoop-0.20.2/bin$ hadoop dfs -cat
outputOfScoreProcessFinal/part-r-00000
```

显示结果的一部分如下所示：

黎十	51
黎升	49
黎卫	48
黎危	49
黎历	48
黎庆	48
黎友	50
黎发	50



(2) 将输出的文件从 HDFS 复制到本地文件系统上, 在本地文件系统上查看命令如下所示:

```
u@master-chukwa:~/hadoop-0.20.2/bin$ hadoop dfs -get outputOfScoreProcessFinal/*  
/home/u/outputOfScoreProcessFinal
```

上述命令的主要功能是将 HDFS 中目录 outputOfScoreProcessFinal 下的所有文件复制到本地文件系统中的目录 /home/u/outputOfScoreProcessFinal 下, 然后就可以方便地查看了。图 5-10 就是在本地文件系统中目录 /home/u/outputOfScoreProcessFinal 下的文件。



图 5-10 复制到本地文件系统中输出文件

可以看到 HDFS 中的输出结果被复制到了本地文件系统中。

另外还可以在命令行中将输出文件 part-r-00000 到 part-r-00005 合并成一个文件, 并复制到本地文件系统中, 下面就是在命令行中进行的操作。

```
u@master-chukwa:~/hadoop-0.20.2/bin$ hadoop dfs -getmerge outputOfScoreProcessFinal  
/home/u/outputScore
```

上述命令的功能就是, 将 HDFS 中 outputOfScoreProcessFinal 目录下的所有文件 (即 part-r-00000 到 part-r-00005) 进行合并, 然后复制到本地文件系统中的 /home/u/outputScore 目录下。

在本地文件系统中可以查看合并后的结果, 图 5-11 就显示了合并后的结果。

(3) 通过 Web 界面查看输出的结果

在浏览器中输入 <http://localhost:50070> 会显示结果, 其中的一部分如图 5-12 所示。

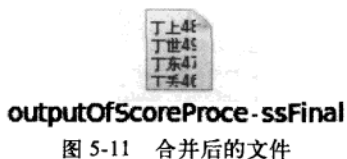


图 5-11 合并后的文件

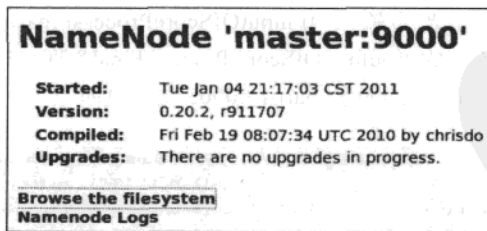


图 5-12 NameNode

点击上图中的 “Browse the filesystem” 即可看到 HDFS 中的内容, 图 5-13 显示的是点击后所显示的 HDFS 中的信息。

Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
chukwa	dir				2011-01-04 00:16	rw-r--r--	u	supergroup
home	dir				2011-01-03 20:33	rw-r--r--	u	supergroup
output07	dir				2011-01-12 20:56	rw-r--r--	u	supergroup
sort18	dir				2011-01-11 21:34	rw-r--r--	u	supergroup
sort19	dir				2011-01-11 21:35	rw-r--r--	u	supergroup
sort20	dir				2011-01-12 19:45	rw-r--r--	u	supergroup
sort21	dir				2011-01-12 19:47	rw-r--r--	u	supergroup
sort22	dir				2011-01-12 19:49	rw-r--r--	u	supergroup
sort23	dir				2011-01-12 19:51	rw-r--r--	u	supergroup
sort24	dir				2011-01-12 20:04	rw-r--r--	u	supergroup
sort25	dir				2011-01-12 20:05	rw-r--r--	u	supergroup
sort26	dir				2011-01-12 20:29	rw-r--r--	u	supergroup
sort27	dir				2011-01-12 20:31	rw-r--r--	u	supergroup

图 5-13 HDFS 中信息

然后再点击图 5-13 中的 **home**，会显示图 5-14 中的内容。

Go to parent directory								
Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
u	dir				2011-01-14 20:01	rw-r--r--	u	supergroup

图 5-14 home 下面的信息

再点击图 5-14 中的 **u**，会显示出最近 HDFS 中存储的内容（如图 5-15 所示）。

Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
MySequenceOutput	file	9.74 MB	3	64 MB	2011-01-14 16:36	rw-r--r--	u	supergroup
MySequenceOutput-uncom	file	131.69 MB	3	64 MB	2011-01-14 16:41	rw-r--r--	u	supergroup
inputOfScoreProcessFinal	dir				2011-01-14 19:53	rw-r--r--	u	supergroup
outputOfScoreProcessFinal	dir				2011-01-14 20:34	rw-r--r--	u	supergroup

图 5-15 u 下面的文件

从上图可以看出，笔者的输入为 inputOfScoreProcessFinal，输出为 outputOfScoreProcessFinal，点击图 5-15 中的“output OfScore Process Final”即可查看输出结果（如图 5-16 所示），它们分别存储在 part-r-00000 到 part-r-00005 文件里面。

Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
logs	dir				2011-01-14 20:01	rw-r--r--	u	supergroup
part-r-00000	file	344.62 KB	2	64 MB	2011-01-14 20:23	rw-r--r--	u	supergroup
part-r-00001	file	349.99 KB	2	64 MB	2011-01-14 20:23	rw-r--r--	u	supergroup
part-r-00002	file	346.9 KB	2	64 MB	2011-01-14 20:34	rw-r--r--	u	supergroup
part-r-00003	file	345.94 KB	2	64 MB	2011-01-14 20:34	rw-r--r--	u	supergroup
part-r-00004	file	351.03 KB	2	64 MB	2011-01-14 20:31	rw-r--r--	u	supergroup
part-r-00005	file	346.46 KB	2	64 MB	2011-01-14 20:31	rw-r--r--	u	supergroup

图 5-16 所要查看的输出文件

如果需要查看每个文件中的详细内容，点击相应的按钮即可查看。

到此，就已经将查看输出结果的三种方式介绍完了。

5.6.4 任务页面

任务页面中有一些链接可以查看该工作里任务的详细信息。例如，点击“map”链接，将看到一个页面，所有的 map 任务信息都列在这一页上。当然，也可以只看已经完成的任务。图 5-17 显示的是带有调试信息的任务页面中的一小部分。表中的每一行都表示一个任务，并且提供了诸如开始时间、结束时间之类的信息，以及由 TaskTracker 提供的错误信息和查看单个任务的计数器的链接。

Hadoop map task list for job_201101042117_0035 on master						
Completed Tasks						
Task	Complete	Status	Start Time	Finish Time	Errors	Counters
task_201101042117_0035_m_000000	100.00%		14-Jan-2011 20:01:26	14-Jan-2011 20:01:35 (9sec)		8
task_201101042117_0035_m_000001	100.00%		14-Jan-2011 20:01:27	14-Jan-2011 20:01:33 (6sec)		8
task_201101042117_0035_m_000002	100.00%		14-Jan-2011 20:01:27	14-Jan-2011 20:01:33 (6sec)		8
task_201101042117_0035_m_000003	100.00%		14-Jan-2011 20:01:29	14-Jan-2011 20:01:38 (9sec)		8
task_201101042117_0035_m_000004	100.00%		14-Jan-2011 20:01:30	14-Jan-2011 20:01:36 (6sec)		8
task_201101042117_0035_m_000005	100.00%		14-Jan-2011 20:01:30	14-Jan-2011 20:01:36 (6sec)		8
task_201101042117_0035_m_000006	100.00%		14-Jan-2011 20:01:33	14-Jan-2011 20:01:39 (6sec)		8
task_201101042117_0035_m_000007	100.00%		14-Jan-2011 20:01:33	14-Jan-2011 20:01:39 (6sec)		8
task_201101042117_0035_m_000008	100.00%		14-Jan-2011 20:01:35	14-Jan-2011 20:01:41 (6sec)		8

图 5-17 map 任务列表

同样，点击“reduce”链接也可以看到一个页面，所有的 reduce 任务信息都列在这一页上。同样可以只看已经完成的任务。图 5-18 显示的是带有调试信息的任务页面的一小部分。和图 5-17 一样，表中的每一行都表示一个任务，它提供了诸如开始时间、结束时间之类的信息，以及由 TaskTracker 提供的错误信息和查看单个任务的计数器的链接。

5.6.5 任务细节页面

从任务页面上可以点击任何任务得到关于它的详细信息。图 5-19 的任务细节页面显示了每个任务的尝试情况。在这里，只有一个任务尝试并且成功完成。这张表还提供了更多的有用数据，比如任务尝试是在哪个节点上运行的，同时还可以查看任务日志文件和计数器的链接。

“Actions”列可终止一个任务尝试的链接。默认情况下，这项功能是没有启用的，网络用户界面只是一个只读接口。将 webinterface.private.actions 设为 true 即可启用这项功能。

Hadoop reduce task list for job_201101042117_0035 on master						
All Tasks						
Task	Complete	Status	Start Time	Finish Time	Errors	Counters
task_201101042117_0035_r_000000	100.00%	reduce > reduce	14-Jan-2011 20:02:34	14-Jan-2011 20:24:07 (21mins, 33sec)		10
task_201101042117_0035_r_000001	100.00%	reduce > reduce	14-Jan-2011 20:02:38	14-Jan-2011 20:24:27 (21mins, 49sec)		10
task_201101042117_0035_r_000002	100.00%	reduce > reduce	14-Jan-2011 20:22:20	14-Jan-2011 20:34:46 (12mins, 26sec)		10
task_201101042117_0035_r_000003	100.00%	reduce > reduce	14-Jan-2011 20:22:27	14-Jan-2011 20:34:29 (12mins, 1sec)		10
task_201101042117_0035_r_000004	100.00%	reduce > reduce	14-Jan-2011 20:23:34	14-Jan-2011 20:31:39 (8mins, 4sec)		10
task_201101042117_0035_r_000005	100.00%	reduce > reduce	14-Jan-2011 20:23:36	14-Jan-2011 20:32:12 (8mins, 35sec)		10
Go back to JobTracker						
Hadoop, 2011.						

图 5-18 reduce 任务列表

对于 map 任务，有一个部分（即图 5-19 中的“Input Split Locations”的那一部分内容）显示了输入的片段被分配到了哪个节点上。

Job Job_201101042117_0035									
All Task Attempts									
Task Attempts	Machine	Status	Progress	Start Time	Finish Time	Errors	Task Logs	Counters	Actions
attempt_201101042117_0035_m_000986_0	/default-rack/slave3	SUCCEEDED	100.00%	14-Jan-2011 20:24:10	14-Jan-2011 20:24:30 (19sec)		Last 4KB Last 8KB All	0	
attempt_201101042117_0035_m_000986_1	/default-rack/slave2	KILLED	100.00%	14-Jan-2011 20:26:20	14-Jan-2011 20:26:24 (4sec)		Last 4KB Last 8KB All	0	
Input Split Locations									
<input type="text" value="/default-rack/slave2"/> <input type="text" value="/default-rack/slave3"/>									
Go back to the job Go back to JobTracker									

图 5-19 任务尝试页面

5.7 性能调优

一个程序可以完成基本功能其实还不够，还有一些具有实际意义的问题需要我们考虑，那就是性能是不是足够好，性能方面有没有提高的空间等。具体来讲包括两个方面的内容：一个是时间性能；另一个是空间性能。衡量性能的指标就是，能够在正确完成功能的基础上，使执行的时间尽量短，占用的空间尽量小。

前面只是实现了程序应该实现的功能，对性能问题还没有加以考虑，下面就从不同的角度来简单地介绍提高性能的方法。

(1) 输入的文件尽量采用大文件，避免使用小文件

在前面的例子中，笔者的实验数据包含 1000 个文件，在 HDFS 中一共占用了 1000 个文件块，而每一个文件的大小都是 2.3MB，相对于 HDFS 默认块的大小 64MB 来说算是比较小的了。如果 MapReduce 在处理数据的 map 阶段，输入的文件较小而数量众多的话，就会产生很多的 map 任务，以前面的输入为例，一共产生了 1000 个 map 任务，每次新的 map 任务操作都会造成一定的性能损失。针对上述 2.2GB 大小的数据，在实验环境中运行的时间大概为 33 分钟。

为了尽量使用大文件的数据，笔者对这 1000 个文件进行了一次预处理，也就是将这些数量众多的小文件合并成了大一些的文件，最终将它们合并成了一个大小为 2.2GB 的大文件。然后再以这个大文件作为输入，在同样的环境中进行测试，结果运行的时间大概为 4 分钟。图 5-20 和图 5-21 是合并文件后和没有合并文件时的运行效果图。

```
Hadoop job_201101192023_0448 on master
User: u
Job Name: ScoreProcessFinal
Job File: hdfs://master:9000/home/u/tmp/mapred/system/job_201101192023_0448/job.xml
Job Setup: Successful
Status: Succeeded
Started at: Fri Jan 21 13:04:27 CST 2011
Finished at: Fri Jan 21 13:08:25 CST 2011
Finished in: 3mins, 57sec
Job Cleanup: Successful
```

图 5-20 合并文件后的运行效果图

```
Hadoop job_201101042117_0035 on master
User: u
Job Name: ScoreProcessFinal
Job File: hdfs://master:9000/home/u/tmp/mapred/system/job_201101042117_0035/job.xml
Job Setup: Successful
Status: Succeeded
Started at: Fri Jan 14 20:01:20 CST 2011
Finished at: Fri Jan 14 20:34:53 CST 2011
Finished in: 33mins, 32sec
Job Cleanup: Successful
```

图 5-21 没有合并文件时的运行效果图

上面两幅图可以很明显地看出二者执行时间上的差别非常大。所以，为了提高性能，应该将小文件做一些合理的预处理，变小为大，从而缩短执行的时间。不仅如此，合并前的众多文件在 HDFS 中占用了 1000 个块，而合并后的文件在 HDFS 中只占用 36 个块（64MB 为一块），也就是说占用空间也相应地变小了，可谓一举两得。

另外，如果不对小文件做合并的预处理，也可以借用 Hadoop 中的 CombineFileInputFormat。它可以将多个文件打包到一个输入单元中，从而每次的 map 操作就会处理更多的数

据。同时，CombineFileInputFormat 会考虑节点和集群的位置信息，以决定哪些文件被打包到一个单元之中，所以使用 CombineFileInputFormat 也会使性能得到相应的提高。

(2) 考虑压缩文件

在分布式系统中，不同节点的数据交换是影响整体性能的一个重要因素。另外在 Hadoop 的 map 阶段所处理的输出大小也会影响整个 MapReduce 程序的执行时间，这是因为 map 阶段的输出首先是存储在一定大小的内存缓冲区中的，如果 map 输出的大小超出一定限度，它就会将结果写入磁盘，等 map 任务结束后再将它们复制到 reduce 任务的节点上，因此如果数据量大，中间的数据交换会占用很多的时间。

一个好的改善性能的方法就是对 map 的输出进行压缩，这样的话会带来两个方面的好处：减少存储文件的空间；加快数据在网络上（不同节点间）的传输速度，以及减少数据在内存和磁盘间交换的时间。可以通过将 `mapred.compress.map.output` 属性设置为 `true` 来对 map 的输出数据进行压缩，同时还可以设置 map 输出数据的压缩格式，通过设置 `mapred.map.output.compression.codec` 属性即可进行压缩格式的设置。

(3) 修改配置文件中的相关属性

属性 `mapred.tasktracker.map.tasks.maximum` 的默认值是 2，属性 `mapred.tasktracker.reduce.tasks.maximum` 的默认值也是 2，因此每个节点上实际处于运行状态的 map 和 reduce 任务数最多为 2，而较为理想的数值应在 10 到 100 之间。因此，可以在 `conf` 目录下修改属性 `mapred.tasktracker.map.tasks.maximum` 和 `mapred.tasktracker.reduce.tasks.maximum` 的取值，将它们设置为一个较大的值，使得每个节点上同时运行的 map 和 reduce 任务数增加，从而缩短运行的时间，提高整体的性能。

例如，下面的修改：

```
<property>
  <name>mapred.tasktracker.map.tasks.maximum</name>
  <value>10</value>
  <description>The maximum number of map tasks that will be run
    simultaneously by a task tracker.
  </description>
</property>

<property>
  <name>mapred.tasktracker.reduce.tasks.maximum</name>
  <value>10</value>
  <description>The maximum number of reduce tasks that will be run
    simultaneously by a task tracker.
  </description>
</property>
```

5.8 MapReduce 工作流

到目前为止，已经讲述了使用 MapReduce 编写程序的机制。不过还没有讨论如何将数

据处理问题转化为 MapReduce 模型。

数据处理问题只不过是解决一些非常简单的问题而已。如果处理过程变得更加复杂，这种复杂性会体现为更多的 MapReduce 工作，而不是更加复杂的 map 函数和 reduce 函数。

对于这些更加复杂的问题，有必要考虑比 MapReduce 更加高级的语言，比如 Pig、Hive 或 Cascading。一个直接好处就是，不必专注于如何将问题转化成 MapReduce 工作，而是更加专注于所分析的问题上。

5.8.1 将问题分解成 MapReduce 工作

假如现在将人体的每个组织（医学词汇）用一个文件存储，作为所有人的健康指标，其中有一个指标包括两项内容，最好状况和最差状况，存储的格式如下所示：

```
赵一 46 94 #赵峰 67 89#
```

而在考察一个人该组织的健康状况时，便用最好状况和最差状况的平均得分作为考察的指标。例如赵一该组织的健康得分为 $(46+94)/2=70$ （分）。

在其他内容和要求上与之前的举例一样，下面要给出每个人所有组织的综合健康状况。如果用 MapReduce 来处理，很明显需要分为两个阶段：

- 计算每个人每个组织的健康得分，即求两项数据的平均数。
- 以上一个阶段的输出结果作为该阶段的 Mapper 输入，之后的处理和前面所讲的处理健康状况的程序一样，最终 Reducer 得到最后的结果。

Mapper 通常用来处理输入格式转化、投影（选择相关的字段）、过滤（去掉那些不感兴趣的记录）等。有时候，应用程序会很复杂，这时候可能会有很多的 Mapper，为了完成一个作业则需要让这些 Mapper 合作，为此可以利用 Hadoop 自带的 ChainMapper 类库来将它们链成一个 Mapper。使用 ChainReducer，可以运行一系列的 Mapper，再运行一个 Reducer 和该 MapReduce 工作中的其他 Mapper 链。

5.8.2 运行相互依赖的工作

当一个 MapReduce 工作流中不止一个工作时，问题就出现了：怎样管理这些工作，保证它们能够有序的运行呢？有多种方法可以解决这个问题，最主要是要考虑这些工作是否构成了一个线性序列，或者能否形成一个有向无环图。

对于线性链，最简单的方法就是执行完一个工作再执行另外一个工作，在一个工作执行完之前另外的工作必须一直等待，代码如下所示。

```
JobClient.runJob(conf1);
JobClient.runJob(conf2);
```

如果工作运行失败，runJob() 方法会抛出一个 IOException，管道中后面的工作就无法执行了。在具体的应用程序中，可能需要抓住异常并将前面的工作所产生的中间数据全部清除。

对于比线性链更加复杂的问题，会有相关的类库帮助合理地安排 workflow。最简单的是 `org.apache.hadoop.mapred.jobcontrol` 包中的 `JobControl` 类。一个 `JobControl` 实例表示一个将要执行的工作图。在添加各个工作的工作配置后，告之 `JobControl` 实例各个工作之间的依赖关系。在一个线程中运行 `JobControl`，然后它再按照依赖关系运行工作。可以查看进程，并且在工作完成后查看工作的状态和每个与失败相关的错误信息。如果有一个工作运行失败，与之有依赖关系的其他工作就不会运行了。

不同于在客户端运行并提交作业的 `JobControl`，Hadoop 工作流调度器 (HWS) 作为一个服务器，允许客户端提交一个工作流给调度器。当工作流完成后调度器发起一个 HTTP 回复给客户端，以通知工作的状态。在同个工作流当中 HWS 可以运行不同类型的工作：例如可以先运行一个 Pig 作业，然后再运行一个 Java MapReduce 作业。

5.9 小结

在本章中，主要介绍了开发 MapReduce 程序的一般框架。


在单节点上完成 `map` 函数和 `reduce` 函数，并且对它们进行测试。待 `map` 和 `reduce` 都能够成功运行后，再在单节点的大数据集进行测试。在进行程序的编写和编译时，最好在集成环境下进行，因为这样便于程序的修改和调试，建议在 Eclipse 下进行编程。

程序可以在集成环境中运行，也可以在命令行中编译打包，然后在命令中执行，最终的结果也有 3 种不同的查看方式：命令行中直接查看；拷贝到本地文件系统中查看；在 Web 用户界面上查看。

对于已经能够完成功能性要求的 MapReduce 程序，还可以从多个方面进行性能上的优化。比如从几个常见的方面入手：变小文件为大文件，减少 `map` 的数量；压缩最终的输出数据或 `map` 的中间输出结果；在 Hadoop 安装路径下的 `conf` 目录下修改属性，使能够同时运行的 `map` 和 `reduce` 任务数增多，从而提高性能。

对于一个大的问题，则可能需要考虑将它们进行分解。这时可以编写多个 Mapper 和 Reducer，然后将它们串起来。





第6章

MapReduce 应用案例

本章内容

- ☐ 单词计数
- ☐ 数据去重
- ☐ 排序
- ☐ 单表关联
- ☐ 多表关联
- ☐ 小结

资源分享

PDG

前面已经介绍了很多关于 MapReduce 的基础知识，比如 Hadoop 集群的配置方法，以及如何开发 MapReduce 应用程序等。本章将从同本书配套的云计算在线监测平台 (<http://cloudcomputing.ruc.edu.cn/>) 上的 MapReduce 编程题目出发，向读者介绍如何挖掘实际问题的并行处理可能性，以及如何设计编写 MapReduce 程序。需要说明的是本章所有给出的代码均在伪分布集群默认的设置下运行通过了，其 Hadoop 版本为 0.20.2，JDK 的版本是 1.6。本章旨在帮助刚接触 MapReduce 的读者入门。

6.1 单词计数

进入云计算在线监测平台后的第一个编程题目是 WordCount，也就是文本中的单词计数。这个例子如同 Java 中的“HelloWorld”经典程序一样是 MapReduce 的入门程序。虽然此例在本书中的其他章节也有涉及，但本章主要从如何挖掘此问题中的并行处理可能性角度出发，让读者了解设计 MapReduce 程序的过程。

6.1.1 实例描述

计算出文件中各个单词的频数。要求输出结果按照单词的字母顺序进行排序。每个单词和其频数占一行，单词和频数之间有间隔。

比如，输入一个文件，其内容如下：

```
hello world
hello hadoop
hello mapreduce
```

对应上面给出的输入样例，其输出样例为：

```
hadoop      1
hello       3
mapreduce   1
world       1
```

6.1.2 设计思路

这个应用实例的解决方案很直接，就是将文件内容切分成单词，然后将所有相同的单词聚集在一块儿，最后计算单词出现的次数并输出。针对 MapReduce 并行程序设计原则可知，解决方案中的内容切分步骤和数据不相关，可以并行化处理，每个拿到原始数据的机器只要将输入数据切分成单词就可以了。所以可以在 map 阶段完成单词切分任务。另外，相同单词的频数计算也可以并行化处理。根据实例要求来看，不同单词之间的频数不相关，所以可以将相同的单词交给一台机器来计算频数，然后输出最终结果。这个过程可以交给 reduce 阶段完成。至于将中间结果根据不同单词分组再分发给 reduce 机器，这正好是 MapReduce 过程中的 shuffle 能够完成的。至此，这个实例的 MapReduce 程序就设计出来了。map 阶段完成由输入数据到单词切分的工作，shuffle 阶段完成相同单词的聚集和分发工作（这个过程是

MapReduce 的默认过程，不用具体配置），reduce 阶段完成接收所有单词并计算其频数的工作。由于 MapReduce 中传递的数据都是 <key,value> 形式的，并且 shuffle 排序聚集分发都是按照 key 值进行的，所以将 map 的输出设计成由 word 作为 key，1 作为 value 的形式，它表示单词 word 出现了一次（map 的输入采用 Hadoop 默认的输入方式：文件一行作为 value，行号作为 key）。reduce 的输入为 map 输出聚集后的结果，即 <key,value-list>，具体到这个实例就是 <word,{1,1,1,1...}>，reduce 的输出会设计成与 map 输出相同的形式，只是后面的数字不再固定是 1，而是具体算出的 word 所对应的频数。下面给出官网中的 WordCount 代码。

6.1.3 程序代码

官网中的 WordCount 代码如下：

```
import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;

public class WordCount {
    // 继承 Mapper 接口，设置 map 的输入类型为 <Object,Text>
    // 输出类型为 <Text, IntWritable>
    public static class TokenizerMapper
        extends Mapper<Object, Text, Text, IntWritable>{
        //one 表示单词出现一次
        private final static IntWritable one = new IntWritable(1);
        //word 存储切下的单词
        private Text word = new Text();

        public void map(Object key, Text value, Context context ) throws IOException,
            InterruptedException {
            StringTokenizer itr = new StringTokenizer(value.toString()); // 对输入的行切词
            while (itr.hasMoreTokens()) {
                word.set(itr.nextToken()); // 切下的单词存入 word
                context.write(word, one);
            }
        }
    }

    // 继承 Reducer 接口，设置 reduce 的输入类型为 <Text,IntWritable>
    // 输出类型为 <Text,IntWritable>
}
```

```

public static class IntSumReducer extends Reducer<Text,IntWritable,Text,IntWr
itable> {
//result 记录单词的频数
private IntWritable result = new IntWritable();

public void reduce(Text key, Iterable<IntWritable> values, Context context )
throws IOException, InterruptedException {
    int sum = 0;
    // 对获取的 <key,value-list> 计算 value 的和
    for (IntWritable val : values) {
        sum += val.get();
    }
    // 将频数设置到 result 中
    result.set(sum);
    // 收集结果
    context.write(key, result);
}

}

public static void main(String[] args) throws Exception {
Configuration conf = new Configuration();
// 检查运行命令
String[] otherArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
if (otherArgs.length != 2) {
    System.err.println("Usage: wordcount <in> <out>");
    System.exit(2);
}
// 配置作业名
Job job = new Job(conf, "word count");
// 配置作业各个类
job.setJarByClass(WordCount.class);
job.setMapperClass(TokenizerMapper.class);
job.setCombinerClass(IntSumReducer.class);
job.setReducerClass(IntSumReducer.class);
job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);
FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}

```

6.1.4 代码解读

WordCount 程序在 map 阶段接收输入的 <key,value> (key 是当前输入的行号, value 是对应行的内容), 然后对此行内容进行切词, 每切一个词下来就将其组织成 <word,1> 的形式输出, 表示 word 出现了一次。

在 reduce 阶段, TaskTracker 会接收到 <word,{1,1,1,1...}> 形式的数据, 也就是特定单词及其出现次数的情况, 其中“1”表示 word 的频数。所以 reduce 每接收一个

<word,{1,1,1,1...}>, 就会在 word 的频数上加 1, 最后组织成 <word,sum> 直接输出。

6.1.5 程序执行

运行条件：将 WordCount.java 文件放在 Hadoop 安装目录下，并在目录下创建输入目录 input，目录下有输入文件 file1、file2，其中：

file1 的内容是：

```
hello world
```

file2 的内容是：

```
hello hadoop
hello mapreduce
```

准备好之后在命令行输入命令运行。

运行命令如下所示：

1) 在集群上创建输入文件夹：

```
bin/hadoop fs -mkdir input
```

2) 上传本地目录 input 下前四个名为 file 的文件到集群上的 input 目录下：

```
bin/hadoop fs -put input/file* input
```

3) 编译 WordCount.java 程序，将结果放入当前目录的 WordCount 目录下：

```
javac -classpath hadoop-0.20.2-core.jar:
lib/commons-cli-1.2.jar -d WordCount WordCount.java
```

4) 将编译结果打成 jar 包：

```
jar -cvf wordcount.jar -C WordCount.
```

5) 在集群上运行 WordCount 程序，以 input 目录作为输入目录，output 目录作为输出目录：

```
bin/hadoop jar wordcount.jar WordCount input output
```

6) 查看输出结果：

```
bin/hadoop fs -cat output/part-r-00000
```

6.1.6 代码结果

运行结果如下所示：

```
hadoop      1
hello       3
mapreduce   1
world       1
```



6.2 数据去重

数据去重这个实例主要是为了让读者掌握并利用并行化思想来对数据进行有意义的筛选。统计大数据集上的数据种类个数、从网站日志中计算访问地等这些看似庞杂的任务都会涉及数据去重。下面就进入这个实例的 MapReduce 程序设计。

6.2.1 实例描述

对数据文件中的数据进行去重。数据文件中的每行都是一个数据。

样例输入如下所示。

file1 :

```
2006-6-9 a
2006-6-10 b
2006-6-11 c
2006-6-12 d
2006-6-13 a
2006-6-14 b
2006-6-15 c
2006-6-11 c
```

file2 :

```
2006-6-9 b
2006-6-10 a
2006-6-11 b
2006-6-12 d
2006-6-13 a
2006-6-14 c
2006-6-15 d
2006-6-11 c
```

样例输出如下所示：

```
2006-6-9 a
2006-6-9 b
2006-6-10 a
2006-6-10 b
2006-6-11 b
2006-6-11 c
2006-6-12 d
2006-6-13 a
2006-6-14 b
2006-6-14 c
2006-6-15 c
2006-6-15 d
```



6.2.2 设计思路

数据去重实例的最终目标是让原始数据中出现次数超过一次的数据在输出文件中只出现一次。我们自然而然会想到将同一个数据的所有记录都交给一台 reduce 机器，无论这个数据出现多少次，只要在最终结果中输出一次就可以了。具体就是 reduce 的输入应该以数据作为 key，而对 value-list 则没有要求。当 reduce 接收到一个 <key,value-list> 时就直接将 key 复制到输出的 key 中，并将 value 设置成空值。在 MapReduce 流程中，map 的输出 <key,value> 经过 shuffle 过程聚集成 <key,value-list> 后会交给 reduce。所以从设计好的 reduce 输入可以反推出 map 的输出 key 应为数据，value 任意。继续反推，map 输出数据的 key 为数据，而在这个实例中每个数据代表输入文件中的一行内容，所以 map 阶段要完成的任务就是在采用 Hadoop 默认的作业输入方式之后，将 value 设置成 key，并直接输出（输出中的 value 任意）。map 中的结果经过 shuffle 过程之后交给 reduce。reduce 阶段不会管每个 key 有多少个 value，它直接将输入的 key 复制为输出的 key，并输出就可以了（输出中的 value 被设置成空了）。

由于此程序简单且执行步骤与单词计数实例完全相同，所以不再赘述，下面只给出程序。

6.2.3 程序代码

程序代码如下所示：

```
import java.io.IOException;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;

public class Dedup {
    //map 将输入中的 value 复制到输出数据的 key 上，并直接输出
    public static class Map extends Mapper<Object, Text, Text, Text>{
        private static Text line = new Text();
        public void map(Object key, Text value, Context context) throws IOException,
            InterruptedException {
            line = value;
            context.write(line, new Text(""));
        }
    }
    //reduce 将输入中的 key 复制到输出数据的 key 上，并直接输出
    public static class Reduce extends Reducer<Text,Text,Text,Text> {
```

```

    public void reduce(Text key, Iterable<Text> values, Context context ) throws
        IOException, InterruptedException {
        context.write(key, new Text(""));
    }
}

public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    String[] otherArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
    if (otherArgs.length != 2) {
        System.err.println("Usage: wordcount <in> <out>");
        System.exit(2);
    }
    Job job = new Job(conf, "Data Deduplication");
    job.setJarByClass(Dedup.class);
    job.setMapperClass(Map.class);
    job.setCombinerClass(Reduce.class);
    job.setReducerClass(Reduce.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(Text.class);
    FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
    FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}

```

6.3 排序

数据排序是许多实际任务执行时要完成的第一项工作，比如学生成绩评比、数据建立索引等。这个实例和数据去重类似，都是先对原始数据进行初步处理，为进一步的数据操作打好基础。下面进入这个实例。

6.3.1 实例描述

对输入文件中的数据进行排序。输入文件中的每行内容均为一个数字，即一个数据。要求在输出中每行有两个间隔的数字，其中，第二个代表原始数据，第一个代表这个原始数据在原始数据集中的位次。

样例输入：

file1 :

```

2
32
654
32
15
756
65223

```

file2 :

```
5956
22
650
92
```

file3 :

```
26
54
6
```

样例输出:

```
1 2
2 6
3 15
4 22
5 26
6 32
7 32
8 54
9 92
10 650
11 654
12 756
13 5956
14 65223
```

6.3.2 设计思路

这个实例仅仅要求对输入数据进行排序,熟悉 MapReduce 过程的读者会很快想到在 MapReduce 过程中就有排序,是否可以利用这个默认的排序,而不需要自己再实现具体的排序呢?答案是肯定的。但是在使用之前首先需要了解它的默认排序规则。它是按照 key 值进行排序的,如果 key 为封装 int 的 IntWritable 类型,那么 MapReduce 按照数字大小对 key 排序,如果 key 为封装 String 的 Text 类型,那么 MapReduce 按照字典顺序对字符串排序。了解了这个细节,我们就知道应该使用封装 int 的 IntWritable 型数据结构了。也就是在 map 中将读入的数据转化成 IntWritable 型,然后作为 key 值输出 (value 任意)。reduce 拿到 <key,value-list> 之后,将输入的 key 作为 value 输出,并根据 value-list 中元素的个数决定输出的次数。输出的 key (即代码中的 linenum) 是一个全局变量,它统计当前 key 的位次。需要注意的是这个程序中没有配置 Combiner,也就是在 MapReduce 过程中不使用 Combiner。这主要是因为使用 map 和 reduce 就已经能够完成任务了。

由于此程序简单且执行步骤与单词计数实例完全相同,所以不再赘述,下面只给出程序。

6.3.3 程序代码

程序代码如下所示:

```

import java.io.IOException;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;

public class Sort {

    //map 将输入中的 value 化成 IntWritable 类型, 作为输出的 key
    public static class Map extends Mapper<Object, Text, IntWritable, IntWritable>{

        private static IntWritable data = new IntWritable();

        public void map(Object key, Text value, Context context) throws IOException,
            InterruptedException {
            String line = value.toString();
            data.set(Integer.parseInt(line));
            context.write(data, new IntWritable(1));
        }
    }

    //reduce 将输入的 key 复制到输出的 value 上, 然后根据输入的
    //value-list 中元素的个数决定 key 的输出次数
    //用全局 linenum 来代表 key 的位次
    public static class Reduce extends Reducer<IntWritable, IntWritable, IntWritable,
        IntWritable> {
        private static IntWritable linenum = new IntWritable(1);

        public void reduce(IntWritable key, Iterable<IntWritable> values, Context
            context) throws IOException, InterruptedException {
            for (IntWritable val : values) {
                context.write(linenum, key);
                linenum = new IntWritable(linenum.get() + 1);
            }
        }
    }

    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        String[] otherArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
        if (otherArgs.length != 2) {
            System.err.println("Usage: wordcount <in> <out>");
            System.exit(2);
        }
    }
}

```

```

Job job = new Job(conf, "Sort");
job.setJarByClass(Sort.class);
job.setMapperClass(Map.class);
job.setReducerClass(Reduce.class);
job.setOutputKeyClass(IntWritable.class);
job.setOutputValueClass(IntWritable.class);
FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}

```

6.4 单表关联

前面的实例都是在数据上进行一些简单的处理，为进一步的操作打基础。单表关联这个实例要求从给出的数据中寻找出所关心的数据，它是对原始数据所包含信息的挖掘。下面进入这个实例。

6.4.1 实例描述

实例中给出 child-parent 表，要求输出 grandchild-grandparent 表。

样例输入如下所示。

file:

```

child parent
Tom Lucy
Tom Jack
Jone Lucy
Jone Jack
Lucy Mary
Lucy Ben
Jack Alice
Jack Jesse
Terry Alice
Terry Jesse
Philip Terry
Philip Alma
Mark Terry
Mark Alma

```

样例输出如下所示。

file:

grandchild	grandparent
Tom	Alice
Tom	Jesse
Jone	Alice
Jone	Jesse



Tom	Mary
Tom	Ben
Jone	Mary
Jone	Ben
Philip	Alice
Philip	Jesse
Mark	Alice
Mark	Jesse

6.4.2 设计思路

分析这个实例，显然需要进行单表连接，连接的是左表的 parent 列和右表的 child 列，且左表和右表是同一个表。连接结果中除去连接的两列就是所需要的结果——grandchild-grandparent 表。要用 MapReduce 解决这个实例，首先应该考虑如何实现表的自连接；其次就是连接列的设置；最后是结果的整理。考虑到 MapReduce 的 shuffle 过程会将相同的 key 值放在一起，所以可以将 map 结果的 key 值设置成待连接的列，然后列中相同的值就自然会连接在一起了。再与最开始的分析联系起来：要连接的是左表的 parent 列和右表的 child 列，且左表和右表是同一个表，所以在 map 阶段将读入数据分割成 child 和 parent 之后，会将 parent 设置成 key，child 设置成 value 进行输出，并作为左表；再将同一对 child 和 parent 中的 child 设置成 key，parent 设置成 value 进行输出，作为右表。为了区分输出中的左右表，需要在输出的 value 中再加上左右表的信息，比如在 value 的 String 最开始处加上字符 1 表示左表，加上字符 2 表示右表。这样在 map 的结果中就形成了左表和右表，然后在 shuffle 过程中完成连接。reduce 接收到连接的结果，其中每个 key 的 value-list 就包含了 grandchild 和 grandparent 关系。取出每个 key 的 value-list 进行解析，将左表中的 child 放入一个数组，右表中的 parent 放入一个数组，然后对两个数组求笛卡儿积就是最后的结果了。

在设计思路中已经包含了对程序的分析，而其程序执行步骤也与单词计数实例完全相同，所以代码解读和程序执行不再赘述，下面只给出代码。

6.4.3 程序代码

程序代码如下所示：

```
import java.io.IOException;
import java.util.*;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```




```

import org.apache.hadoop.util.GenericOptionsParser;

public class STjoin {
    public static int time = 0;
    //map 将输入分割成 child 和 parent, 然后正序输出一次作为右表, 反
    //序输出一次作为左表, 需要注意的是在输出的 value 中必须加上左右表
    //区别标志
    public static class Map extends Mapper<Object, Text, Text, Text>{

        public void map(Object key, Text value, Context context) throws IOException,
            InterruptedException {
            String childname = new String();
            String parentname = new String();
            String relationtype = new String();
            String line = value.toString();
            int i = 0;
            while(line.charAt(i)!=' '){
                i++;
            }

            String[] values = {line.substring(0,i),line.substring(i+1)};
            if(values[0].compareTo("child") != 0)
            {
                childname = values[0];
                parentname = values[1];
                relationtype = "1"; //左右表区分标志
                context.write(new Text(values[1]), new Text(relationtype + "+" +
                    childname + "+" + parentname));
                //左表
                relationtype = "2";
                context.write(new Text(values[0]), new Text(relationtype + "+" +
                    childname + "+" + parentname));
                //右表
            }
        }
    }

    public static class Reduce extends Reducer<Text,Text,Text,Text> {

        public void reduce(Text key, Iterable<Text> values,Context context) throws
            IOException, InterruptedException {

            if(time == 0){ //输出表头
                context.write(new Text("grandchild"),new Text("grandparent"));
                time++;
            }

            int grandchildnum = 0;
            String grandchild[] = new String[10];
            int grandparentnum = 0;
            String grandparent[] = new String[10];
            Iterator ite = values.iterator();

```

```

while(ite.hasNext())
{
    String record = ite.next().toString();
    int len = record.length();
    int i = 2;
    if(len == 0) continue;
    char relationtype = record.charAt(0);
    String childname = new String();
    String parentname = new String();
    // 获取 value-list 中 value 的 child
    while(record.charAt(i) != '+')
    {
        childname = childname + record.charAt(i);
        i++;
    }
    i = i+1;
    // 获取 value-list 中 value 的 parent
    while(i < len)
    {
        parentname = parentname + record.charAt(i);
        i++;
    }
    // 左表, 取出 child 放入 grandchild
    if(relationtype == '1'){
        grandchild[grandchildnum] = childname;
        grandchildnum++;
    }
    else{// 右表, 取出 parent 放入 grandparent
        grandparent[grandparentnum] = parentname;
        grandparentnum++;
    }
}

//grandchild 和 grandparent 数组求笛卡儿积
if(grandparentnum != 0 && grandchildnum != 0){
    for(int m = 0; m < grandchildnum; m++){
        for(int n = 0; n < grandparentnum; n++){
            context.write(new Text(grandchild[m]), new Text(grandparent[n])); // 输出结果
        }
    }
}

}
}

public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    String[] otherArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
    if (otherArgs.length != 2) {
        System.err.println("Usage: wordcount <in> <out>");
        System.exit(2);
    }
}

```

```

    }
    Job job = new Job(conf, "single table join");
    job.setJarByClass(STJoin.class);
    job.setMapperClass(Map.class);
    job.setReducerClass(Reduce.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(Text.class);
    FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
    FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}

```

6.5 多表关联

6.5.1 实例描述

多表关联和单表关联类似，它也是通过对原始数据进行一定的处理，从其中挖掘出关心的信息。下面进入这个实例。

输入是两个文件，一个代表工厂表，包含工厂名列和地址编号列；另一个代表地址表，包含地址名列和地址编号列。要求从输入数据中找出工厂名和地址名的对应关系，输出工厂名 - 地址名表。

样例输入如下所示。

factory :

```

factoryname addressed
Beijing Red Star 1
Shenzhen Thunder 3
Guangzhou Honda 2
Beijing Rising 1
Guangzhou Development Bank 2
Tencent 3
Bank of Beijing 1

```

address :

```

addressID addressname
1 Beijing
2 Guangzhou
3 Shenzhen
4 Xian

```

样例输出如下所示。

```

factoryname addressname
Bank of Beijing Beijing
Beijing Red Star Beijing

```



Beijing Rising Beijing
 Guangzhou Development Bank Guangzhou
 Guangzhou Honda Guangzhou
 Shenzhen Thunder Shenzhen
 Tencent Shenzhen

6.5.2 设计思路

多表关联和单表关联相似，都类似于数据库中的自然连接。相比单表关联，多表关联的左右表和连接列更加清楚。所以可以采用和单表关联相同的处理方式，map 识别出输入的行属于哪个表之后，对其进行分割，将连接的列值保存在 key 中，另一列和左右表标志保存在 value 中，然后输出。reduce 拿到连接结果之后，解析 value 内容，根据标志将左右表内容分开存放，然后求笛卡儿积，最后直接输出。

这个实例的具体分析参考单表关联实例。下面给出代码。

6.5.3 程序代码

程序代码如下所示：

```
import java.io.IOException;
import java.util.*;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;

public class MTjoin {
    public static int time = 0;

    public static class Map extends Mapper<Object, Text, Text, Text>{
        // 在 map 中先区分输入行属于左表还是右表，然后对两列值进行分割，
        // 保存连接列在 key 值，剩余列和左右表标志在 value 中，最后输出
        public void map(Object key, Text value, Context context) throws IOException,
            InterruptedException {
            String line = value.toString();
            int i = 0;
            // 输入文件首行，不处理
            if(line.contains("factoryname") == true || line.contains("addressID") ==
                true){
                return;
            }
        }
    }
}
```

```

// 找出数据中的分割点
while(line.charAt(i) >= '9' || line.charAt(i) <= '0'){
    i++;
}

if(line.charAt(0) >= '9' || line.charAt(0) <= '0') {
// 左表
    int j = i-1;
    while(line.charAt(j) != ' ') j--;
    String[] values = {line.substring(0,j),line.substring(i)};
    context.write(new Text(values[1]), new Text("1+" + values[0]));
}
else{ // 右表
    int j = i + 1;
    while(line.charAt(j) != ' ') j++;
    String[] values = {line.substring(0,i+1),line.substring(j)};
    context.write(new Text(values[0]), new Text("2+" + values[1]));
}
}

public static class Reduce extends Reducer<Text,Text,Text,Text> {
//reduce 解析 map 输出, 将 value 中数据按照左右表分别保存, 然后求
// 笛卡儿积, 输出
    public void reduce(Text key, Iterable<Text> values,Context context) throws
        IOException, InterruptedException {

        if(time == 0){ // 输出文件第一行
            context.write(new Text("factoryname"),new Text("addressname"));
            time++;
        }

        int factorynum = 0;
        String factory[] = new String[10];
        int addressnum = 0;
        String address[] = new String[10];
        Iterator ite = values.iterator();
        while(ite.hasNext())
        {
            String record = ite.next().toString();
            int len = record.length();
            int i = 2;
            char type = record.charAt(0);
            String factoryname = new String();
            String addressname = new String();
            if(type == '1'){ // 左表
                factory[factorynum] = record.substring(2);
                factorynum++;
            }

```

```

        else{ //右表
            address[addressnum] = record.substring(2);
            addressnum++;
        }
    }
    if(factorynum != 0 && addressnum != 0){ //求笛卡儿积
        for(int m = 0; m < factorynum; m++){
            for(int n = 0; n < addressnum; n++){
                context.write(new Text(factory[m]), new Text(address[n]));
            }
        }
    }
}

public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    String[] otherArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
    if (otherArgs.length != 2) {
        System.err.println("Usage: wordcount <in> <out>");
        System.exit(2);
    }
    Job job = new Job(conf, "multiple table join");
    job.setJarByClass(MTJoin.class);
    job.setMapperClass(Map.class);
    job.setReducerClass(Reduce.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(Text.class);
    FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
    FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}

```

6.6 小结

本章通过五个实例，向读者呈现了如何使用 MapReduce 程序解决实际问题。其中第一个 WordCount 实例是 MapReduce 的入门程序，它能统计出数据文件中单词的频数。实例二数据去重和实例三数据排序，都是对原始数据的初步操作，为进一步的数据分析打下基础。实例四单表关联和实例五多表关联是对数据的进一步操作，从中挖掘有用的信息。虽然五个实例相对简单普通，但是都能利用 Hadoop 平台对大数据集进行并行处理，展示了 MapReduce 编程框架的魅力所在。



第7章

MapReduce 工作机制

本章内容

- ☐ MapReduce 作业的执行流程
- ☐ 错误处理机制
- ☐ 作业调度机制
- ☐ shuffle 和排序
- ☐ 任务执行
- ☐ 小结

资源如书
PDG

关于 MapReduce 的准备知识和应用案例在本书前面的章节中已经介绍得很详细了，本章将从 MapReduce 作业的执行情况、作业运行过程中的错误机制、作业的调度策略、shuffle 和排序、任务的执行等几个方面详细讲解 MapReduce，让读者更加深入地了解 MapReduce 的运行机制，为深入学习使用 Hadoop 和 Hadoop 子项目打下基础。

7.1 MapReduce 作业的执行流程

从第 6 章的 MapReduce 编程实例中可以看出，只要在 main() 函数中写入 JobClient.runJob(conf)，然后将程序提交到 Hadoop 上，MapReduce 作业就可以在 Hadoop 上运行。另外，在前面的其他章节中也从 task 运行角度介绍了 Map 和 Reduce 的过程，但是从运行“Hadoop jar”到看到作业运行结果，这中间实际上还涉及很多其他细节。那么 Hadoop 运行 MapReduce 作业的完整步骤是什么呢？每一步又是如何具体实现的呢？本节将详细介绍。

7.1.1 MapReduce 任务的执行总流程

通过前面的知识我们知道，一个 MapReduce 作业的执行流程是：代码编写→作业配置→作业提交→Map 任务的分配和执行→处理中间结果→Reduce 任务的分配和执行→作业完成，而在每个任务的执行过程中，又包含输入准备→任务执行→输出结果。

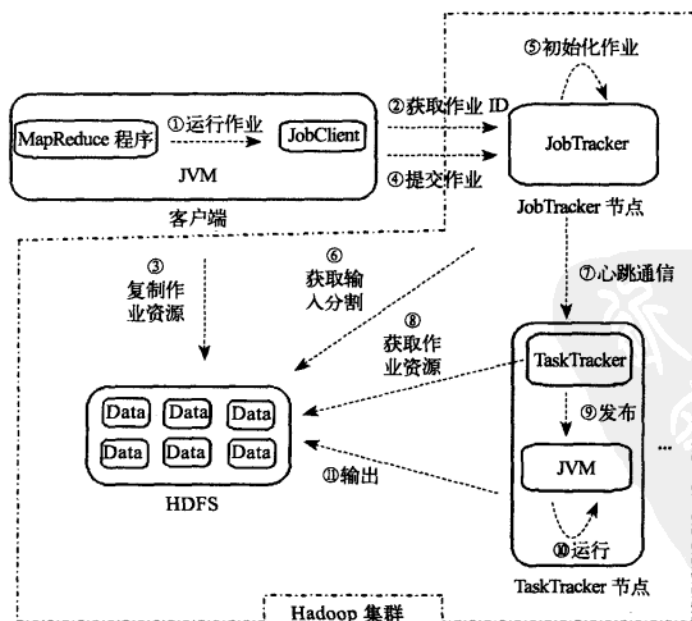


图 7-1 MapReduce 作业执行的流程图

图 7-1 给出了 MapReduce 作业详细的执行流程图。从图中可以看出 MapReduce 作业的执行可以分为 11 个步骤，涉及 4 个独立的实体。它们在 MapReduce 执行过程中的主要作用是：

- ❑ 客户端 (client)：编写 MapReduce 代码，配置作业，提交作业；
- ❑ JobTracker：初始化作业，分配作业，与 TaskTracker 通信，协调整个作业的执行；
- ❑ TaskTracker：保持 JobTracker 的通信，在分配的数据片段上执行 Map 或 Reduce 任务。需要注意的是图 7-1 中 TaskTracker 节点后的省略号表示 Hadoop 集群中可以包含多个 TaskTracker；
- ❑ HDFS：保存作业的数据、配置信息等，保存作业结果。

下面按照图 7-1 中 MapReduce 作业的执行流程结合代码详细介绍各个步骤。

7.1.2 提交作业

一个 MapReduce 作业在提交到 Hadoop 上之后，会进入完全地自动化执行过程。在这个过程中，用户除了监控程序的执行情况和强制中止作业之外，不能对作业的执行过程进行任何干扰。所以在作业提交之前，用户需要将所有应该配置的参数按照自己的意愿配置完毕。需要配置的主要内容有：

- ❑ 程序代码：这里主要是指 map 函数和 reduce 函数的具体代码，这是一个 MapReduce 作业对应的程序必不可少的部分，并且这部分代码的逻辑正确与否与运行结果直接相关。
- ❑ Map 接口和 Reduce 接口的配置：在 MapReduce 中，Map 接口需要派生自 Mapper<k1,v1,k2,v2> 接口，Reduce 接口则要派生自 Reducer<k2,v2,k3,v3>。它们都对应唯一一个方法，分别是 map 函数和 reduce 函数，也就是在上一点中所写的代码。在调用这两个方法时需要配置它们的四个参数，分别是输入 key 的数据类型、输入 value 的数据类型、输出 key-value 对的数据类型和 Reporter 实例，其中输入输出的数据类型要与继承时所设置的数据类型相同，还有一个要求是 Map 接口的输出 key-value 类型和 Reduce 接口的输入 key-value 类型要对应，因为 map 输出组合 value 之后，它们会成为 reduce 的输入内容（入门者请特别注意，很多入门者编写的 MapReduce 程序中会忽视这个问题）。
- ❑ 输入输出路径：作业提交之前还需要在主函数中配置 MapReduce 作业在 Hadoop 集群上的输入路径和输出路径（必须保证输出路径不存在，如果存在程序会报错，这也是初学者经常犯的错误）。具体的代码是：

```
FileInputFormat.setInputPaths(conf, new Path(args[0]));
FileOutputFormat.setOutputPath(conf, new Path(args[1]));
```

- ❑ 其他类型设置，比如调用 runJob 方法：先要在主函数中配置如 Output 的 key 和 value 类型、作业名称、InputFormat 和 OutputFormat 等，最后再调用 JobClient 的 runJob 方法。

配置完作业的所有内容并确认无误之后就可以运行作业了，也就是执行图 7-1 中的步

骤①（具体提交方法不再赘述，请参考本书的第5章）。

用户程序调用 JobClient 的 runJob 方法，在提交 JobConf 对象之后，runJob 方法会先行调用 JobSubmissionProtocol 接口所定义的 submitJob 方法，并将作业提交给 JobTracker。紧接着，runJob 不断循环，并在循环中调用 JobSubmissionProtocol 的 getTaskCompletionEvents 方法，获取 TaskCompletionEvent 类的对象实例，了解作业的实时执行情况。如果发现作业运行状态有更新，就将状态报告给 JobTracker。作业完成后，如果成功就显示作业计数器，否则，将导致作业失败的错误记录到控制台。

上面介绍了作业提交的过程，可以看出，最关键的是 JobClient 对象中 submitJob() 方法的调用执行，那么 submitJob() 方法具体是怎么做的呢？下面从 submitJob() 方法的代码出发介绍作业提交的详细过程。

```
public RunningJob submitJob(JobConf job) throws FileNotFoundException,
    InvalidJobConfException, IOException {
    // 从 JobTracker 得到当前任务的 ID
    JobID jobId = jobSubmitClient.getNewJobId();
    // 获取 HDFS 路径:
    Path submitJobDir = new Path(getSystemDir(), jobId.toString());
    // 获取提交作业的 JAR 文件目录
    Path submitJarFile = new Path(submitJobDir, "job.jar");
    // 获取提交输入文件分割信息的路径
    Path submitSplitFile = new Path(submitJobDir, "job.split");
    // 此处将 -libjars 命令行指定的 jar 上传至 HDFS
    configureCommandLineOptions(job, submitJobDir, submitJarFile);
    // 获取作业配置文档的路径
    Path submitJobFile = new Path(submitJobDir, "job.xml");
    .....
    // 通过 input format 的格式获得相应的 input split, 默认类型为
    //FileSplit InputSplit[] splits = job.getInputFormat().getSplits(job, job.
        getNumMapTasks());
    // 生成一个写入流，将 input split 的信息写入 job.split 文件
    FSDataOutputStream out = FileSystem.create(fs, submitSplitFile, new
        FsPermission(JOB_FILE_PERMISSION));
    try {
        writeSplitsFile(splits, out);
    } finally {
        out.close();
    }
    job.set("mapred.job.split.file", submitSplitFile.toString());
    // 根据 split 的个数设定 map task 的个数
    job.setNumMapTasks(splits.length);
    // 将 job 的配置信息写入 job.xml 文件
    out = FileSystem.create(fs, submitJobFile,
        new FsPermission(JOB_FILE_PERMISSION));
    try {
        job.writeXml(out);
    } finally {
```

```

        out.close();
    }
    // 真正地调用 JobTracker 来提交任务
    JobStatus status = jobSubmitClient.submitJob(jobId);
    .....
}

```

从上面的代码可以看出，整个提交过程包含以下步骤：

- 1) 通过调用 JobTracker 对象的 `getNewJobId()` 方法从 JobTracker 处获取当前作业 ID 号（见图 7-1 中的步骤②）。
- 2) 检查作业相关路径。在代码中获取各个路径信息的时候会对作业的对应路径进行检查。比如，如果没有指定输出目录或它已经存在，作业就不会被提交，并且会给 MapReduce 程序返回错误信息，再比如输入目录不存在也会返回错误等。
- 3) 计算作业的输入划分，并将划分信息写入 `job.split` 文件，如果写入失败就会返回错误。`split` 文件的信息主要包括：`split` 文件头、`split` 文件版本号、`split` 的个数。这些信息中每一条都会包括以下内容：`split` 类型名（默认 `FileSplit`）、`split` 的大小、`split` 的内容（对于 `FileSplit` 来说是写入的文件名，此 `split` 在文件中的起始位置上）、`split` 的 location 信息（即在哪个 `DataNode` 上）。
- 4) 将运行作业所需要的资源——包括作业 JAR 文件、配置文件和计算所得的输入划分等——复制到作业对应的 HDFS 上（见图 7-1 的步骤③）。
- 5) 调用 JobTracker 对象的 `submitJob()` 方法来真正提交作业，告诉 JobTracker 作业准备执行（见图 7-1 的步骤④）。

7.1.3 初始化作业

在客户端用户作业调用 JobTracker 对象的 `submitJob()` 方法后，JobTracker 会把此调用放入内部的 `TaskScheduler` 变量中，然后进行调度，默认的调度方法是 `JobQueueTaskScheduler`，也就是 FIFO 调度方式。当客户作业被调度执行时，JobTracker 会创建一个代表这个作业的 `JobInProgress` 对象，并将任务和记录信息封装到这个对象中，以便跟踪任务的状态和进程。接下来 `JobInProgress` 对象的 `initTasks` 函数会对任务进行初始化操作（见图 7-1 的步骤⑤）。下面仍然从 `initTasks` 函数的代码出发详细讲解初始化过程。

```

public synchronized void initTasks() throws IOException {
    .....
    // 从 HDFS 中作业对应的路径读取 job.split 文件，生成 input
    // splits 为下面 map 的划分做好准备
    String jobFile = profile.getJobFile();
    Path sysDir = new Path(this.jobtracker.getSystemDir());
    FileSystem fs = sysDir.getFileSystem(conf);
    DataInputStream splitFile =
        fs.open(new Path(conf.get("mapred.job.split.file")));
    JobClient.RawSplit[] splits;
}

```

```

try {
    splits = JobClient.readSplitFile(splitFile);
} finally {
    splitFile.close();
}
// 根据 input split 设置 map task 个数
numMapTasks = splits.length;
// 为每个 map tasks 生成一个 TaskInProgress 来处理一个 input split
maps = new TaskInProgress[numMapTasks];
for(int i=0; i < numMapTasks; ++i) {
    inputLength += splits[i].getDataLength();
    maps[i] = new TaskInProgress(jobId, jobFile, splits[i], jobtracker, conf,
        this, i); }
if (numMapTasks > 0) {
    //map task 放入 nonRunningMapCache, 其将在 JobTracker 向
    //TaskTracker 分配 map task 的时候使用
    nonRunningMapCache = createCache(splits, maxLevel);
}
// 创建 reduce task
this.reduces = new TaskInProgress[numReduceTasks];
for (int i = 0; i < numReduceTasks; i++) {
    reduces[i] = new TaskInProgress(jobId, jobFile, numMapTasks, i, jobtracker,
        conf, this);
//reduce task 放入 nonRunningReduces, 其将在 JobTracker 向
//TaskTracker 分配 reduce task 的时候使用
    nonRunningReduces.add(reduces[i]);
}

// 清理 map 和 reduce
cleanup = new TaskInProgress[2];
cleanup[0] = new TaskInProgress(jobId, jobFile, splits[0], jobtracker, conf,
    this, numMapTasks);
cleanup[0].setJobCleanupTask();
cleanup[1] = new TaskInProgress(jobId, jobFile, numMapTasks, numReduceTasks,
    jobtracker, conf, this);
cleanup[1].setJobCleanupTask();
// 创建两个初始化 task, 一个初始化 map, 一个初始化 reduce
setup = new TaskInProgress[2];
setup[0] = new TaskInProgress(jobId, jobFile, splits[0], jobtracker, conf, this,
    numMapTasks + 1 );
setup[0].setJobSetupTask();
setup[1] = new TaskInProgress(jobId, jobFile, numMapTasks, numReduceTasks + 1,
    jobtracker, conf, this);
setup[1].setJobSetupTask();
tasksInitiated.set(true); // 初始化完毕
.....
}

```

从上面的代码可以看出在初始化过程中主要有以下步骤:

1) 从 HDFS 中读取作业对应的 job.split (见图 7-1 的步骤⑥)。JobTracker 从 HDFS 中

作业对应的路径获取 JobClient 在步骤③中写入的 job.split 文件，得到输入数据的划分信息。为后面初始化过程中 map 任务的分配做好准备。

2) 创建并初始化 map 任务和 reduce 任务。initTasks 先根据输入数据划分信息中的个数设定 map task 的个数，然后为每个 map tasks 生成一个 TaskInProgress 来处理 input split，并将 map task 放入 nonRunningMapCache 中，以便在 JobTracker 向 TaskTracker 分配 map task 的时候使用。接下来根据 JobConf 中的 mapred.reduce.tasks 属性利用 setNumReduceTasks() 方法来设置 reduce task 的个数，然后采用类似 map task 的方式将 reduce task 放入 nonRunning-Reduces 中，以便在向 TaskTracker 分配 reduce task 的时候使用。

3) 最后就是创建两个初始化 task，根据个数和输入划分已经配置的信息，并分别初始化 map 和 reduce。

7.1.4 分配任务

在前面的介绍中已经知道，TaskTracker 和 JobTracker 之间的通信与任务的分配是通过心跳机制完成的。TaskTracker 作为一个单独的 JVM 执行一个简单的循环，主要实现每隔一段时间向 JobTracker 发送心跳 (heartbeat)：告诉 JobTracker，此 TaskTracker 是否存活，是否准备执行新的任务。在 JobTracker 接收到心跳信息后，如果有待分配的任务，它就会为 TaskTracker 分配一个任务，并将分配信息封装在心跳通信的返回值中返回给 TaskTracker，TaskTracker 从心跳方法的 Response 中得知此 TaskTracker 需要做的事情，如果是一个新的 task 则将 task 加入本机的任务队列中（见图 7-1 的步骤④）。

下面从 TaskTracker 中的 transmitHeartBeat() 方法和 JobTracker 中的 heartbeat() 方法的主要代码出发，介绍任务分配的详细过程，以及在此过程中 TaskTracker 和 JobTracker 的通信。

TaskTracker 中 transmitHeartBeat() 方法的主要代码如下所示：

```
// 向 JobTracker 报告 TaskTracker 的当前状态
if (status == null) {
    synchronized (this) {
        status = new TaskTrackerStatus(taskTrackerName, localHostname,
            httpPort, cloneAndResetRunningTaskStatuses(sendCounters), failures,
            maxCurrentMapTasks, maxCurrentReduceTasks);
    }
}
.....
// 根据条件是否满足来确定此 TaskTracker 是否请求 JobTracker
// 为其分配新的 Task
boolean askForNewTask;
long localMinSpaceStart;
synchronized (this) {
    askForNewTask = (status.countMapTasks() < maxCurrentMapTasks ||
        status.countReduceTasks() < maxCurrentReduceTasks) &&
        acceptNewTasks;
```


7.1.5 执行任务

在 TaskTracker 申请到新的任务之后，就要在本地运行任务了。运行任务的第一步是将任务本地化（将任务运行所必需的数据、配置信息、程序代码从 HDFS 复制到 TaskTracker 本地，见图 7-1 的步骤⑧）。这主要是通过调用 `localizeJob()` 方法来完成的（此方法的具体代码并不复杂，不再列出）。这个方法主要通过下面几个步骤来完成任务的本地化：

- 1) 将 `job.split` 拷贝到本地；
- 2) 将 `job.jar` 拷贝到本地；
- 3) 将 `job` 的配置信息写入 `job.xml`；
- 4) 创建本地任务目录，解压 `job.jar`；
- 5) 调用 `launchTaskForJob()` 方法发布任务（见图 7-1 的步骤⑨）。

任务本地化之后，就可通过调用 `launchTaskForJob()` 真正启动起来。接下来 `launchTaskForJob()` 又会调用 `launchTask()` 方法启动任务。`launchTask()` 方法的主要代码如下：

```
...
// 创建 task 本地运行目录
localizeTask(task);
if (this.taskStatus.getRunState() == TaskStatus.State.UNASSIGNED) {
    this.taskStatus.setRunState(TaskStatus.State.RUNNING);
}
// 创建并启动 TaskRunner
this.runner = task.createRunner(TaskTracker.this, this);
this.runner.start();
this.taskStatus.setStartTime(System.currentTimeMillis());
...
```

从代码中可以看出 `launchTask()` 方法先会为任务创建本地目录，然后启动 `TaskRunner`。在启动 `TaskRunner` 后，对于 `map` 任务，会启动 `MapTaskRunner`；对于 `reduce` 任务则启动 `ReduceTaskRunner`。

这之后，`TaskRunner` 又会启动新的 Java 虚拟机来运行每个任务（见图 7-1 的步骤⑩）。以 `map` 任务为例，任务执行的简单流程是：

- 1) 配置任务的执行参数（获取 Java 程序的执行环境和配置参数等）；
- 2) 在 `Child` 临时文件表中添加 `map` 任务信息（运行 `map` 和 `reduce` 任务的主进程是 `Child` 类）；
- 3) 配置 `log` 文件夹，然后配置 `map` 任务的通信和输出参数；
- 4) 读取 `input split`，生成 `RecordReader` 读取数据；
- 5) 为 `map` 任务生成 `MapRunnable`，依次从 `RecordReader` 中接收数据，并调用 `Mapper` 的 `map` 函数进行处理；
- 6) 最后将 `map` 函数的输出调用 `collect` 收集到 `MapOutputBuffer` 中（见图 7-1 的步骤⑪）。

7.1.6 更新任务执行进度和状态

在本章的作业提交过程中曾介绍：一个 MapReduce 作业在提交到 Hadoop 上之后，会进入完全地自动化执行过程，用户只能监控程序的执行状态和强制中止作业。但是 MapReduce 作业是一个长时间运行的批量作业，有时候可能需要运行数小时。所以对于用户而言，能够得知作业的运行状态是非常重要的。在 Linux 终端运行 MapReduce 作业时，可以看到在作业执行过程中有一些简单的作业执行状态报告，这能让用户大致了解作业的运行情况，并通过与预期运行情况进行对比来确定作业是否按照预定方式运行。

在 MapReduce 作业中，作业的进度主要由一些可衡量可计数的小操作组成。比如在 map 任务中，其任务进度就是已处理输入的百分比，比如完成 100 条记录中的 50 条，那么 map 任务的进度就是 50%（这里只是针对一个 map 任务举例，并不是指在 Linux 终端中执行 MapReduce 任务时出现的 map 50%，在终端中出现的 50% 是总体 map 任务的进度，这是将所有 map 任务的进度组合起来的结果）。总体来讲，MapReduce 作业的进度由下面几项组成：mapper（或 reducer）读入或写出一条记录，在报告中设置状态描述，增加计数器，调用 Reporter 对象的 progress() 方法。

由 MapReduce 作业分割的每个任务中都有一组计数器，它们对任务执行过程中的进度组成事件进行计数。如果任务要报告进度，它便会设置一个标志以表明状态变化将会发送到 TaskTracker 上。另一个监听线程检查到这标志后，会告知 TaskTracker 当前的任务状态。具体代码如下（这是 MapTask 中 run 函数的部分代码）：

```
// 同 TaskTracker 通信，汇报任务执行进度
final Reporter reporter = getReporter(umbilical);
startCommunicationThread(umbilical);
initialize(job, reporter);
```

同时，TaskTracker 每隔 5 秒在发送给 JobTracker 的心跳中封装任务状态，报告自己的任务执行状态。具体代码如下（这是 TaskTracker 中 transmitHeartBeat() 方法的部分代码）：

```
// 每隔一段时间，向 JobTracker 返回一些统计信息
boolean sendCounters;
if (now > (previousUpdate + COUNTER_UPDATE_INTERVAL)) {
    sendCounters = true;    previousUpdate = now;
}
else {
    sendCounters = false;
}
```

通过心跳通信机制，所有 TaskTracker 的统计信息都会汇总到 JobTracker 处。JobTracker 将这些统计信息合并起来，产生一个全局作业进度统计信息，用来表明正在运行的所有作业，以及其中所含任务的状态。最后，JobClient 通过每秒查看 JobTracker 来接收作业进度的最新状态。具体代码如下（这是 JobClient 中用来提交作业的 runJob() 方法的部分代码）：


```
// 首先生成一个 JobClient 对象
JobClient jc = new JobClient(job);
.....
// 调用 submitJob 来提交一个任务
running = jc.submitJob(job);
JobID jobId = running.getID();
.....
//while 循环中不断查看 JobTracker 来获得作业进度
while (true) {
    ...
}
```

7.1.7 完成作业

所有 TaskTracker 任务的执行进度信息都会汇总到 JobTracker 处，当 JobTracker 接收到最后一个任务的已完成通知后，便把作业的状态设置为“成功”。然后，JobClient 也将及时得知任务已成功完成，它便会显示一条信息告知用户作业已完成，最后从 runJob() 方法处返回（在返回后 JobTracker 会清空作业的工作状态，并指示 TaskTracker 也清空作业的工作状态，比如删除中间输出等）。

7.2 错误处理机制

众所周知，Hadoop 有很强的容错性。这主要是针对由成千上万台普通机器组成的集群中常态化的硬件故障的，Hadoop 能够利用冗余数据方式来解决硬件故障，以保证数据安全和任务执行。那么 MapReduce 在具体执行作业过程中遇到硬件故障会如何处理呢？对于用户代码的缺陷或进程崩溃引起的错误又会如何处理呢？本节将从硬件故障和任务失败两个方面说明 MapReduce 的错误处理机制。

7.2.1 硬件故障

从 MapReduce 任务的执行角度出发，所涉及的硬件主要是 JobTracker 和 TaskTracker（对应从 HDFS 出发就是 NameNode 和 DataNode）。显然硬件故障就是 JobTracker 机器故障和 TaskTracker 机器故障。

在 Hadoop 集群中，任何时候都只有唯一一个 JobTracker。所以 JobTracker 故障就是单点故障，这是所有错误中最严重的错误。到目前为止，在 Hadoop 中还没有相应的解决办法。能够想到的是通过创建多个备用 JobTracker 节点，在主 JobTracker 失败之后采用领导选举算法（Hadoop 中常用的一种确定 master 的算法）来重新确定 JobTracker 节点。在一些企业使用 Hadoop 提供服务时，就采用了这样的方法来避免 JobTracker 错误。

机器故障除了 JobTracker 错误外就是 TaskTracker 错误。TaskTracker 故障相对较为常见，并且 MapReduce 也有相应的解决办法，主要是重新执行任务。下面将详细介绍当作业遇到 TaskTracker 错误时，MapReduce 所采取的解决步骤。

在 Hadoop 中, 正常情况下, TaskTracker 会不断地与系统 JobTracker 通过心跳机制进行通信。如果某 TaskTracker 出现故障或运行缓慢, 它会停止或很少向 JobTracker 发送心跳。如果一个 TaskTracker 在一定时间内(默认是 1 分钟)没有与 JobTracker 通信, 那么 JobTracker 会将此 TaskTracker 从等待任务调度的 TaskTracker 集合中移除。同时 JobTracker 会要求此 TaskTracker 上的任务立刻返回, 如果此 TaskTracker 任务是仍然在 mapping 阶段的 map 任务, 那么 JobTracker 会要求其他的 TaskTracker 重新执行所有原本由故障 TaskTracker 执行的 map 任务。如果任务是在 reduce 阶段的 reduce 任务, 那么 JobTracker 会要求其他 TaskTracker 重新执行故障 TaskTracker 未完成的 reduce 任务。比如, 一个 TaskTracker 已经完成被分配的 3 个 reduce 任务中的 2 个, 由于 reduce 任务一旦完成会将数据写到 HDFS 上, 所以只有第三个未完成的 reduce 需要重新执行。但是对于 map 任务来说, 即使 TaskTracker 完成了部分 map, 但是 reduce 仍可能无法获取此节点上所有 map 的所有输出。所以无论 map 任务完成与否, 故障 TaskTracker 上的 map 任务都必须重新执行。

7.2.2 任务失败

在实际任务中, MapReduce 作业还会遇到用户代码缺陷或进程崩溃引起的任务失败。用户代码缺陷会导致它在执行过程中抛出异常。此时, 任务 JVM 进程会自动退出, 并向 TaskTracker 父进程发送错误消息, 同时错误消息也会写入 log 文件, 最后 TaskTracker 将此次任务尝试标记失败。对于进程崩溃引起的任务失败, TaskTracker 的监听程序会发现进程退出, 此时 TaskTracker 也会将此次任务尝试标记为失败。对于死循环程序或执行时间太长的程序, 由于 TaskTracker 没有接收到进度更新, 它也会将此次任务尝试标记为失败, 并杀死程序对应的进程。

在以上情况中, TaskTracker 将任务尝试标记为失败之后会将 TaskTracker 自身的任务计数器减 1, 以便向 JobTracker 申请新的任务。TaskTracker 也会通过心跳机制告诉 JobTracker 本地的一个任务尝试失败。JobTracker 接到任务失败的通知后, 通过重置任务状态, 将其加入调度队列来重新分配该任务执行(JobTracker 会尝试避免将失败的任务再次分配给运行失败的 TaskTracker)。如果此任务尝试了 4 次(次数可以进行设置)仍没有完成, 就不会再被重试, 此时整个作业也就失败了。

7.3 作业调度机制

在 0.19.0 版本之前, Hadoop 集群上的用户作业采用先进先出(FIFO, First Input First Output)的调度算法, 即按照作业提交的顺序来运行。同时每个作业都会使用整个集群, 因此它们只有轮到自己运行时才能享受整个集群的服务。虽然 FIFO 调度器最后又支持设置了优先级的功能, 但是由于不支持优先级抢占, 所以这种单用户的调度算法仍然不符合云计算中采用并行计算来提供服务的宗旨。从 0.19.0 版本开始, Hadoop 除了默认的 FIFO 调度器外, 还提供了支持多用户同时服务和集群资源公平共享的调度器, 即公平调度器(Fair Scheduler)

Guide) 和容量调度器 (Capacity Scheduler Guide)。下面主要介绍公平调度器。

公平调度是为作业分配资源的方法, 其目的是随着时间的推移, 让提交的作业获取等量的集群共享资源, 让用户公平地共享集群。具体做法是: 当集群上只有一个作业在运行时, 它将使用整个集群, 当有其他作业提交时, 系统会将 TaskTracker 节点空闲的时间片分配给这些新的作业, 并保证每一个作业都得到大概等量的 CPU 时间。

公平调度器按作业池来组织作业, 它会按照提交作业的用户数目将资源公平地分到这些作业池里。默认情况下, 每一个用户拥有一个独立的作业池, 以使每个用户都能获得一份等量的集群资源而不会管它们提交了多少作业。在每一个资源池内, 会使用公平共享的方法在运行作业之间共享容量。除了提供公平共享方法外, 公平调度器还允许为作业池设置最小的共享资源, 以确保特定用户、群组或生产应用程序总能获取到足够的资源。对于设置了最小共享资源的作业池来说, 如果它包含了作业, 它至少能获取到最小的共享资源。但是如果最小共享资源超过作业需要的资源时, 额外的资源会在其他作业池间进行切分。

在常规操作中, 当提交了一个新作业时, 公平调度器会等待已运行作业中的任务完成, 以释放时间片给新的作业。但公平调度器也支持作业抢占。如果新的作业在一定时间 (即超时时间, 可以配置) 内还未获取公平的资源分配, 公平调度器就会允许这个作业抢占已运行作业中的任务, 以获取运行所需要的资源。另外, 如果作业在超时时间内获取的资源不到公平共享资源的一半时也允许对任务进行抢占。而在选择时, 公平调度器会在所有运行任务中选择最近运行起来的任务, 这样浪费的计算相对较少。由于 Hadoop 作业能容忍丢失任务, 抢占不会导致被抢占的作业失败, 只是让被抢占作业的运行时间更长。

最后, 公平调度器还可以限制每个用户和每个作业池并发运行的作业数量。这个限制可以在一个用户一次性提交数百个作业或当大量作业并发执行时来确保中间数据不会塞满集群上的磁盘空间。超出限制的作业会被列入调度器的队列中进行等待, 直到早期作业运行完毕。公平调度器再根据作业优先权和提交时间的排列情况从等待作业中调度即将运行的作业。

7.4 shuffle 和排序

从前面的介绍中, 我们得知 map 的输出会经过一个名为 shuffle 的过程交给 reduce 处理 (在 “MapReduce 数据流” 图中也可以看出), 当然也有 map 的结果经过 sort-merge 交给 reduce 处理的。其实在 MapReduce 流程中, 为了让 reduce 可以并行处理 map 结果, 必须对 map 的输出进行一定的排序和分割, 然后再交给对应的 reduce, 而这个将 map 输出进行进一步整理并交给 reduce 的过程就成为了 shuffle。从 shuffle 的过程中可以看出, 它是 MapReduce 的核心所在, shuffle 过程的性能与整个 MapReduce 的性能直接相关。

总体来说, shuffle 过程包含在 map 和 reduce 两端中。在 map 端的 shuffle 过程是对 map 的结果进行划分 (partition)、排序 (sort) 和分割 (spill), 然后将属于同一个划分的输出合并在一起 (merge), 并写在磁盘上, 同时按照不同的划分将结果发送给对应的 reduce (map 输出的划分与 reduce 的对应关系由 JobTracker 确定)。reduce 端又会将各个 map 送来的属于

同一个划分的输出进行合并 (merge)，然后对 merge 的结果进行排序，最后交给 reduce 处理。下面将从 map 和 reduce 两端详细介绍 shuffle 过程。

7.4.1 map 端

从 MapReduce 的程序中可以看出，map 的输出结果是由 collector 处理的，所以 map 端的 shuffle 过程包含在 collect 函数对 map 输出结果的处理过程中。下面从具体的代码来分析 map 端的 shuffle 过程。

首先从 collect 函数的代码入手。从下面的代码段可以看出 map 函数的输出内存缓冲区是一个环形结构。

```
final int kvnext = (kvindex + 1) % kvoffsets.length;
```

当输出内存缓冲区内容达到设定的阈值时，就需要把缓冲区内容分割 (spill) 到磁盘中了。但是在分割的时候 map 并不会阻止继续向缓冲区中写入结果，如果 map 结果生成的速度快于写出速度，那么缓冲区会写满，这时 map 任务必须等待，直到分割写出过程结束。这个过程可以参考下面的代码。

```
do {
    // 在环形缓冲区中，如果下一个空闲位置同起始位置相等，那么缓冲区
    // 已满
    kvfull = kvnext == kvstart;
    // 环形缓冲区的内容是否达到写出的阈值
    final boolean kvsoftlimit = ((kvnext > kvend)
        ? kvnext - kvend > softRecordLimit
        : kvend - kvnext <= kvoffsets.length - softRecordLimit);
    // 达到阈值，写出缓冲区内容，形成 spill 文件
    if (kvstart == kvend && kvsoftlimit) {
        startSpill();
    }
    // 如果缓冲区满，则 map 任务等待写出过程结束
    if (kvfull) {
        while (kvstart != kvend) {
            reporter.progress();
            spillDone.await();
        }
    }
} while (kvfull);
```

在 collect 函数中将缓冲区中的内容写出时会调用 sortAndSpill 函数。sortAndSpill 函数每被调用一次就会创建一个 spill 文件，然后按照 key 值对需要写出的数据进行排序，最后按照划分的顺序将所有需要写出的结果写入这个 spill 文件中。如果用户作业配置了 combiner 类，那么在写出过程中会先调用 combineAndSpill() 再写出，对结果进行进一步的合并 (combine) 是为了让 map 的输出数据更加紧凑。sortAndSpill 函数的执行过程可以参考下面 sortAndSpill 函数的代码。

```

// 创建 spill 文件
Path filename = mapOutputFile.getSpillFileForWrite(getTaskID(), numSpills,
    size);
out = rfs.create(filename);
.....
// 按照 key 值对待写出数据进行排序
sorter.sort(MapOutputBuffer.this, kvstart, endPosition, reporter);
...
// 按照划分将数据写入文件
for (int i = 0; i < partitions; ++i) {
    IFile.Writer<K, V> writer = null;
    long segmentStart = out.getPos();
    writer = new Writer<K, V>(job, out, keyClass, valClass, codec);
    // 如果没有配置 combiner 类, 数据直接写入文件
    if (null == combinerClass) {
        .....
    }
    else {
        .....
    }
    // 如果配置了 combiner 类, 则先调用 combineAndSpill 函数后再写入文件
    combineAndSpill(kvIter, combineInputCounter);
} }

```

显然, 直接将每个 map 生成的众多 spill 文件 (因为 map 过程中, 每一次缓冲区写出都会产生一个 spill 文件) 交给 reduce 处理不现实。所以在每个 map 任务结束之后在 map 的 TaskTracker 上还会执行合并操作 (merge), 这个操作的主要目的是将 map 生成的众多 spill 文件中的数据按照划分重新组织, 以便于 reduce 处理。主要做法是针对指定的分区, 从各个 spill 文件中拿出属于同一个分区的所有数据, 然后将它们合并在一起, 并写入一个已分区且已排序的 map 输出文件中。这个过程的具体情况请参考 mergeParts() 函数的代码, 这里不再列出。

待唯一的已分区且已排序的 map 输出文件写入最后一条记录后, map 端的 shuffle 阶段就结束了。下面就进入 reduce 端的 shuffle 阶段。

7.4.2 reduce 端

在 reduce 端, shuffle 阶段可以分成三个阶段: 复制 map 输出、排序合并、reduce 处理。下面按照这三个阶段进行详细介绍。

如前文所述, map 任务成功完成后, 会通知父 TaskTracker 状态已更新, 进而 TaskTracker 通知 JobTracker (这些通知在心跳机制中进行)。所以, 对于指定作业来说, JobTracker 能够记录 map 输出和 TaskTracker 的映射关系。reduce 会定期向 JobTracker 获取 map 的输出位置。一旦拿到输出位置, reduce 任务就会从此输出对应的 TaskTracker 上复制输出到本地 (如果 map 的输出很小, 则会被复制到执行 reduce 任务的 TaskTracker 节点的内存中, 便于进一步的处理, 否则会放入磁盘), 而不会等到所有的 map 任务结束。这就是

reduce 任务的复制阶段。

在 reduce 复制 map 的输出结果的同时，reduce 任务就进入了合并（merge）阶段。这一阶段主要的任务是将各个 map TaskTracker 上复制的 map 输出文件（无论在内存还是在磁盘）进行整合，并维持数据原来的顺序。

reduce 端的最后阶段就是对合并的文件进行 reduce 处理。reduce TaskTracker 从合并的文件中按照顺序先拿出一条数据，交给 reduce 函数处理，然后将结果输出到本地的 HDFS 上（因为在 Hadoop 集群上，TaskTracker 节点一般也是 DataNode 节点），接着继续拿出下一条数据，再进行处理。下面是 reduce Task 上 run 函数的部分代码，从这个函数可以看出整个 reduce 端的三个步骤。

```
// 复制阶段，从 map TaskTracker 处获取 map 输出
boolean isLocal = "local".equals(job.get("mapred.job.tracker", "local"));
if (!isLocal) {
    reduceCopier = new ReduceCopier(umbilical, job);
    if (!reduceCopier.fetchOutputs()) {
        .....
    }
}
// 复制阶段结束
copyPhase.complete();
// 合并阶段，将得到的 map 输出合并
setPhase(TaskStatus.Phase.SORT);
.....
// 合并阶段结束
sortPhase.complete();
// reduce 阶段
setPhase(TaskStatus.Phase.REDUCE);
.....
Reducer reducer = ReflectionUtils.newInstance(job.getReducerClass(), job);

// 逐条读出每一条记录，然后调用 Reducer 的 reduce 函数
while (values.more()) {
    reduceInputKeyCounter.increment(1);
    reducer.reduce(values.getKey(), values, collector, reporter); values.
    nextKey(); values.informReduceProgress();
}
values.informReduceProgress();
}
reducer.close();
out.close(reporter);
done(umbilical);
}
```

7.4.3 shuffle 过程的优化

熟悉了上面介绍的 shuffle 过程，可能有读者会说：这个 shuffle 过程不是最优的。是的，Hadoop 采用的 shuffle 过程并不是最优的。举个简单的例子，如果现在需要 Hadoop 集群完

成两个集合的并操作，事实上并操作只需要让两个集群中重复的元素在最后的結果中出现一次就可以了，并不要求結果的元素是按顺序排列的。但是如果使用 Hadoop 默认的 shuffle 过程，那么結果势必是排好序的，显然这个处理就不是必須的了。在这里简单介绍从 Hadoop 参数的配置出发来优化 shuffle 过程。在一个任务中，完成单位任务使用时间最多的一般都是 I/O 操作。在 map 端，主要就是 shuffle 阶段中缓冲区内容超过阈值后的写出操作。可以通过合理地设置 `ip.sort.*` 属性来减少这种情况下的写出次数，具体来说就是增加 `io.sort.mb` 的值。在 reduce 端，直接在复制 map 输出的时候将复制的結果放在内存中同样能够提升性能，这样可以让部分数据少做两次 I/O 操作（前提是留下的内存足够 reduce 任务执行）。所以在 reduce 函数的内存需求很小的情况下，将 `mapred.inmem.merge.threshold` 设置为 0，将 `mapred.job.reduce.input.buffer.percent` 设置为 1.0（或者一个更低的值）能够让 I/O 操作更少，提升 shuffle 的性能。

7.5 任务执行

本章前面详细介绍了 MapReduce 作业的执行流程，也简单介绍了基于 Hadoop 自身的一些参数优化。本节再介绍一些 Hadoop 在任务执行时的具体策略，让读者进一步了解 MapReduce 任务的执行细节，以便控制细节。

7.5.1 推测式执行

所谓推测式执行，是指当作业的所有任务都开始运行时，JobTracker 会统计所有任务的平均进度，如果某个任务所在的 TaskTracker 节点由于配置比较低或 CPU 负载过高，导致任务执行的速度比总体任务的平均速度要慢，此时 JobTracker 就会启动一个新的备份任务，原有任务和新任务哪个先执行完就把另外一个 kill 掉，这也是经常在 JobTracker 页面看到任务执行成功，但是总有些任务被 kill 的原因。

MapReduce 将待执行作业分割成一些小任务，然后并行运行这些任务，提高作业运行的效率，使作业的整体执行时间少于顺序执行的时间。但很明显，运行缓慢的任务（可能因为配置问题、硬件问题或 CPU 负载过高）将成为 MapReduce 的性能瓶颈。因为只要有一个运行缓慢的任务，整个作业的完成时间将被大大延长。这个时候就需要采用推测式执行来避免出现这种情况。当 JobTracker 检测到所有任务中存在运行过于缓慢的任务时，就会启动另一个相同的任务作为备份。原始任务和备份任务中只要有一个完成，另一个就会被中止。推测式执行的任务只有在在一个作业的所有任务开始执行之后才会启动，并且只针对运行一段时间之后，其执行速度慢于整个作业的平均执行速度的情况。

推测式执行在默认情况下是启用的。这种执行方式有一个很明显的缺陷：对于因为代码缺陷导致的任务执行速度过慢，它所启用的备份任务并不会解决问题。除此之外，由于推测式执行会启动新的任务，所以这种执行方式无可避免地会增加集群的负担。所以在利用 Hadoop 集群运行作业的时候可以根据具体情况选择开启或关闭推测式执行策略（通过设置

`mapred.map.tasks.speculative.execution` 和 `mapred.reduce.tasks.speculative.execution` 属性的值为 `map` 任务和 `reduce` 任务开启或关闭推测式执行策略)。

7.5.2 任务 JVM 重用

从本章图 7-1 中可以看出,不论是 `map` 任务还是 `reduce` 任务,都是在 `TaskTracker` 节点上的 Java 虚拟机 (JVM) 中运行的。当 `TaskTracker` 被分配一个任务时,就会在本地启动一个新的 Java 虚拟机来运行这个任务。对于有大量零碎输入文件的 `map` 任务而言,为每一个 `map` 任务启动一个 Java 虚拟机这种做法显然还有很大的改善空间。如果在一个非常短的任务结束之后让后续的任务重用此 Java 虚拟机,这样就可以省下新任务启动新的 Java 虚拟机的时间,这就是所谓的任务 JVM 重用。需要注意的是,虽然一个 `TaskTracker` 上可能会有多个任务在同时运行,但这些正在执行的任务都是在相互独立的 JVM 上的。`TaskTracker` 上的其他任务必须等待,因为即使启用 JVM 重用, JVM 也只能按顺序执行任务。

控制 JVM 重用的属性是 `mapred.job.reuse.jvm.num.tasks`。这个属性定义了单个 JVM 上运行任务的最大数目,默认情况下是 1,意味着每个 JVM 上运行一个任务。可以将这个属性设置为一个大于 1 的值来启用 JVM 重用,也可以将此属性设为 -1,表明共享此 JVM 的任务数目不受限制。

7.5.3 跳过坏记录

MapReduce 作业处理的数据集非常庞大,用户在基于 MapReduce 编写处理程序时可能并不会考虑到数据集中的每一种数据格式和字段(特别是某些坏的记录)。所以,用户代码在处理数据集中的某个特定记录时可能会崩溃。这个时候即使 MapReduce 有错误处理机制,但是因为存在这种代码缺陷,即使重新执行 4 次(默认的最大重新执行次数),这个任务仍然会失败,最终也会导致整个作业失败。所以针对这种由于坏数据而导致任务抛出的异常,重新运行任务是无济于事的。可是,如果想要在庞大的数据集中找出这个坏记录,然后在程序中添加相应的处理代码或直接除去这条坏记录显然也是很困难的一件事情,况且并不能保证没有其他坏记录。所以最好的办法就是在当前代码对应的任务执行期间,遇到坏记录时就直接跳过去(由于数据集巨大,忽略这种极少数的坏记录是可以接受的),然后继续执行,这就是 Hadoop 中的忽略模式(skipling 模式)。当忽略模式启动时,如果任务失败两次之后,它会将正在处理的记录告诉 `TaskTracker`,然后 `TaskTracker` 会重新运行该任务并在运行到先前任务报告的记录地方时直接跳过。从忽略模式的工作方式可以看出,忽略模式只能检测并忽略一个错误记录,因此这种机制仅适用于检测个别错误记录。如果增加任务尝试次数最大值(这由 `mapred.map.max.attempts` 和 `mapred.reduce.max.attempts` 两个属性决定),可以增加忽略模式能够检测并忽略的错误记录数目。默认情况下忽略模式是关闭的,可以使用 `SkipBadRedcord` 类单独为 `map` 和 `reduce` 任务启用它。

7.5.4 任务执行环境

Hadoop 能够为执行任务的 TaskTracker 提供执行所需要的环境信息。例如，map 任务可以知道自己所处理文件的名称，自己在作业任务群中的 ID 号等。JobTracker 分配任务给 TaskTracker 时，就会将作业的配置信息发送给 TaskTracker，TaskTracker 将此信息保存在本地。从本章前面的介绍中知道 TaskTracker 是在本节点单独的 JVM 上以子进程的形式执行 map 任务或 reduce 任务的。所以启动 map 或 reduce task 时，会直接从父 TaskTracker 处继承任务的执行环境。图 7-2 列出了每个 task 执行时使用的本地参数（从作业配置中获取，返回给 task 的是配置信息）。

名称	类型	描述
mapred.job.id	String	job id
mapred.jar	String	job目录下job.jar的位置
job.local.dir	String	job指定的共享存储空间
mapred.tip.id	String	task id
mapred.task.id	String	task尝试id
mapred.task.is.map	boolean	是否是map task
mapred.task.partition	int	task在job中的id
map.input.file	String	map读取的文件名
map.input.start	long	map输入的数据块的起始位置偏移
map.input.length	long	map输入的数据块的字节数
mapred.work.output.dir	String	task临时输出目录

图 7-2 task 的本地参数表

当 Job 启动时，TaskTracker 会根据配置文件创建 Job 和本地缓存。TaskTracker 的本地目录是 `${mapred.local.dir}/taskTracker/`。在这个目录下有两个子目录：一个是作业的分布式缓存目录，路径是在本地目录后面加上 `archive/`；一个是本地 Job 目录，路径是在本地目录后面加上 `jobcache/$jobid/`，在这个目录下保存了 Job 执行的共享目录（各个任务可以使用这个空间作为暂存空间，用于任务之间的文件共享，此目录通过 `job.local.dir` 参数暴露给用户）、存放 Jar 包的目录（保存作业的 jar 文件和展开的 jar 文件）、一个 xml 文件（此 xml 文件是本地通用的作业配置文件）和根据任务 ID 分配的任务目录（每个任务都有一个这样的目录，目录中包含本地化的任务作业配置文件，存放中间结果的输出文件目录、任务当前工作目录和任务临时目录）。

关于任务的输出文件需要注意的是，应该确保同一个任务的多个实例不会尝试向同一个文件进行写操作。因为这可能会存在两个问题，第一个问题是，如果任务失败并被重试，那么会先删除第一个任务的旧文件。第二个问题是，在推测式执行的情况下同一任务的两个实例会向同一个文件进行写操作。Hadoop 通过将输出写到任务的临时文件夹来解决上面的两个问题。这个临时目录是 `{mapred.out.put.dir}/_temporary/${mapred.task.id}`。如果任务执行成功，目录的内容（任务输出）就会被复制到此作业的输出目录（`${mapred.out.put.dir}`）。因此，如果一个任务失败并重试，第一个任务尝试的部分输出就会被消除。同时推测式执行时的备份任务和原始任务位于不同的工作目录，它们的临时输出文件夹并不相同，只有先完成的任务才会把其工作目录中的输出内容传到输出目录中，而另外一个任务的工作目录就会


被丢弃。

7.6 小结

本章以 MapReduce 程序中的 `JobClient.runJob(conf)` 开始，给出了 MapReduce 执行的流程图，并分析了流程图中的四个核心实体，结合实际代码介绍了 MapReduce 执行的详细流程。MapReduce 的执行流程简单概括如下：用户作业执行 `JobClient.runJob(conf)` 代码会在 Hadoop 集群上将其启动。启动之后 `JobClient` 实例会向 `JobTracker` 获取 `JobId`，而且客户端会将作业执行需要的作业资源复制到 HDFS 上，然后将作业提交给 `JobTracker`。`JobTracker` 在本地初始化作业，再从 HDFS 作业资源中获取作业输入的分割信息，根据这些信息 `JobTracker` 将作业分割成了多个任务，然后分配给在 `JobTracker` 心跳通信中请求任务的 `TaskTracker`。`TaskTracker` 接收到新的任务之后会先从 HDFS 上获取作业资源，包括作业配置信息和本作业分片的输入，然后在本地启动一个 JVM 并执行任务。任务结束之后将结果写回 HDFS。

介绍完 MapReduce 作业的详细流程后，本章还重点介绍了 MapReduce 中采用的两种机制，分别是错误处理机制和作业调度机制。在错误处理机制中，如果遇到硬件故障，MapReduce 会将故障节点上的任务分配给其他节点处理。如果遇到任务失败，则会重新执行。在作业调度机制中，主要介绍了公平调度器。这种调度策略能够按照提交作业的用户数目将资源公平地分到用户的作业池中，以达到用户公平共享整个集群的目的。

本章最后介绍了 MapReduce 中两个流程的细节，分别是 shuffle 和任务执行。在 shuffle 中，从代码入手介绍了 map 端和 reduce 端的 shuffle 过程及 shuffle 的优化。shuffle 的过程可以概括为：在 map 端，当缓冲区内容达到阈值时 map 写出内容。写出时按照 key 值对数据排序，再按照划分将数据写入文件，然后进行 merge 并将结果交给 reduce。在 reduce 端，`TaskTracker` 先从执行 map 的 `TaskTracker` 节点上复制 map 输出，然后对排序合并，最后进行 reduce 处理。关于任务执行则主要介绍了三个任务执行的细节，分别是推测式执行、JVM 重用和执行环境。推测式执行是指 `JobTracker` 在作业执行过程中，发现某个作业执行速度过慢，为了不影响整个作业的完成进度，会启动和这个作业完全相同的备份作业让 `TaskTracker` 执行，最后保留二者中较快完成的结果。JVM 重用主要是针对比较零碎的任务，对于新任务不是启动新的 JVM，而是在先前任务执行完毕的 JVM 上直接执行，这样节省了 JVM 启动的时间。在任务执行环境中主要介绍了任务执行参数的内容和任务目录结构，以及任务临时文件夹的使用情况。



第 8 章

Hadoop I/O 操作

本章内容

- ☐ I/O 操作中的数据检查
- ☐ 数据的压缩
- ☐ 数据的 I/O 中序列化操作
- ☐ 针对 MapReduce 的文件类
- ☐ 小结

资源分享

PDG

Hadoop 提供了如下一些与 I/O 相关的类：

- ❑ org.apache.hadoop.io
- ❑ org.apache.hadoop.io.compress
- ❑ org.apache.hadoop.io.file.tfile
- ❑ org.apache.hadoop.io.serializer
- ❑ org.apache.hadoop.io.serializer.avro

除了 org.apache.hadoop.io.serializer.avro 用于为 Avro（与 Hadoop 相关的 Apache 的另一个顶级项目）提供数据序列化操作外，其余都是用于 Hadoop 的 I/O 操作。

除此以外，部分 fs 类中的内容也与本章有关，所以本章也会提及一些，不过大都是一些通用的东西，由于对 HDFS 的介绍不是本章的重点，因此不再详述。

可以说，Hadoop 的 I/O 由传统的 I/O 操作而来，但是又有些不同。第一，在我们常见的计算机系统中，数据是集成的，无论多少电影、音乐，或者 Word 文档，它只会存在于一台主机中。而 Hadoop 则不同，Hadoop 系统中的数据经常是分散在多个计算机系统中的；第二，一般而言，传统计算机系统中的数据量相对较小，大多在 GB 级别。而 Hadoop 处理的数据经常是 PB 级别的。

变化就会带来问题，这两个变化带给我们的问题就是 Hadoop 的 I/O 操作不仅要考虑本地主机的 I/O 操作成本，还要考虑数据在不同主机之间的传输成本。同时 Hadoop 的数据寻址方式也要改变，才能应对庞大数据带来的寻址压力。

虽说 Hadoop 的 I/O 操作与传统方式已经有了一些变化，但是仍未脱离传统的数据 I/O 操作，因此如果你熟悉传统的 I/O 操作，你会发现本章的内容非常简单。

8.1 I/O 操作中的数据检查

Apache 的 Hadoop 官网上有一个名为 Sort900 的具体的 Hadoop 配置实例，所谓 Sort900 就是在 900 台主机上对 9TB 的数据进行排序。一般而言，在 Hadoop 集群的实际应用中，主机的数目是很大的，Sort900 使用了 900 台主机，而淘宝目前则使用了 1100 台主机来存储它们的数据（据说计划扩充到 1500 台）。在这么多的主机同时运行时，你会发现主机损坏是非常常见的，这就会涉及很多程序上的预处理了。对于本章而言，也就体现于在 Hadoop 中进行数据完整性检查的重要性上。

校验和方式是检查数据完整性的重要方式。一般会通过对比新旧校验和来确定数据情况，如果两者不同则说明数据已经损坏。比如，在传输数据前生成了一个校验和，将数据传输到目的主机时再次计算校验和，如果两次的校验和不同，则说明数据已经损坏。或者在系统启动时计算校验和，如果其值和硬盘上已经存在的校验和不同，那么也说明数据已经损坏。校验和不能恢复数据，只能检测错误。

Hadoop 采用 CRC-32（Cyclic Redundancy Check——循环冗余校验，其中的 32 指生成的校验和是 32 位的）的方式检查数据完整性。这是一种非常常见的校验和验证方式，检错

能力强, 开销小, 易于实现。如果读者有兴趣可以自行查阅资料了解。

因为 Hadoop 采用 HDFS 作为默认的文件系统, 因此我们需要讨论两方面的数据完整性:

- 本地文件系统的数据完整性;
- HDFS 的数据完整性。

1. 对本地文件 I/O 的检查

在 Hadoop 中, 本地文件系统的数据完整性由客户端负责。重点是在存储和读取文件时进行校验和的处理。

具体做法是, 每当 Hadoop 创建文件 a 时, Hadoop 就会同时在同一文件夹下创建隐藏文件 .a.crc, 这个文件记录了文件 a 的校验和。针对数据文件的大小, 每 512 个字节 Hadoop 就会生成一个 32 位的校验和 (4 字节), 你可以在 src/core/core-default.xml 中通过修改 io.bytes.per.checksum 的大小来修改每个校验和所针对的文件大小。如下所示:

```
<property>
<name>io.bytes.per.checksum</name>
<value>512</value>
<description>The number of bytes per checksum. Must not be larger than io.file.
    buffer.size.</description>
</property>
```

一般来说, 主流的文件系统都能在一定程度上保证数据的完整性, 因此有可能你并不需要 Hadoop 的这部分功能。如果你不需要这部分功能, 你可以通过修改文件 src/core/core-default.xml 中 fs.file.impl 的值来禁用校验和机制, 如下所示:

```
<property>
<name>fs.file.impl</name>
<value>org.apache.hadoop.fs.LocalFileSystem</value>
<description>The FileSystem for file: uris.</description>
</property>
```

把值修改为 org.apache.hadoop.fs.RawLocalFileSystem 即可禁用校验和机制。

如果你只想在程序中对某些读取禁用校验和检验, 那么你可以声明 RawLocalFileSystem 实例。例如:

```
FileSystem fs = new RawFileSystem();
Fs.initialize( null , conf );
```

在 Hadoop 中, 校验和系统单独为一类——org.apache.hadoop.fs.ChecksumFileSystem, 当你需要校验和机制时, 你可以很方便地调用它来为你服务。

引用方法为:

```
FileSystem rawFS = ...;
FileSystem checksumFS = new ChecksumFileSystem(rawFS);
```

事实上, org.apache.hadoop.fs.ChecksumFileSystem 是 org.apache.hadoop.fs.FileSystem 子类

的子类，其继承关系如下：

```
java.lang.Object
  -org.apache.hadoop.conf.Configured
    -org.apache.hadoop.fs.FileSystem
      -org.apache.hadoop.fs.FilterFileSystem
        -org.apache.hadoop.fs.ChecksumFileSystem
          -org.apache.hadoop.fs.LocalFileSystem
```

如果你对这些类的作用感兴趣，你可以查阅 Hadoop 的 app 文档，地址为 <http://hadoop.apache.org/common/docs/current/api/index.html>。读取文件时，如果 ChecksumFileSystem 检测到错误，便会调用 reportChecksumFailure。这是一个布尔类型的函数，此时，LocalFileSystem 会把这些问题文件及其校验和一起移动到同一台主机的次级目录下，命名为 bad_files。一般而言，使用者需要经常处理这些文件。

2. 对 HDFS 的 I/O 数据进行检查

一般来说，HDFS 会在三种情况下检验校验和：

(1) DataNode 接收数据后，存储数据前

要了解这种情况，读者先要了解 DataNode 一般会在什么时候接收数据。它接收数据一般有两种情况：一是用户从客户端上传数据；二是 DataNode 从其他 DataNode 上接收数据。一般来说，客户端往往也是 DataNode，不过有时候客户端仅仅是客户端而已，并不是 Hadoop 集群中的节点。当客户端上传数据时，Hadoop 会根据预定规则形成一条数据管线。图 8-1 就是一个典型的副本管线（数据备份为 3）。原数据是数据 0，数据 1、数据 2、数据 3 是备份。

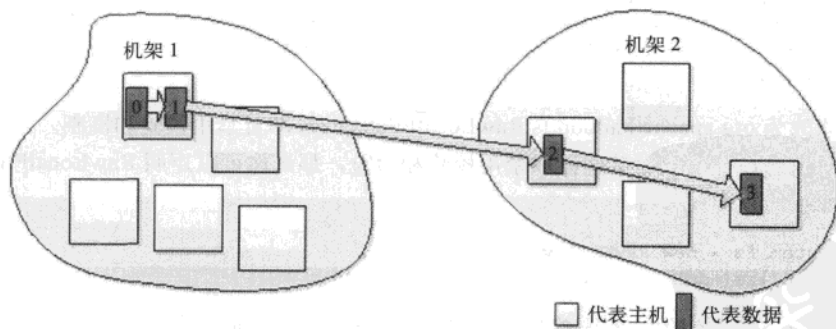


图 8-1 数据管线及数据备份流程图

让数据按管线流动以完成数据的上传及备份过程，图 8-1 中的顺序就是先在客户端这个节点上保存数据（在这幅图上，客户端也是 Hadoop 集群中的一个节点）。注意这个流动的过程，备份 1 在接收数据的同时也会把接收到的数据发送给备份 2 所在的机器，因此如果过程执行顺利，三个备份形成的时间相差不多（相对依次备份而言）。这里面涉及一个负载均衡

的问题，不过这个问题不是本章的重点，这里不再详述，我们在这里只关心数据完整性的问题。Hadoop 不会在数据每流动到一个 DataNode 时都检查校验和，它只会在数据流动到最后一个节点时才检验校验和。也就是说 Hadoop 会在备份 3 所在的 DataNode 接收完数据后检查校验和。这就是从客户端上传数据时 Hadoop 对数据完整性检测进行的相关处理。

DataNode 从其他 DataNode 上接收数据时也是同样的处理过程。

(2) 客户端读取 DataNode 上的数据时

Hadoop 会在客户端读取 DataNode 上的数据时检查校验和。

(3) DataNode 后台守护进程的定期检测

DataNode 会在后台运行 DataBlockScanner，这个程序会定期检测此 DataNode 上的所有数据块。

3. 数据恢复策略

当 Hadoop 发现某数据块已经失效时，它可以采用复制完整的备份方式来恢复数据。

与读取本地文件的情况相同，用户也可以使用命令来禁用检验和。有两种方法可以达到这个目的。

一个是在使用 open() 读取文件前，设置 FileSystem 中的 setVerifyChecksum 值为 false。

```
FileSystem fs = new FileSystem();
fs.setVerifyChecksum(false);
```

另一个是使用 shell 命令，比如 get 命令和 copyToLocal 命令。

get 命令的使用方法如下所示：

```
hadoop fs -get [-ignoreCrc] [-crc] <src> <localdst>
```

举个例子：

```
hadoop fs -get -ignoreCrc input ~/Desktop/
```

get 命令会复制文件到本地文件系统。可用 -ignoreCrc 选项复制 CRC 校验失败的文件，或者使用 -crc 选项复制文件，以及 CRC 信息。

copyToLocal 命令的使用方法如下所示：

```
hadoop fs -copyToLocal [-ignorecrc] [-crc] URI <localdst>
```

再举个例子：

```
hadoop fs -copyToLocal -ignoreCrc input ~/Desktop
```

除了要限定目标路径是一个本地文件外，其他和 get 命令类似。

禁用校验和的最主要目的并不是节约时间，用于检验校验和的开销一般情况都是可以接受的，禁用校验和的主要原因是，如果你不禁用校验和，你无法下载那些已经损坏的文件来查看是否可以挽救，而有时候即使是只能挽救一小部分文件也是很值得的。

8.2 数据的压缩

对于任何大容量的分布式存储系统而言，文件压缩都是必须的，文件压缩带来了两个好处：

- 减少了文件所需的存储空间；
- 加快了文件在网络上或磁盘间的传输速度。

Hadoop 关于文件压缩的代码几乎都在 package org.apache.hadoop.io.compress 中。本节的内容将会主要围绕这一部分展开。

8.2.1 Hadoop 对压缩工具的选择

有许多压缩格式和压缩算法是可以应用到 Hadoop 中的，但是不同的算法也有各自不同的特点。表 8-1 是一些 Hadoop 中使用的压缩算法。

表 8-1 压缩格式及编码解码器

压缩格式	Hadoop 压缩编码 / 解码器
DEFLATE	org.apache.hadoop.io.compress.DefaultCodec
Gzip	org.apache.hadoop.io.compress.GzipCodec
bzip2	org.apache.hadoop.io.compress.BZip2Codec
LZO	com.hadoop.compression.lzo.LzopCodec

LZO 是基于 GPL 许可的，不能通过 Apache 发布许可，因此 LZO 的 Hadoop 编码 / 解码器必须单独下载，下载地址为 <http://code.google.com/p/hadoop-gpl-compression/>。

表 8-2 是它们的特点。

表 8-2 压缩格式及其特点

压缩格式	工 具	算 法	文件扩展名	多文件	可分割性
DEFLATE*	无	DEFLATE	.deflate	否	否
Gzip	gzip	DEFLATE	.gz	否	否
bzip2	bzip2	bzip2-	.bz2	否	是
LZO	Lzop	LZO	.lzo	否	否

Hadoop 目前尚不支持 Zip 压缩。

压缩一般都是在时间和空间上的一种权衡，一般来说，更长的压缩时间会节省更多的空间。不同的压缩算法之间也会有一定的区别，而同样的压缩算法在压缩不同类型的文件时表现也不同。在 jeff 的实验比较报告中，包含了面对不同文件时在各种要求（最佳压缩、最快速度等）下的最佳压缩工具。如果读者感兴趣可以自行查阅，地址为 <http://compression.ca/act/act-summary.html>（这个地址是总体评价，此网站还有不同压缩工具面对不同类型文件时的具体表现）。

8.2.2 压缩分割和输入分割

压缩分割和输入分割是很重要的内容，比如，如果读者需要处理经 LZO 压缩后的 5GB 大小的文件，按前文介绍过的分割方式，Hadoop 会将其分割为 80 块（每块 64MB，这是默认值，可以根据需要修改）。但是这是没有意义的，因为程序无法分开读取每块的内容，那么也就无法创建多个 map 程序分别来处理每块内容。因为在这种情况下，Hadoop 不会分割存储 LZO 压缩的文件。

而 bzip2 的情况就不一样了，因为 bzip2 支持文件分割，你可以分开读取每块内容并分别处理之，因此 bzip2 压缩的文件可分割存储。

8.2.3 在 MapReduce 程序中使用压缩

在 MapReduce 程序中使用压缩非常简单，你只需在它进行 Job 配置时配置好 conf 就可以了。

设置 map 处理后数据的压缩代码示例如下：

```
JobConf conf = new Jobconf();
conf.setBoolean("mapred.compress.map.output", true);
```

设置 output 输出压缩的代码示例如下：

```
JobConf conf = new Jobconf();
conf.setBoolean("mapred.output.compress", true);
conf.setClass("mapred.output.compression.codec", GzipCodec.class, CompressionCodec.class);
```

对于一般情况而言，压缩总是好的，无论是对最终结果的压缩还是对 map 处理后的中间数据进行压缩。对于 map 而言，它处理后的数据都要输出到硬盘上并经过网络传输的，使用数据压缩一般都会加快这一过程。最终结果的压缩不单会加快数据存储的速度，也会节省硬盘空间。

下面我们做一个实验来看看在 MapReduce 中使用压缩与不使用压缩时的效率差别。

先来叙述一下我们的实验环境，这是由六台主机组成的一个小集群（一台 master，三台 slave）。输入文件为未压缩的大约为 300MB 的文件，它是由随机的英文符号字符串组成的，每个字符串都是 6 位的英文字母（大小写被认为是不同的），形如 “AdEfr”，以空格隔开，每 50 个一行，共 50 000 000 个字符串。对这个文件进行 WordCount。Map 的输出压缩采用默认的压缩算法，output 的输出采用 Gzip 压缩方法，我们关注的内容是程序执行的速度差别。

执行压缩操作的 WordCount 程序与基本的 WordCount 程序相似，只是在 conf 设置时写入了以下几行代码：

```
conf.setBoolean("mapred.compress.map.output", true);
conf.setBoolean("mapred.output.compress", true);
conf.setIfUnset("mapred.output.compression.type", "BLOCK");
conf.setClass("mapred.output.compression.codec", GzipCodec.class,
    CompressionCodec.class);
```

下面分别执行编译打包两个程序，在运行时用 `time` 命令记录程序的执行时间，如下所示：

```
time bin/hadoop jar WordCount.jar WordCount XWTInput xwtOutput

real          12m41.308s

time bin/hadoop jar CompressionWordCount.jar CompressionWordCount XWTInput
xwtOutput2

real          8m9.714s
```

`CompressionWordCount.jar` 是带压缩的 `WordCount` 程序的打包，从上面可以看出执行压缩的程序要比不执行压缩的程序快 4 分钟，或者说，在这个实验环境下，使用压缩会使 `WordCount` 的效率提高 1/3。

8.3 数据的 I/O 中序列化操作

序列化是将对象转化为字节流的方法，或者说用字节流描述对象的方法。与序列化相对的是反序列化，反序列化是将字节流转化为对象的方法。序列化有两个目的：

- ❑ 进程间通信；
 - ❑ 数据持久性存储。
- Hadoop 采用 RPC 来实现进程间通信。一般而言，RPC 的序列化机制有以下特点：
- ❑ 紧凑：紧凑的格式可以充分利用带宽，加快传输速度；
 - ❑ 快速：能减少序列化和反序列化的开销，这会有效减少进程间通信的时间；
 - ❑ 可扩展：可以逐步改变，是客户端与服务器端直接相关的。例如，可以随时加入一个新的参数方法调用；
 - ❑ 互操作性：支持不同语言编写的客户端与服务器端交换数据。

同时 Hadoop 也希望数据持久性存储同样具有以上这些优点，因此它的数据序列化机制也是依照以上这些目的而设计的（或者说是希望设计成这样）。

在 Hadoop 中，序列化处于核心地位。因为无论是存储文件还是在计算中传输数据，都需要执行序列化的过程。序列化与反序列化的速度、序列化后的数据大小等都会影响数据传输的速度，以致影响计算的效率。正是因为这些原因，Hadoop 并没有采用 Java 提供的序列化机制（Java Object Serialization），而是自己重新写了一个序列化机制 `Writable`。Writable 具有紧凑、快速的优点（但不易扩展，也不利于不同语言的互操作），同时也允许你对自己定义的类加入序列化与反序列化方法，而且很方便。

8.3.1 Writable 类

Writable 是 Hadoop 的核心，Hadoop 通过它定义了 Hadoop 中基本的数据类型及其操作。一般来说，无论是上传下载数据还是运行 MapReduce 程序，你无时无刻不需要使用 Writable

类，因此 Hadoop 中具有庞大的一类 Writable 类（见图 8-2），不过 Writable 类本身却很简单。

Writable 类中只定义了两个方法：

```
void write(DataOutput out)
    throws IOException
void readFields(DataInput in)
    throws IOException
```

Hadoop 还有很多其他的 Writable 类。比如 WritableComparable、ArrayWritable、TwoDArrayWritable 及 AbstractMapWritable，它们直接继承自 Writable 类。还有一些类，如 BooleanWritable、ByteWritable 等，它们不是直接继承于 Writable 类，而是继承自 WritableComparable 类。Hadoop 的基本数据类型就是由这些类构成的。这些类构成了以下的层次关系（如图 8-2 所示）。

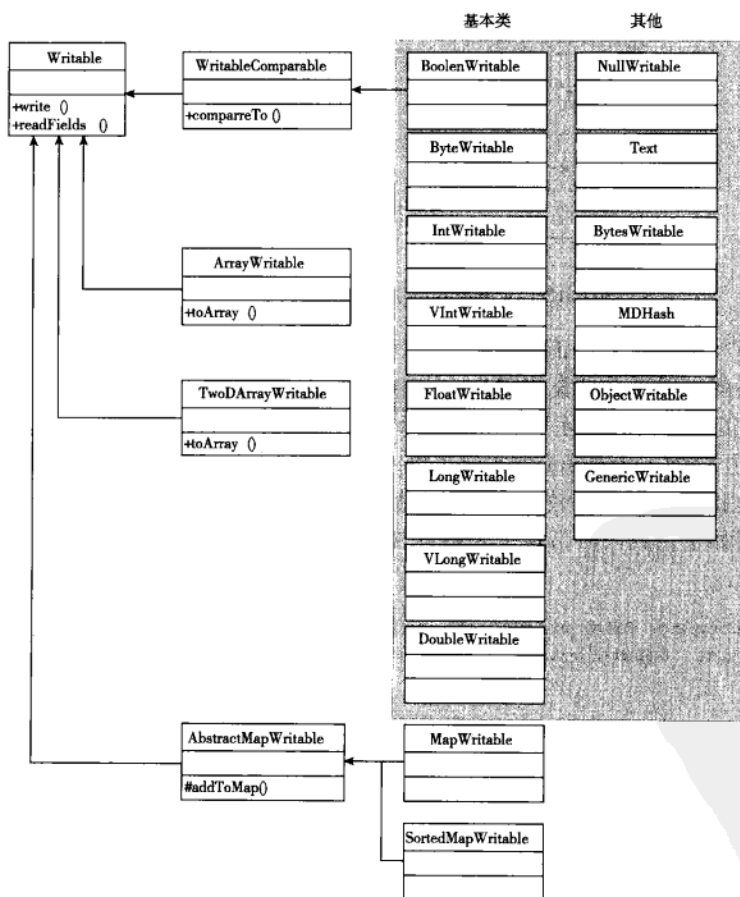


图 8-2 Writable 类层次关系图

1. Hadoop 的比较器

WritableComparable 是 Hadoop 中非常重要的接口类。它继承自 org.apache.hadoop.io.Writable 类和 java.lang.Comparable 类。WritableComparator 是 WritableComparable 的比较器，它是 RawComparator 针对 WritableComparable 类的一个通用实现，而 RawComparator 继承自 java.util.Comparator，它们之间的关系如图 8-3 所示。

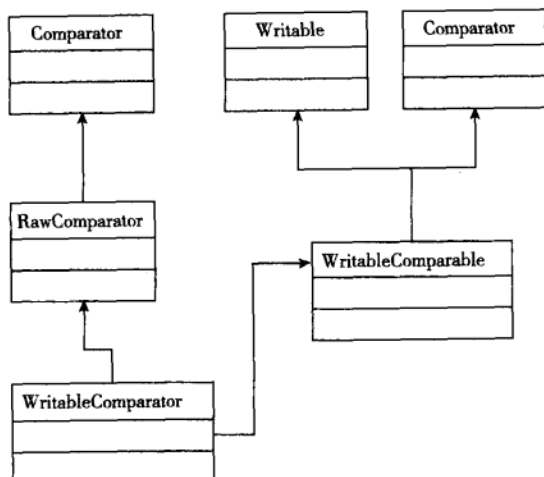


图 8-3 WritableComparable 和 WritableComparator 类层次关系图

这两个类对 MapReduce 而言至关重要，大家都知道，MapReduce 执行时，Reducer（执行 reduce 任务的机器）会搜集相同 key 值的 key/value 对，并且在 reduce 之前会有一个排序过程，这些键值的比较都是对 WritableComparator 类型进行的比较。

Hadoop 在 RawComparator 中实现了对未反序列化对象的读取。这样做的好处是，可以不必创建对象就比较想要比较的内容（多是 key 值），从而省去了创建对象的开销。例如，读者可以使用如下函数，对指定了开始位置（s1 和 s2）及固定长度（l1 和 l2）的数组进行比较：

```
public interface RawComparator<T> extends Comparator<T> {
    public int compare(byte[] b1, int s1, int l1, byte[] b2, int s2, int l2);
}
```

WritableComparator 是 RawComparator 的子类，在这里，添加了一个默认的对象进行反序列化，并调用了比较函数 compare() 进行比较。下面是在 WritableComparator 中对固定字节反序列化的执行情况，以及比较的实现过程：

```
public int compare(byte[] b1, int s1, int l1, byte[] b2, int s2, int l2) {
    try {
        buffer.reset(b1, s1, l1);           // parse key1
```

```

key1.readFields(buffer);

buffer.reset(b2, s2, l2);           // parse key2
key2.readFields(buffer);

} catch (IOException e) {
    throw new RuntimeException(e);
}
return compare(key1, key2);          // compare them
}

```

2. Writable 类中的数据类型

(1) 基本类

Writable 中封装有很多 Java 的基本类，如表 8-3 所示。

表 8-3 Writable 中的 Java 基本类

Java 基本类型	Writable 中的类型	序列化后字节数
boolean	BooleanWritable	1
byte	ByteWritable	1
int	IntWritable	4
	VIntWritable	1-5
float	FloatWritable	4
long	LongWritable	8
	VLongWritable	1-9
double	DoubleWritable	8

其中最简单的要数 Hadoop 中对 Boolean 的实现，如下所示：

```

import java.io.*;

public class BooleanWritable implements WritableComparable {
    private boolean value;
    public BooleanWritable() {};
    public BooleanWritable(boolean value) {
        set(value);
    }
    public void set(boolean value) {
        this.value = value;
    }
    public boolean get() {
        return value;
    }
    public void readFields(DataInput in) throws IOException {
        value = in.readBoolean();
    }
    public void write(DataOutput out) throws IOException {

```

```

        out.writeBoolean(value);
    }

    public boolean equals(Object o) {
        if (!(o instanceof BooleanWritable)) {
            return false;
        }
        BooleanWritable other = (BooleanWritable) o;
        return this.value == other.value;
    }

    public int hashCode() {
        return value ? 0 : 1;
    }

    public int compareTo(Object o) {
        boolean a = this.value;
        boolean b = ((BooleanWritable) o).value;
        return ((a == b) ? 0 : (a == false) ? -1 : 1);
    }

    public String toString() {
        return Boolean.toString(get());
    }

    public static class Comparator extends WritableComparator {
        public Comparator() {
            super(BooleanWritable.class);
        }

        public int compare(byte[] b1, int s1, int l1,
                           byte[] b2, int s2, int l2) {
            boolean a = (readInt(b1, s1) == 1) ? true : false;
            boolean b = (readInt(b2, s2) == 1) ? true : false;
            return ((a == b) ? 0 : (a == false) ? -1 : 1);
        }
    }

    static {
        WritableComparator.define(BooleanWritable.class, new Comparator());
    }
}

```

可以看到 Hadoop 直接将 boolean 写入到字节流 (out.writeBoolean (value)) 中了, 并没有采用 Java 的序列化机制。同时, 除了构造函数、set() 函数、get() 函数等外, Hadoop 还定义了三个用于比较的函数: equals()、compareTo()、compare()。前两个很简单, 第三个就是前文中重点介绍的比较器。Hadoop 中封装定义的其他 Java 基本数据类型 (如 Boolean、byte、int、float、long、double) 都是相似的。

如果读者对 Java 流处理比较了解的话可能会知道, Java 流处理中并没有 DataOutput.writeVInt()。实际上, 这是 Hadoop 自己定义的变长类型 (VInt, VLong), 而且 VInt 和 VLong 的处理方式实际上是一样的。

```

public static void writeVInt(DataOutput stream, int i) throws IOException {
    writeVLong(stream, i);
}

```

Hadoop 对 VLong 类型的处理方法如下。

```
public static void writeVLong(DataOutput stream, long i) throws IOException {
    if (i >= -112 && i <= 127) {
        stream.writeByte((byte)i);
        return;
    }
    int len = -112;
    if (i < 0) {
        i ^= -1L; // take one's complement
        len = -120;
    }
    long tmp = i;
    while (tmp != 0) {
        tmp = tmp >> 8;
        len--;
    }
    stream.writeByte((byte)len);
    len = (len < -120) ? -(len + 120) : -(len + 112);
    for (int idx = len; idx != 0; idx--) {
        int shiftbits = (idx - 1) * 8;
        long mask = 0xFFL << shiftbits;
        stream.writeByte((byte)((i & mask) >> shiftbits));
    }
}
```

上面代码的意思是如果数值较小（在 -112 和 127 之间），那么就直接将这个数值写入数据流内（stream.writeByte((byte)i)）。如果不是，则先用 len 表示字节长度与正负，并写入数据流中，然后在其后写入这个数值。

(2) 其他类

下面将按照先易后难的顺序一一讲解。

❑ NullWritable

这是一个占位符，它的序列化长度为零，没有数值从流中读出或是写入流中。

```
public void readFields(DataInput in) throws IOException {}
public void write(DataOutput out) throws IOException {}
```

在任何编程语言或编程框架中，占位符都是很有用的，这个类型不可以和其他类型比较，在 MapReduce，你可以将任何键或值设为空值。

❑ BytesWritable

BytesWritable 是一个二进制数据数组的封装。它对输出流的处理如下所示：

```
public BytesWritable(byte[] bytes) {
    this.bytes = bytes;
    this.size = bytes.length;
}
public void write(DataOutput out) throws IOException {
```

```

    out.writeInt(size);
    out.write(bytes, 0, size);
}

```

可以看到，它首先会把这个二进制数据数组的长度写入输入流中，这个长度一般是在声明时所获得的二进制数据数组的实际长度。当然这个值也可以人为设定。如果你要把长度为 3，位置为 129 的字节数组序列化，根据程序可知，结果应为：

```
Size=00000003 bytes[]={ (01), (02), (09) }
```

数据流中的值就是：

```
00000003010209
```

□ Text

这可能是这几个自定义类型中相对复杂的一个了。实际上，这是 Hadoop 中对 string 类型的重写，但是又与其有一些不同。Text 使用标准的 UTF-8 编码，同时 Hadoop 使用变长类型 VInt 来存储字符串，其存储上限是 2GB。

Text 类型与 String 类型的主要差别如下：

- String 的长度定义为 String 包含的字符个数；Text 的长度定义为 UTF-8 编码的字节数。
- String 内的 indexOf() 方法返回的是 char 类型字符的索引，比如字符串 (1234)，字符 3 的位置就是 2（字符 1 的位置是字符 0）；而 Text 的 find() 方法返回的是字节偏移量。
- String 的 charAt() 方法返回的是指定位置的 char 字符；而 Text 的 charAT() 方法需要指定偏移量。

另外，在 Text 内定义了一个方法 toString()，它用于将 Text 类型转化为 String 类型。

看如下这个例子：

```

import java.io.*;
import org.apache.hadoop.io.*;

public class MyMapre {
    public static void strings(){
        String s="\u0041\u00DF\u6771\uD801\uDC00";
        System.out.println(s.length());
        System.out.println(s.indexOf("\u0041"));
        System.out.println(s.indexOf("\u00DF"));
        System.out.println(s.indexOf("\u6771"));
        System.out.println(s.indexOf("\uD801\uDC00"));
    }
    public static void texts(){
        Text t = new Text("\u0041\u00DF\u6771\uD801\uDC00");
        System.out.println(t.getLength());
        System.out.println(t.find("\u0041"));
        System.out.println(t.find("\u00DF"));
        System.out.println(t.find("\u6771"));
        System.out.println(t.find("\uD801\uDC00"));
    }
}

```




```

    }
    public static void main(String args[]){
        strings();
        texts();
    }
}

```

输出结果为：

```

5
0
1
2
3
10
0
1
3
6

```

上面这个例子可以验证前面所列的那些差别。

❑ ObjectWritable

ObjectWritable 是一种多类型的封装。可以适用于 Java 的基本类型、字符串等。不过，这并不是一个好方法，因为 Java 在每次被序列化时，都要写入被封装类型的类名中。但是如果类型过多，使用静态数组难以表示时，采用这个类仍是不错的做法。

❑ ArrayWritable 和 TwoDArrayWritable

ArrayWritable 和 TwoDArrayWritable，顾名思义，是针对数组和二维数组构建的数据类型。这两个类型声明的变量需要在使用时指定类型，因为 ArrayWritable 和 TwoDArrayWritable 并没有空值的构造函数。

```
ArrayWritable a = new ArrayWritable(IntWritable.class)
```

同样，在声明它们的子类时，你必须使用 super() 来指定 ArrayWritable 和 TwoDArrayWritable 的数据类型。

```

public class IntArrayWritable extends ArrayWritable{
    public IntArrayWritable(){
        super(IntWritable.class);
    }
}

```

一般地，ArrayWritable 和 TwoDArrayWritable 都有 set() 函数与 get() 函数，在将 Text 转化为 String 时，它们也都提供了一个转化函数 toArray()。但是它们没有提供比较器 comparator，这点需要注意。

❑ MapWritable 和 SortedMapWritable

MapWritable 和 SortedMapWritable 分别是 java.util.Map() 与 java.util.SortedMap() 的实现。

这两个实例是按照如下格式声明的：

```
private Map<Writable, Writable> instance;
private SortedMap<WritableComparable, Writable> instance;
```

你可以用 Hadoop 定义的 Writable 类型来填充 key 或 value，你也可以使用自己定义的 Writable 类型来填充。

在 java.util.Map() 和 java.util.SortedMap() 中定义的功能，如 getKey()、getValue()、keySet() 等，在这两个类中均有实现。Map 的使用也很简单，见如下程序，需要注意的是，不同 key 值对应的 value 数据类型可以不同：

```
import java.io.*;
import java.util.*;
import org.apache.hadoop.io.*;

public class MyMapre {
    public static void main(String args[]) throws IOException{
        MapWritable a = new MapWritable();
        a.put(new IntWritable(1), new Text("Hello"));
        a.put(new IntWritable(2), new Text("World"));

        MapWritable b = new MapWritable();
        WritableUtils.cloneInto(b, a);
        System.out.println(b.get(new IntWritable(1)));
        System.out.println(b.get(new IntWritable(2)));
    }
}
```

显示结果为：

```
Hello
World
```

8.3.2 实现自己的 Hadoop 数据类型

实现自己的 Hadoop 数据类型具有非常重要的意义，虽然 Hadoop 已经定义了很多有用的数据类型，但在实际应用中，你总是需要定义自己的数据类型以满足程序的需要。

我们定义一个简单的整数对 <LongWritable, LongWritable>，这个类可以用来记录文章中单词出现的位置，第一个 LongWritable 代表行数，第二个 LongWritable 代表它是该行的第几个单词。定义 NumPair，如下所示：

```
import java.io.*;
import org.apache.hadoop.io.*;

public class NumPair implements WritableComparable<NumPair> {
    private LongWritable line;
    private LongWritable location;
    public NumPair() {
        set(new LongWritable(0), new LongWritable(0));
    }
}
```

```

    }
    public void set(LongWritable first, LongWritable second)
    {
        this.line=first;
        this.location=second;
    }
    public NumPair(LongWritable first,LongWritable second){
        set(first,second);
    }
    public NumPair(int first, int second){
        set(new LongWritable(first),new LongWritable(second));
    }
    public LongWritable getLine(){
        return line;
    }
    public LongWritable getLocation(){
        return location;
    }
    @Override
    public void readFields(DataInput in) throws IOException
    {
        line.readFields(in);
        location.readFields(in);
    }
    @Override
    public void write(DataOutput out) throws IOException {
        line.write(out);
        location.write(out);
    }
    public boolean equals(NumPair o){
        if((this.line==o.line)&&(this.location==o.location))
            return true;
        return false;
    }
    @Override
    public int hashCode(){
        return line.hashCode()*13+location.hashCode();
    }
    @Override
    public int compareTo(NumPair o) {
        if((this.line==o.line)&&(this.location==o.location))
            return 0;
        return -1;
    }
}

```

8.4 针对 MapReduce 的文件类

Hadoop 定义了两种数据类型以适应 MapReduce 编程框架的需要。这两种类型非常重

要，其中 map 输出的中间结果就是由它们表示的。它们指的是 SequenceFile 和 MapFile。其中，MapFile 是经过排序并带有索引的 SequenceFile。

8.4.1 SequenceFile 类

SequenceFile 记录的是 key/value 对的列表，是序列化之后的二进制文件，因此是不能直接查看的，你可以通过如下命令来查看这个文件的内容。

```
hadoop fs -text MySequenceFile( 你的 SequenceFile 文件 )
```

Sequence 有三种不同类型的结构：

- ❑ 未压缩的 key/value 对；
- ❑ 记录压缩的 key/value 对（这种情况下只有 value 被压缩）；
- ❑ Block 压缩的 key/value 对（在这种情况下，key 与 value 被分别记录到块中并压缩）。

下面详细介绍它们的结构。

1. 未压缩和只压缩 value 的 SequenceFile 数据格式

未压缩和只压缩 value 的 SequenceFile 数据格式基本是相同的。

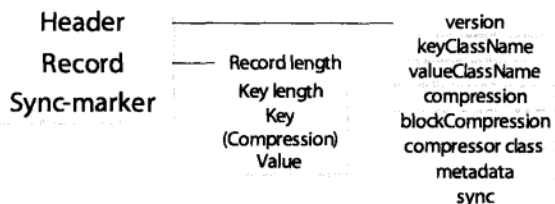


图 8-4 SequenceFile 数据格式（未压缩和 Record 压缩格式）

Header 是头，它记录的内容如图 8-4 所示，现在一一对其进行解释：

- ❑ version（版本号）：这是一个形如 SEQ4 或 SEQ5 的字节数组，一共占 4 字节；
- ❑ keyClassName（key 类名）和 valueClassName（value 类名）：这两个都是 string 类型，记录的是 key 和 value 的数据类型；
- ❑ compression（压缩）：这是一个布尔类型，它记录的是在这个文件中压缩是否启用；
- ❑ blockCompression（Block 压缩）：布尔类型，记录 Block 压缩是否启用；
- ❑ compressor class（压缩类）：这是 Hadoop 内封装的用于压缩 key 和 value 的代码；
- ❑ metadata（元数据）：用于记录文件的元数据，文件的元数据是一个 < 属性名，值 > 对的列表；
- ❑ Record：它是数据内容，其内容简单明了，相信读者看图就很容易明白；
- ❑ Sync-marker：它是一个标记，可以允许程序快速找到文件中随机的一个点。它可以使 MapReduce 程序更有效率地分割大文件。

需要注意的是，Sync-marker 会每隔几百个字节出现一次，因此最后的 SequenceFile 会

是形成如图 8-5 所示的序列文件。

```
Header Recorder Recorder Recorder Sync Recorder Recorder Sync
```

图 8-5 SequenceFile 数据存储示例

Sync 出现的位置取决于字节数，而不是间隔的 Recorder 的个数。

从上面的内容可以知道，未压缩与只压缩 value 的 SequenceFile 数据格式有两点不同，一是 compression（是否压缩）的值不同；二是 value 存储的数据是否经过了压缩不同。

2. Block 压缩的 SequenceFile 数据格式

Block 压缩的 SequenceFile 数据格式与上面两种格式也很相似，它们的头与其是一样的，同时也会标记一个 Sync-marker。不过它们的 recorder 格式是不同的，并且 Sync-marker 是标记在每个块前面的。下面是 Block 压缩的 SequenceFile 的 Recorder 格式（如图 8-6 所示）。

```
Compressed key-lengths block-size
Compressed key-lengths block
Compressed keys block-size
Compressed keys block
Compressed value-lengths block-size
Compressed value-lengths block
Compressed values block-size
Compressed values block
```

图 8-6 SequenceFile 数据格式 Recorder 部分（Block 压缩）

Block 压缩一次会压缩多个 Recorder，Recorder 在达到一个值时被记录，这个值是由 io.seqfile.compress.blocksize 定义的。Block 压缩的 SequenceFile 是形成如图 8-7 所示的序列文件。

```
Header Sync Recorder Sync Recorder Sync Recorder
```

图 8-7 SequenceFile 数据存储示例（Block 压缩）

我们可以通过编写程序生成读取 SequenceFile 文件来实践一下。

程序如代码清单 8-1 所示（注意这个程序生成的数据大概会有 150MB，需要的话可以减少循环次数以减少运行时间）：

代码清单 8-1 SequenceFileWriteFile.java

```
import java.io.IOException;
import java.net.URI;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.*;
import org.apache.hadoop.io.*;
```

```

public class SequenceFileWriteDemo {
    private static String[] myValue = {
        "hello world",
        "bye world",
        "hello hadoop",
        "bye hadoop"
    };
    public static void main(String[] args) throws IOException {
        String uri = "你想要生成的 SequenceFile 的位置";
        Configuration conf = new Configuration();
        FileSystem fs = FileSystem.get(URI.create(uri), conf);
        Path path = new Path(uri);
        IntWritable key = new IntWritable();
        Text value = new Text();
        SequenceFile.Writer writer = null;
        try {
            writer = SequenceFile.createWriter(fs, conf, path, key.getClass(),
                value.getClass());
            for (int i = 0; i < 5000000; i++) {
                key.set(5000000 - i);
                value.set(myValue[i%myValue.length]);
                writer.append(key, value);
            }
        } finally {
            IOUtils.closeStream(writer);
        }
    }
}

```

程序结果是生成了一个 SequenceFile 文件，你可以使用前文提到的命令，Hadoop fs-text 你的 SequenceFile 文件名，来查看这个文件。因为内容太多只展示一部分，其内容如下：

```

5000000    hello world
4999999    bye world
4999998    hello hadoop
4999997    bye hadoop
4999996    hello world
4999995    bye world
4999994    hello hadoop
4999993    bye hadoop
4999992    hello world
4999991    bye world
.....
10 hello hadoop
9  bye hadoop
8  hello world
7  bye world
6  hello hadoop
5  bye hadoop
4  hello world

```



```

3  bye world
2  hello hadoop
1  bye hadoop

```

这个程序的关键是下面这两句话：

```

SequenceFile.Writer writer = null;

writer = SequenceFile.createWriter(fs, conf, path, key.getClass(), value.
    getClass());

writer.append(key, value);

```

你需要声明 `SequenceFile.Writer` 类并使用函数 `SequenceFile.createWriter()` 来给它赋值。这个函数中你至少要指定四个参数，即输出流（fs）、conf 对象（conf）、key 的类型、value 的类型，同时它还有很多重构函数，你可以设置压缩等。然后你就可以使用 `writer.append()` 来向流中写入 key/value 对了。

读取 `SequenceFile` 文件内容的程序也很简单，如代码清单 8-2 所示。

代码清单 8-2 SequenceFileReadFile

```

import java.io.IOException;
import java.net.URI;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.FileSystem;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IOUtils;
import org.apache.hadoop.io.SequenceFile;
import org.apache.hadoop.io.Writable;
import org.apache.hadoop.util.ReflectionUtils;

public class SequenceFileReadFile {
    public static void main(String[] args) throws IOException {
        String uri = "你想要读取的 SequenceFile 所在位置";
        Configuration conf = new Configuration();
        FileSystem fs = FileSystem.get(URI.create(uri), conf);
        Path path = new Path(uri);
        SequenceFile.Reader reader = null;
        try {
            reader = new SequenceFile.Reader(fs, path, conf);
            Writable key = (Writable)ReflectionUtils.newInstance(
                reader.getKeyClass(), conf);
            Writable value = (Writable)ReflectionUtils.newInstance(
                reader.getValueClass(), conf);
            long position = reader.getPosition();
            while (reader.next(key, value)) {
                String syncSeen = reader.syncSeen() ? "*" : "";
                System.out.printf("[%s%s]\t%s\t%s\n", position, syncSeen, key, value);
                position = reader.getPosition(); // beginning of next record
            }
        }
    }
}

```

```

    }
} finally {
    IOUtils.closeStream(reader);
}
}
}
}

```

读取 SequenceFile 文件的程序关键是以下几行代码：

```

SequenceFile.Reader reader = null;
reader = new SequenceFile.Reader(fs, path, conf);
reader.next(key, value);
Writable key = (Writable)ReflectionUtils.newInstance(reader.
getKeyClass(), conf);
Writable value = (Writable)ReflectionUtils.newInstance(reader.
getValueClass(), conf);

```

很简单，声明 reader 并赋值之后，你可以通过 getKeyClass() 和 getValueClass() 得到 key 与 value 的类型，并通过 ReflectionUtils 直接实例化对象，然后你就可以通过 reader.next() 跳到下一个 key/value 值了，以遍历文件中所有的 key/value 对。

根据前文所述，生成 SequenceFile 文件时是可以采用压缩方式的，下面就采用 Block 压缩方式生成 SequenceFile 文件，此程序与生成不压缩 SequenceFile 文件的程序基本相同，只是在 SequenceFile.createWrite() 时修改了一下设置，如下所示：

```

SequenceFile.createWriter(fs, conf, path, key.getClass(), value.
getClass(), CompressionType.BLOCK)

```

然后查看生成的两个文件大小：

```

-rwxrwxrwx 1 u u 10214801 2011-01-14 16:31 MySequenceOutput
-rwxrwxrwx 1 u u 159062628 2011-01-14 16:25 MySequenceOutput2

```

文件大小是以 byte 显示的，可以看到，采用 Block 压缩的文件是不压缩的 1/16 左右。我们可以将这个 Java 文件编译打包，在运行时使用 time 函数记录这两个 jar 包的执行时间，如下所示：

```

// 这是不使用压缩的程序
time hadoop jar UnComSequenceFileWriteFile.jar UnComSequence
FileWriteFile
real    0m47.668s
// 这是使用压缩的程序
time hadoop jar ComSequenceFileWriteFile.jar ComSequenceFile
WriteFile
real    0m7.539s

```

上面记录了程序具体运行的时间，以毫秒为单位。可以看出，使用压缩的程序其执行效率要远远高于不使用压缩的程序。我们推测这个时间的差距主要是受硬盘写入时间

的影响。再加上传输 10MB 的数据所花的时间明显要远远少于 159MB 的数据。这就能很好地解释为什么在 MapReduce 程序中采用压缩会提高效率了（因为一般而言，这是 map 的输出文件）。

8.4.2 MapFile 类

MapFile 的使用与 SequenceFile 类似，建立 MapFile 文件的程序，如代码清单 8-3 所示。

代码清单 8-3 MapFileWriteFile.java

```
import java.io.IOException;
import java.net.URI;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.*;
import org.apache.hadoop.io.*;

public class MapFileWriteFile {
    private static final String[] myValue = {
        "hello world",
        "bye world",
        "hello hadoop",
        "bye hadoop"
    };

    public static void main(String[] args) throws IOException {
        String uri = "你想要生成 SequenceFile 的位置";
        Configuration conf = new Configuration();
        FileSystem fs = FileSystem.get(URI.create(uri), conf);
        IntWritable key = new IntWritable();
        Text value = new Text();
        MapFile.Writer writer = null;
        try {
            writer = new MapFile.Writer(conf, fs, uri, key.getClass(),
                value.getClass());
            for (int i = 0; i < 500; i++) {
                key.set(i);
                value.set(myValue[i % myValue.length]);
                writer.append(key, value);
            }
        } finally {
            IOUtils.closeStream(writer);
        }
    }
}
```

这个程序与建立 SequenceFile 文件的程序极其类似，这里就不详述了。与 SequenceFile 只生成一个文件不同的是，这个程序生成的是一个文件夹。如下所示：

```
-rw-r--r-- * * supergroup 16018 * /user/root/MapFileOutput/data
-rw-r--r-- * * supergroup 227 * /user/root/MapFileOutput/index
```

其中 data 是存储的数据，即 MapFile 文件（经过排序 SequenceFile 文件），index 就是索引了，在这个程序中，其内容如下：

0	128
128	4200
256	8272
384	12344

可以看出，索引是按每 128 个键建立的，这个值可以通过修改 io.map.index.interval 的大小来修改。key 值后面是偏移量，用于记录 key 的位置。

读取 MapFile 文件的程序也很简单，其内容如下所示：

```
import java.io.IOException;
import java.net.URI;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.FileSystem;
import org.apache.hadoop.io.IOUtils;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.MapFile;
import org.apache.hadoop.io.Writable;
import org.apache.hadoop.io.WritableComparable;
import org.apache.hadoop.util.ReflectionUtils;

public class MapFileReadFile {
    public static void main(String[] args) throws IOException {
        String uri = "你想要读取的 MapFile 文件位置";
        Configuration conf = new Configuration();
        FileSystem fs = FileSystem.get(URI.create(uri), conf);
        MapFile.Reader reader = null;
        try {
            reader = new MapFile.Reader(fs, uri, conf);
            WritableComparable key = (WritableComparable)
                ReflectionUtils.newInstance(reader.getKeyClass(), conf);
            Writable value = (Writable)ReflectionUtils.
                newInstance(reader.getValueClass(), conf);
            while (reader.next(key, value)) {
                System.out.printf("%s\t%s\n", key, value);
            }
            reader.get(new IntWritable(7), value);
            System.out.printf("%s\n", value);
        } finally {
            IOUtils.closeStream(reader);
        }
    }
}
```

其特别之处是，MapFile 可以查找单个键所对应的 value 值，见下面这段话：

执行这个操作时，MapFile.Reader() 需要先把 index 读入内存中，然后执行一个简单的

二叉搜索找到数据，`MapFile.Reader()` 在查找时，会先在索引文件中找到小于你想要找的 key 值的索引 key 值，然后再到 data 文件中向后查找。

大型 `MapFile` 文件的索引通常会占用很大的内存，这时我们可以通过重设索引，增加索引间隔的方法来降低索引文件的大小，但是重设索引是一个很麻烦的事情。Hadoop 提供了另一个非常有效的方法，就是读取索引文件时，可以隔几个索引 key 再读取索引 key 值，这样就可以有效地降低读入内存的索引文件的大小。至于跳过 key 的个数是通过 `io.map.index.skip` 来设置的。

8.5 小结

本章介绍了 Hadoop 的 I/O 操作，主要有以下几个内容，数据完整性、压缩、序列化和基于文件的数据结构。数据完整性主要介绍了 Hadoop 是如何通过校验和机制保证数据完整性的；压缩介绍了目前 Hadoop 开发的几种压缩算法及它们的优缺点，其中压缩分割和输入分割是我们编写 MapReduce 程序时经常要用到的，要理解清楚；序列化主要介绍了 Hadoop 自己的序列化机制，这个序列化机制非常简单直接，并不像 Java 的序列化机制那样面面俱到，这样可以使数据更加紧凑，同时也可以加快序列化和反序列化的速度。最后介绍了 Hadoop 自己定义的两类数据结构（也可以看成一类），这是一种非常常用的基于文件数据结构，MapReduce 程序中 map 程序生成的中间结果就是用这种基于文件的数据结构表示的，它也是本章中非常重要的一个内容。





第 9 章

HDFS 详解

本章内容

- ☐ Hadoop 的文件系统
- ☐ HDFS 简介
- ☐ HDFS 体系结构
- ☐ HDFS 的基本操作
- ☐ HDFS 常用 Java API 详解
- ☐ HDFS 中的读写数据流
- ☐ HDFS 命令详解
- ☐ 小结

资源分享

PDG

HDFS (Hadoop Distributed File System) 是 Hadoop 项目的核心子项目, 是 Hadoop 主要应用的一个分布式文件系统, 本章将对它进行详细介绍。实际上, 在 Hadoop 中有一个综合性的文件系统抽象, 它提供了文件系统实现的各类接口, HDFS 只是这个抽象文件系统的—个实例。

在本章中, 笔者首先会对 Hadoop 的文件系统给予一个总体的介绍, 然后对 HDFS 的相关内容给予重点的讲解, 包括 HDFS 的特点、基本操作、常用 API 及读写数据流等, 并在 9.9 节中分小节介绍。

9.1 Hadoop 的文件系统

Hadoop 整合了众多文件系统, 它首先提供了一个高层的文件系统抽象类 `org.apache.hadoop.fs.FileSystem`, 这个抽象类展示了一个分布式文件系统, 并有几个具体实现, 如表 9-1 所示。

表 9-1 Hadoop 的文件系统

文件系统	URI 方案	Java 实现 (<code>org.apache.hadoop</code>)	定义
Local	file	<code>fs.LocalFileSystem</code>	支持有客户端校验和的本地文件系统。带有校验和的本地文件系统在 <code>fs.RawLocalFileSystem</code> 中实现
HDFS	hdfs	<code>hdfs.DistributedFileSystem</code>	Hadoop 的分布式文件系统
HFTP	hftp	<code>hdfs.HftpFileSystem</code>	支持通过 HTTP 方式以只读的方式访问 HDFS, <code>distcp</code> 经常用在不同的 HDFS 集群间复制数据
HSFTP	hsftp	<code>hdfs.HsftpFileSystem</code>	支持通过 HTTPS 方式以只读的方式访问 HDFS
HAR	har	<code>fs.HarFileSystem</code>	构建在 Hadoop 文件系统之上, 对文件进行归档。Hadoop 归档文件主要用来减少 NameNode 的内存使用
KFS	kfs	<code>fs.kfs.KosmosFileSystem</code>	Cloudstore (其前身是 Kosmos 文件系统) 文件系统是类似于 HDFS 和 Google 的 GFS 文件系统, 使用 C++ 编写
FTP	ftp	<code>fs.ftp.FtpFileSystem</code>	由 FTP 服务器支持的文件系统
S3 (本地)	s3n	<code>fs.s3native.NativeS3FileSystem</code>	基于 Amazon S3 的文件系统
S3 (基于块)	s3	<code>fs.s3.NativeS3FileSystem</code>	基于 Amazon S3 的文件系统, 以块格式存储解决了 S3 的 5GB 文件大小的限制

Hadoop 提供了许多文件系统的接口, 用户可使用 URI 方案选取合适的文件系统来实现交互。比如, 可以使用 9.4.1 节介绍的文件系统命令行接口来进行 Hadoop 文件系统的操作。如果想列出本地文件系统的目录, 那么执行以下 shell 命令即可:

```
hadoop fs -ls file:///
```

(1) 接口

Hadoop 是使用 Java 编写的，而 Hadoop 中不同文件系统之间的交互是由 Java API 进行调解的。事实上，前面使用的文件系统的 shell 就是一个 Java 应用，它使用 Java 文件系统类来提供文件系统操作。即使其他文件系统比如 FTP、S3 都有自己的访问工具，这些接口在 HDFS 中还是被广泛使用，主要用来进行 Hadoop 文件系统之间的协作。

(2) Thrift

上面提到可以通过 Java API 与 Hadoop 的文件系统进行交互，而其他非 Java 应用访问 Hadoop 文件系统则比较麻烦。Thriftfs 分类单元中的 Thrift API 可通过将 Hadoop 文件系统展示为一个 Apache Thrift 服务来填补这个不足，让任何有 Thrift 绑定的语言都能轻松地与 Hadoop 文件系统进行交互。Thrift 是由 Facebook 公司开发的一种可伸缩的跨语言服务的发展软件框架。Thrift 解决了各系统间大数据量的传输通信问题，以及系统之间因语言环境不同而需要跨平台的问题。在多种不同的语言之间通信时，Thrift 可以作为二进制的高性能的通信中间件，它支持数据（对象）序列化和多种类型的 RPC 服务。

下面来看如何使用 Thrift API。要使用 Thrift API，首先要运行提供 Thrift 服务的 Java 服务器，并以代理的方式访问 Hadoop 文件系统。Thrift API 包含很多其他语言生成的 stub，包括 C++、Perl、PHP、Python 等。Thrift 支持不同的版本，因此可以从同一个客户代码中访问不同版本的 Hadoop 文件系统，但要运行针对不同版本的代理。

关于安装与使用教程，可以参考 src/contrib/thriftfs 目录中关于 Hadoop 分布的参考文档。

(3) C 语言库

Hadoop 提供了映射 Java 文件系统接口的 C 语言库——libhdfs。libhdfs 可以编写为一个访问 HDFS 的 C 语言库，实际上，它可以访问任意的 Hadoop 文件系统，它也可以使用 JNI（Java Native Interface）来调用 Java 文件系统的客户端。

这里的 C 语言的接口和 Java 的使用非常相似，只是稍滞后于 Java，目前还不支持一些新特性。相关资料可参见 libhdfs/docs/api 目录中关于 Hadoop 分布的 C API 文档。

(4) FUSE

FUSE（Filesystem in Userspace）允许将文件系统整合为一个 Unix 文件系统并在用户空间中执行。通过使用 Hadoop Fuse-DFS 的 contrib 模块来支持任意的 Hadoop 文件系统作为一个标准的文件系统挂载，并且可以使用 Unix 的工具（像 ls、cat）和文件系统进行交互，还可以通过任意一种编程语言使用 POSIX 库来访问文件系统。

Fuse-DFS 是用 C 语言实现的，可使用 libhdfs 作为与 HDFS 的接口，关于如何编译和运行 Fuse-DFS，可以参见 src/contrib/fuse-dfs 中的相关文档。

(5) WebDAV

WebDAV 是一系列支持编辑和更新文件的 HTTP 的扩展，在大部分操作系统中，WebDAV 共享都可以作为文件系统挂载，因此，可通过 WebDAV 来对外提供 HDFS 或其他 Hadoop 文件系统，可以将 HDFS 作为一个标准的文件系统进行访问。

(6) 其他 HDFS 接口

HDFS 接口还提供了以下其他两种特定的接口。

- HTTP。HDFS 定义了一个只读接口，用来在 HTTP 上检索目录列表和数据。NameNode 的嵌入式 Web 服务器运行在 50070 端口上，以 XML 格式提供服务，文件数据由 DataNode 通过它们的 Web 服务器 50075 端口向 NameNode 提供。这个协议并不拘泥于某个 HDFS 版本，所以用户可以自己编写使用 HTTP 从运行不同版本的 Hadoop 的 HDFS 中读取数据。HftpFileSystem 就是其中的一种实现，它是一个通过 HTTP 和 HDFS 交流的 Hadoop 文件系统，是 HTTPS 的变体。
- FTP。Hadoop 接口中还有一个 HDFS 的 FTP 接口，它允许使用 FTP 协议和 HDFS 交互，即使用 FTP 客户端和 HDFS 进行交互。

9.2 HDFS 简介

HDFS 是基于流数据模式访问和处理超大文件的需求而开发的，它可以运行于廉价的商用服务器上。总的来说，可以将 HDFS 的主要特点概括为以下几点。

(1) 处理超大文件

这里的超大文件通常是指数百 MB、甚至数百 TB 大小的文件。目前在实际应用中，HDFS 已经能用来存储管理 PB (PeteBytes) 级的数据了。在 Yahoo!，Hadoop 集群也已经扩展到了 4 000 个节点。

(2) 流式地访问数据

HDFS 的设计建立在更多地响应“一次写入、多次读取”任务的基础之上。这意味着一个数据集一旦由数据源生成，就会被复制分发到不同的存储节点中，然后响应各种各样的数据分析任务请求。在多数情况下，分析任务都会涉及数据集中的大部分数据，也就是说，对 HDFS 来说，请求读取整个数据集要比读取一条记录更加高效。

(3) 运行于廉价的商用机器集群上

Hadoop 设计对硬件需求比较低，只须运行在廉价的商用硬件集群上，而无须昂贵的高可用性机器上。廉价的商用机也就意味着大型集群中出现节点故障情况的概率非常高。这就要求在设计 HDFS 时要充分考虑数据的可靠性、安全性及高可用性。

正是由于以上的种种考虑，我们会发现现在的 HDFS 在处理一些特定问题时不但没有优势，而且有一定的局限性，主要表现在以下几方面。

(1) 不适合低延迟数据访问

如果要处理一些用户要求时间比较短的低延迟应用请求，则 HDFS 不适合。HDFS 是为了处理大型数据集分析任务的，主要是为达到高的数据吞吐量而设计的，这就可能要求以高延迟作为代价。目前有一些补充的方案，比如使用 HBase，通过上层数据管理项目来尽可能地弥补这个不足。

(2) 无法高效存储大量小文件

在 Hadoop 中需要用 NameNode (名称节点) 来管理文件系统的元数据, 以响应客户端请求返回文件位置等, 因此文件数量大小的限制要由 NameNode 来决定。例如, 每个文件、索引目录及块大约占 100 字节, 如果有 100 万个文件, 每个文件占一个块, 那么至少要消耗 200MB 内存, 这似乎还可以接受。但如果有更多文件, 那么 NameNode 的工作压力更大, 检索处理元数据的时间就不可接受了。

(3) 不支持多用户写入及任意修改文件

在 HDFS 的一个文件中只有一个写入者, 而且写操作只能在文件末尾完成, 即只能执行追加操作。目前 HDFS 还不支持多个用户对同一文件的写操作, 以及在文件任意位置进行修改。

当然, 以上几点都是当前的问题, 相信随着研究者的努力, HDFS 会更加成熟, 可满足更多的应用需要。以下链接是 Hadoop 的一些热点研究方向, 读者可以自行参考: <http://wiki.apache.org/hadoop/ProjectSuggestions>。

9.3 HDFS 体系结构

想要了解 HDFS 的体系结构, 首先要从 HDFS 的相关概念入手, 下面将介绍 HDFS 中的几个重要概念。

9.3.1 HDFS 的相关概念

1. 块 (Block)

我们知道, 在操作系统中都有一个文件块的概念, 文件以块的形式存储在磁盘中, 此处块的大小代表系统读/写可操作的最小文件大小。也就是说, 文件系统每次只能操作磁盘块大小的整数倍数据。通常来说, 一个文件系统块为几千字节, 而磁盘块大小为 512 字节。文件的操作都由系统完成, 这些对用户来说都是透明的。

这里, 我们所要介绍的 HDFS 中的块是一个抽象的概念, 它比上面操作系统中所说的块要大得多。在配置 Hadoop 系统时会看到, 它的默认块为 64MB。和单机上的文件系统相同, HDFS 分布式文件系统中的文件也被分成块进行存储, 它是文件存储处理的逻辑单元 (如果没有特别指出, 后文中所描述的块都是指 HDFS 中的块)。

HDFS 作为一个分布式文件系统, 是设计用来处理大文件的, 使用抽象的块会带来很多好处。一个好处是可以存储任意大的文件, 而又不会受到网络中任一单个节点磁盘大小的限制。可以想象一下, 单个节点存储 100TB 的数据是不可能的, 但是由于逻辑块的设计, HDFS 可以将这个超大的文件分成众多块, 分别存储在集群的各台机器上。另外一个好处是使用抽象块作为操作的单元可简化存储子系统。这里之所以提到简化, 是因为这是所有系统的追求, 而对故障出现频繁和种类繁多的分布式系统来说, 简化就显得尤为重要。在 HDFS

中块的大小固定，这样它就简化了存储系统的管理，特别是元数据信息可以和文件块内容分开存储。不仅如此，块更有利于分布式文件系统中复制容错的实现。在 HDFS 中为了处理节点故障，默认将文件块副本数设定为 3 份，分别存储在集群的不同节点上。当一个块损坏时，系统会通过 NameNode 获取元数据信息，在另外的机器上读取一个副本并进行存储，这个过程对用户来说都是透明的。当然，这里的文件块副本冗余量可以通过文件进行配置，比如在有些应用中，可能会为操作频率较高的文件块设置较高的副本数量以提高集群的吞吐量。

在 HDFS 中，可以通过终端命令直接获得文件和块信息，比如以下命令可以列出文件系统中组成各个文件的块（有关 HDFS 的命令，将会在 9.4 节中详细讲解）：

```
hadoop fsck / -files -blocks
```

2. NameNode 和 DataNode

HDFS 体系结构中有两类节点，一类是 NameNode；另一类是 DataNode。这两类节点分别承担 Master 和 Worker 的任务。NameNode 就是 Master 管理集群中的执行调度，DataNode 就是 Worker 具体任务的执行节点。NameNode 管理文件系统的命名空间，维护整个文件系统的文件目录树及这些文件的索引目录。这些信息以两种形式存储在本地文件系统中，一种是命名空间镜像（Namespace image）；一种是编辑日志（Edit log）。从 NameNode 中你可以获得每个文件的每个块所在的 DataNode。有一点需要注意的是，这些信息不是永久保存的，NameNode 会在每次启动系统时动态地重建这些信息。当运行任务时，客户端通过 NameNode 获取元数据信息，和 DataNode 进行交互以访问整个文件系统。系统会提供一个类似于 POSIX 的文件接口，这样用户在编程时无须考虑 NameNode 和 DataNode 的具体功能。

DataNode 是文件系统 Worker 中的节点，用来执行具体的任务：存储文件块，被客户端和 NameNode 调用。同时，它会通过心跳（heartbeat）定时向 NameNode 发送所存储的文件块信息。

9.3.2 HDFS 的体系结构

如图 9-1 所示，HDFS 采用 master/slave 架构对文件系统进行管理。一个 HDFS 集群是由一个 NameNode 和一定数目的 DataNodes 组成的。NameNode 是一个中心服务器，负责管理文件系统的名字空间（Namespace）及客户端对文件的访问。集群中的 DataNode 一般是一个节点运行一个 DataNode 进程，负责管理它所在节点上的存储。HDFS 展示了文件系统的名字空间，用户能够以文件的形式在上面存储数据。从内部看，一个文件其实被分成了一个或多个数据块，这些块存储在一组 DataNode 上。NameNode 执行文件系统的名字空间操作，比如打开、关闭、重命名文件或目录。它也负责确定数据块到具体 DataNode 节点的映射。DataNode 负责处理文件系统客户端的读/写请求。在 NameNode 的统一调度下进行数据块的创建、删除和复制。

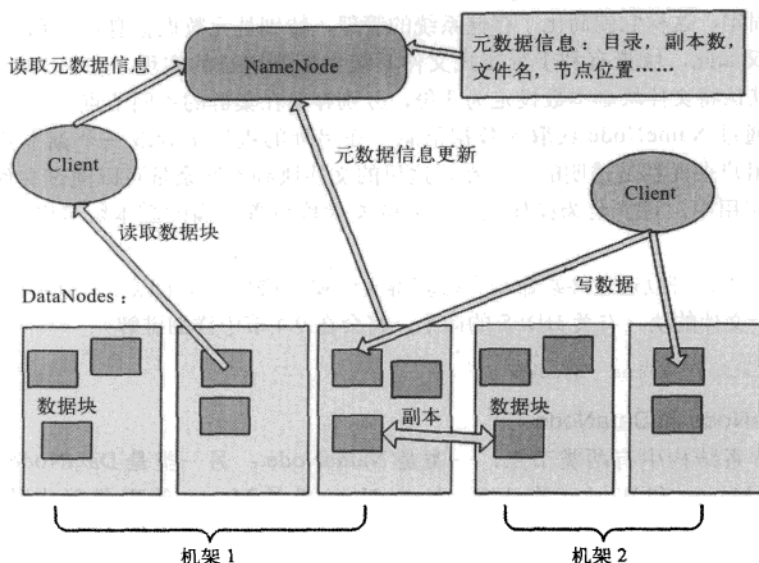


图 9-1 HDFS 的体系结构

1. 副本存放与读取策略

副本的存放是 HDFS 可靠性和性能的关键，优化的副本存放策略也正是 HDFS 区别于其他大部分分布式文件系统的重要特性。HDFS 采用一种称为机架感知（rack-aware）的策略来改进数据的可靠性、可用性和网络带宽的利用率。大型 HDFS 实例一般运行在跨越多个机架的计算机组成的集群上，不同机架上的两台机器之间的通信需要经过交换机，这样会增加数据传输的成本。在大多数情况下，同一个机架内的两台机器间的带宽会比不同机架的两台机器间的带宽大。

一方面，通过一个机架感知的过程，NameNode 可以确定每个 DataNode 所属的机架 ID。目前 HDFS 采用的策略就是将副本存放在不同的机架架上。这样可以有效防止当整个机架失效时数据的丢失，并且允许读数据的时候充分利用多个机架的带宽。这种策略设置可以将副本均匀地分布在集群中，有利于在组件失效情况下的负载均衡。但是，因为这种策略的一个写操作需要传输数据块到多个机架，这增加了写操作的成本。

举例来看，在大多数情况下，副本系数是 3，HDFS 的存放策略是将一个副本存放在本地机架的节点上，另一个副本放在同一机架的另一个节点上，最后一个副本放在不同机架的节点上。这种策略减少了机架间的数据传输，这就提高了写操作的效率。机架的错误远比节点的错误少，所以这个策略不会影响数据的可靠性和可用性。同时，因为数据块只放在两个不同的机架架上，所以此策略减少了读取数据时需要的网络传输总带宽。这一策略在不损害数据可靠性和读取性能的情况下改进了写的性能。

另一方面，在读取数据时，为了减少整体的带宽消耗和降低整体的带宽延时，HDFS 会

尽量让读取程序读取离客户端最近的副本。如果在读取程序的同一个机架上有有一个副本，那么就读取该副本。如果一个 HDFS 集群跨越了多个数据中心，那么客户端也将首先读取本地数据中心的副本。

2. 安全模式

NameNode 启动后会进入一个称为安全模式的特殊状态。处于安全模式的 NameNode 是不会进行数据块的复制的。NameNode 从所有的 DataNode 上接收心跳信号和块状态报告。块状态报告包括了某个 DataNode 所有的数据块列表。每个数据块都有一个指定的最小副本数。当 NameNode 检测确认某个数据块的副本数目达到最小值时，那么该数据块就会被认为是副本安全的；在一定百分比（这个参数可配置）的数据块被 NameNode 检测确认安全之后（加上一个额外的 30 秒等待时间），NameNode 将退出安全模式状态。接下来它会确定还有哪些数据块的副本没有达到指定数目，并将这些数据块复制到其他 DataNode 上。9.7 节中将详细介绍安全模式的相关命令。

3. 文件安全

NameNode 的重要性是显而易见的，没有它客户端将无法获得文件块的位置。在实际应用中，如果集群的 NameNode 出现故障，就意味着整个文件系统中全部的文件会丢失，因为我们无法再通过 DataNode 上的文件块来重构文件。下面简单介绍 Hadoop 是采用哪种机制来确保 NameNode 的安全的。

第一种方法是，备份 NameNode 上持久化存储的元数据文件，然后将其转储到其他文件系统中，这种转储是同步的、原子的操作。通常的实现方法是，将 NameNode 中的元数据转储到远程的 NFS 文件系统中。

第二种方法是，系统中同步运行一个 Secondary NameNode（二级 NameNode）。这个节点的主要作用就是周期性地合并编辑日志中的命名空间镜像，以避免编辑日志过大。Secondary NameNode 的运行通常需要大量的 CPU 和内存去做合并操作，这就要求其运行在一台单独的机器上。在这台机器上会存储合并过的命名空间镜像，这些镜像文件会在 NameNode 宕机后做替补使用，以便最大限度地减少文件的损失。但是，需要注意的是，Secondary NameNode 的同步备份总会滞后于 NameNode，所以损失是必然的。有关文件系统镜像和编辑日志的详细介绍请参见第 10 章。

9.4 HDFS 的基本操作

本节将对 HDFS 的命令行操作及其 Web 界面进行介绍。

9.4.1 HDFS 的命令行操作

可以通过命令行接口来和 HDFS 进行交互。当然，命令行接口只是 HDFS 的访问接口之一，它的特点是更加简单直观，便于使用，可以进行一些基本操作。下面会简要介绍在单机

上运行 Hadoop，执行单机伪分布（我的环境为 Windows 下的单节点情况，与其他情况下命令一样，读者可自行参考），具体的安装与配置可以参看本书第 2 章，随后会告诉你如何运行在集群机器上，以支持可扩展性和容错。

在单机伪分布的配置中需要修改两个配置属性。第一个需要修改的配置文件属性为 `fs.default.name`，将其设置为 `hdfs://localhost/`，用来设定一个默认的 Hadoop 文件系统，再使用一个 `hdfsURI` 来配置说明，Hadoop 默认使用 HDFS 文件系统。HDFS 的守护进程会通过这个属性来为 NameNode 定义 HDFS 中的主机和端口。这里在本机 `localhost` 上运行 HDFS，其端口采用默认的 8020。HDFS 的客户端可以通过这个属性访问各个节点。第二个需要修改的配置文件属性为 `dfs.replication`，因为采用单机伪分布，所以不支持副本，HDFS 不可能将副本存储到其他两个节点，因此要将配置文件中默认的副本数 3 改为 1。

下面就具体介绍如何通过命令行访问 HDFS 文件系统。本节主要讨论一些基本的文件操作，比如读文件、创建文件存储路径、转移文件、删除文件、列出文件列表等操作。在终端中你可以通过输入 `hadoop fs -help` 获得 HDFS 操作的详细帮助信息。

首先，我们将本地的一个文件复制到 HDFS 中，操作命令如下：

```
hadoop fs -copyFromLocal testInput/hello.txt hdfs://localhost/user/admin/In/hello.txt
```

这条命令调用了 Hadoop 的终端命令 `fs`。`fs` 支持很多子命令，这里使用 `-copyFromLocal` 命令将本地的文件 `hello.txt` 复制到 HDFS 中的 `/user/ admin/In/hello.txt` 下。事实上，使用 `fs` 命令可以省略 URI 中的访问协议和主机名，而直接使用配置文件 `core-site.xml` 中的默认属性值 `hdfs://localhost`，即将命令改为如下形式亦可：

```
hadoop fs -copyFromLocal testInput/hello.txt /user/admin/In/hello.txt
```

其次，看如何将 HDFS 中的文件拷贝到本机，操作命令如下：

```
hadoop fs -copyToLocal /user/admin/In/hello.txt testInput/hello.copy.txt
```

命令执行后，用户可查看根目录 `testInput` 文件夹下的 `hello.copy.txt` 文件以验证完成从 HDFS 到本机的文件拷贝。

下面查看创建文件夹的方法：

```
hadoop fs -mkdir testDir
```

最后，用命令行查看 HDFS 文件列表：

```
Admin@Admin ~/hadoop-0.20.2
hadoop fs -lsr In
```

```
-rw-r--r-- 1 admin supergroup 348624 2010-10-11 11:34 /user/admin/In/CHANGES.txt
-rw-r--r-- 1 admin supergroup 13366 2010-10-11 11:34 /user/admin/In/LICENSE.txt
-rw-r--r-- 1 admin supergroup 101 2010-10-11 11:34 /user/admin/In/NOTICE.txt
-rw-r--r-- 1 admin supergroup 1366 2010-10-11 11:34 /user/admin/In/README.txt
-rw-r--r-- 1 admin supergroup 13 2010-10-17 15:14 /user/admin/In/hello.txt
```

从以上文件列表中可以看到，命令返回的结果和 Linux 下 `ls -l` 命令返回的结果很相似。返回结果的第一列是文件属性，第二列是文件的副本因子，而这是传统的 Linux 系统没有的。为了方便，我将副本因子设置为 1，所以这里显示为 1，我们也看到了从本地拷贝到 In 文件夹下的 `hello.txt` 文件。

9.4.2 HDFS 的 Web 界面

在部署好 Hadoop 集群之后，便可以直接通过 `http://NameNodeIP:50070` 访问 HDFS 的 Web 界面，如图 9-2 所示。

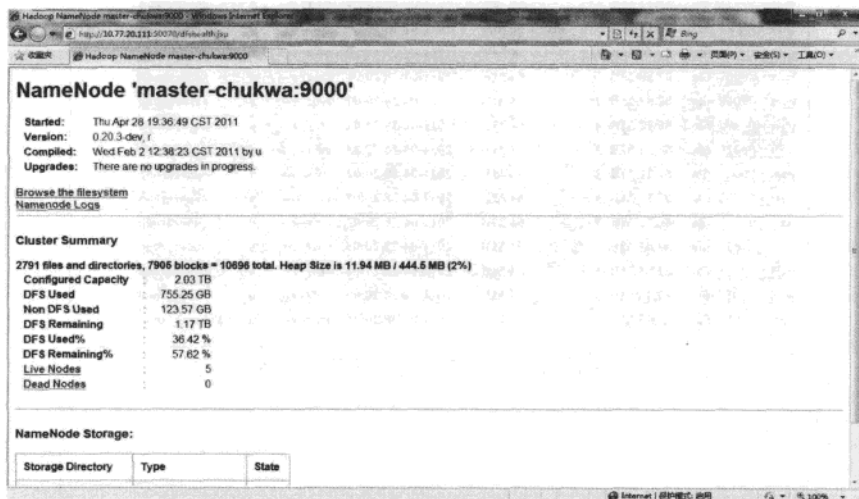


图 9-2 HDFS 的 Web 界面

从图 9-2 中可以看到，HDFS 的 Web 界面提供了基本的文件系统信息，其中包括集群启动时间、版本号、编译时间及是否又升级：

```
Started: Thu Apr 28 19:36:49 CST 2011
Version: 0.20.3-dev, r
Compiled: Wed Feb 2 12:38:23 CST 2011 by u
```

从图 9-2 中的 Upgrades : There are no upgrades in progress. 可以看出，HDFS 的 Web 界面还提供了文件系统的基本功能：Browse the filesystem（浏览文件系统），点击链接即可看到，它将 HDFS 的文件结构通过目录的形式展现出来，增加了对文件系统的可读性，如图 9-3 所示。

此外，可以直接通过 Web 界面访问文件内容，如图 9-4 所示。同时，HDFS 的 Web 界面还将该文件块所在的节点位置展现出来。可以通过设置 Chunk size to view 来设置一次读取并展示的文件块大小。

Name	Type	Size	Replication	Block Size	Modification Time	Permission	Owner	Group
logs	dir				2011-04-28 19:44	rxr-xr-x	u	supergroup
part-r-00000	file	87.65 MB	3	64 MB	2011-04-28 19:58	rw-r--r--	u	supergroup
part-r-00001	file	87.01 MB	3	64 MB	2011-04-28 19:58	rw-r--r--	u	supergroup
part-r-00002	file	86.65 MB	3	64 MB	2011-04-28 19:58	rw-r--r--	u	supergroup
part-r-00003	file	87.68 MB	3	64 MB	2011-04-28 19:58	rw-r--r--	u	supergroup
part-r-00004	file	87.12 MB	3	64 MB	2011-04-28 19:58	rw-r--r--	u	supergroup
part-r-00005	file	87.15 MB	3	64 MB	2011-04-28 19:58	rw-r--r--	u	supergroup
part-r-00006	file	87.46 MB	3	64 MB	2011-04-28 19:58	rw-r--r--	u	supergroup
part-r-00007	file	87.88 MB	3	64 MB	2011-04-28 20:01	rw-r--r--	u	supergroup
part-r-00008	file	87.57 MB	3	64 MB	2011-04-28 20:02	rw-r--r--	u	supergroup
part-r-00009	file	86.72 MB	3	64 MB	2011-04-28 20:01	rw-r--r--	u	supergroup
part-r-00010	file	88.32 MB	3	64 MB	2011-04-28 20:03	rw-r--r--	u	supergroup
part-r-00011	file	87.31 MB	3	64 MB	2011-04-28 20:02	rw-r--r--	u	supergroup
part-r-00012	file	87.12 MB	3	64 MB	2011-04-28 20:02	rw-r--r--	u	supergroup
part-r-00013	file	86.71 MB	3	64 MB	2011-04-28 20:02	rw-r--r--	u	supergroup
part-r-00014	file	87.43 MB	3	64 MB	2011-04-28 20:03	rw-r--r--	u	supergroup
part-r-00015	file	87.55 MB	3	64 MB	2011-04-28 20:04	rw-r--r--	u	supergroup
part-r-00016	file	87.36 MB	3	64 MB	2011-04-28 20:05	rw-r--r--	u	supergroup
part-r-00017	file	87.35 MB	3	64 MB	2011-04-28 20:06	rw-r--r--	u	supergroup

图 9-3 通过 HDFS 的 Web 界面访问文件系统目录

Go back to dir listing

Advanced view/download options

View Next chunk

*,on 20

*2.5G/622M/155M 20

5BGS 20

*9. 20

*ASP方法*对新产品的开发顺序进行排序,利用企业现有的资源进行重点突破,保证新产品<WP-3>品开发顺利进行,从根本上改变本企业产品结构,业的产品从供给家用电器市场拓展到包括家用电器、计算机、信息通信、汽车电器等终端。 20

*Activity 20

All-In-One-Card 20

*Base 40

Bi-syndrome, 20

*Borneol 20

CANE-COOKED-RICE 20

*City 60

*Classical 20

*Component-based 20

*Computer 40

*Cool 20

*Daniel 20

*East 40

*Eligih 20

*Electromagnetism 20

ErYa 20

*Every 20

*Exploit 20

*Genre 20

Download this file

Tail this file

Chunk size to view (in bytes, up to file's DFS block size): 32768 Refresh

Total number of blocks: 2

6846332361295758414: 10.77.20.202-50010 10.77.20.201-50010 10.77.20.112-50010

528472438790713098: 10.77.20.201-50010 10.77.20.203-50010 10.77.20.112-50010

图 9-4 通过 HDFS 的 Web 界面访问文件内容

除了在本节中展示的信息之外，HDFS 的 Web 界面还提供了 NameNode 的日志列表、运行中的节点列表及宕机的节点列表等信息。

9.5 HDFS 常用 Java API 详解

在 9.1 节中已经了解了 Java API 的重要性，本节将深入介绍 Hadoop 的 Filesystem 类与 Hadoop 文件系统进行交互的 API。

9.5.1 使用 Hadoop URL 读取数据

如果想从 Hadoop 中读取数据，最简单的办法就是使用 `java.net.URL` 对象打开一个数据流，并从中读取数据，一般的调用格式如下：

```
InputStream in = null;
try {
    in = new URL("hdfs://NameNodeIP/path").openStream();
    // process in
} finally {
    IOUtils.closeStream(in);
}
```

这里要进行的处理是，通过 `FsUrlStreamHandlerFactory` 实例来调用在 URL 中的 `setURLStreamHandlerFactory` 方法。这种方法在一个 Java 虚拟机中只能调用一次，因此放在一个静态方法中执行。这意味着如果程序的其他部分也设置了一个 `URLStreamHandlerFactory`，那么会导致无法再从 Hadoop 中读取数据。

读取文件系统中路径为 `hdfs://NameNodeIP/user/admin/In/ hello.txt` 的文件 `hello.txt`，如代码清单 9-1 所示。这里假设 `hello.txt` 的文件内容为“Hello Hadoop！”。

代码清单 9-1 使用 `URLStreamHandler` 以标准输出显示 Hadoop 文件系统文件

```
import java.io.*;
import java.net.URL;
import org.apache.hadoop.fs.FsUrlStreamHandlerFactory;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.filecache.DistributedCache;
import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;

public class URLCat {
    static {
        URL.setURLStreamHandlerFactory(new FsUrlStreamHandlerFactory());
    }

    public static void main(String[] args) throws Exception {
        InputStream in = null;
```

```

try {
    in = new URL(args[0]).openStream();
    IOUtils.copyBytes(in, System.out, 4096, false);
} finally {
    IOUtils.closeStream(in);
}
}
}

```

然后在 Eclipse 下设置程序运行参数为: `hdfs://NameNodeIP/user/admin/In/hello.txt`, 运行程序即可看到 `hello.txt` 中的文本内容。

需要说明的是, 这里使用了 Hadoop 中简洁的 `IOUtils` 类来关闭 `finally` 子句中的数据流, 同时复制了输出流之间的字节 (`System.out`)。代码清单 9-1 中用到的 `IOUtils.copyBytes()` 方法, 其中的两个参数, 前者表示复制缓冲区的大小, 后者表示复制后关闭数据流。

9.5.2 使用 FileSystem API 读取数据

9.5.1 节提到在应用中会出现不能使用 `URLStreamHandlerFactory` 的情况, 这时就需要使用 `FileSystem` 的 API 打开一个文件的输入流了。

文件在 Hadoop 文件系统中被视为一个 `Hadoop Path` 对象。我们可以把一个路径视为 Hadoop 的文件系统 URI, 比如上文中的 `hdfs://localhost/user/admin/In/hello.txt`。

`FileSystemAPI` 是一个高层抽象的文件系统 API, 所以, 首先要找到这里的文件系统实例 `HDFS`。取得 `FileSystem` 实例有两种静态工厂方法:

```

public static FileSystem get(Configuration conf) throws IOException
public static FileSystem get(URI uri, Configuration conf) throws IOException

```

`Configuration` 对象封装了一个客户端或服务器的配置, 这是用路径读取的配置文件设置的, 一般为 `conf/core-site.xml`。第一个方法返回的是默认文件系统, 如果没有设置, 则为默认的本地文件系统。第二个方法使用指定的 URI 方案决定文件系统的权限, 如果指定的 URI 中没有指定方案, 则退回默认的文件系统。

有了 `FileSystem` 实例后, 可通过 `open()` 方法得到一个文件的输入流:

```

public FSDataInputStream open(Path f) throws IOException
public abstract FSDataInputStream open(Path f, int bufferSize) throws IOException

```

第一个方法直接使用默认的 4KB 的缓冲区, 如代码清单 9-2 所示。

代码清单 9-2 使用 `FileSystem` API 显示 Hadoop 文件系统中的文件

```

public class FileSystemCat {
    public static void main(String[] args) throws Exception {
        String uri = args[0];
        Configuration conf = new Configuration();
        FileSystem fs = FileSystem.get(URI.create(uri), conf);
    }
}

```



```

InputStream in = null;
try {
    in = fs.open(new Path(uri));
    IOUtils.copyBytes(in, System.out, 4096, false);
} finally {
    IOUtils.closeStream(in);
}
}
}
}

```

然后设置程序运行参数为 `hdfs://localhost/user/admin/In/hello.txt`，运行程序即可看到 `hello.txt` 中的文本内容“Hello Hadoop!”。

下面对代码清单 9-2 中的程序进行扩展，重点关注 `FSDDataInputStream`。

`FileSystem` 中的 `open` 方法实际上返回的是一个 `FSDDataInputStream`，而不是标准的 `java.io` 类。这个类是 `java.io.DataInputStream` 的一个子类，支持随机访问，并可以从流的任意位置读取，代码如下：

```

public class FSDDataInputStream extends DataInputStream
implements Seekable, PositionedReadable {
    // implementation elided
}

```

`Seekable` 接口允许在文件中定位并提供一个查询方法用于查询当前位置相对于文件开始处的偏移量 (`getPos()`)，代码如下：

```

public interface Seekable {
    void seek(long pos) throws IOException;
    long getPos() throws IOException;
    boolean seekToNewSource(long targetPos) throws IOException;
}

```

其中，调用 `seek()` 来定位大于文件长度的位置会导致 `IOException` 异常。开发人员并不常用 `seekToNewSource()` 方法，此方法倾向于切换到数据的另一个副本，并在新的副本中寻找 `targetPos` 制定的位置。HDFS 就采用这样的方法在数据节点出现故障时为客户端提供可靠的数据流访问，如代码清单 9-3 所示。

代码清单 9-3 扩展代码清单 9-2，通过使用 `seek` 读取一次后，重新定位到文件头第三位，再次显示 Hadoop 文件系统中的文件内容

```

import java.io.*;
import java.net.URI;
import java.net.URL;
import java.util.*;
import org.apache.hadoop.fs.FSDDataInputStream;
import org.apache.hadoop.fs.FsUrlStreamHandlerFactory;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.fs.FileSystem;

```

```
import org.apache.hadoop.filecache.DistributedCache;
import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapred.*;
import org.apache.hadoop.util.*;

public class DoubleCat {
    public static void main(String[] args) throws Exception {
        String uri = args[0];
        Configuration conf = new Configuration();
        FileSystem fs = FileSystem.get(URI.create(uri), conf);
        FSDataInputStream in = null;
        try {
            in = fs.open(new Path(uri));
            IOUtils.copyBytes(in, System.out, 4096, false);
            in.seek(3); // go back to pos 3 of the file
            IOUtils.copyBytes(in, System.out, 4096, false);
        } finally {
            IOUtils.closeStream(in);
        }
    }
}
```

然后设置程序的运行参数为 `hdfs://localhost/user/admin/In/hello.txt`, 运行程序即可看到 `hello.txt` 中的文本内容 “Hello Hadoop!lo Hadoop!”。

同时, `FSDataInputStream` 也实现了 `PositionedReadable` 接口, 从一个指定位置读取一部分数据。这里不再详细介绍, 大家可以参考以下源代码。

```
public interface PositionedReadable {
    public int read(long position, byte[] buffer, int offset, int length)
        throws IOException;
    public void readFully(long position, byte[] buffer, int offset, int length)
        throws IOException;
    public void readFully(long position, byte[] buffer) throws IOException;
}
```

需要注意的是, `seek()` 是一个高开销的操作, 需要慎重。通常我们是依靠流数据 MapReduce 构建应用访问模式, 而不是大量地执行 `seek()` 操作。

9.5.3 创建目录

`FileSystem` 显然也提供了创建目录的方法, 代码如下:

```
public boolean mkdirs(Path f) throws IOException
```

这个方法会按照客户端请求创建未存在的父目录, 就像 `java.io.File` 的 `mkdirs()` 一样。如果目录包括所有父目录且创建成功, 那么它会返回 `true`。事实上, 一般不需要特别地创建一个目录, 因为调用 `creat()` 时写入文件会自动生成所有的父目录。

9.5.4 写数据

FileSystem 还有一系列创建文件的方法，最简单的就是给拟创建的文件指定一个路径对象，然后返回一个写输出流，代码如下：

```
public FSDataOutputStream create(Path f) throws IOException
```

这个方法有很多重载方法，例如，可以设定是否强制覆盖原文件、设定文件副本数量、设置写入文件缓冲区大小、文件块大小及设置文件许可等。

还有一个用于传递回调接口的重载方法 Progressable，通过这个方法就可以获得数据节点写入进度，代码如下：

```
package org.apache.hadoop.util;
public interface Progressable {
    public void progress();
}
```

新建文件也可以使用 append() 在一个已有文件中追加内容，这个方法也有重载，代码如下：

```
public FSDataOutputStream append(Path f) throws IOException
```

这个方法对于写入日志文件很有用，比如重启后可以在之前的日志中继续添加内容，但并不是所有的 Hadoop 文件系统都支持此方法，比如 HDFS 支持，但 S3 不支持。

代码清单 9-4 展示了如何将本地文件复制到 Hadoop 的文件系统，当 Hadoop 调用 progress() 方法时，也就是在每 64KB 大小的数据包写入数据节点的管道之后，打印一个星号来展示整个过程。

代码清单 9-4 将本地文件复制到 Hadoop 文件系统并显示进度

```
public class FileCopyWithProgress {
    public static void main(String[] args) throws Exception {
        String localSrc = args[0];
        String dst = args[1];
        InputStream in = new BufferedInputStream(new FileInputStream(localSrc));
        Configuration conf = new Configuration();
        FileSystem fs = FileSystem.get(URI.create(dst), conf);
        OutputStream out = fs.create(new Path(dst), new Progressable() {
            public void progress() {
                System.out.print("*");
            }
        });
        IOUtils.copyBytes(in, out, 4096, true);
    }
}
```

然后配置应用参数，可以看到控制台输出 “*****”，即上传显示进度，每写入 64KB

即输出一个*。目前其他文件系统写入时都不会调用 progress()。

9.5.3 节在介绍读数据时曾提到 FSDataInputStream, 这里 FileSystem 中的 creat() 方法也返回一个 FSDataOutputStream, 它也有一个查询文件当前位置的方法, 代码如下:

```
package org.apache.hadoop.fs;
public class FSDataOutputStream extends DataOutputStream implements Syncable {
    public long getPos() throws IOException {
        // implementation elided
    }
    // implementation elided
}
```

但是它与 FSDataInputStream 不同, FSDataOutputStream 不允许定位。这是因为 HDFS 只对一个打开的文件顺序写入, 或者向一个已有的文件添加。换句话说, 它不支持对除文件尾部以外的其他位置的写入操作, 这样, 写入时的定位就没有意义了。

9.5.5 删除数据

使用 FileSystem 的 delete() 可以永久删除 Hadoop 中的文件或目录。

```
public boolean delete(Path f, boolean recursive) throws IOException
```

如果传入的 f 为空文件或空目录, 那么 recursive 值会被忽略。只有当 recursive 的值为 true 时, 非空的文件或目录才会被删除, 否则抛出异常。

9.5.6 文件系统查询

同样, Java API 提供了文件系统的基本查询接口。通过这个接口, 可以查询系统的元数据信息和文件目录结构, 并可以执行更复杂的目录匹配等操作。下面将一一进行介绍。

1. 文件元数据: FileStatus

任何文件系统都要具备的重要功能就是定位其目录结构以及检索器存储的文件和目录信息。FileStatus 类封装了文件系统中文件和目录的元数据, 其中包括文件长度、块大小、副本、修改时间、所有者和许可信息等。

FileSystem 的 getFileStatus() 方法提供了获取一个文件或目录的状态对象的方法, 如代码清单 9-5 所示。

代码清单 9-5 获取文件状态信息

```
public class ShowFileStatusTest {
    private MiniDFSCluster cluster; // use an in-process HDFS cluster for testing
    private FileSystem fs;
    @Before
    public void setUp() throws IOException {
        Configuration conf = new Configuration();
        if (System.getProperty("test.build.data") == null) {
```

```

System.setProperty("test.build.data", "/tmp");
}
cluster = new MiniDFSCluster(conf, 1, true, null);
fs = cluster.getFileSystem();
OutputStream out = fs.create(new Path("/dir/file"));
out.write("content".getBytes("UTF-8"));
out.close();
}
@After
public void tearDown() throws IOException {
    if (fs != null) { fs.close(); }
    if (cluster != null) { cluster.shutdown(); }
}
@Test(expected = FileNotFoundException.class)
public void throwsFileNotFoundExceptionForNonExistentFile() throws IOException {
    fs.getFileStatus(new Path("no-such-file"));
}
@Test
public void fileStatusForFile() throws IOException {
    Path file = new Path("/dir/file");
    FileStatus stat = fs.getFileStatus(file);
    assertEquals(stat.getPath().toUri().getPath(), is("/dir/file"));
    assertTrue(stat.isDir(), is(false));
    assertEquals(stat.getLen(), is(7L));
    assertEquals(stat.getModificationTime(),
        is(lessThanOrEqualTo(System.currentTimeMillis())));
    assertEquals(stat.getReplication(), is((short) 1));
    assertEquals(stat.getBlockSize(), is(64 * 1024 * 1024L));
    assertEquals(stat.getOwner(), is("tom"));
    assertEquals(stat.getGroup(), is("supergroup"));
    assertEquals(stat.getPermission().toString(), is("rw-r--r--"));
}
@Test
public void fileStatusForDirectory() throws IOException {
    Path dir = new Path("/dir");
    FileStatus stat = fs.getFileStatus(dir);
    assertEquals(stat.getPath().toUri().getPath(), is("/dir"));
    assertTrue(stat.isDir(), is(true));
    assertEquals(stat.getLen(), is(0L));
    assertEquals(stat.getModificationTime(),
        is(lessThanOrEqualTo(System.currentTimeMillis())));
    assertEquals(stat.getReplication(), is((short) 0));
    assertEquals(stat.getBlockSize(), is(0L));
    assertEquals(stat.getOwner(), is("tom"));
    assertEquals(stat.getGroup(), is("supergroup"));
    assertEquals(stat.getPermission().toString(), is("rwxr-xr-x"));
}
}

```

如果文件或者目录不存在，就会抛出 `FileNotFoundException` 异常；如果只对文件或目

录是否存在感兴趣，那么用 `exists()` 方法更方便：

```
public boolean exists(Path f) throws IOException
```

2. 列出目录文件信息

查找文件或者目录信息很有用，但是，有时需要列出目录的内容，这需要使用 `listStatus()` 方法，代码如下：

```
public FileStatus[] listStatus(Path f) throws IOException
public FileStatus[] listStatus(Path f, PathFilter filter) throws IOException
public FileStatus[] listStatus(Path[] files) throws IOException
public FileStatus[] listStatus(Path[] files, PathFilter filter) throws IOException
```

当传入参数是一个文件时，它会简单地返回长度为 1 的 `FileStatus` 对象的一个数组。当传入参数为一个目录时，它会返回 0 个或多个 `FileStatus` 对象，代表该目录所包含的文件和子目录。

我们看到 `listStatus()` 有很多重载方法，可以使用 `PathFilter` 来限制匹配的文件和目录。如果把路径数组作为参数来调用 `listStatus()` 方法，其结果与一次对多个目录进行查询，再将 `FileStatus` 对象数组收集到一个单一的数组的结果是相同的。当然我们可以感受到，前者更为方便。代码清单 9-6 是一个简单的示范。

代码清单 9-6 显示 Hadoop 文件系统中的—个目录的文件信息

```
import java.util.*;
import org.apache.hadoop.fs.FSDataInputStream;
import org.apache.hadoop.fs.FileStatus;
import org.apache.hadoop.fs.FileUtil;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.fs.FileSystem;

public class ListStatus {
    public static void main(String[] args) throws Exception {
        String uri = args[0];
        Configuration conf = new Configuration();
        FileSystem fs = FileSystem.get(URI.create(uri), conf);
        Path[] paths = new Path[args.length];
        for (int i = 0; i < paths.length; i++) {
            paths[i] = new Path(args[i]);
        }
        FileStatus[] status = fs.listStatus(paths);
        Path[] listedPaths = FileUtil.stat2Paths(status);
        for (Path p : listedPaths) {
            System.out.println(p);
        }
    }
}
```

配置应用参数可以查看文件系统的目录，可以查看 HDFS 中对应文件目录下的文件信息。

3. 通过通配符实现目录筛选

有时候，我们需要批量处理文件，比如处理日志文件，这时可能要求 MapReduce 任务分析一个月的文件。这些文件包含在大量目录中，这就要求我们执行一个通配符操作，并使用通配符核对多个文件。Hadoop 为通配符提供了两个方法，可以在 FileSystem 中找到：

```
public FileStatus[] globStatus(Path pathPattern) throws IOException
public FileStatus[] globStatus(Path pathPattern, PathFilter filter) throws
    IOException
```

globStatus() 返回了其路径匹配所提供的 FileStatus 对象数组，再按路径进行排序，其中可选的 PathFilter 命令可以进一步限定匹配。

表 9-2 显示了 Hadoop 支持的一系列通配符。

表 9-2 Hadoop 支持的通配符及其作用

通配符	名称	匹配功能
*	星号	匹配 0 个或多个字符
?	问号	匹配一个字符
[ab]	字符类别	匹配 {a, b} 中的一个字符
[^ab]	非此字符类别	匹配不属于 {a, b} 中的一个字符
[a-b]	字符范围	匹配在 {a, b} 范围内的字符（包括 a, b），a 在字典顺序上要小于等于 b
[^a-b]	非此字符范围	匹配不在 {a, b} 范围内的字符（包括 a, b），a 在字典顺序上要小于等于 b
{a, b}	或选择	匹配包含 a 或 b 中的语句
\c	转义字符	匹配元字符 c

下面通过例子进行详细说明，假设一个日志文件的存储目录是分层组织的，其中目录格式为年/月/日：/2009/12/30、/2009/12/31、/2010/01/01、/2010/01/02。表 9-3 展示了通配符的部分样例。

表 9-3 通配符使用样例

通配符	匹配结果
/*	/2009 /2010
/*/*	/2009/12 /2010/01
/*/12/*	/2009/12/30 /2009/12/31
/200*	/2009
/200[9-10]	/2009 /2010
/200[^012345678]	/2009
/*/*/{31,01}	/2009/12/31 /2010/01/01
/*/{12/31, 01/01}	/2009/12/31 /2010/01/01

4. PathFilter 对象

使用通配符有时也不一定能够精确地定位到要访问的文件集合，比如排除一个特定的文件，这时可以使用 FileSystem 中的 listStatus() 方法和 globStatus() 方法提供可选的 PathFilter 对象来通过编程的方法控制匹配结果，如下面的代码所示。

```
package org.apache.hadoop.fs;
public interface PathFilter {
    boolean accept(Path path);
}
```

下面来看一个 PathFilter 的应用，如代码清单 9-7 所示。

代码清单 9-7 使用 PathFilter 排除匹配正则表达式的目录

```
public class RegexExcludePathFilter implements PathFilter {
    private final String regex;
    public RegexExcludePathFilter(String regex) {
        this.regex = regex;
    }
    public boolean accept(Path path) {
        return !path.toString().matches(regex);
    }
}
```

这个过滤器将留下与正则表达式不匹配的文件。

9.6 HDFS 中的读写数据流

在本节中，我们将对 HDFS 的读写数据流进行详细介绍，以帮助读者理解 HDFS 具体是如何工作的。

9.6.1 文件的读取

本节将详细介绍在执行读操作时客户端与 HDFS 交互过程的实现，以及 NameNode 和各 DataNode 之间的数据流是什么。下面将围绕图 9-5 进行具体的讲解。

首先，客户端通过调用 FileSystem 对象中的 open() 函数来读取它需要的数据。FileSystem 是 HDFS 中 DistributedFileSystem 的一个实例（参见图 9-5 第①步）。DistributedFileSystem 会通过 RPC 协议调用 NameNode 来确定请求文件块所在的位置。这里需要注意的是，NameNode 只会返回所调用文件中开始的几个块而不是全部返回（参见图 9-5 第②步）。对于每个返回的块，都包含块所在的 DataNode 地址。随后，这些返回的 DataNode 会按照 Hadoop 定义的集群拓扑结构得出客户端的距离，然后再进行排序（有关拓扑结构的相关知识可参见本书第 10 章）。如果客户端本身就是一个 DataNode，那么它将从本地读取文件。

其次，DistributedFileSystem 会向客户端返回一个支持文件定位的输入流对象

`FSDDataInputStream`，用于给客户端读取数据。`FSDDataInputStream` 包含一个 `DFSInputStream` 对象，这个对象用来管理 `DataNode` 和 `NameNode` 之间的 I/O。

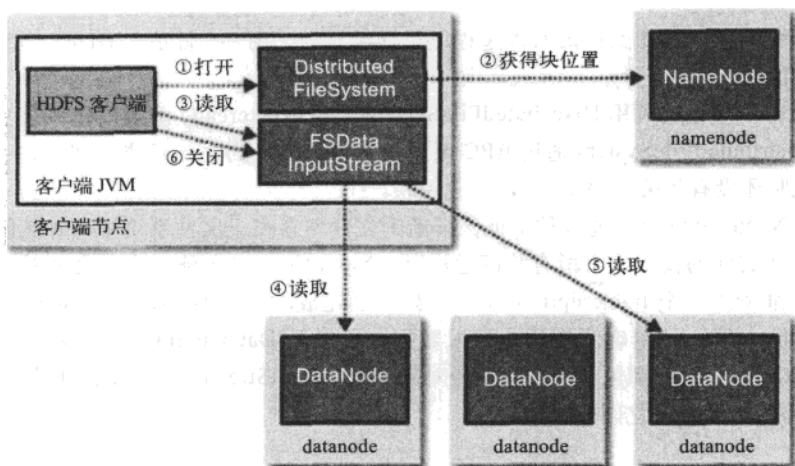


图 9-5 客户端从 HDFS 中读取数据

当以上步骤完成时，客户端便会在这个输入流之上调用 `read()` 函数（参见图 9-5 第③步）。`DFSInputStream` 对象中包含文件开始部分的数据块所在的 `DataNode` 地址，首先它会连接包含文件第一个块最近的 `DataNode`。随后，在数据流中重复调用 `read()` 函数，直到这个块全部读完为止（参见图 9-5 第④步）。当最后一个块读取完毕时，`DFSInputStream` 会关闭连接，并查找存储下一个数据块距离客户端最近的 `DataNode`（参见图 9-5 第⑤步）。以上这些步骤对客户端来说都是透明的。

客户端按照 `DFSInputStream` 打开和 `DataNode` 连接返回的数据流的顺序读取该块，它也会调用 `NameNode` 来检索下一组块所在的 `DataNode` 的位置信息。当客户端完成所有文件的读取时，则会在 `FSDDataInputStream` 中调用 `close()` 函数（参见图 9-5 第⑥步）。

当然，HDFS 会考虑读取中节点出现故障的情况。目前 HDFS 是这样处理的，如果客户端和所连接的 `DataNode` 在读取时出现故障，那么它就会去尝试连接存储这个块的下一个最近的 `DataNode`，同时它会记录这个节点的故障，这样它就不会再去尝试连接和读取块。客户端还会验证从 `DataNode` 传送过来的数据校验和。如果发现一个损坏的块，那么客户端将会再尝试从别的 `DataNode` 读取数据块，向 `NameNode` 报告这个信息，`NameNode` 也会更新保存的文件信息。

这里要关注的一个设计要点是，客户端通过 `NameNode` 引导获取最合适的 `DataNode` 地址，然后直接连接 `DataNode` 读取数据。这种设计的好处是，可以使 HDFS 扩展到更大规模的客户端并行处理，这是因为数据的流动是在所有 `DataNode` 之间分散进行的。同时 `NameNode` 的压力也变小了，使得 `NameNode` 只用提供请求块所在的位置信息就可以了，而

不用通过它提供数据，这样就避免了 NameNode 随着客户端数量的增长而成为系统瓶颈。

9.6.2 文件的写入

本小节将对 HDFS 中文件的写入过程进行详细介绍。图 9-6 就是在 HDFS 中写入一个新文件的数据流图。

第一，客户端通过调用 `DistributedFileSystem` 对象中的 `creat()` 函数创建一个文件（参见图 9-6）。`DistributedFileSystem` 通过 RPC 调用在 NameNode 的文件系统命名空间中创建一个新文件，此时还没有相关的 DataNode 与之关联。

第二，NameNode 会通过多种验证保证新的文件不存在于文件系统中，并且确保请求客户端拥有创建文件的权限。当所有验证通过时，NameNode 会创建一个新文件的记录，如果创建失败，则抛出一个 `IOException` 异常；如果创建成功，则 `DistributedFileSystem` 返回一个 `FSDDataOutputStream` 给客户端用来写入数据。这里 `FSDDataOutputStream` 和读取数据时的 `FSDDataInputStream` 一样都包含一个数据流对象 `DFSOutputStream`，客户端将使用它来处理与 DataNode 和 NameNode 之间的通信。

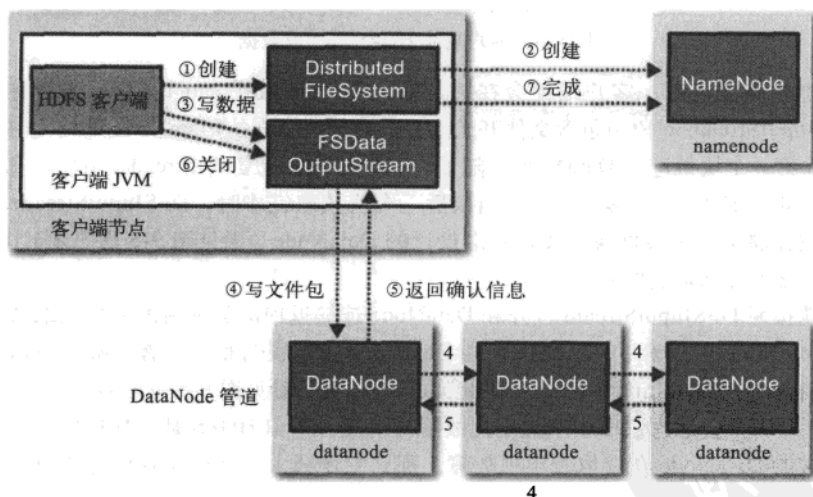


图 9-6 客户端在 HDFS 中写入数据

第三，当客户端写入数据时，`DFSOutputStream` 会将文件分割成包，然后放入一个内部队列，我们称为“数据队列”。`DataStreamer` 会将这些小的文件包放入数据流中，`DataStreamer` 的作用是请求 NameNode 为新的文件包分配合适的 DataNode 存放副本。返回的 DataNode 列表形成一个“管道”，假设这里的副本数是 3，那么这个管道中就会有 3 个 DataNode。`DataStreamer` 将文件包以流的方式传送给队列中的第一个 DataNode，第一个 DataNode 会存储这个包，然后将它推送到第二个 DataNode 中，随后照这样进行，直到管道

中的最后一个 DataNode。

第四，DFSOutputStream 同时也会保存一个包的内部队列，用来等待管道中的 DataNode 返回确认信息，这个队列被称为确认队列（ack queue）。只有当所有管道中的 DataNode 都返回了写入成功的返回信息文件包，才会从确认队列中删除。

当然，HDFS 会考虑写入失败的情况，当数据写入节点失败时，HDFS 会做出以下反应。首先管道会被关闭，任何在确认通知队列中的文件包都会被添加到数据队列的前端，这样管道中失败的 DataNode 都不会丢失数据。当前存放在正常工作的 DataNode 之上的文件块会被赋予一个新的身份，并且和 NameNode 进行关联，这样，如果失败的 DataNode 过段时间后会从故障中恢复出来，其中的部分数据块就会被删除。然后，管道会把失败的 DataNode 删除，文件会继续被写到管道中的另外两个 DataNode 中。最后，NameNode 会注意到现在的文件块副本数没有达到配置属性要求，会在另外的 DataNode 上重新安排创建一个副本，随后的文件会正常执行写入操作。

当然，在文件块写入期间，多个 DataNode 同时出现故障的可能性存在，但是很小。只要 dfs.replication.min 的属性值（默认为 1）成功写入，这个文件块就会被异步复制到集群的其他 DataNode 中，直到满足 dfs.replication.min 的属性值（默认为 3）。

客户端成功完成数据写入的操作后，就会调用 6 种 close() 函数关闭数据流（参见图 9-6 第⑥步）。这步操作会在连接 NameNode 确认文件写入完全之前将所有剩下的文件包放入 DataNode 管道，等待通知确认信息。NameNode 会知道哪些块组成一个文件（通过 DataStreamer 获得块位置信息），这样 NameNode 只要在返回成功标志前等待块被复制（要求复制数量不小于最小量 dfs.replication.min）即可。

9.6.3 一致性模型

文件系统的一致性模型描述了文件读写的可见性。HDFS 牺牲了一些 POSIX 的需求来补偿性能，所以有些操作可能会和传统的文件系统不同。

当创建一个文件时，它在文件系统的命名空间中是可见的，代码如下：

```
Path p = new Path("p");
fs.create(p);
assertThat(fs.exists(p), is(true));
```

但是对这个文件的任何写操作不保证是可见的，即使在数据流已经刷新的情况下，文件的长度很长时间也会显示为 0：

```
Path p = new Path("p");
OutputStream out = fs.create(p);
out.write("content".getBytes("UTF-8"));
out.flush();
assertThat(fs.getFileStatus(p).getLen(), is(0L));
```

一旦一个数据块写入成功，那么大家提出的新请求就可以看到这个块，但大家看不见当

前写入的块。HDFS 提供了使所有缓存和 DataNodes 之间的数据强制同步的方法，这个方法就是 FSDataOutputStream 中的 sync() 函数。当 sync() 函数返回成功时，HDFS 就可以保证此时写入的文件数据是一致的并且对于所有新的用户都是可见的。即使 HDFS 客户端之间发生冲突，也会导致数据丢失，代码如下：

```
Path p = new Path("p");
FSDataOutputStream out = fs.create(p);
out.write("content".getBytes("UTF-8"));
out.flush();
out.sync();
assertThat(fs.getFileStatus(p).getLen(), is(((long) "content".length())));
```

这个操作类似于 Unix 系统中的 fsync 系统调用，为一个文件描述符提交缓存数据，利用 Java API 写入本地数据，这样就可以保证看到刷新流并且同步之后的数据，代码如下：

```
FileOutputStream out = new FileOutputStream(localFile);
out.write("content".getBytes("UTF-8"));
out.flush(); // flush to operating system
out.getFD().sync(); // sync to disk
assertThat(localFile.length(), is(((long) "content".length())));
```

在 HDFS 中关闭一个文件也隐式地执行了 sync() 函数，代码如下：

```
Path p = new Path("p");
OutputStream out = fs.create(p);
out.write("content".getBytes("UTF-8"));
out.close();
assertThat(fs.getFileStatus(p).getLen(), is(((long) "content".length())));
```

下面来了解一致性模型对应用设计的重要性。文件系统的一致性与设计应用程序的方法有关。如果不调用 sync()，那么需要为因客户端或者系统发生故障而丢失部分数据做好准备。对大多数应用程序来说，这是不可接受的，所以需要在合适的地方调用 sync()，比如在写入一定量的数据之后。尽管 sync() 被设计用来最大限度地减少 HDFS 的负担，但是它仍然有不可忽视的开销，所以需要在数据健壮性和吞吐量之间做好权衡，其中一个较好的参考平衡点就是：通过“测试应用程序”来选择不同 sync() 频率间的最佳平衡点。

9.7 HDFS 命令详解

Hadoop 提供了一组 shell 命令在命令行终端对 Hadoop 进行操作。这些操作包括诸如格式化文件系统、上传和下载文件、启动 DataNode、查看文件系统使用情况、运行 jar 包等几乎所有与 Hadoop 相关的操作。本节将具体介绍 HDFS 的相关命令和操作。

9.7.1 通过 distcp 进行并行复制

Java API 等多种接口对 HDFS 访问模型都集中于单线程的存取，如果要对一个文件

集进行操作,就需要编写一个程序来执行并行操作。HDFS 提供了一个非常实用的程序 `distcp`,用来在 Hadoop 文件系统中并行地复制大数据量文件。`distcp` 一般适用于在两个 HDFS 集群间传送数据的情况。如果两个集群都运行在同一个 Hadoop 版本上,那么可以使用 HDFS 模式:

```
hadoop distcp hdfs://NameNode1/foo hdfs://NameNode2/bar
```

这条命令会将第一个集群中的 `/foo` 文件夹以及文件夹下的文件复制到第二个集群中的 `/bar` 目录下,即在第二个集群中会以 `/bar/foo` 的形式出现。如果 `/bar` 目录不存在,则系统会新建一个。也可以指定多个数据源,并且所有的内容都会被复制到目标路径。需要注意的是,源路径必须是绝对路径。

默认情况下,虽然 `distcp` 会跳过在目标路径上已经存在的文件,但是通过 `-overwrite` 选项可以选择对这些文件进行覆盖重写,也可以使用 `-update` 选项仅对更新过的文件进行重写。

`distcp` 操作有很多选项可以设置,比如忽略失败、限制文件或者复制的数据量等。若选择直接输入指令或者不附加选项则可以查看此操作的使用说明。具体实现时,`distcp` 操作会被解析为一个 MapReduce 操作来执行。当没有 reducer 操作时,“复制操作”被作为“map 操作”并行地在集群节点中运行。因此,每个文件都可被当做一个“map 操作”来执行“复制操作”。另外,`distcp` 会通过执行文件聚集捆绑操作,从而尽可能地保证每个 map 操作执行了相同数量的数据。那么,执行 `distcp` 时,map 操作如何确定呢?

系统需要保证每个 map 操作执行的数据量是合理的,以尽可能地减少 map 执行的开销,而按规定,每个 map 最少要执行 256MB 的数据量(除非复制的全部数据量小于 256MB),比如,如果要复制 1GB 的数据,那么系统就会分配 4 个 map 任务,当数据量非常大时,就需要对执行的 map 任务数进行限制,以限制网络带宽和集群的使用率。默认情况下,每个集群的一个节点最多执行 20 个 map 任务。比如,要复制 1000GB 数据到 100 节点的集群中,那么系统就会分配 2000 个 map 任务(每个节点 20 个),也就是说,每个节点会平均复制 512MB。还可以通过调整 `distcp` 的 `-m` 参数来减少 map 任务量,比如 `-m 1000` 就意味着分配 1000 个 maps,每个节点分配 1GB 数据量。

如果尝试使用 `distcp` 进行 HDFS 集群间的复制,使用 HDFS 模式之后,HDFS 运行在不同的 Hadoop 版本之上,复制将会因为 RPC 系统的不匹配而失败。为了纠正这个错误,可以使用基于 HTTP 的 HFTP 进行访问。因为任务要在目标集群中执行,所以 HDFS 的 RPC 版本需要匹配,在 HFTP 模式下运行的代码如下:

```
hadoop distcp hftp://NameNode1:50070/foo hdfs://NameNode2/bar
```

需要注意的是,要定义访问源的 URI 中 NameNode 的网络接口,这个接口会通过 `dfs.address` 的属性值设定,默认值为 50070。

9.7.2 HDFS 的平衡

当复制大规模数据到 HDFS 时,要考虑的一个重要因素是文件系统的平衡。当系统中

的文件块能够很好地均衡分布到集群的各个节点时，HDFS 才能够更好地工作，所以要保证 `distcp` 操作不会打破这个平衡。回到前面复制 1000GB 数据的例子，当设定 `-m` 为 1，就意味着 1 个 `map` 操作可以完成 1000GB 的操作。这样不仅会让复制操作非常慢，而且不能充分利用集群的性能。最重要的是，复制文件的第一个块都要存储在执行 `map` 任务的那个节点上，直到这个节点的磁盘被写满，显然这个节点是不平衡的。通常我们通过设置更多的超过集群节点的 `map` 任务数来避免不平衡情况的发生，所以最好的选择是刚开始并且还是使用的默认属性值，每个节点分配 20 个 `map` 任务。

当然，我们不能保证集群总能够保持平衡，有时可能会限制 `map` 的数量以便节点可以被其他任务使用，这样 HDFS 还提供了一个工具 `balancer`（参见第 10 章）来改变集群中的文件块存储的平衡。

9.7.3 使用 Hadoop 归档文件

在 9.2 节中介绍过，HDFS 采用块的方式来存储文件，在系统运行时，文件块的元数据信息会被存储在 `NameNode` 的内存中，因此，对 HDFS 来说，大规模存储小文件显然是低效的，很多小文件会耗尽 `NameNode` 的大部分内存。

Hadoop 归档文件和 `HAR` 文件是可以将文件高效地放入 HDFS 块中的文件存档设备，在减少 `NameNode` 内存使用的同时，仍然允许对文件进行透明访问。具体来说，Hadoop 归档文件可以作为 `MapReduce` 的输入。这里需要注意的是，小文件并不会占用太多的磁盘空间，比如设定一个大小为 128MB 的块文件来存储 1MB 的文件，实际上存储这个文件只需要 1MB 磁盘空间，而不是 128MB。

Hadoop 归档文件是通过 `archive` 命令工具根据文件集合创建的。因为这个工具需要运行一个 `MapReduce` 来并行处理输入文件，所以需要有一个运行 `MapReduce` 的集群。而 HDFS 中有些文件是需要进行归档的，例如：

```
hadoop fs -lsr /user/admin/In/
-rw-r--r--   3 admin\admin supergroup      13 2010-10-18 20:15 /user/admin/In/hello.c.txt
-rw-r--r--   1 admin\admin supergroup      13 2010-10-17 15:13 /user/admin/In/hello.txt
```

运行 `archive` 命令如下：

```
hadoop archive -archiveName files.har /user/admin/In/ /user/admin/
```

```
10/10/18 20:46:47 INFO mapred.JobClient: Running job: job_201010182044_0001
10/10/18 20:46:48 INFO mapred.JobClient: map 0% reduce 0%
10/10/18 20:47:21 INFO mapred.JobClient: map 100% reduce 0%
10/10/18 20:47:39 INFO mapred.JobClient: map 100% reduce 100%
10/10/18 20:47:41 INFO mapred.JobClient: Job complete: job_201010182044_0001
10/10/18 20:47:41 INFO mapred.JobClient: Counters: 17
10/10/18 20:47:41 INFO mapred.JobClient: Job Counters
10/10/18 20:47:41 INFO mapred.JobClient: Launched reduce tasks=1
10/10/18 20:47:41 INFO mapred.JobClient: Launched map tasks=1
```

```

10/10/18 20:47:41 INFO mapred.JobClient: FileSystemCounters
10/10/18 20:47:41 INFO mapred.JobClient: FILE_BYTES_READ=540
10/10/18 20:47:41 INFO mapred.JobClient: HDFS_BYTES_READ=531
10/10/18 20:47:41 INFO mapred.JobClient: FILE_BYTES_WRITTEN=870
10/10/18 20:47:41 INFO mapred.JobClient: HDFS_BYTES_WRITTEN=305
10/10/18 20:47:41 INFO mapred.JobClient: Map-Reduce Framework
10/10/18 20:47:41 INFO mapred.JobClient: Reduce input groups=6
10/10/18 20:47:41 INFO mapred.JobClient: Combine output records=0
10/10/18 20:47:41 INFO mapred.JobClient: Map input records=6
10/10/18 20:47:41 INFO mapred.JobClient: Reduce shuffle bytes=0
10/10/18 20:47:41 INFO mapred.JobClient: Reduce output records=0
10/10/18 20:47:41 INFO mapred.JobClient: Spilled Records=12
10/10/18 20:47:41 INFO mapred.JobClient: Map output bytes=280
10/10/18 20:47:41 INFO mapred.JobClient: Map input bytes=399
10/10/18 20:47:41 INFO mapred.JobClient: Combine input records=0
10/10/18 20:47:41 INFO mapred.JobClient: Map output records=6
10/10/18 20:47:41 INFO mapred.JobClient: Reduce input records=6

```

在命令行中，第一个参数是归档文件的名称，这里是 file.har 文件；第二个参数是要归档的文件源，这里我们只归档一个源文件夹，即 HDFS 下 /user/admin/In/ 中的文件，但事实上，archive 命令可以接收多个文件源；最后一个参数，即本例中的 /user/admin/ 是 HAR 文件的输出目录。可以看到这个命令的执行流程为一个 MapReduce 任务。

下面我们来看这个归档文件是怎么创建的：

```

hadoop fs -ls /user/admin/In/ /user/admin/
Found 2 items
-rw-r--r--  3 admin\admin supergroup    13 2010-10-18 20:15 /user/admin/In/hello.c.txt
-rw-r--r--  1 admin\admin supergroup    13 2010-10-17 15:13 /user/admin/In/hello.txt
Found 3 items
drwxr-xr-x  - admin\admin supergroup    0 2010-10-18 20:15 /user/admin/In
drwxr-xr-x  - admin\admin supergroup    0 2010-10-18 18:53 /user/admin/admin
drwxr-xr-x  - admin\admin supergroup    0 2010-10-18 20:47 /user/admin/files.har

```

这个目录列表展示了一个 HAR 文件的组成：两个索引文件和部分文件（part file）的集合。这里的部分文件包含已经连接在一起的大量源文件的内容，同时索引文件可以检索部分文件中的归档文件，包括它的长度起始位置等。但是，这些细节在使用 har URI 模式访问 HAR 文件时多数都是隐藏的。HAR 文件系统是建立在底层文件系统上的（此处是 HDFS），以下命令以递归的方式列出了归档文件中的文件：

```

hadoop fs -lsr har:///user/admin/files.har

drw-r--r--  -admin\admin supergroup    0 2010-10-18 20:47 /user/admin/files.har/
user
drw-r--r--  -admin\admin supergroup    0 2010-10-18 20:47 /user/admin/files.har/
user/admin
drw-r--r--  -admin\admin supergroup    0 2010-10-18 20:47 /user/admin/files.har/
user/admin/In

```

```

-rw-r--r-- 10 admin\admin supergroup          13 2010-10-18 20:47 /user/admin/
files.har/user/admin/In/hello.c.txt
-rw-r--r-- 10 admin\admin supergroup          13 2010-10-18 20:47 /user/admin/
files.har/user/admin/In/hello.txt

```

如果 HAR 文件所在的文件系统是默认的文件系统，那么这里的内容就非常直观和易懂。但是，如果你想要在其他文件系统中使用 HAR 文件，就需要使用不同格式的 URI 路径。下面两个命令即具有相同的作用：

```

hadoop fs -lsr har:/// user/admin/files.har/my/files/dir
hadoop fs -lsr har://hdfs-localhost:8020/ user/admin/files.har/my/files/dir

```

第二个命令，它仍然使用 HAR 模式描述一个 HAR 文件系统，但是使用 HDFS 作为底层的文件系统模式，它在 HAR 模式之后紧跟 HDFS 系统的主机和端口号。并且，HAR 文件系统会将 har URI 转换为底层的文件系统访问 URI。在本例中即为 hdfs://localhost:8020/user/admin/archive/files.har，文件的剩余部分路径即为文件归档部分的路径 /my/files/dir。

如果想要删除 HAR 文件，需要使用删除的递归格式，这是因为底层的文件系统 HAR 文件是一个目录，删除命令为 `hadoop fs -rmr /user/admin/files.har`。

关于 HAR 文件，我们还需要了解它的一些不足。当创建一个归档文件时，还会创建原始文件的一个副本，这样就需要额外的磁盘空间（尽管归档完成后会删除原始文件）。而且当前还没有针对归档文件的压缩方法，只能对写入归档文件的原始文件进行压缩。归档文件一旦创建就不能改变，如果要增加或者删除文件，就要重新创建。事实上，这对于那些写后不能更改的文件不是问题，因为可以按日或者按周定期进行成批归档。

如前所述，HAR 文件可以作为 MapReduce 的一个输入文件，然而，没有一个基于归档的 InputFormat 可以将多个文件打包到一个单一的 MapReduce 中去。所以，即使是 HAR 文件，处理小文件时的效率仍然不高。

9.7.4 其他命令

其他相关命令还包括以下这些：

- ❑ `namenode -format`：格式化 DFS 文件系统
- ❑ `secondaryNameNode`：运行 DFS 的 Secondarynamenode 进程
- ❑ `namenode`：运行 DFS 的 namenode 进程
- ❑ `datanode`：运行 DFS 的 datanode 进程
- ❑ `dfsadmin`：运行 DFS 的管理客户端
- ❑ `mradmin`：运行 MapReduce 的管理客户端
- ❑ `fsck`：运行 HDFS 的检测进程
- ❑ `fs`：运行一个文件系统工具
- ❑ `balancer`：运行一个文件系统平衡进程
- ❑ `jobtracker`：运行一个 jobtracker 进程

- ❑ pipes：运行一个 Pipes 任务
- ❑ tasktracker：运行一个 tasktracker 进程
- ❑ job：管理运行中的 MapReduce 任务
- ❑ queue：获得运行中的 MapReduce 队列的信息
- ❑ version：打印版本号
- ❑ jar <jar>：运行一个 jar 文件
- ❑ daemonlog：读取 / 设置守护进程的日志记录级别

相信大家已经对这些命令中的一部分很熟悉了，比如，在命令行终端中 jar 是用来运行 Java 程序的，version 命令可以查看 Hadoop 的当前版本，或者在安装时必须运行的 namenode-format 命令。这一小节我们介绍的是与 HDFS 有关的命令，其中与 HDFS 相关的命令有如下几个：secondarynamenode、namenode、datanode、dfsadmin、fsck、fs、balancer、distcp 和 archives。

它们的统一格式如下：

```
bin/hadoop command [genericOptions] [commandOptions]
```

其中只有 dfsadmin、fsck、fs 具有选项 genericOptions 和 commandOptions，其余的命令只有 commandOptions。下面先介绍只有 commandOptions 选项的命令。

❑ distcp。distcp 命令用于 DistCp（即 Dist 分布式，Cp 拷贝）分布式拷贝。用于在集群内部及集群之间拷贝数据。

❑ archives。archives 命令是 Hadoop 定义的档案格式。archive 对应一个文件系统，它的扩展名是 .har，包含元数据和数据文件。

这两个命令在前文中已有介绍，这里就不再赘述了。

❑ datanode。datanode 命令要简单一些。你可以使用如下命令将 Hadoop 回滚到前一个版本，它的用法如下：

```
hadoop datanode [-rollback]
```

❑ namenode。namenode 命令稍稍复杂一些，它的用法如下：

```
hadoop namenode
[-format] // 格式化 namenode
[-upgrade] // 在 Hadoop 升级后，应该使用这个命令启动 namenode
[-rollback] // 使用 namenode 回滚前一个版本
[-finalize] // 删除文件系统的前一个状态，这会导致系统不能回滚到前一个状态
[-importCheckpoint] // 复制和备份 checkpoint 的状态到当前 checkpoint
```

❑ secondarynamenode。secondarynamenode 命令的用法如下：

```
hadoop secondarynamenode
[-checkpoint [force]]
// 当 EditLog 超过规定大小（默认 64MB）时，启动检查 secondarynamenode 的 checkpoint 过程；如果启用 force 选项，则强制执行 checkpoint 过程
```

```
[-geteditssize]
// 在终端上显示 EditLog 文件的大小
```

- ❑ balancer。balancer 命令如解释中所说，用于分担负载。很多原因都会造成数据在集群内分布不均衡，一般来说，当集群中添加新的 datanode 时，可以使用这个命令来进行负载均衡。其用法如下：

```
hadoop balancer
```

接下来的 dfsadmin、fsck、fs 这三个命令有一个共同的选项 genericOptions，这个选项一般与系统相关，其用法如下：

```
-conf <configuration file>      // 指定配置文件
-D <property=value>              // 指定某属性的属性值
-fs <local|namenode:port>        // 指定 datanode 及其端口
```

- ❑ dfsadmin。在 dfsadmin 命令中可以执行一些类似 Windows 中高级用户才能执行的命令，比如升级和回滚等。其用法如下：

```
hadoop dfsadmin [GENERIC_OPTIONS]
[-report] // 在终端上显示文件系统的基本信息
[-safemode enter | leave | get | wait]/* Hadoop 的安全模式及相关维护；在安全模式中系统是只读的，数据块也不可以删除或复制
[-refreshNodes] [-finalizeUpgrade]/* 重新读取 hosts 和 exclude 文件，将新的被允许加入到集群中的 datanode 连入，同时断开与那些从集群出去的 datanode 的连接 */
[-upgradeProgress status | details | force]// 获得当前系统的升级状态和细节，或者强制进行升级过程
[-metasave filename]// 将 namenode 的主要数据结构保存到指定目录下
[-setQuota <quota> <dirname>...<dirname>]// 为每个目录设定配额
[-clrQuota <dirname>...<dirname>]// 清除这些目录的配额
[-setSpaceQuota <quota> <dirname>...<dirname>]// 为每个目录设置配额空间
[-clrSpaceQuota <dirname>...<dirname>]// 清除这些目录的配额空间
[-help [cmd]]// 显示命令的帮助信息
```

- ❑ fsck。fsck 在 HDFS 中被用来检查系统中的不一致情况，比如某文件只有目录，但数据块已经丢失或副本数目不足。与 Linux 不同，这个命令只用于检测，不能进行修复。其使用方法如下：

```
hadoop fsck [GENERIC_OPTIONS] <path> [-move | -delete | -openforwrite] [-files [-blocks
[-locations | -racks]]]
//<path> 检查的起始目录
//--move 移动受损文件到 /lost+found
//--delete 删除受损文件
//--openforwrite 在终端上显示为写打开的文件
//--files 在终端上显示正被检查的文件
//--blocks 在终端上显示块信息
//--location 在终端上显示每个块的位置
//--rack 显示 datanode 的网络拓扑结构图
```

□ fs. fs 可以说是 HDFS 最常用的命令，这是一个高度类似 Linux 文件系统的命令集。你可以使用这些命令查看 HDFS 上的目录结构文件、上传和下载文件、创建文件夹、复制文件等。其使用方法如下：

```
hadoop fs [genericOptions]
[-ls <path>] // 显示目标路径中当前目录下的所有文件
[-lsr <path>] // 递归显示目标路径下的所有目录及文件（深度优先）
[-du <path>] // 以字节为单位显示目录中所有文件的大小，或该文件的大小（如果目标为文件）
[-dus <path>] // 以字节为单位显示目标文件大小（用于查看文件夹大小）
[-count [-q] <path>] // 将目录的大小、包含文件（包括文件）个数的信息输出到屏幕（标准 stdout）
[-mv <src> <dst>] // 把文件或目录移动到目标路径，这个命令允许同时移动多个文件，但是只允许
                    移动到一个目标路径中，参数中的最后一个文件夹即为目标路径
[-cp <src> <dst>] // 复制文件或目录到目标路径，这个命令允许同时复制多个文件，如果复制多个
                    文件，目标路径必须是文件夹
[-rm [-skipTrash] <path>] // 删除文件，这个命令不能删除文件夹
[-rmr [-skipTrash] <path>] // 删除文件夹及其下的所有文件
[-expunge]
[-put <localsrc> ... <dst>] // 从本地文件系统上传文件到 HDFS 中
[-copyFromLocal <localsrc> ... <dst>] // 与 put 相同，在 0.18 版中，只允许上传一个文件，
                    但是在 0.20 中，这两个命令已经完全一样了
[-moveFromLocal <localsrc> ... <dst>] // 与 put 相同，但是文件上传之后会从本地文件系统中删除
[-get [-ignoreCrc] [-crc] <src> <localdst>] // 复制文件到本地文件系统。这个命令可以选择是否忽视校验和，忽视校验和下载主要用于挽救那些已经发生错误的文件
[-getmerge <src> <localdst> [addnl]] // 对源目录中的所有文件进行排序并写入目标文件中，
                    文件之间以换行符分隔
[-cat <src>] // 在终端显示（标准输出 stdout）文件中的内容，类似 Linux 系统中的 cat
[-text <src>]
[-copyToLocal [-ignoreCrc] [-crc] <src> <localdst>] // 与 get 相同
[-moveToLocal [-crc] <src> <localdst>] // 这个命令在 0.20 版本中还没有完成
[-mkdir <path>] // 创建文件夹
[-setrep [-R] [-w] <rep> <path/file>] // 改变一个文件的副本个数。参数 -R 可以递归地对该
                    目录下的所有文件执行统一操作
[-touchz <path>] // 类似 Linux 中的 touch，创建一个空文件
[-test [-ezd] <path>] // 将源文件输出为文本格式并显示到终端上，通过这个命令可以查看
                    TextRecordInputStream（SequenceFile 等）或 zip 文件
[-stat [format] <path>] // 以指定格式返回路径的信息
[-tail [-f] <file>] // 在终端上显示（标准输出 stdout）文件的最后 1kb 内容。-f 选项的行为
                    与 Linux 中一致，会持续检测新添加到文件中的内容，这在查看日志文件
                    时会显得非常方便
[-chmod [-R] <MODE[,MODE]... | OCTALMODE> PATH...] // 改变文件的权限，只有文件的所有者或是超级用户才能使用这个命令。-R 可以递归地改变文件夹内所有文件的权限
[-chown [-R] [OWNER][:[GROUP]] PATH...] // 改变文件的拥有者，-R 可以递归地改变文件夹内所有文件的拥有者。同样，这个命令只有超级用户才能使用
[-chgrp [-R] GROUP PATH...] // 改变文件所属的组，-R 可以递归地改变文件夹内所有文件所属的组。这个命令只有超级用户才能使用
[-help [cmd]] // 这是命令的帮助信息
```

在这些命令中，参数 <path> 的完整格式是 `hdfs://NameNodeIP:port/`，比如你的 NameNode 地址是 192.168.0.1，端口是 9000，那么，如果你想访问 HDFS 上路径为 `/user/root/hello` 的文件，则需要输入的地址是 `hdfs://192.168.0.1:9000/user/root/hello`。在 Hadoop 中，如果参数 <path> 没有 NameNodeIP，那么会默认按照 `core-site.xml` 中属性 `fs.default.name` 的设置，然后附加 `“/user/ 你的用户名”` 作为路径，这是为了方便使用以及对不同的用户进行区分。

9.8 小结

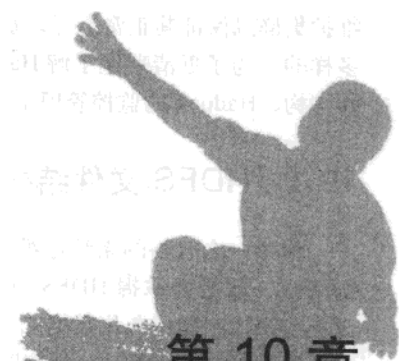
本章深入介绍了分布式文件系统 HDFS。

首先，对 Hadoop 的文件系统进行了总体的概括，随后针对 HDFS 进行了简单介绍，分析了它的研究背景和设计基础。有了这样的背景知识，就可以在随后的章节中更好地理解 HDFS 的功能和实现。本章还从结构上对 HDFS 进行了描述，给出了 HDFS 的相关概念，包括块、NameNode、DataNode 等。通过对 HDFS 概念的学习，还可以了解 HDFS 的体系结构。

其次，在基本概念的基础上，我们还介绍了 HDFS 的基本操作及其操作接口。HDFS 为开发者提供了丰富的接口，包括命令行接口和各种方便使用的 Java 接口，可以通过 Java API 对 HDFS 中的文件执行常规的文件操作。不仅如此，在使用 API 对 HDFS 文件系统进行管理的基础上，还对 HDFS 中文件流的读写进行了详细的介绍，这对更深入地了解 HDFS 有很大帮助。

最后，本章对 HDFS 的命令进行了详细的讲解，并对其中特有的 `distcp` 操作和归档文件进行了具体的说明，理解它们可以更好地帮助大家了解 Hadoop 的文件系统。





第 10 章

Hadoop 的管理

本章内容

- ☐ HDFS 文件结构
- ☐ Hadoop 的状态监视和管理工具
- ☐ Hadoop 集群的维护
- ☐ 小结



在第2章我们已经详细介绍了如何安装和部署 Hadoop 集群，本章我们将具体介绍如何维护集群以保证其正常运行。毋庸置疑，维护一个大型集群的稳定运行是必要的，手段也是多样的。为了更清晰地了解 Hadoop 的集群管理相关内容，本章我们主要从 HDFS 本身的文件结构、Hadoop 的监控管理工具，以及集群常用的维护功能三方面进行讲解。

10.1 HDFS 文件结构

作为一名合格的系统运维人员，首先要全面掌握系统的文件组织目录，对于 Hadoop 系统来说，就是要掌握 HDFS 中的 NameNode、DataNode、Secondary NameNode 是如何在磁盘上组织、存储持久化数据的。只有这样，当遇到问题时，管理员才能借助系统本身的文件存储机制来快速诊断、分析问题。下面我们就从 HDFS 的几个方面来分别介绍。

1. NameNode 的文件结构

最新格式化的 NameNode 会创建以下的目录结构：

```
${dfs.name.dir}/current/VERSION
                        /edits
                        /fsimage
                        /fstime
```

其中，dfs.name.dir 属性是一个目录列表，是每个目录的镜像。VERSION 文件是 Java 属性文件，其中包含运行 HDFS 的版本信息。下面是一个典型的 VERSION 文件所包含的内容：

```
#Wed Mar 23 16:03:27 CST 2011
namespaceID=1064465394
cTime=0
storageType=NAME_NODE
layoutVersion=-18
```

其中，namespaceID 是文件系统的唯一标识符，当文件系统第一次格式化时便会被创建，这个标识符也要求各 DataNode 节点和 NameNode 保持一致。NameNode 会使用它识别新的 DataNode，DataNode 只有在向 NameNode 注册后才会获得此 namespaceID。cTime 属性标记了 NameNode 存储空间创建的时间。对于新格式化的存储空间，这里的属性值虽为 0，但是只要文件系统被更新，它就会更新到一个新的时间戳上。storageType 指出此存储目录包含一个 NameNode 的数据结构，在 DataNode 中它的属性值为 DATA_NODE。

layoutVersion 是一个负的整数，定义了 HDFS 持久数据结构的版本。请注意，该版本号与 Hadoop 的发行版本号无关。每次 HDFS 的布局发生改变，该版本号就会递减（比如 -18 版本号之后是 -19），在这种情况下，HDFS 就需要更新升级了，因为一个新的 NameNode 或 DataNode 如果还处在旧版本上，那么系统就无法正常运行，各节点的版本号需要一致。

在 NameNode 的存储目录中包含 edits、fsimage、fstime 三个文件。它们都是二进制的

文件，可以通过 HadoopWritable 对象进行序列化。下面将深入介绍 NameNode 的工作原理以便更清晰地理解这三个文件的作用。

2. 编辑日志 (edit log) 及文件系统映像 (filesystem image)

当客户端执行写操作时，首先 NameNode 会在编辑日志中写下记录，并在内存中保存一个文件系统元数据，这个描述符会在编辑日志有了改动后进行更新。内存中的元数据用来提供读数据请求服务。

编辑日志会在每次成功操作之后，且成功代码尚未返回给客户端之前进行刷新和同步。对于要写入多个目录的操作，该写入流要刷新和同步到所有的副本上，这就保证了操作不会因故障而丢失数据。

fsimage 文件是文件系统元数据的持久性检查点，和编辑日志不同，它不会在每个文件系统的写操作后进行更新，因为写出 fsimage 文件会非常慢 (fsimage 可能增长到 GB 大小)。这种设计并不会影响系统的恢复力，因为如果 NameNode 失败，那么元数据的最新状态可以通过从磁盘中读出 fsimage 文件加载到内存中来进行重建恢复，然后重新执行编辑日志中的操作。事实上，这也正是 NameNode 启动时要做的事情。一个 fsimage 文件中包含以序列化格式存储的文件系统目录和文件 inodes。每个 inodes 表征一个文件或目录的元数据信息，以及文件的副本数、修改和访问时间等信息。

正如上面所描述，Hadoop 文件系统会出现编辑日志的不断增长情况。尽管在 NameNode 运行期间不会对系统造成影响，但是，如果 NameNode 重新启动，它将会花很长时间来运行编辑日志中的每个操作。在此期间 (即安全模式时间)，文件系统还是不可用的，通常来说这是不符合应用需求的。

为了解决这个问题，Hadoop 在 NameNode 之外的节点上运行了一个 Secondary NameNode 进程，它的任务就是为原 NameNode 内存中的文件系统元数据产生检查点。其实 Secondary NameNode 是一个辅助 NameNode 处理 fsimage 和编辑日志的节点，它从 NameNode 中拷贝 fsimage 和编辑日志到临时目录并定期合并成一个新的 fsimage。随后它会将新的 fsimage 上传到 NameNode，这样，NameNode 便可更新 fsimage 并删除原来的编辑日志了。下面我们看图 10-1 对检查点处理过程进行描述。

检查点处理过程如下：

- 1) Secondary NameNode 首先请求原 NameNode 进行 edits 的滚动，这样新的编辑操作就能够进入一个新的文件中了。
- 2) Secondary NameNode 通过 HTTP 方式读取原 NameNode 中的 fsimage 及 edits。
- 3) Secondary NameNode 读取 fsimage 到内存中，然后执行 edits 中的每个操作，并创建一个新的统一的 fsimage 文件。
- 4) Secondary NameNode (通过 HTTP 方式) 将新的 fsimage 发送到原 NameNode。
- 5) 原 NameNode 用新的 fsimage 替换旧的 fsimage，旧的 edits 文件用步骤①中的 edits 进行替换，同时系统会更新 fsimage 文件到文件的记录检查点时间。

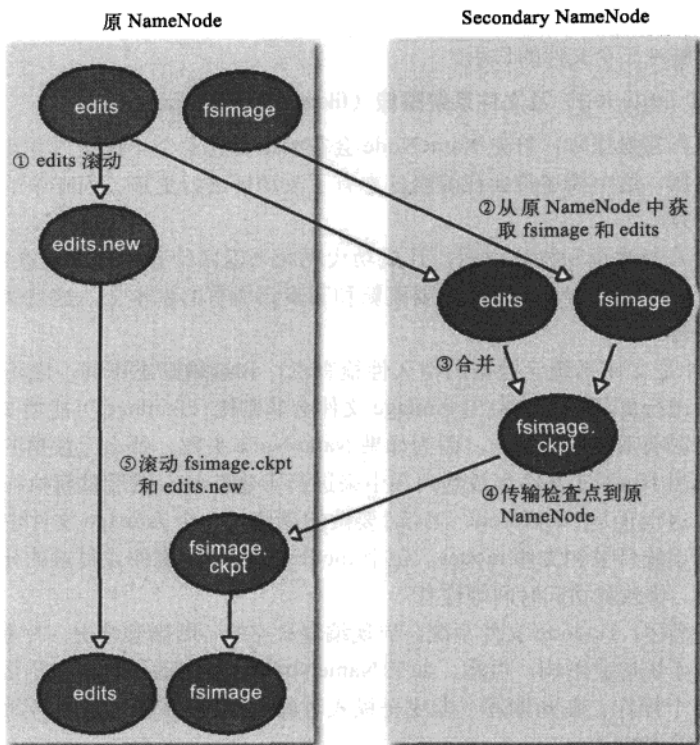


图 10-1 检查点处理过程

在这个过程结束后, NameNode 就有了最新的 `fsimage` 文件和更小的 `edits` 文件。事实上, NameNode 在安全模式下, 管理员可以通过以下命令运行这个过程:

```
hadoop dfsadmin -saveNamespace
```

这个过程清晰地表明了为什么 Secondary NameNode 要有和原 NameNode 一样的内存需求 (它要把 `fsimage` 加载到内存中), 因此 Secondary NameNode 在集群中也需要有专用机器。

有关检查点的时间表由两个配置参数决定。Secondary NameNode 每隔一个小时会插入一个检查点 (`fs.checkpoint.period`, 以秒为单位), 如果编辑日志达到 64MB (`fs.checkpoint.size`, 以字节为单位), 则间隔时间更短, 每隔 5 分钟会检查一次。

3. SecondaryNameNode 的目录结构

Secondary NameNode 在每次处理过程结束后都有一个检查点, 这个检查点可以在一个子目录 `/previous.checkpoint` 中找到。这个检查点可以用来作为 NameNode 的元数据备份源, 其目录如下:

```
${fs.checkpoint.dir}/current/VERSION
```



```

        /edits
        /fsimage
        /fstime
    /previous.checkpoint/VERSION
        /edits
        /fsimage
        /fstime

```

以上目录和 Secondary NameNode 的 /current 目录结构是完全相同的。这样设计的目的是万一整个 NameNode 发生故障，并且没有用于恢复的备份，甚至 NFS 中也没有，就可以直接从 Secondary NameNode 恢复。具体方式有两种，第一是直接复制相关的目录到新的 NameNode 中。第二是启动 NameNode 的守护进程时，Secondary NameNode 可以使用 -importCheckpoint 选项，接管并作为新的 NameNode 继续运行任务。-importCheckpoint 选项将加载 fs.checkpoint.dir 属性定义目录中的最新检查点的 NameNode 数据，但这种操作只有在 dfs.name.dir 所指定的目录下没有元数据的情况下才进行，这样就避免了重写之前元数据的风险。

4. DataNode 的目录结构

DataNode 不需要进行格式化，它们会在启动时自己创建存储目录，其中关键的文件和目录如下：

```

${dfs.data.dir}/current/VERSION
    /blk_<id_1>
    /blk_<id_1>.meta
    /blk_<id_2>
    /blk_<id_2>.meta
    /...

    /subdir0/
    /subdir1/
    /...
    /subdir63/

```

DataNode 的 VERSION 文件和 NameNode 的非常类似，内容如下：

```

#Tue Mar 10 21:32:31 GMT 2010
namespaceID=134368441
storageID=DS-547717739-172.16.85.1-50010-1236720751627
cTime=0
storageType=DATA_NODE
layoutVersion=-18

```

其中，namespaceID 和 cTime 与 layoutVersion 及 NameNode 的值都一样，namespaceID 在第一次连接 NameNode 时就会从中获取此值。storageID 相对于 DataNode 来说是唯一的，用于 NameNode 标识 DataNode。storageType 将这个目录标志为 DataNode 数据存储目录。

DataNode 中 current 目录下的其他文件都有 blk_refix 前缀，它有两种类型：

- 1) HDFS 中的文件块本身，存储的是原始文件内容；

2) 块的元数据信息 (使用 .meta 后缀标识)。一个块文件由存储的原始文件字节组成, 元数据文件由一个包含版本和类型信息的头文件与一系列块的区域校验和组成。

当目录中存储的块数量增加到一定规模时, DataNode 会创建一个新的目录, 用于保存新的块及元数据。目录块中的块数量达到 64 (可由 dfs.edatanode.numblocks 属性值确定) 时, 便会新建一个子目录, 这样就会形成一个更宽的文件树结构, 避免了系统由于存储大量数据块而导致目录很深, 影响检索性能。通过这样的措施, 数据节点可以确保每个目录中的文件块数是可控的, 这也避免了一个目录中存在过多的文件。

10.2 Hadoop 的状态监视和管理工具

对于一个系统运维管理员来说, 进行系统监控是必须的。监控的目的就是让我们知道系统何时出现问题, 并找到问题出在哪里, 从而做出相应的处理。管理守护进程对监控 NameNode、DataNode 和 JobTracker 是非常重要的。实际运行中, 因为 DataNode 及 TaskTracker 的故障可能随时出现, 所以集群需要提供额外的功能以应对少部分节点所出现的故障。管理员也要隔一段时间执行一些监测任务, 以获知当前集群的运行状态。本节我们将详细介绍 Hadoop 是如何实现系统监控的。

10.2.1 审计日志

HDFS 通过审计日志可以实现记录文件系统中所有文件访问请求的功能, 其审计日志功能通过 log4j 实现, 但是在默认配置下这个功能是关闭的: log 的记录等级在 log4j.properties 中被设置为 WARN:

```
log4j.logger.org.apache.hadoop.fs.FSNamesystem.audit=WARN
```

此处将 WARN 修改为 INFO, 便可打开审计日志功能。这样每个 HDFS 事件后, 系统就会在 NameNode 的 log 文件中写入一行记录。下面是一个请求 /usr/hadoop 文件的例子:

```
2010-03-13 07:11:22,982 INFO org.apache.hadoop.hdfs.server.NameNode.FSNamesystem.audit:
ugi=admin,staff,admin ip=/127.0.0.1 cmd=listStatus src=/user/admin=null
perm=null
```

关于 log4j 还有很多其他配置可改, 比如可以将审计日志从 NameNode 的日志文件中分离出来等。具体操作可查看 Hadoop 的 Wiki: <http://wiki.apache.org/hadoop/HowToConfigure>。

10.2.2 监控日志

所有 Hadoop 守护进程都会产生一个日志文件, 这对管理员来说非常重要, 下面我们就介绍如何使用这些日志文件。

1. 设置日志级别

当进行故障调试排除时, 很有必要临时调整日志的级别, 以获得系统不同类型的信

息。log4j 日志一般包含以下几个级别：OFF、FATAL、ERROR、WARN、INFO、DEBUG、ALL，或者用户自定义的级别。

Hadoop 守护进程有一个网络页面可以用来调整任何 log4j 日志的级别，在守护进程的网络 UI 附后缀 /logLevel 即可访问。按照规定，日志的名称和它所对应的执行日志记录的类名必须是一样的，你可以通过查找源代码找到日志名称。例如，为了调试 JobTracker 类的日志，可以访问 JobTracker 的网络 UI：<http://jobtracker-host:50030/logLevel>，同时设置日志名称 org.apache.hadoop.mapred.JobTracker 到层级 DEBUG。当然也可以通过命令行进行调整，代码如下：

```
hadoop daemonlog -setlevel jobtracker-host:50030/
org.apache.hadoop.mapred.JobTracker DEBUG
```

通过命令行修改的日志级别需要在守护进程重启时重置，如果想要持久化地改变日志级别，那么只要改变 log4j.properties 文件内容即可。我们可以在文件中加入以下行：

```
log4j.logger.org.apache.hadoop.mapred.JobTracker=DEBUG
```

2. 获取堆栈信息

有关系统的堆栈信息，Hadoop 守护进程为我们提供了一个网络页面（在网络 UI 附后缀 /stacks 访问），该网络页面可以为管理员提供所有守护进程 JVM 中运行的线程信息。你可以通过以下链接访问：<http://jobtracker-host:50030/stacks>。

10.2.3 Metrics

事实上，除了 Hadoop 自带的日志功能以外，还有很多其他可以扩展的 Hadoop 监控程序供管理员使用。在介绍这些监控工具之前，我们首先对系统的可度量信息（Metrics）进行简单的讲解。

HDFS 及 MapReduce 的守护进程会按照一定的规则来收集系统的度量信息。这种度量规则我们称为 Metrics。例如，DataNode 会收集如下度量信息：写入的字节数、被复制的文件块数及来自客户端的请求数等。

Metrics 属于一个上下文，当前 Hadoop 拥有 dfs、mapred、rpc、jvm 等上下文。Hadoop 守护进程会收集多个上下文的度量信息。所谓上下文即应用程序进入系统执行时，系统为用户提供的完整的运行时环境。进程的运行环境是由它的程序代码和程序运行所需要的数据结构及硬件环境组成的。

这里我们认为，一个上下文定义了一个单元，比如你可以选择获取 dfs 上下文或 jvm 上下文。我们可以通过配置 conf/hadoopmetrics.properties 文件设定 Metrics。在默认情况下，所有上下文都会配置为 NullContext 类，这代表它们不会发布任何 Metrics。下面是配置文件的默认配置情况：

```
dfs.class=org.apache.hadoop.metrics.spi.NullContext
```

```
mapred.class=org.apache.hadoop.metrics.spi.NullContext
jvm.class=org.apache.hadoop.metrics.spi.NullContext
rpc.class=org.apache.hadoop.metrics.spi.NullContext
```

其中每一行都针对一个不同的上下文单元，同时每一行定义了处理此上下文 Metrics 的类。这里的类必须是 MetricsContext 接口的一个实现。在上面的例子中，这些 NullContext 类正如其名，也就是说它什么都不做，既不会发布也不会更新它们的 Metrics。

下面我们来介绍 MetricsContext 接口的实现。

1. FileContext

FileContext 可将 Metrics 写入本地文件，FileContext 拥有两个属性：fileName——定义文件的名称；period——指定文件更新的间隔。这两个属性都是可选的，如果不进行设置，那么 Metrics 就会每隔 5 秒写入标准输出。

配置属性将应用于指定的上下文中，并通过在上下文名称后附加点“.”及属性名进行标示。比如，为了将 jvm 导出到一个文件中，我们会通过以下方法调整它的配置：

```
jvm.class=org.apache.hadoop.metrics.file.FileContext
jvm.fileName=/tmp/jvm_metrics.log
```

其中，第一行使用 FileContext 来改变 JVM 的上下文，第二行将 JVM 上下文导出到临时文件中。

需要注意的是，FileContext 非常适合于本地系统的调试，但是它并不适合在大型集群中使用，因为它的输出文件会被分散到集群中，这使得分析的时间成本变得很高。

2. GangliaContext

Ganglia (<http://ganglia.info/>) 是一个开源的分布式监控系统，主要应用于大型分布式集群的监控。通过它可以更好地监控和调整集群中每个机器节点的资源分配。Ganglia 本身会收集一些监控信息，包括 CPU 和内存使用率等。通过使用 GangliaContext 我们可以非常方便地将 Hadoop 的一些测量内容注入 Ganglia 中。此外，GangliaContext 有一个必需的属性——servers，它的属性值是 Ganglia 服务器主机地址：端口。我们将在 10.2.5 节中进行详细讲解。

3. NullContextWithUpdateThread

通过前面的介绍，我们会发现 FileContext 和 GangliaContext 都是将 Metrics 推广到外部系统。而 Hadoop 内部度量信息的获取则需要另外的工具，比如著名的 Java 管理扩展 (JMX, Java Management Extensions)，JMX 中的 NullContextWithUpdateThread 就是用来解决这个问题的（我们将在后面进行详细讲解）。和 NullContext 相似，它不会发布任何 Metrics，但是它会运行一个定时器周期性地更新内存中的 Metrics，以保证另外的系统可以获得最新的 Metrics。

除了 NullContextWithUpdateThread 以外，所有的 MetricsContext 都会执行这种在内存中定时更新的方法，所以只有当你不使用其他输出进行 Metrics 收集时，才需要使用

NullContextWithUpdateThread。举例来说，如果你之前正在使用 GangliaContext，那么随后只要确认 Metrics 是否被更新，而且只需要使用 JMX，不用进一步对 Metrics 系统进行配置。

4. CompositeContext

CompositeContext 允许我们输出多个上下文中相同的 Metrics，比如下面这个例子：

```
jvm.class=org.apache.hadoop.metrics.spi.CompositeContext
jvm.arity=2
jvm.sub1.class=org.apache.hadoop.metrics.file.FileContext
jvm.fileName=/tmp/jvm_metrics.log
jvm.sub2.class=org.apache.hadoop.metrics.ganglia.GangliaContext
jvm.servers=ip-10-70-20-111.ec2.internal:8699
```

其中 arity 属性用来定义子上下文数量，在这里它的值为 2。所有子上下文的属性名称都可以使用下面的句子设置：jvm.sub1.class=org.apache.hadoop.metrics.file.FileContext。

10.2.4 Java 管理扩展

Java 管理扩展 (JMX) 是一个为应用程序、设备、系统等植入管理功能的框架。JMX 可以跨越一系列异构操作系统平台、系统体系结构和网络传输协议，灵活地开发无缝集成的系统、网络和服务管理应用。Hadoop 包含多个 MBean (Managed Bean, 管理服务，它描述一个可管理的资源)，MBean 可以将 Hadoop 的 Metrics 应用到基于 JMX 的应用程序中。当前 MBeans 可以将 Metrics 展示到 dfs 和 rpc 上下文中，但不能在 mapred 及 JVM 上下文上实现。表 10-1 是 MBeans 的列表。

表 10-1 Hadoop 的 MBeans

MBean 类	后台进程	说 明
NameNodeActivityMBean	名称节点	名称节点活动的度量，比如创建文件操作的数量
FSNamesystemMbean	名称节点	名称节点状态的度量，比如已连接的数据节点数量
DataNodeActivityMbean	数据节点	数据节点活动度量，比如读入的字节数
FSDatasetMbean	数据节点	数据节点储存度量，比如容量、空闲存储空间
RpcActivityMbean	所有使用 RPC 的守护进程：名称节点、数据节点、JobTracker 和 TaskTracker	RPC 统计数据，比如平均处理时间

JDK 中的 Jconsole 工具可以帮助我们查看 JVM 中运行的 MBeans 信息，使我们方便地浏览 Hadoop 中的监控信息。很多第三方监控和调整系统 (Nagios 和 Hyperic 等) 可以查询 MBeans，这样 JMX 自然就成为了我们监控 Hadoop 系统的最好工具。但是，你需要设置支持远程访问的 JMX，并且设置一定的安全级别，包括密码权限、SSL 链接及 SSL 客户端权限设置等。为了使系统支持远程访问，JMX 需要我们更改一些选项，其中包括设置 Java 系统的属性 (可以通过编辑 Hadoop 的 conf/hadoop-env.sh 文件实现)。下面的例子展示了如何通过密码远程访问 NameNode 中的 JMX (SSL 不可用条件下)：

```
export HADOOP_NAMENODE_OPTS="-Dcom.sun.management.jmxremote
-Dcom.sun.management.jmxremote.ssl=false
-Dcom.sun.management.jmxremote.password.file=$HADOOP_CONF_DIR/jmxremote.password
-Dcom.sun.management.jmxremote.port=8004 $HADOOP_NAMENODE_OPTS"
```

jmxremote.password 文件以纯文本的格式列出了所有的用户名和密码。JMX 文档有关于 jmxremote.password 文件的更进一步的格式信息。

通过以上的配置，我们可以使用 JConsole 工具浏览远程 NameNode 中的 MBean 监控信息。事实上，我们还有其他很多方法可以实现这个功能，比如通过 jmxquery（一个命令行工具，具体信息可查看 <http://code.google.com/p/jmxquery/>）来检索低于副本要求的块：

```
./check_jmx -U service:jmx:rmi:///jndi/rmi://namenode-host:8004/jmxrmi -O \
hadoop:service=NameNode,name=FSNamesystemState -A UnderReplicatedBlocks \
-w 100 -c 1000 -username monitorRole -password secret
JMX OK - UnderReplicatedBlocks is 0
```

jmxquery 命令创建一个 JMX RMI 链接，链接到 NameNode 的主机地址上，端口号为 8004。它会读取对象名为 `hadoop:service=NameNode, name=FSNamesystemState` 的 `UnderReplicatedBlocks` 属性，并将读出的值写入终端。`-w`、`-c` 选项定义了警告和数值的临界值，这个临界值的选定要在我们运维集群一段时间以后才能选出比较合适的经验值。

这里需要注意的是，尽管你可以通过 JMX 的默认监控信息配置看到 Hadoop 的监控信息，但是它们不会自动更新，除非你更改了 `MetricContext` 的具体实现。比如，如果 JMX 是你唯一使用的监控系统信息的方法，那么你就可以把 `MetricContext` 的实现更改为 `NullContextWithUpdateThread`。

通常大多数人会使用 Ganglia 和另外一个可选的系统（比如 Nagios）来进行 Hadoop 集群的检测工作。Ganglia 可以很好地完成大数据量监控信息的收集和图形化工作，而 Nagios 及类似的系统则更擅长处理小规模的监控数据，并在监控信息超出设定的监控阈值时发出警告。管理者可以根据需求选择合适的工具。下一节我们就对 Ganglia 的使用配置进行详细的讲解。

10.2.5 Ganglia

Ganglia 是 UC Berkeley 发起的一个开源集群监视项目，用于测量数以千计的节点集群。Ganglia 的核心包含两个 Daemon（分别是客户端 Ganglia Monitoring Daemon (gmond) 和服务端 Ganglia Meta Daemon (gmetad)），以及一个 Web 前端。它主要是用来监控系统性能的，如：CPU、memory、硬盘利用率、I/O 负载、网络流量情况等，通过 Web 前端页面我们可以获得曲线描述的各个节点工作状态。通过 Ganglia 可以帮助我们合理调整、分配系统资源，为提高系统整体性能起到了重要作用。

监控下的每台节点计算机都需要运行一个收集和发送度量数据的名叫 `gmond` 的守护进程。接收所有度量数据的主机可以显示这些数据并且可以将这些数据传递到监控主机中。`gmond` 带来的系统负载非常少，它的运行不会影响用户应用进程的性能。所有这些数据多次

收集则会影响节点性能。网络中的“抖动”发生在大量小消息同时出现时，可以通过将节点时钟保持一致，来避免这个问题。

gmetad 可以部署在集群内任一节点或通过网络连接到集群的独立主机上，它通过单播路由的方式与 gmond 通信，收集区域内节点的状态信息，并以 XML 数据的形式保存到数据库中。最终由 RRDTool 工具处理数据，并生成相应的图形显示，以 Web 方式直观地提供给客户端。

1. 服务器端的安装与配置

这个服务器可以看做是一个信息收集的装置，它能同时监控多个客户端的系统状况，并把信息显示在 Web 界面上。通过 Web 端连接这个服务器，就可以看到它所监控的所有机器状态。

首先我们需要在服务器端安装下列包：ganglia-gmetad-3.0.3-1.fc4.i386.rpm（从各个网段获取汇总监控信息）；rrdtool-1.2.18-1.el4.rf.i386.rpm（显示图像的工具）；rrdtool-devel-1.2.18-1.el4.rf.i386.rpm；ganglia-web-3.0.3-1.noarch.rpm（Ganglia 的 Web 程序）；perl-rrdtool-1.2.18-1.el4.rf.i386.rpm。使用 #rpm -ivh 软件包 .rpm 可以安装这些包。

安装完成之后，找到 Ganglia 服务器端的配置文档：/etc/gmetad.conf。可以根据不同的需要来进行配置。在这里只说一下如何添加或修改要监控的系统。首先通过 #vi /etc/gmetad.conf 命令（进入编辑）找到：data_source “Login FARM” 10.77.20.111:8651 10.77.20.111:8699（后面的这些 IP 地址就是进行监控的主机，冒号后跟的是要监听的端口号，这个端口号将在介绍客户端的配置时提到）。其他属性保持默认配置即可。

配置完成后要重启 gmetad 服务：#service gmetad restart。下面我们来配置虚拟主机，设置路径为 DocumentRoot “/var/www/html/ganglia”：

```
# 配置虚拟主机
<Directory "/var/www/html/ganglia">
    Options Indexes FollowSymLinks
    AllowOverride None
    Order allow,deny
    Allow from all
</Directory>
```

然后重启 httpd 服务：service httpd restart，这样即完成了服务器端的安装和配置。

2. 客户端的安装与配置

在客户端安装 Ganglia，是为了收集本机的信息，并通过设置好的端口把信息传给服务器端，我们需要在所有节点上进行相应的安装与配置。下面来讲解如何进行安装。

首先在客户端安装软件包：ganglia-gmond-3.0.3-1.fc4.i386.rpm。安装完成后，找到它的配置文档 /etc/gmond.conf 并打开编辑 #vi /etc/gmond.conf。找到配置文件中如下部分并按照所给出的例子进行配置：

```
/* You can specify as many tcp_accept_channels as you like to share
   an xml description of the state of the cluster */
```

```

tcp_accept_channel {
    port = 8699 /* 注释: 这个是端口, 通过它来传递系统信息。注意要和服务器端监听的端口一致 */
    acl {
        default = "deny"
        access {
            ip = 10.77.20.111 /* 注释: 这里是服务器的 IP 地址 */
            mask = 32
            action = "allow"
        }
    }
}

```

配置完成后, 重启 gmond 服务 #service gmond restart 即可。至此 Ganglia 在服务器端和客户端的安装完成, 我们可以查看它的运行状态 (如图 10-2 所示)。

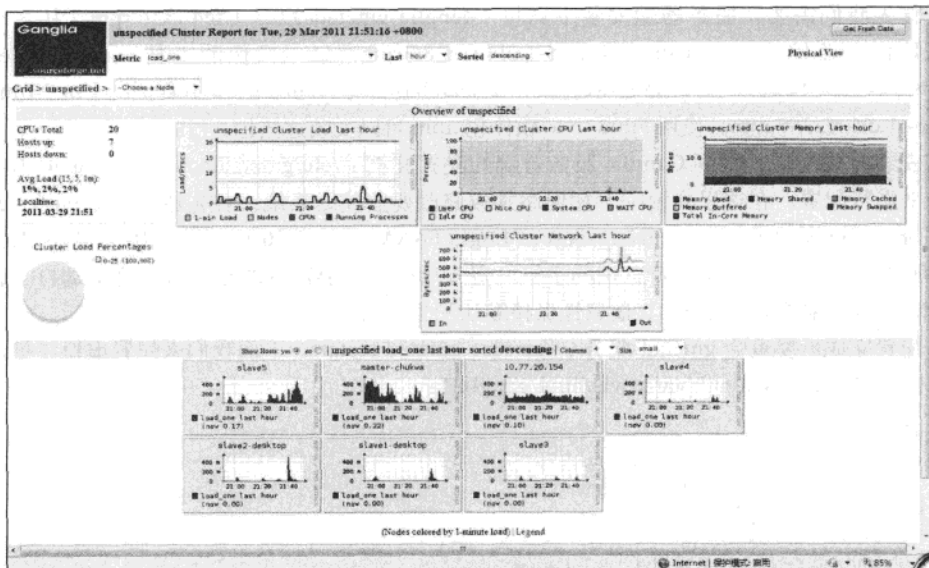


图 10-2 Ganglia 的监控页面

事实上, 有很多其他可以扩展 Hadoop 监控能力的工具, 比如在本书第 17 章中介绍的 Chukwa, 它是一个数据收集和监控系统, 构建于 HDFS 和 MapReduce 之上, 也是可供管理员选择的监控工具。Chukwa 可以统计分析日志文件, 从而提供给管理员想要的信息。

10.2.6 Hadoop 管理命令

我们在介绍扩展的监控管理工具时, 也不能忘记 Hadoop 本身为我们提供了相应的系统管理工具, 本节我们就对相关的工具进行介绍。

1. dfsadmin

dfsadmin 是一个多任务的工具, 我们可以使用它来获取 HDFS 的状态信息, 以及在

HDFS 上执行的管理操作。管理员可以在终端中通过 Hadoop dfsadmin 命令调用它，这里需要使用超级用户权限。dfsadmin 相关的命令如表 10-2 所示。

表 10-2 dfsadmin 命令的解析

命令选项	描 述
-report	报告文件系统的基本信息和统计信息
-safemode enter leave get wait	安全模式维护命令。安全模式是 NameNode 的一种状态，这种状态下，NameNode： 1) 不接受对名字空间的更改（只读） 2) 不复制或删除块 NameNode 会在启动时自动进入安全模式，当配置块的最小百分数满足最小副本数的条件时，会自动离开安全模式。可以手动进入安全模式，但是也必须手动关闭之
-refreshNodes	重新读取 hosts 和 exclude 文件，使新的节点或需要退出集群的节点能够重新被 NameNode 识别
-finalizeUpgrade	终结 HDFS 的升级操作。DataNode 删除前一个版本的工作目录，之后 NameNode 也会这样做。该操作完结整个升级过程
-upgradeProgress status details force	请求当前系统的升级状态、升级状态的细节，或者进行强制升级操作
-metasave filename	保存 NameNode 的主要数据结构到 hadoop.log.dir 属性指定目录下的 <filename> 文件上。对于下面的每一项，<filename> 中都有一行内容与之相对应： 1) NameNode 收到的 DataNode 的心跳信号 2) 等待被复制的块 3) 正在被复制的块 4) 等待被删除的块
-setQuota <quota> <dirname>...<dirname>	为每个目录 <dirname> 设定配额 <quota>。目录配额是一个长整型整数，强制限定目录树下的名字个数。以下情况会报错： 1) N 不是一个正整数 2) 用户不是管理员 3) 这个目录不存在或为文件 4) 目录会马上超出新设定的配额
-clrQuota <dirname>...<dirname>	为每个目录 <dirname> 清除配额设定。以下情况会报错： 1) 该目录不存在或为文件 2) 用户不是管理员 如果目录原来没有配额，则不会报错
-help [cmd]	显示给定命令的帮助信息，如果没有给定命令，则显示所有命令的帮助信息

2. 文件系统验证 (fsck)

Hadoop 提供了 fsck 工具来验证 HDFS 中的文件是否正常可用。这个工具可以检测文件块是否在 DataNode 中丢失，是否低于或高于文件副本要求。下面给出使用的例子：

```
hadoop fsck /
.....Status: HEALTHY
Total size: 511799225 B
```

```

Total dirs: 10
Total files: 22
Total blocks (validated): 22 (avg. block size 23263601 B)
Minimally replicated blocks: 22 (100.0 %)
Over-replicated blocks: 0 (0.0 %)
Under-replicated blocks: 0 (0.0 %)
Mis-replicated blocks: 0 (0.0 %)
Default replication factor: 3
Average block replication: 3.0
Corrupt blocks: 0
Missing replicas: 0 (0.0 %)
Number of data-nodes: 4
Number of racks: 1
The filesystem under path '/' is HEALTHY

```

fsck 会递归遍历文件系统的 Namespace，从文件系统的根目录开始检测它所找到的全部文件，并在它验证过的文件上标记一个点。要检查一个文件，fsck 首先会检索元数据中文件的块，然后查看是否有问题或是否一致。这里需要注意的是，fsck 验证只和 NameNode 通信而不和 DataNode 通信。

以下是几个 fsck 的输出情况，在这里进行说明：

(1) Over-replicated blocks

Over-replicated blocks 用来指明一些文件块副本数超出了它所属文件的限定。通常来说，过量的副本数存在并不是问题，HDFS 会自动删除多余的副本。

(2) Under-replicated blocks

Under-replicated blocks 用来指明文件块数未达到所属文件要求的副本数量。HDFS 也会自动创建新的块直到该块的副本数能够达到要求。你可以通过 `hadoop dfsadmin -metasave` 命令获得正在被复制的块信息。

(3) Misreplicated blocks

Misreplicated blocks 用来指明不满足块副本存储位置策略的块。例如，假设副本因子为 3，如果一个块的所有副本都存在于一个机架中，那么这个块就是 Misreplicated blocks。针对这个问题，HDFS 不会自动调整，我们只能通过手动设置来提高该文件的副本数，然后再将它的副本数设置为正常值来解决。

(4) Corrupt blocks

Corrupt blocks 指所有的块副本全部出现问题，只要块存在的副本可用，它就不会被报告为 Corrupt blocks。NameNode 会使用没有出现问题的块进行复制操作，直到达到目标值。

(5) Missing replicas

Missing replicas 用来表明集群中不存在副本的文件块。

Missing replicas 及 Corrupt 块是受关注最多的情况，因为这意味着数据的丢失。fsck 默认不去处理那些丢失或出现问题的文件块，但是可以通过命令使其执行以下操作：

❑ 通过 `-move`，将出现问题的文件放入 HDFS 的 `/lost+found` 文件夹下。

❑ 通过 `-delete` 选项将出现问题的文件删除，删除后即不可恢复。

`fsck` 提供了一种简单的方法去查找属于某个文件的所有块。代码如下：

```
hadoop fsck /user/admin/In/hello.txt -files -blocks -racks
/user/admin/In/hello.txt 13 bytes, 1 block(s): OK
0. blk_-8114668855310504639_1056 len=13 repl=1 [/default-rack/127.0.0.1:50010]
```

```
Status: HEALTHY
Total size:      13 B
Total dirs:      0
Total files:     1
Total blocks (validated):      1 (avg. block size 13 B)
Minimally replicated blocks:  1 (100.0 %)
Over-replicated blocks:       0 (0.0 %)
Under-replicated blocks:      0 (0.0 %)
Mis-replicated blocks:        0 (0.0 %)
Default replication factor:    1
Average block replication:     1.0
Corrupt blocks:                0
Missing replicas:              0 (0.0 %)
Number of data-nodes:         1
Number of racks:              1

The filesystem under path '/user/admin/In/hello.txt' is HEALTHY
```

从以上输出中可以看到文件 `hello.txt` 由一个块组成，并且命令也返回了它所在的 `DataNode`。`fsck` 的选项用法如下：

❑ `-files`：显示文件的名称、大小、块数量及是否可用（是否存在丢失的块）；

❑ `-blocks`：显示每个块在文件中的信息，一个块用一行显示；

❑ `-racks`：展示了每个块所处的机架位置及 `DataNode` 的位置。

运行 `fsck` 命令时如果不加参数选项，则执行以上所有指令。

3. `DataNode` 块扫描任务

每个 `DataNode` 都会执行一个块扫描任务，它会周期性地验证它所存储的块。这就允许有问题的块在客户端读取时被删除或修整。`DataBlockScanner` 可维护一个块列表，它会一个一个地扫描这些块，并进行校验和验证。

进行块验证的周期可以通过 `dfs.datanode.scan.period.hours` 属性值进行设定，默认为 504 小时，即 3 周。出现问题的块将会被报告给 `NameNode` 进行处理。

你也可以通过访问 `DataNode` 的 Web 接口获得块验证的信息：`http://datanodeIP:50075/block ScannerReport`，下面是一个报告的样本。

```
Total Blocks      :      32
Verified in last hour :      1
Verified in last day  :      1
Verified in last week :     12
```

```

Verified in last four weeks      :      31
Verified in SCAN_PERIOD         :      31

Not yet verified                 :        1
Verified since restart           :        2
Scans since restart              :        2
Scan errors since restart        :        0
Transient scan errors            :        0
Current scan rate limit KBps     :    1024

Progress this period             :      8%
Time left in cur period          :    99.96%

```

通过后缀 `listblocks` 参数 (<http://datanodeIP:50075/blockScannerReport?listblocks>), 报告会在前面这个 `DataNode` 中加入所有块的最新验证状态信息。

4. 均衡器 (balancer)

HDFS 不间断地运行, 隔一段时间可能就会出现文件在集群中分布不均匀的情况, 一个不平衡的集群会影响系统资源的充分利用, 所以我们要想办法避免它。

`balancer` 程序是 Hadoop 的守护进程, 它会通过将文件块从高负载的 `DataNode` 转移到低使用率的 `DataNode` 上, 即进行文件块的重新分布, 达到集群的平衡。同时还要考虑 HDFS 的块副本分配策略。`balancer` 的目的是使集群达到相对平衡, 这里的相对平衡是指每个 `DataNode` 的磁盘使用率和整个集群的资源使用率的差值小于给定的阈值。我们可以通过以下命令来运行 `balancer` 程序: `start-balancer.sh`。-threshold 参数设定了多个可以接受的集群平衡点。超过这个平衡预置, 就要进行平衡调整, 对文件块进行重分布。这个参数值在大多数情况下为 10%, 当然也可通过命令行设置。`balancer` 被设计运行于集群后台中, 不会增加集群运行负担。我们可以通过参数设置, 限制 `balancer` 在执行 `DataNode` 之间的数据转移时占用的带宽资源。这个属性值可以通过 `hdfs-site.xml` 配置文件中的 `dfs.balance.bandwidthPerSec` 属性进行修改, 默认为 1MB。

10.3 Hadoop 集群的维护

10.3.1 安全模式

当 `NameNode` 启动时, 第一件要做的事情就是将映像文件 `fsimage` 加载到内存, 并应用 `edits` 文件记录编辑日志。一旦成功重构和之前文件系统一致且居于内存的文件系统元数据, `NameNode` 就会创建一个新的 `fsimage` 文件 (这样就可以更高效地记录检查点, 而不用依赖于 `Secondary NameNode`) 和一个空的编辑日志文件。只有全部完成了这些工作, `NameNode` 才会监听 `RPC` 和 `HTTP` 请求。然而, 如果 `NameNode` 运行于安全模式, 那么文件系统只能对客户端提供只读模式的视图了。

文件块的位置信息并没有持久化地存储在 NameNode 中，这些信息都存储在各 DataNode 中。在文件系统的常规操作期间，NameNode 会在内存中存储一个块位置的映射。在安全模式下，需要给 DataNode 一定的时间向 NameNode 上传它们存储块的列表，这样 NameNode 才能获得充足的块位置信息，才会使文件系统更加高效。如果 NameNode 没有足够的时间来等待获取这些信息，它就会认为该块没有足够的副本；进而安排其他 DataNode 复制，这在很多情况下显然是没有必要的，并且还会浪费系统资源。在安全模式下，NameNode 不会处理任何块复制和删除指令。

当最小副本条件达到要求时，系统就会退出安全模式，这需要延期 30 秒（这个时间由 `dfs.safemode.extension` 属性值确定，默认为 30，一些小的集群比如 10 个节点，可以设置为 0）。这里我们所说的最小副本条件是指系统中 99.9%（这个值由 `dfs.safemode.threshold.pct` 属性确定，默认为 0.999）的文件块达到 `dfs.replication.min` 属性值所设置的副本数（默认为 1）。

当格式化一个新的 HDFS 时，NameNode 不会进入安全模式，因为此时系统中还没有任何文件块。

使用以下命令可以查看 NameNode 是否已进入安全模式：

```
hadoop dfsadmin -safemode get
Safe mode is ON
```

在有些情况下，需要在运行命令前等待名称节点退出安全模式，这时你可以使用以下命令：

```
hadoop dfsadmin -safemode wait
# command to read or write a file
```

作为管理员，也应具有使 NameNode 进入或退出安全模式的方法，这些操作有时是必需的，比如在升级完集群后需要确认数据是否仍然可读等。这时我们可以使用以下命令：

```
hadoop dfsadmin -safemode enter
Safe mode is ON
```

当 NameNode 仍处于安全模式时，也可以使用以上命令来保证 NameNode 没有退出安全模式。要使系统退出安全模式可执行以下命令：

```
hadoop dfsadmin -safemode leave
Safe mode is OFF
```

10.3.2 Hadoop 的备份

1. 元数据的备份

如果 NameNode 中存储的持久化元数据信息丢失或遭到破坏，那么整个文件系统就不可用了。因此元数据的备份至关重要，需要备份不同时期的元数据信息（1 小时、1 天、1 周……）以避免突然的宕机带来的破坏。

备份一个最直接的办法就是编写一个脚本程序，然后周期性地将 Secondary NameNode 中 `previous.checkpoint` 子目录（该目录由 `fs.checkpoint.dir` 属性值确定）下的文件归档到另外的机器上。该脚本需要额外验证所拷贝备份文件的完整性。这个验证可以通过在 NameNode 的守护进程中运行一个验证程序来实现，验证其是否成功地从内存中读取了 `fsimage` 及 `edits` 文件。

2. 数据的备份

HDFS 的设计目标之一就是能够可靠地在分布式集群中储存数据，由于 HDFS 允许数据丢失的发生，所以数据的备份就显得至关重要了。由于 Hadoop 可以存储大规模的数据，所以备份哪些数据，备份到哪里就成为一个挑战。在备份过程中，最优先备份的应该是那些不能再生的数据和对商业应用最关键的数据。而对于那些可以通过其他手段再生的数据或对于商业应用价值不是很大的数据，可以考虑不进行备份。

这里需要强调的是，不要认为 HDFS 的副本机制可以代替数据的备份。HDFS 中的 Bug 也会导致副本丢失，同样硬件也会出现故障。尽管 Hadoop 可以承受集群中廉价商用机器的故障，但是有些极端情况不能排除在外，特别是系统有时还会出现软件 Bug 和人为失误的情况。

通常 Hadoop 会设置用户目录的策略，比如，每个用户都有一个空间配额，每天晚上都可进行备份工作。但是不管策略如何，都需要通知用户知晓，以避免客户反映问题。

前面介绍的 `distcp` 工具（参见第 9 章“HDFS 详解”）是在不同 HDFS 之间或不同 Hadoop 文件系统之间转储、进行数据备份的好工具，因为 `distcp` 可以并行运行数据复制。

10.3.3 Hadoop 的节点管理

作为 Hadoop 集群的管理员，可能随时都要处理增加机器节点、移除机器节点的任务。例如，要增加集群的存储容量，就要增加新的节点。相反，如果希望缩小集群的规模，那么需要撤销已存在的节点。如果一个节点频繁地发生故障或运行缓慢，那么也要考虑撤销已存在的节点。节点一般承担 `DataNode` 和 `TaskTracker` 的任务，Hadoop 支持对它们的添加和撤销操作。

1. 添加新的节点

在第 2 章，我们介绍了如何部署 Hadoop 集群，可以看到添加一个新的节点虽然只用配置 `hdfs-site.xml` 文件和 `mapred-site.xml` 文件，但最好还是配置一个授权节点列表。

如果允许任何机器都可以链接到 NameNode 上并充当 `DataNode` 是存在安全隐患的，因为这样的机器可能能够获得未授权文件的访问权限。此外这样的机器并不是真正的 `DataNode`，但它可以存储数据，而又不在于集群的控制之下，并且任何时候都有可能停止运行，从而造成数据丢失。由于其配置简单或存在配置错误，即使在防火墙内这样处理也可能存在风险，在集群中也要对 `DataNode` 进行明确的管理。

在 `dfs.hosts` 文件中指定可以链接到 NameNode 的 `DataNode` 列表。`dfs.hosts` 文件存储在 NameNode 的本地文件系统之上，包含每个 `DataNode` 的网络地址，一行表示一个

DataNode。如果要为一个 DataNode 设置多个网络地址，则把它们写到一行中，中间用空格分开。类似的，TaskTracker 是在 `mapred.hosts` 中设置的。一般来说，DataNode 和 TaskTracker 列表中都存在一个共享文件，名为 `include file`，它被 `dfs.hosts` 及 `mapred.hosts` 两者引用，因为在大多数情况下，集群中的机器会同时运行 DataNode 及 TaskTracker 守护进程。

这里需要注意的是，`dfs.hosts` 和 `mapred.hosts` 这两个文件与 `slaves` 文件不同，`slaves` 文件被 Hadoop 的执行脚本用来执行集群范围的操作，例如集群的重启等，但它从来不会被 Hadoop 的守护进程使用。

向集群添加新的节点，要执行以下步骤：

- 1) 向 `include` 文件中添加新节点的网络地址；
- 2) 使用以下命令更新 NameNode 中具有连接权限的 DataNode 集合：

```
hadoop dfsadmin -refreshNodes
```

3) 更新带有新节点的 `slaves` 文件，以便 Hadoop 控制脚本执行后续操作时可以使用更新后的 `slaves` 文件中的所有节点；

- 4) 启动新的数据节点；
- 5) 重新启动 MapReduce 集群；
- 6) 检查在网页用户界面是否有新的 DataNode 和 TaskTracker。

需要注意的是，HDFS 不会自动将旧的数据转移到新的 DataNode 中，但你可以运行平衡器命令进行集群均衡。

2. 撤销节点

撤销数据节点时要避免数据的丢失，在撤销前，需先通知 NameNode 要撤销的节点，然后使其在此节点撤销前将其上的数据块转移出去。而如果关闭正在运行的 TaskTracker，那么 JobTracker 会意识到错误并将任务分配到其他 TaskTracker 中去。

撤销节点过程由 `exclude` 文件控制，对于 HDFS 来说，可以通过 `dfs.hosts.exclude` 属性来控制；对于 MapReduce 来说，可以由 `mapred.hosts.exclude` 来设置。

要看 TaskTracker 是否可以连接到 JobTracker 规则上很简单，只有 `include` 文件中包含，但 `exclude` 文件不包含时，TaskTracker 才可以连接到 JobTracker 执行任务。没有定义或空的 `include` 文件意味着所有节点都在 `include` 文件中。

对于 HDFS 来说规则则有些许不同，表 10-3 总结了 `include` 和 `exclude` 存放节点的情况。对于 TaskTracker 来说，一个未定义或空的 `include` 文件意味着所有的节点都包含其中。

表 10-3 HDFS 的 `include` 和 `exclude` 的文件优先级

include 文件中包含	exclude 文件中包含	解 释
否	否	节点可以连接
否	是	节点不可以连接
是	否	节点可以连接
是	是	节点可以连接和撤销

要想从集群中移除节点，需要执行以下步骤：

- 1) 将需要撤销的节点网络地址增加到 `exculde` 文件中，但请注意不要在此时更新 `include` 文件；
- 2) 重新启动 MapReduce 集群来终止已撤销节点的 TaskTracker；
- 3) 用以下命令更新具有新的许可 DataNode 节点集的 NameNode：

```
hadoop dfsadmin -refreshNodes
```

4) 进入网络用户界面，先检查已撤销的 DataNode 的管理状态是否变为 “DecommissionIn Progress”，然后把数据块复制到集群的其他 DataNode 中；

5) 当所有 DataNode 报告其状态为 “Decommissioned” 时，所有数据块也都会被复制，此时可以关闭已撤销的节点；

- 6) 从 `include` 中删除节点网络地址，然后再次运行命令：

```
hadoop dfsadmin -refreshNodes
```

- 7) 从 `slaves` 文件中删除节点。

10.3.4 系统升级

升级 HDFS 和 MapReduce 集群需要一个合理的操作步骤，这里我们主要讲解 HDFS 的升级。如果文件系统升级后文件格局发生了变化，那么升级时会将文件系统的数据和元数据迁移到与新版本一致的格式上。由于任何涉及数据迁移的操作都会导致数据丢失，所以必须保证数据和元数据都有备份（具体操作参看 10.3.2 节）。进行升级时，可以先在小型集群中进行测试，以便正式运行时可以解决所有问题。

Hadoop 对自身的兼容性要求非常高，所有 Hadoop 1.0 之前版本的兼容性要求最严格，只有来自相同发布版本的组件才能保证相互的兼容性，这就意味着整个系统从守护进程到客户端都要同时更新，还需要集群停机一段时间。后期发布的版本则支持回滚升级，允许集群守护进程分阶段升级，以便在更新期间可以运行客户端。

如果文件系统的布局不改变，那么集群升级就非常简单了。首先在集群中安装新的 HDFS 和 MapReduce（同时在客户端也要安装），然后关闭旧的守护进程，升级配置文件，启动新的守护进程和客户端更新库。这个过程是可逆的，因此升级后的版本回滚到之前版本也很简单。

每次成功升级后都要执行一系列的最终清除步骤：

- 1) 从集群上删除旧的安装和配置文件；
- 2) 修复代码和配置中的每个错误警告。

以上的系统升级非常简单，但是如果我们升级文件系统，就需要更进一步的操作了。

如果使用以上方法进行升级，并且 HDFS 是一个不同的布局版本，那么 NameNode 就不会正常运行了。NameNode 的日志会产生以下信息：

```
File system image contains an old layout version -15.
An upgrade to version -18 is required.
Please restart NameNode with -upgrade option.
```


想确定是否需要升级文件系统，最好的办法就是在一个小集群上进行测试。

HDFS 升级将复制以前版本的元数据和数据。升级并不需要两倍的集群存储空间，因为 DataNode 使用硬链接来保留对同一个数据块的两个引用，这样就可以在需要的时候轻松实现回滚到以前版本的文件系统上。

需要注意的是，升级后只能保留前一个版本的文件系统，而不能回滚到多个文件系统，因此执行另一个对 HDFS 的升级时需要删除以前的版本，这个过程被称为确定更新（finalizing the upgrade）。一旦更新被确定，那它就不会回滚到以前的版本了。

需要说明的是，只有可以正常运行的健康的系统才能被正确升级。所以在进行升级之前，必须进行一个全面的 fsck 操作。要防止意外，可以将系统中的所有文件及块的列表（fsck 的输出）进行备份。这样就可以在升级后将运行的输出与之进行对比，检测是否全部正确升级，有没有数据丢失。

还需要注意的是，在升级之前要删除临时文件。这些临时文件包括 HDFS 上 MapReduce 系统目录中的文件和本地临时文件。

完成以上这些工作后就可以进行集群的升级和文件系统的迁移了，具体步骤如下：

- 1) 确保之前的升级操作全部完成，不会影响此次升级；
- 2) 关闭 MapReduce，终止 TaskTracker 上的所有任务进程；
- 3) 关闭 HDFS 并备份 NameNode 目录；
- 4) 在集群和客户端上安装新版本的 Hadoop HDFS 和同步的 MapReduce；
- 5) 使用 -upgrade 选项启动 HDFS；
- 6) 等待操作完成；
- 7) 在 HDFS 上进行健康检查；
- 8) 启动 MapReduce；
- 9) 回滚或确定升级。

运行升级程序时，最好能从 PATH 环境变量中删除 Hadoop 脚本，这样可以避免运行不确定版本的脚本程序。安装目录定义两个环境变量是很方便的，在以下指令中已经定义了 OLD_HADOOP_INSTALL 和 NEW_HADOOP_INSTALL。以上步骤中的第 5) 步我们要运行以下指令：

```
$NEW_HADOOP_INSTALL/bin/start-dfs.sh -upgrade
```

NameNode 升级它的元数据，并将以前的版本放入新建的目录 previous 中：

```
$ {dfs.name.dir}/current/VERSION
    /edits
    /fsimage
    /fstime
                                /previous/VERSION
    /edits
    /fsimage
    /fstime
```

采用类似的方式，DataNode 升级它的存储目录，将旧的目录拷贝到 previous 目录中去。

升级过程需要一段时间才能完成。可以使用 dfsadmin 命令来检查升级的进度，升级的事件同样会记录在守护进程的日志文件中。第 6) 步执行以下命令：

```
$NEW_HADOOP_INSTALL/bin/hadoop dfsadmin -upgradeProgress status
Upgrade for version -18 has been completed.
Upgrade is not finalized.
```

以上代码表明升级已经完成，在这个阶段必须在文件系统上进行一些健康检查（即第 7) 步，比如使用 fsck 进行文件和块的检查）。当进行检查（只读模式）时，可以让 HDFS 进入安全模式，以防止其他检查对文件进行更改。

第 9) 步是可选操作，如果在升级后发现问题；则可以回滚到之前版本。

首先，关闭新的守护进程：

```
$NEW_HADOOP_INSTALL/bin/stop-dfs.sh
```

再次，用 -rollback 选项启动旧版本的 HDFS：

```
$OLD_HADOOP_INSTALL/bin/start-dfs.sh -rollback
```

这个命令会使用 NameNode 和 DataNode 以前的副本替换它们当前存储目录下的内容，文件系统即返回到原始状态。

如果对新升级的版本感到满意，那么可以执行确定升级（即第 9) 步，可选），并删除以前的存储目录。需要注意的是在升级确定后，就不能回滚到前面的版本了。

需要执行以下步骤后才能进行另一次升级：

```
$NEW_HADOOP_INSTALL/bin/hadoop dfsadmin -finalizeUpgrade
$NEW_HADOOP_INSTALL/bin/hadoop dfsadmin -upgradeProgress status
There are no upgrades in progress.
```

至此，HDFS 升级到了最新的版本。

10.4 小结

本章我们重点介绍了 Hadoop 监控和管理方面的相关内容。

首先，从 HDFS 的文件结构开始进行相关介绍。HDFS 作为 Hadoop 的核心分布式文件系统，许多应用都构建在其核心分布式文件系统上。核心分布式文件系统作为基础的架构，管理员要给予更多的关注。

其次，本章从整体上对 Hadoop 的监控机制和相关的监控工具进行了分析，着重分析了 Hadoop 监控的支持基础、日志和度量，同时提出了诸多系统监控的解决方案，并着重介绍了 Ganglia 监控软件。

最后，对实际应用中经常遇到的维护要求，比如增删节点、数据备份、系统升级等进行了介绍。



第 11 章

Hive 详解

本章内容

- ☐ Hive 简介
- ☐ Hive 的基本操作
- ☐ Hive QL 详解
- ☐ Hive 的网络 (WebUI) 接口
- ☐ Hive 的 JDBC 接口
- ☐ Hive 的优化
- ☐ 小结



Hive 是 Hadoop 中的一个重要子项目，它利用 MapReduce 编程技术，实现了部分 SQL 语句，提供了类 SQL 的编程接口。Hive 的出现极大地推进了 Hadoop 在数据仓库方面的发展。事实上，目前业界对何谓大规模数据分析最佳方法仍在进行着辩论。由于传统应用的惯性，业界保守派依然青睐于关系型数据库和 SQL 语言。而在学术界中，互联网阵营则更集中于支持 MapReduce 的开发模式。本章我们将对基于 Hive 的数据仓库解决方案进行介绍。

11.1 Hive 简介

Hive 是一个基于 Hadoop 文件系统上的数据仓库架构。它为数据仓库的管理提供了许多功能：数据 ETL（抽取、转换和加载）工具、数据存储管理和大型数据集的查询与分析能力。同时 Hive 还定义了类 SQL 的语言——Hive QL，Hive QL 允许用户进行和 SQL 相似的操作。Hive QL 还允许开发人员方便地使用 mapper 和 reducer 操作，这样对 MapReduce 框架是一个强有力的支持。

由于 Hadoop 是批量处理系统，任务是高延迟性的，所以在任务提交和处理过程中会消耗一些时间成本。同样，即使 Hive 处理的数据集非常小（比如几百 MB），在执行时也会出现延迟现象。这样，Hive 的性能就不可能很好地和传统的 Oracle 数据库进行比较了。Hive 不能提供数据排序和查询 cache 功能，也不提供在线事务处理，不提供实时的查询和记录级的更新，但 Hive 能更好地处理不变的大规模数据集（例如网络日志）上的批量任务。所以，Hive 最大的价值是可扩展性（基于 Hadoop 平台，可以自动适应机器数目和数据量的动态变化）、可延展性（结合 MapReduce 和用户定义的函数库），并且拥有良好的容错性和低约束的数据输入格式。

Hive 本身建立在 Hadoop 的体系架构上，提供了一个 SQL 解析过程，并从外部接口中获取命令，以对用户指令进行解析。Hive 可将外部命令解析成一个 MapReduce 可执行计划，并按照该计划生成 MapReduce 任务后交给 Hadoop 集群处理，Hive 的体系结构如图 11-1 所示。

11.1.1 Hive 的数据存储

Hive 的存储是建立在 Hadoop 文件系统之上的。Hive 本身没有专门的数据存储格式，也不能为数据建立索引，用户可以非常自由地组织 Hive 中的表，只需要在创建表的时候告诉 Hive 数据中的列分隔符和行分隔符就可以解析数据了。

Hive 中主要包含四类数据模型：表（Table）、外部表（External Table）、分区（Partition）和桶（Bucket）。

Hive 中的表和数据库中的表在概念上是类似的，每个表在 Hive 中都有一个对应的存储目录。例如，一个表 htable 在 HDFS 中的路径为 /datawarehouse/htable，其中，/datawarehouse 是 hive-site.xml 配置文件中由 \${hive.metastore.warehouse.dir} 指定的数据仓库的目录，所有的表数据（除了外部表）都保存在这个目录中。

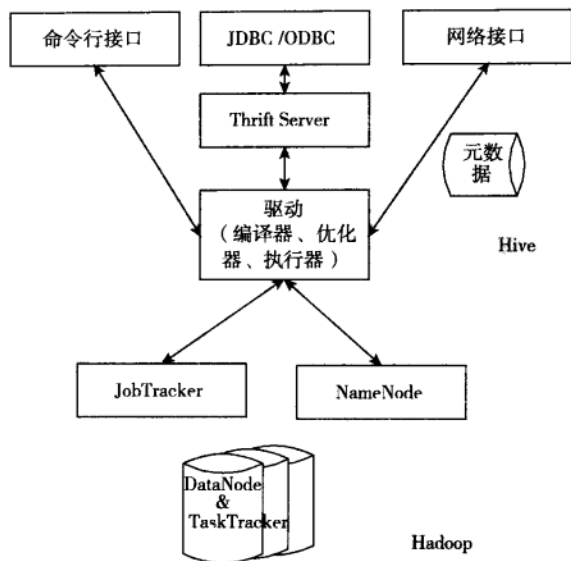


图 11-1 Hive 的体系结构

Hive 中每个分区都对应数据库中相应分区列的一个索引，但是分区的组织方式和传统关系型数据库不同。在 Hive 中，表中的一个分区对应表下的一个目录，所有分区的数据都存储在对应的目录中。例如，htable 表中包含的 ds 和 city 两个分区，分别对应两个目录：对应于 ds=20100301, city=Beijing 的 HDFS 子目录为：/datawarehouse/htable/ds=20100301/city=Beijing；对应于 ds=20100301, city=Shanghai 的 HDFS 子目录为：/datawarehouse/htable/ds=20100301/city=Shanghai。

桶对指定列进行哈希 (hash) 计算时，根据哈希值切分数据，每个桶对应一个文件。例如，将属性列 user 列分散到 32 个桶中，首先要对 user 列的值进行 hash 计算，对应哈希值为 0 的桶写入 HDFS 的目录为：/datawarehouse/htable/ds=20100301/city=Beijing/part-00000；对应哈希值为 10 的 HDFS 目录为：/datawarehouse/htable/ds=20100301/city=Beijing/part-00010，依此类推。

外部表指向已经在 HDFS 中存在的数据库，也可以创建分区。它和表在元数据的组织上是相同的，而实际数据的存储则存在较大差异，主要表现在以下两点上。

1) 创建表的操作 (CREATE Table) 包含两个步骤：表创建过程和数据加载步骤（这两个过程可以在同一语句中完成）。在数据加载过程中，实际数据会移动到数据仓库目录中。之后的数据访问将会直接在数据仓库目录中完成。删除表时，表中的数据和元数据将会被同时删除。

2) 外部表的创建只有一个步骤，加载数据和创建表同时完成，实际数据存储在创建语句 LOCATION 指定的 HDFS 路径中，并不会移动到数据仓库目录中。如果删除一个外部表，

仅会删除元数据，表中的数据不会被删除。

11.1.2 Hive 的元数据存储

由于 Hive 的元数据可能要面临不断的更新、修改和读取，所以它显然不适合使用 Hadoop 文件系统进行存储。目前 Hive 将元数据存储于 RDBMS 中，比如 MySQL、Derby 中，Hive 有三种模式可以连接到 Derby 数据库：1) Single User Mode，此模式连接到一个 In-memory（内存）数据库 Derby，一般用于单元测试；2) Multi User Mode，通过网络连接到一个数据库中，是最常用的模式；3) Remote Server Mode，用于非 Java 客户端访问元数据库，在服务器端启动一个 MetaStoreServer，客户端利用 Thrift 协议通过 MetaStoreServer 访问元数据库。

关于 Hive 元数据的使用配置，我们将在 11.5 节“Hive 的 JDBC 接口”中进行详细介绍。

11.2 Hive 的基本操作

本节中我们将介绍 Hive 的基本操作，包括 Hive 在集群上的安装配置，以及 Hive 的 Web UI 的使用。

11.2.1 在集群上安装 Hive

Hadoop-0.19.0 之后的版本都集成了 Hive，这样安装起来非常方便。当然读者也可自行下载安装，下面我们具体介绍如何下载、安装、配置 Hive。

读者通过以下命令下载 Hive 安装包：

```
wget http://www.apache.org/dist/hadoop/hive/hive-0.5.0/hive-0.5.0-bin.tar.gz
tar xzf hive-0.5.0-bin.tar.gz
cd hive-0.5.0
```

然后，修改解压文件夹下 bin/hive-config.sh 的文件内容，添加如下两行记录以便告诉 Hive，Hadoop 的安装路径和 Hive 的安装目录：

```
export HADOOP=/usr/local/hadoop-0.20.2 //Hadoop 安装目录
export HIVE_HOME=/usr/local/hadoop-0.20.2/contrib/hive //Hive 安装目录
```

以如上方式安装 Hive 的前提是你已正确安装了 Hadoop，接下来继续配置 Hive 安装目录下的 conf/hive-site.xml 文件，由于 conf 里没有 hive-site.xml 文件，所以需要复制一份 hive-default.xml，主要配置项如下：

- 1) hive.metastore.warehouse.dir：(HDFS 上的) 数据目录。
- 2) hive.exec.scratchdir：(HDFS 上的) 临时文件目录。
- 3) 连接数据库配置。

默认 meta 数据库为 Derby，当然你也可以配置元数据库为 MySQL，这里我们需要修改 hive-site.xml，配置以下参数：

- `javax.jdo.option.ConnectionURL` : 元数据连接字符串。
 - `javax.jdo.option.ConnectionDriverName` : DB 连接引擎, MySQL 为 `com.mysql.jdbc.Driver`。
 - `javax.jdo.option.ConnectionUserName` : DB 连接用户名。
 - `javax.jdo.option.ConnectionPassword` : DB 连接密码。
- 下面是 `hive-site.xml` 文件的配置示例:

```
<property>
<name>javax.jdo.option.ConnectionURL</name>
<value>jdbc:mysql://localhost:3306/shicheng_hive?useUnicode=true&characterEncoding=UTF-8&createDatabaseIfNotExist=true</value>
<description>JDBC connect string for a JDBC metastore</description>
</property>
<property>
<name>javax.jdo.option.ConnectionDriverName</name>
<value>com.mysql.jdbc.Driver</value>
<description>Driver class name for a JDBC metastore</description>
</property>
<property>
<name>javax.jdo.option.ConnectionUserName</name>
<value>hive</value>
<description>username to use against metastore database</description>
</property>
<property>
<name>javax.jdo.option.ConnectionPassword</name>
<value>hive123</value>
<description>password to use against metastore database</description>
</property>
```

此外, 我们还要在 `~/.bashrc` 里设置环境变量, 将 Hive 的配置目录指向用户主目录:

```
export HADOOP_HOME=/home/hadoop/hadoop/hadoop-current
export HIVE_HOME=/home/hive/hive (你自己的hive安装路径)
export PATH=$PATH:$HADOOP_HOME/bin:$HIVE_HOME/bin
```

这样就可以直接在用户目录下运行 Hadoop 和 hive 命令了。

下面在终端输入 `bin/hive`, 即可启动类似 MySQL 的 shell, `hive -f` 参数直接执行某个文件里的查询语句。

当你看到 `hive>`, 那么恭喜你, Hive 已经正确安装可以运行了。

另外, Hive 还提供了丰富的 Wiki 文档, 读者可以参考以下链接。

- Hive 的 Wiki 页面: <http://wiki.apache.org/hadoop/Hive>
- Hive 入门指南: <http://wiki.apache.org/hadoop/Hive/GettingStarted>
- HQL 查询语言指南: <http://wiki.apache.org/hadoop/Hive/HiveQL>
- 演示文稿: <http://wiki.apache.org/hadoop/Hive/Presentations>

由于 Hive 本身还处在不断的发展中, 很多时候文档更新的速度还赶不上 Hive 本身的

更新速度，所以，如果你想了解 Hive 最新的发展动态或想和研究者进行交流，那么可以加入 Hive 的邮件列表，用户：hive-user@hadoop.apache.org，开发者：hive-dev@hadoop.apache.org。

11.2.2 配置 Hive

按照 11.2.1 节的描述安装好 Hive 后，就可以进行简单的数据操作了。但在实际应用中，不可避免地要进行参数的配置和调优，本节我们将对 Hive 参数的设置进行介绍。

首先，在进行操作前要确保目录权限配置正确：将 /tmp 目录配置成所有用户都有 write 权限，Table 对应目录的 owner 必须是 Hive 启动用户。

其次，可以通过调整 Hive 的参数来调优 HQL 代码的执行效率或帮助管理员进行定位。参数设置可以通过配置文件、命令行参数或参数声明的方式进行设定。下面具体进行介绍。

1. 配置文件

Hive 的配置文件包括：

□ 用户自定义配置文件，即 \$HIVE_CONF_DIR/hive-site.xml；

□ 默认配置文件，即 \$HIVE_CONF_DIR/hive-default.xml。

需要注意的是，用户自定义配置会覆盖默认配置。另外，Hive 也会读入 Hadoop 的配置，因为 Hive 是作为 Hadoop 的客户端启动的，Hadoop 的配置文件包括：

□ \$HADOOP_CONF_DIR/hive-site.xml；

□ \$HADOOP_CONF_DIR/hive-default.xml。

同样，Hive 的配置文件会覆盖 Hadoop 的配置，这里配置文件的修改对所有启动的 Hive 都有效。

2. 运行时配置

当运行 Hive QL 时可以进行参数声明。Hive 的查询可通过执行 MapReduce 任务来实现，而有些查询可以通过控制 Hadoop 的配置参数来实现。在命令行接口（CLI）中可以通过命令 SET 来设置参数，例如：

```
hive> SET mapred.job.tracker=myhost.mycompany.com:50030
hive> SET mapred.reduce.tasks=100;
hive> SET -v
```

通过 SET -v 命令可以查看当前设定的所有信息。需要指出的是，通过 CLI 的 SET 命令设定的作用域是 Session 级的，只对本次操作有作用。此外，SerDe 参数必须写在建表语句中。例如：

```
create table if not exists t_Student (
name    string
)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe'
WITH SERDEPROPERTIES (
```



```
'field.delim'='\t',
'escape.delim'='\\',
'serialization.null.format'=' '
) STORED AS TEXTFILE;
```

类似 `serialization.null.format` 这样的参数，必须和某个表或分区关联。在 DDL 外部声明不起作用。

3. 设置本地模式

对于大多数查询 Query，Hive 编译器会产生 MapReduce 任务，这些任务会被提交到 MapReduce 集群，集群可以用参数 `mapred.job.tracker` 指明。

需要说明的是，Hadoop 支持在本地或集群中运行 Hive 提交的查询，这对于小数据集查询的运行是非常有用的，可以避免将任务分布到大型集群中而降低了效率。在 MapReduce 任务提交给 Hadoop 之后，HDFS 中的文件访问对用户来说是透明的。相反，如果是大数据集的查询，那么需要设定 Hive 的查询交给集群运行，这样就可以利用集群的并行性来提高效率了。我们可以通过以下参数设定 Hive 查询在本地运行：

```
hive> SET mapred.job.tracker=local;
```

最新的 Hive 版本都支持在本地自动运行 MapReduce 任务：

```
hive> SET hive.exec.mode.local.auto=false;
```

可以看到该属性默认是关闭的。如果设定为开启 `enable`，Hive 就会先分析查询中的每个 MapReduce 任务，当任务的输入数据规模小于 `hive.exec.mode.local.auto.inputbytes.max` 属性值（默认为 128MB），并且全部的 map 数少于 `hive.exec.mode.local.auto.tasks.max` 的属性值（默认为 4），全部的 Reduce 任务数为 1 或 0 时，则任务会自动选择在本地模式运行。

4. Error Logs 错误日志

Hive 使用 `log4j` 记录日志。在默认情况下，日志文件的记录等级是 `WARN`（即存储紧急程度为 `WARN` 及以上的错误信息），存储在 `/tmp/{user.name}/hive.log` 文件夹下。如果用户想要在终端看到日志内容，则可以通过设置以下参数达到目的：

```
bin/hive -hiveconf hive.root.logger=INFO,console
```

同样，用户也可以改变日志记录等级：

```
bin/hive -hiveconf hive.root.logger=INFO,DRFA
```

Hive 在 Hadoop 执行阶段的日志由 Hadoop 配置文件配置。通常来说，Hadoop 会对每个 map 和 reduce 任务对应的执行节点生成一个日志文件。这个日志文件可以通过 JobTracker 的 Web UI 获得。错误日志对调试错误非常有用，当运行过程中遇到 Bug 时可以向 `hive-d ev@hadoop.apache.org` 提交。

11.3 Hive QL 详解

11.3.1 数据定义（DDL）操作

数据定义语句主要包括以下语句。

1. 创建表

创建表语法：

```
CREATE [EXTERNAL] TABLE [IF NOT EXISTS] table_name
[(col_name data_type [COMMENT col_comment], ...)]
[COMMENT table_comment]
[PARTITIONED BY (col_name data_type [col_comment], col_name data_type [COMMENT
col_comment], ...)]
[CLUSTERED BY (col_name, col_name, ...) [SORTED BY (col_name, ...)] INTO num_
buckets BUCKETS]
[ROW FORMAT row_format]
[STORED AS file_format]
[LOCATION hdfs_path]
[AS select_statement] (Note: this feature is only available on the latest trunk
or versions higher than 0.4.0.)

CREATE [EXTERNAL] TABLE [IF NOT EXISTS] table_name
LIKE existing_table_name
[LOCATION hdfs_path]

data_type
: primitive_type
| array_type
| map_type
primitive_type
: TINYINT
| SMALLINT
| INT
| BIGINT
| BOOLEAN
| FLOAT
| DOUBLE
| STRING
array_type
: ARRAY < primitive_type >
map_type
: MAP < primitive_type, primitive_type >
row_format
: DELIMITED [FIELDS TERMINATED BY char] [COLLECTION ITEMS TERMINATED BY char]
[MAP KEYS TERMINATED BY char]
| SERDE serde_name [WITH SERDEPROPERTIES property_name=property_value, property_
name=property_value, ...]
file_format:
: SEQUENCEFILE
| TEXTFILE
```

```
| INPUTFORMAT input_format_classname OUTPUTFORMAT output_format_classname
```

下面是相关的说明。

- ❑ **CREATE TABLE** 为创建一个指定名字的表。如果相同名字的表已经存在，则抛出异常，用户可以用 **IF NOT EXIST** 选项来忽略这个异常。
- ❑ **EXTERNAL** 关键字可以让用户创建一个外部表，在创建表的同时指定一个指向实际数据的路径（**LOCATION**），Hive 创建内部表时，会将数据移动到数据仓库指向的路径中；若创建外部表，则仅记录数据所在的路径，不对数据的位置做任何改变。当删除表时，内部表的元数据和数据会被一起删除，而外部表只删除元数据，不删除数据。
- ❑ **LIKE** 格式修饰的 **CREATE TABLE** 命令允许复制一个已存在表的定义，而不复制它的数据内容。

这里还需要说明的是，用户可以使用自定义的 SerDe 或自带的 SerDe 创建表。SerDe 是 **Serialize/Deserialize** 的简称，目的是用于序列化和反序列化。在 Hive 中，序列化和反序列化即是在 **key/value** 和 **hive table** 的每个列值之间的转化。如果没有指定 **ROW FORMAT** 或 **ROW FORMAT DELIMITED**，创建表就使用自带的 SerDe。如果使用自带的 SerDe，则必须指定字段列表。关于字段类型，请参考用户指南的类型部分。定制的 SerDe 字段列表可以是指定的，但是 Hive 将通过查询 SerDe 决定实际的字段列表。

如果数据需要存储为纯文本文件，则请使用 **STORED AS TEXTFILE**。如果数据需要压缩，则使用 **STORED AS SEQUENCEFILE**。**INPUTFORMAT** 和 **OUTPUTFORMAT** 定义一个与 **InputFormat** 和 **OutputFormat** 类相对应的名字作为字符串，例如，将 “**org.apache.hadoop.hive.contrib.fileformat.base64**” 定义为 “**Base64TextInputFormat**”。

Hive 还支持建立带有分区（**Partition**）的表。有分区的表可以在创建的时候使用 **PARTITIONED BY** 语句。一个表可以拥有一个或多个分区，每个分区单独存在于一个目录下。而且，表和分区都可以对某个列进行 **CLUSTERED BY** 操作，将若干个列放入一个桶（**bucket**）中。也可以利用 **SORT BY** 列来存储数据，以提高查询性能。

表名和列名不区分大小写，但 SerDe 和属性名是区分大小写的。表和列的注释分别是单引号的字符串。

下面是一个创建表的例子：

```
CREATE TABLE page_view(viewTime INT, userid BIGINT,
page_url STRING, referrer_url STRING,
ip STRING COMMENT 'IP Address of the User')
COMMENT 'This is the page view table'
PARTITIONED BY(dt STRING, country STRING)
STORED AS SEQUENCEFILE;
```

该语句创建 **page_view** 表，包括 **viewTime**、**userid**、**page_url**、**referrer_url**、**ip** 列及注释，并建立表分区，其中用 **Ctrl+A** 分隔字段，并使用换行分隔符来定义文件中的数据格式。

```
CREATE TABLE page_view(viewTime INT, userid BIGINT,
page_url STRING, referrer_url STRING,
ip STRING COMMENT 'IP Address of the User')
COMMENT 'This is the page view table'
PARTITIONED BY(dt STRING, country STRING)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '\001'
STORED AS SEQUENCEFILE;
```

下面的语句用于创建和上表相同的表 (page_view) :

```
CREATE TABLE page_view(viewTime INT, userid BIGINT,
page_url STRING, referrer_url STRING,
ip STRING COMMENT 'IP Address of the User')
COMMENT 'This is the page view table'
PARTITIONED BY(dt STRING, country STRING)
CLUSTERED BY(userid) SORTED BY(viewTime) INTO 32 BUCKETS
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '\001'
COLLECTION ITEMS TERMINATED BY '\002'
MAP KEYS TERMINATED BY '\003'
STORED AS SEQUENCEFILE;
```

在上述示例中, page_view表按照userid进行分区划分到不同的桶中, 并按照viewTime值的大小进行排序存储。这样的组织结构允许用户通过userid属性高效地对集群列进行采样。

到目前为止的所有例子中, 数据都存储在 <hive.metastore.warehouse.dir> / page_view 中, 它的值在 Hive 配置的文件 hive-site.xml 中设定, 代码如下:

```
CREATE EXTERNAL TABLE page_view(viewTime INT, userid BIGINT,
page_url STRING, referrer_url STRING,
ip STRING COMMENT 'IP Address of the User',
country STRING COMMENT 'country of origination')
COMMENT 'This is the staging page view table'
ROW FORMAT DELIMITED FIELDS TERMINATED BY '\054'
STORED AS TEXTFILE
LOCATION '<hdfs_location>';
```

用上述语句可以创建一个 page_view 表, 并将表存储在 HDFS 的任何目录下。但是必须确保上面的查询数据是分隔开的。

2. 删除表

命令如下所示:

```
DROP TABLE table_name
```

DROP TABLE 用于删除表的元数据和数据。如果配置了 Trash, 数据将删除 Trash/Current 目录, 元数据将完全丢失。当删除 EXTERNAL 定义的表时, 表中的数据不会从文件

系统中删除。

3. 修改表、分区语句

ALTER TABLE 语句用于改变一个已经存在的表结构，比如增加列或分区，改变 SerDe、添加表和 SerDe 的属性或重命名表。

❑ 增加分区，命令如下所示：

```
ALTER TABLE table_name ADD partition_spec [ LOCATION 'location1' ] partition_spec [
    LOCATION 'location2' ] ...
partition_spec:
: PARTITION (partition_col = partition_col_value, partition_col = partiton_col_
    value, ...)
```

用户可以用 ALTER TABLE ADD PARTITION 来对表增加分区。当分区名是字符串时加引号，例如：

```
ALTER TABLE page_view ADD
    PARTITION (dt='2010-08-08', country='us')
        location '/path/to/us/part080808'
    PARTITION (dt='2010-08-09', country='us')
        location '/path/to/us/part080809';
```

❑ 删除分区，命令如下所示：

```
ALTER TABLE table_name DROP
    partition_spec, partition_spec,...
```

用户可以用 ALTER TABLE DROP PARTITION 来删除分区，分区的元数据和数据将被一并删除，例如：

```
ALTER TABLE page_view
    DROP PARTITION (dt='2010-08-08', country='us');
```

❑ 重命名表，命令如下所示：

```
ALTER TABLE table_name RENAME TO new_table_name
```

这个命令可以让用户为表更名。数据所在的位置和分区名并不改变。换言之，老的表名并未“释放”，对老表的更改会改变新表的数据。

❑ 改变列名字 / 类型 / 位置 / 注释，命令如下所示：

```
ALTER TABLE table_name CHANGE [COLUMN]
    col_old_name col_new_name column_type
    [COMMENT col_comment]
    [FIRST|AFTER column_name]
```

这个命令允许用户修改列的名称、数据类型、注释或位置，例如：

```
CREATE TABLE test_change (a int, b int, c int);
```

```
ALTER TABLE test_change CHANGE a a1 INT; // 将 a 列的名字改为 a1
ALTER TABLE test_change CHANGE a a1 STRING AFTER b; // 将 a 列的名字改为 a1, a 列的数据类型改为 string, 并将它放置在列 b 之后
```

修改后, 新的表结构为: b int, a1 string, c int。

```
ALTER TABLE test_change CHANGE b b1 INT FIRST; // 会将 b 列的名字修改为 b1, 并将它放在第一列
```

修改后, 新表的结构为: b1 int, a string, c int。

注意 列的改变只会修改 Hive 的元数据, 而不会改变实际数据。用户应该确保元数据的定义和实际数据结构的一致性。

□ 增加 / 更新列, 命令如下所示:

```
ALTER TABLE table_name ADD|REPLACE
  COLUMNS (col_name data_type [COMMENT col_comment], ...)
```

ADD COLUMNS 允许用户在当前列的末尾、分区列之前增加新的列。REPLACE COLUMNS 删除当前列, 加入新的列。只有在使用 native 的 SerDe (DynamicSerDe 或 MetadataTypeColumnsetSerDe) 时才可以这么做。

□ 增加表属性, 命令如下所示:

```
ALTER TABLE table_name SET TBLPROPERTIES table_properties
table_properties:
  : (property_name = property_value, property_name = property_value, ... )
```

用户可以用这个命令向表中增加 metadata, 目前 last_modified_user、last_modified_time 属性都是由 Hive 自动管理的。用户可以向列表中增加自己的属性, 可以使用 DESCRIBE EXTENDED TABLE 来获得这些信息。

□ 增加 SerDe 属性, 命令如下所示:

```
ALTER TABLE table_name
  SET SERDE serde_class_name
  [WITH SERDEPROPERTIES serde_properties]
```

```
ALTER TABLE table_name
  SET SERDEPROPERTIES serde_properties
```

```
serde_properties:
  : (property_name = property_value,
    property_name = property_value, ... )
```

这个命令允许用户向 SerDe 对象增加用户定义的元数据。Hive 为了序列化和反序列化数据, 将会初始化 SerDe 属性, 并将属性传给表的 SerDe。这样, 用户可以为自定义的 SerDe 存储属性。

□ 改变表文件格式和组织，命令如下所示：

```
ALTER TABLE table_name SET FILEFORMAT file_format
ALTER TABLE table_name CLUSTERED BY (col_name, col_name, ...)
[SORTED BY (col_name, ...)] INTO num_buckets BUCKETS
```

这个命令修改了表的物理存储属性。

注意 这些命令只是修改 Hive 的元数据，不能重组或格式化现有的数据。用户应该确定实际数据的分布符合元数据的定义。

4. 创建 / 删除视图

目前，只有 Hive 0.6 之后的版本才支持视图。

□ 创建表视图，命令如下所示：

```
CREATE VIEW [IF NOT EXISTS] view_name [ (column_name [COMMENT column_comment], ...) ]
[COMMENT view_comment]
AS SELECT ...
```

CREATE VIEW 以指定的名字创建一个表视图。如果表或视图的名字已经存在，则报错，也可以使用 IF NOT EXISTS 忽略这个错误。

如果没有提供表名，则视图列的名字将由定义的 SELECT 表达式自动生成；如果 SELECT 包括如 $x + y$ 无标量的表达式，则视图列的名字将生成 _C0、_C1 等形式。当重命名列时，可选择性地提供列注释。注释不会从底层列自动继承。如果定义 SELECT 表达式的视图是无效的，那么 CREATE VIEW 语句将失败。

注意，没有关联存储的视图是纯粹的逻辑对象。目前在 Hive 中不支持物化视图。当一个查询引用一个视图时，可以评估视图的定义并为下一步查询提供记录集合。这是一种概念的描述，实际上，作为查询优化的一部分，Hive 可以将视图的定义与查询的定义结合起来，例如从查询到视图所使用的过滤器。

在视图创建的同时确定视图的架构，如果随后再改变基本表（如添加一列）将不会在视图的架构中体现。如果基本表被删除或以不兼容的方式被修改，则该无效视图的查询将失败。

视图是只读的，不能用于 LOAD/INSERT/ALTER。

视图可能包含 ORDER BY 和 LIMIT 子句，如果一个引用了视图的查询也包含了这些子句，那么在执行这些子句时首先要查看视图语句，然后返回结果按视图中的语句执行。例如，如果一个视图 v 指定返回记录 LIMIT 为 5，执行查询语句：select * from v LIMIT 10，那么这个查询最多返回 5 行记录。

以下是创建视图的例子：

```
CREATE VIEW onion_referrers(url COMMENT 'URL of Referring page')
COMMENT 'Referrers to The Onion website'
```

```
AS
SELECT DISTINCT referrer_url
FROM page_view
WHERE page_url='http://www.theonion.com';
```

❑ 删除表视图，命令如下所示：

```
DROP VIEW view_name
```

DROP VIEW 为删除指定视图的元数据，在视图中使用 DROP TABLE 是错误的。例如：

```
DROP VIEW onion_referrers;
```

5. 创建 / 删除函数

❑ 创建函数，命令如下所示：

```
CREATE TEMPORARY FUNCTION function_name AS class_name
```

该语句创建了一个由类名实现的函数。在 Hive 中用户可以使用 Hive 类路径中的任何类，用户通过执行 ADD FILES 语句将函数类添加到类路径，并且可持续使用该函数进行操作。请参阅用户指南 CLI 部分了解有关在 Hive 中添加 / 删除的更多信息。使用该语句注册用户定义函数。

❑ 删除函数

注销用户定义函数的格式如下：

```
DROP TEMPORARY FUNCTION function_name
```

6. 展示描述语句

在 Hive 中，该语句提供一种方法对现有的数据和元数据进行查询。

❑ 显示表，命令如下所示：

```
SHOW TABLES identifier_with_wildcards
```

SHOW TABLES 列出了所有基表与给定正则表达式相匹配的视图。正则表达式只能包含 '*' 作为任意字符以及 [s] 或 |。例如 'page_view'、'page_v *'、'*view|page*'，所有这些将匹配 'page_view' 表。匹配表按字母顺序排列。在元存储中，如果没有找到匹配的表，则不提示错误。

❑ 显示分区，命令如下所示：

```
SHOW PARTITIONS table_name
```

SHOW PARTITIONS 列出了给定基表中所有的现有分区，分区按字母顺序排列。

❑ 显示表 / 分区扩展，命令如下所示：

```
SHOW TABLE EXTENDED [IN|FROM database_name] LIKE identifier_with_wildcards
[PARTITION(partition_desc)]
```

SHOW TABLE EXTENDED 为列出所有给定的匹配正则表达式的表信息。如果分区规范存在，那么用户不能使用正则表达式作为表名。该命令的输出包括基本表信息和文件系统信息，

例如，文件总数、文件总大小、最大文件大小、最小文件大小、最新存储时间和最新更新时间。如果分区存在，则它会输出给定分区的文件系统信息，而不是表中的文件系统信息。

作为视图，SHOW TABLE EXTENDED 用于检索视图的定义。

□ 显示函数，命令如下所示：

```
SHOW FUNCTIONS "a.*"
```

SHOW FUNCTIONS 为列出用户定义和建立所有匹配给定的正则表达式的函数。可以给所有函数用 ".*"。

□ 描述表 / 列，命令如下所示：

```
DESCRIBE [EXTENDED] table_name[DOT col_name]
DESCRIBE [EXTENDED] table_name[DOT col_name ( [DOT field_name] | [DOT '$elem$'] |
[DOT '$key$'] | [DOT '$value$'] ) * ]
```

DESCRIBE TABLE 为显示列信息，包括给定表的分区。如果指定 EXTENDED 关键字，则将在序列化形式中显示表的所有元数据。DESCRIBE TABLE 通常只用于调试，而不适用于平常的使用中。

如果表有复杂的列，可以通过指定数组元素 table_name.complex_col_name (和 '\$elem \$' 作为数组元素，'\$key\$' 为图的主键，'\$value\$' 为图的属性) 来检查该列的属性。对于复杂的列类型，可以使用这些定义递归地进行查询。

□ 描述分区，命令如下所示：

```
DESCRIBE [EXTENDED] table_name partition_spec
```

该语句列出了给定分区的元数据，输出和 DESCRIBE TABLE 类似。目前，在查询计划准备阶段不能使用这些列信息。

11.3.2 数据操作 (DML)

下面我们将详细介绍 DML，它是数据操作类语言，其中包括向数据表加载文件，写查询结果等操作。

1. 向数据表中加载文件

当数据被加载至表中时，不会对数据进行任何转换。Load 操作只是将数据复制 / 移动到 Hive 表对应的位置上，代码如下：

```
LOAD DATA [LOCAL] INPATH 'filepath' [OVERWRITE]
INTO TABLE tablename
[PARTITION (partcol1=val1, partcol2=val2 ...)]
```

其中，filepath 可以是相对路径 (例如：project/data1)，filepath 可以是绝对路径 (例如：/user/admin/project/ data1)，或者 filepath 可以是完整的 URI (例如：hdfs://namenodeIP:9000/user/admin/project/data1)。加载的目标可以是一个表或分区。如果表包含分区，则必须指定

每个分区的分区名。filepath 可以引用一个文件（这种情况下，Hive 会将文件移动到表所对应的目录中）或一个目录（这种情况下，Hive 会将目录中的所有文件移动至表所对应的目录中）。如果指定 LOCAL，那么 load 命令会去查找本地文件系统中的 filepath。如果发现是相对路径，则路径会被解释为相对于当前用户的当前路径。用户也可以为本地文件指定一个完整的 URI，比如 file:///user/hive/project/data。load 命令会将 filepath 中的文件复制到目标文件系统中，目标文件系统由表的位置属性决定。被复制的数据文件移动到表的数据对应的位置上。如果没有指定 LOCAL 关键字，filepath 指向一个完整的 URI，Hive 则会直接使用这个 URI。如果没有指定 schema 或 authority，Hive 则会使用在 Hadoop 配置文件中定义的 schema 和 authority，fs.default.name 属性指定 NameNode 的 URI。如果路径不是绝对的，则 Hive 相对于 /user/ 进行解释。Hive 还会将 filepath 中指定的文件内容移动到 table（或者 partition）所指定的路径中。如果使用 OVERWRITE 关键字，则目标表（或者分区）中的内容（如果有）会被删除，并将 filepath 指向的文件 / 目录中的内容添加到表 / 分区中。如果目标表（分区）中已经有文件，并且文件名和 filepath 中的文件名冲突，那么现有的文件会被新文件替代。

2. 将查询结果插入 Hive 表中

查询的结果通过 insert 语法加入到表中，代码如下：

```
INSERT OVERWRITE TABLE tablename1 [PARTITION (partcol1=val1, partcol2=val2 ...)]
  select_statement1 FROM from_statement
Hive extension (multiple inserts):
FROM from_statement
INSERT OVERWRITE TABLE tablename1 [PARTITION (partcol1=val1, partcol2=val2 ...)]
  select_statement1
[INSERT OVERWRITE TABLE tablename2 [PARTITION ...] select_statement2] ...
Hive extension (dynamic partition inserts):
INSERT OVERWRITE TABLE tablename PARTITION (partcol1 [=val1], partcol2 [=val2] ...)
  select_statement FROM from_statement
```

这里需要注意的是，插入可以针对一个表或一个分区进行操作，如果对一个表进行了划分，那么插入时就要指定划分列的属性值以确定分区。每个 Select 语句的结果会被写入选择的表或分区中，OVERWRITE 关键字会强制将输出结果写入。其中输出格式和序列化方式由表的元数据决定。Hive 中的多表插入，可以减少数据扫描的次数，因为 Hive 可以只扫描输入数据一次，而对输入数据进行多个操作命令。

3. 将查询的结果写入文件系统

查询结果可以通过如下命令插入文件系统目录中：

```
INSERT OVERWRITE [LOCAL] DIRECTORY directory1 SELECT ... FROM ...
Hive extension (multiple inserts):
FROM from_statement
INSERT OVERWRITE [LOCAL] DIRECTORY directory1 select_statement1
[INSERT OVERWRITE [LOCAL] DIRECTORY directory2 select_statement2] ...
```

这里需要注意的是，目录可以是完整的 URI。如果 scheme 或 authority 没有定义，那么 Hive 会使用 Hadoop 配置参数 `fs.default.name` 中的 scheme 和 authority 来定义 NameNode 的 URI。如果使用 LOCAL 关键字，那么 Hive 会将数据写入本地文件系统中。

数据写入文件系统时会进行文本序列化，且每列用 ^A 来区分，换行表示一行数据结束。如果任何一列不是原始类型，那么这些列将会被序列化为 JSON 格式。

11.3.3 SQL 操作

下面是一个标准的 Select 语句语法定义：

```
SELECT [ALL | DISTINCT] select_expr, select_expr, ...
FROM table_reference
[WHERE where_condition]
[GROUP BY col_list]
[   CLUSTER BY col_list
  | [DISTRIBUTE BY col_list] [SORT BY col_list]
]
[LIMIT number]
```

下面对其中重要的定义进行说明。

(1) table_reference

table_reference 指明查询的输入，它可以是一个表或一个视图、一个子查询。下面是一个简单的查询，可以检索所有表 t1 中的列和行：

```
SELECT * FROM t1
```

(2) WHERE

where_condition 是一个布尔表达式。比如下面的查询只输出 sales 表中 amount>10 并且 region 属性值为 US 的记录：

```
SELECT * FROM sales WHERE amount > 10 AND region = "US"
```

(3) ALL 和 DISTINCT

ALL 和 DISTINCT 选项可以定义重复的行是否要返回。如果没有定义，那么默认为 ALL，即输出所有的匹配记录而不删除重复的记录，代码如下：

```
hive> SELECT col1, col2 FROM t1
1 3
1 3
1 4
2 5
hive> SELECT DISTINCT col1, col2 FROM t1
1 3
1 4
2 5
hive> SELECT DISTINCT col1 FROM t1
1
2
```

(4) LIMIT

LIMIT 可以控制输出的记录数，并随机选取检索结果中的相应数目输出：

```
SELECT * FROM t1 LIMIT 5
```

下面代码为输出 Top-k, k=5 的查询结果：

```
SET mapred.reduce.tasks = 1
SELECT * FROM sales SORT BY amount DESC LIMIT 5
```

(5) 使用正则表达式

SELECT 声明可以匹配使用一个正则表达式的列，下面的例子会对 sales 表中除了 ds 和 hr 的所有列进行扫描：

```
SELECT '(ds|hr)?+.' FROM sales
```

(6) 基于分区的查询

通常来说，SELECT 查询要扫描全部的表。如果一个表是使用 PARTITIONED BY 语句产生的，那么查询可以对输入进行“剪枝”，只对表的相关部分进行扫描。Hive 现在只对 WHERE 中指定的分区断言进行“剪枝”式的扫描。举例来说，如果一个表 page_view 按照 date 列的值进行了分区，那么下面的查询可以检索出日期为 2010-03-01 的行记录：

```
SELECT page_views.*
FROM page_views
WHERE page_views.date >= '2010-03-01' AND page_views.date <= '2010-03-31'
```

(7) HAVING

Hive 目前不支持 HAVING 语句，但是可以使用子查询实现，以下为例：

```
SELECT col1 FROM t1 GROUP BY col1 HAVING SUM(col2) > 10
```

可以表示为：

```
SELECT col1 FROM (SELECT col1, SUM(col2) AS col2sum FROM t1 GROUP BY col1) t2
WHERE t2.col2sum > 10
```

我们可以将查询的结果写入到目录中：

```
hive> INSERT OVERWRITE DIRECTORY '/tmp/hdfs_out' SELECT a.* FROM invites a WHERE
a.ds='2009-09-01';
```

上面的例子将查询结果写入 /tmp/hdfs_out 目录中，或者写入本地文件路径，代码如下所示：

```
hive> INSERT OVERWRITE LOCAL DIRECTORY '/tmp/local_out' SELECT a.* FROM pokes a;
```

其他例如 GROUP BY JOIN 的作用和 SQL 相同，就不再赘述，下面是使用举例，详细的可以查看 <http://wiki.apache.org/hadoop/Hive/LanguageManual>。

(8) GROUP BY

```
hive> FROM invites a INSERT OVERWRITE TABLE events SELECT a.bar, count(*) WHERE
```

```
a.foo > 0 GROUP BY a.bar;
hive> INSERT OVERWRITE TABLE events SELECT a.bar, count(*) FROM invites a WHERE
a.foo > 0 GROUP BY a.bar;
```

(9) JOIN

```
hive> FROM pokes t1 JOIN invites t2 ON (t1.bar = t2.bar) INSERT OVERWRITE TABLE
events SELECT t1.bar, t1.foo, t2.foo;
```

(10) 多表 INSERT

```
FROM src
INSERT OVERWRITE TABLE dest1 SELECT src.* WHERE src.key < 100
INSERT OVERWRITE TABLE dest2 SELECT src.key, src.value WHERE src.key >= 100 and
src.key < 200
INSERT OVERWRITE TABLE dest3 PARTITION(ds='2010-04-08', hr='12') SELECT src.key
WHERE src.key >= 200 and src.key < 300
INSERT OVERWRITE LOCAL DIRECTORY '/tmp/dest4.out' SELECT src.value WHERE src.
key >= 300;
```

(11) STREAMING

```
hive> FROM invites a INSERT OVERWRITE TABLE events SELECT TRANSFORM(a.foo, a.bar)
AS (oof, rab) USING '/bin/cat' WHERE a.ds > '2010-08-09';
```

这个命令会将数据输入给 map 操作（通过 /bin/cat 命令），同样也可以将数据流式输入给 reduce 操作。

11.3.4 Hive QL 的使用实例

下面用两个例子来对 Hive QL 的使用方法进行介绍，我们可以看到它与传统 SQL 语句的异同点。

1. 电影评分

首先创建表，并使用 tab 空格定义文本格式：

```
CREATE TABLE u_data (
  userid INT,
  movieid INT,
  rating INT,
  unixtime STRING)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '\t'
STORED AS TEXTFILE;
```

然后下载数据文本文件，并解压，代码如下所示：

```
wget http://www.grouplens.org/system/files/ml-data.tar__0.gz
tar xvfz ml-data.tar__0.gz
```



将文件加载到表中，代码如下所示：

```
LOAD DATA LOCAL INPATH 'ml-data/u.data'
OVERWRITE INTO TABLE u_data;
Count the number of rows in table u_data:
SELECT COUNT(*) FROM u_data; // 由于版本问题，如果此处出现错误，你可能需要使用 COUNT(1) 替换
COUNT(*)
```

下面可以基于该表进行一些复杂的数据分析操作，此处我们使用 Python 语言，代码如下所示：

```
# 创建 weekday_mapper.py 文件:
import sys
import datetime
for line in sys.stdin:
    line = line.strip()
    userid, movieid, rating, unixtime = line.split('\t')
    weekday = datetime.datetime.fromtimestamp(float(unixtime)).isoweekday()
    print '\t'.join([userid, movieid, rating, str(weekday)])
# 使用下面的 mapper 脚本:
CREATE TABLE u_data_new (
    userid INT,
    movieid INT,
    rating INT,
    weekday INT)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY '\t';
add FILE weekday_mapper.py;
INSERT OVERWRITE TABLE u_data_new
SELECT
    TRANSFORM (userid, movieid, rating, unixtime)
    USING 'python weekday_mapper.py'
    AS (userid, movieid, rating, weekday)
FROM u_data;
SELECT weekday, COUNT(*)
FROM u_data_new
GROUP BY weekday;
```

2. Apache 网络日志数据 (Weblog)

Apache 网络日志数据格式可以定制，一般管理者都使用默认的。对于默认设置的 Apache Weblog 可以使用以下命令创建表：

```
add jar ../build/contrib/hive_contrib.jar;
CREATE TABLE apachelog (
    host STRING,
    identity STRING,
    user STRING,
    time STRING,
    request STRING,
```




```
description>
</property>
```

某些情况下,读者可能需要自动在服务器模式下安装 Derby,具体内容可参见 <http://wiki.apache.org/hadoop/HiveDerbyServerMode>,这样便可以同时运行多个会话。

随后启动服务即可,在配置文件中,监听端口默认是 9999,也可以自己到 hive-default.xml 中定制,代码如下所示:

```
bin/hive --service hwi
10/03/01 03:32:03 INFO hwi.HWIServer: HWI is starting up
10/03/01 03:32:03 INFO mortbay.log: Logging to org.slf4j.impl.
Log4jLoggerAdapter(org.mortbay.log) via org.mortbay.log.Slf4jLog
10/03/01 03:32:03 INFO mortbay.log: jetty-6.1.14
10/03/01 03:32:03 INFO mortbay.log: Extract jar:file:/home/hadoop/hive-0.3.99.1+0/
lib/hive_hwi.war! to /tmp/Jetty_0_0_0_9999_hive_hwi.war_hwi__bw65n0/webapp
10/03/01 03:32:03 INFO mortbay.log: Started SocketConnector@0.0.0.0:9999
```

这样我们通过浏览器访问网络接口的地址: <http://MsaterIP:9999/hwi> 即可,如图 11-2 所示。

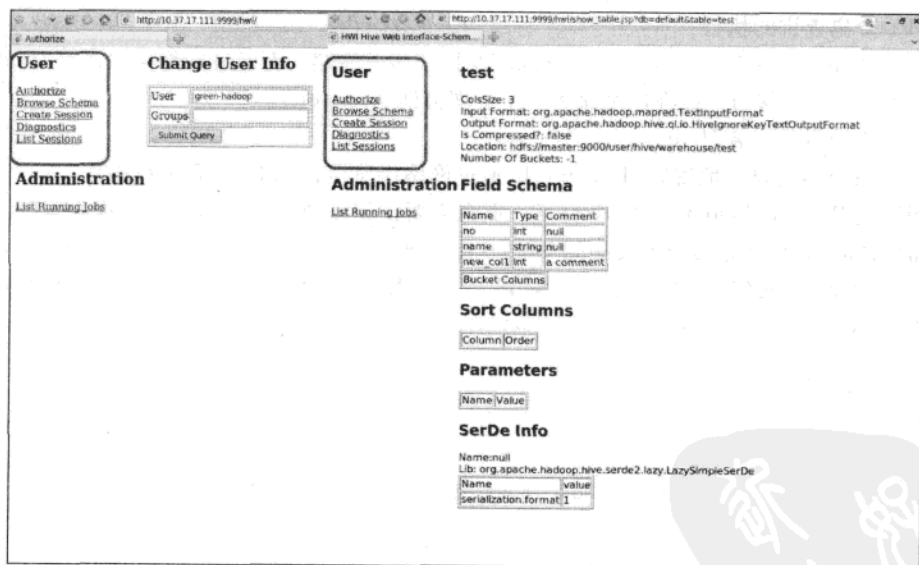


图 11-2 Hive 的网络接口 (WebUI)

可以看到 Hive 的网络接口拉近了用户和系统的距离,我们可以通过网络直接创建会话,并进行查询。用户界面和功能展示非常直观,适合于刚接触到 Hive 的用户。

11.5 Hive 的 JDBC 接口

通过上面的介绍,我们知道用户可以使用命令行接口 (CLI) 和 Hive 进行交互,也可以

使用网络接口（Web UI）和 Hive 进行交互。本节我们将具体介绍 JDBC 接口。如果是集群中的节点作为客户端来访问 Hive，则可以直接使用 jdbc:hive:// 来进行访问。对于一个非集群节点的客户端来说，可以使用 jdbc:hive://host:port/dbname 来进行访问。

使用 JDBC 接口，首先需要配置 Hive 的配置文件 hive-site.xml，代码如下所示：

```
<property>
  <name>javax.jdo.option.ConnectionURL</name>
  <value>jdbc:derby:;databaseName=metastore_db;create=true</value>
  <!-- 此处为 true，表明使用嵌入式的 derby 建立元数据表为 metastore-->
  <!--<value>jdbc:derby://10.37.17.111:7777/metastore;create=true</value>--> <!-- 表示使用客服模式的 derby，metastore 为数据库名，10.37.17.111 为 derby 服务端的 IP 地址，而 7777 为服务端的端口号 -->
  <description>JDBC connect string for a JDBC metastore</description>
</property>
<property>
  <name>javax.jdo.option.ConnectionDriverName</name>
  <value>org.apache.derby.jdbc.EmbeddedDriver</value>
  <!-- 表示使用嵌入式的 derby-->
  <!--<value>org.apache.derby.jdbc.ClientDriver</value>--> 表示使用客服模式的 derby
  <description>Driver class name for a JDBC metastore</description>
</property>
```

上例的配置文件是使用 Derby 数据库作为元数据的存储数据库，当然，也可以使用 MySQL 或 Oracle 的传统关系型数据库进行元数据的存储，下面给出使用 MySQL 存储元数据的配置文件，需要修改 Hive 配置文件 hive-site.xml 的属性值，代码如下所示：

```
<property>
  <name>hive.metastore.local</name>
  <value>true</value>
</property>

<property>
  <name>javax.jdo.option.ConnectionURL</name>
  <value>jdbc:mysql://centos1:3306/hive?createDatabaseIfNotExist=true</value>
</property>

<property>
  <name>javax.jdo.option.ConnectionDriverName</name>
  <value>com.mysql.jdbc.Driver</value>
</property>

<property>
  <name>javax.jdo.option.ConnectionUserName</name>
  <value>hadoop 可以访问的 Mysql 账户 </value>
</property>

<property>
  <name>javax.jdo.option.ConnectionPassword</name>
  <value>密码 </value>
```

```
</property>
```

使用 MySQL 存储元数据，除了需要修改配置文件外，还需要确保你的 Hadoop 用户可以访问 MySQL（如何配置 MySQL 这里不再赘述），同时要将 MySQL 相应的连接驱动 jar 包拷贝到 Hive 的 lib 目录下。

下面是一个使用 Java 编写的 JDBC 客户端访问的代码样例：

```
import java.sql.SQLException;
import java.sql.Connection;
import java.sql.ResultSet;
import java.sql.Statement;
import java.sql.DriverManager;

public class HiveJdbcClient {
    private static String driverName = "org.apache.hadoop.hive.jdbc.HiveDriver";

    /**
     * @param args
     * @throws SQLException
     */
    public static void main(String[] args) throws SQLException {
        try {
            Class.forName(driverName);
        } catch (ClassNotFoundException e) {
            // TODO Auto-generated catch block
            e.printStackTrace();
            System.exit(1);
        }
        Connection con = DriverManager.getConnection("jdbc:hive://localhost:10000/
            default", "", "");
        Statement stmt = con.createStatement();
        String tableName = "testHiveDriverTable";
        stmt.executeQuery("drop table " + tableName);
        ResultSet res = stmt.executeQuery("create table " + tableName + " (key int,
            value string)");
        // show tables
        String sql = "show tables '" + tableName + "'";
        System.out.println("Running: " + sql);
        res = stmt.executeQuery(sql);
        if (res.next()) {
            System.out.println(res.getString(1));
        }
        // describe table
        sql = "describe " + tableName;
        System.out.println("Running: " + sql);
        res = stmt.executeQuery(sql);
        while (res.next()) {
            System.out.println(res.getString(1) + "\t" + res.getString(2));
        }
    }
}
```

```

// load data into table
// NOTE: filepath has to be local to the hive server
// NOTE: /tmp/a.txt is a ctrl-A separated file with two fields per line
String filepath = "/tmp/a.txt";
sql = "load data local inpath '" + filepath + "' into table " + tableName;
System.out.println("Running: " + sql);
res = stmt.executeQuery(sql);

// select * query
sql = "select * from " + tableName;
System.out.println("Running: " + sql);
res = stmt.executeQuery(sql);
while (res.next()) {
    System.out.println(String.valueOf(res.getInt(1)) + "\t" + res.getString(2));
}

// regular hive query
sql = "select count(1) from " + tableName;
System.out.println("Running: " + sql);
res = stmt.executeQuery(sql);
while (res.next()) {
    System.out.println(res.getString(1));
}
}
}

```

当前的 JDBC 接口只支持查询的执行及结果的获取，并支持部分元数据的读取。Hive 支持的接口除了 JDBC 外，还有 Python、PHP、ODBC 等。读者可以访问 <http://wiki.apache.org/hadoop/Hive/HiveClient#JDBC> 查看相关信息。

11.6 Hive 的优化

Hive 是针对不同的查询进行优化的，优化可以通过配置来进行控制。本节我们将介绍部分优化策略及优化控制选项。

1. 列裁剪 (Column Pruning)

读数据时，只读取查询中需要用到列，而忽略其他列。例如如下查询：

```
SELECT a,b FROM t WHERE e < 10;
```

其中，表 t 包含 5 个列 (a,b,c,d,e)，经过列裁剪，列 c 和 d 将会被忽略，执行中只会读取 a, b, e 列，要实现列裁剪，需要设置参数 `hive.optimize.cp = true`。

2. 分区裁剪 (Partition Pruning)

在查询过程中减少不必要的分区。例如如下查询：

```
SELECT * FROM (SELECT c1, COUNT(1)
```

```
FROM T GROUP BY c1) subq
WHERE subq.prtm = 100;
SELECT * FROM T1 JOIN
  (SELECT * FROM T2) subq ON (T1.c1=subq.c2)
WHERE subq.prtm = 100;
```

经过分区裁剪优化的查询，会在子查询中考虑 `subq.prtm = 100` 条件，从而减少读入的分区数目。要实现分区裁剪，须设置 `hive.optimize.prtm=true`。

3. Join 操作

当使用写有 Join 操作的查询语句时，有一条原则：应该将条目少的表 / 子查询放在 Join 操作符的左边。原因是在 Join 操作的 reduce 阶段，位于 Join 操作符左边的表的内容会被加载到内存中，将条目少的表放在左边，这样可以有效减少发生内存溢出（OOM：Out of Memory）的几率。

对于一条语句中有多个 Join 的情况，如果 Join 的条件相同可以进行优化，比如如下查询：

```
INSERT OVERWRITE TABLE pv_users
SELECT pv.pageid, u.age FROM page_view p
JOIN user u ON (pv.userid = u.userid)
JOIN newuser x ON (u.userid = x.userid);
```

我们可以进行的优化是，如果 Join 的 key 相同，那么不管有多少个表，都会合并为一个 MapReduce。如果 Join 的条件不相同，比如：

```
INSERT OVERWRITE TABLE pv_users
SELECT pv.pageid, u.age FROM page_view p
JOIN user u ON (pv.userid = u.userid)
JOIN newuser x ON (u.age = x.age);
```

如果 MapReduce 的任务数目和 Join 操作的数目是对应的，那么上述查询和以下查询是等价的：

```
INSERT OVERWRITE TABLE tmptable
SELECT * FROM page_view p JOIN user u
ON (pv.userid = u.userid);

INSERT OVERWRITE TABLE pv_users
SELECT x.pageid, x.age FROM tmptable x
JOIN newuser y ON (x.age = y.age);
```

4. Map Join 操作

Map Join 操作无须 reduce 操作就可以在 map 阶段全部完成，前提是在 map 过程中可以访问到全部需要的数据。比如如下查询：

```
INSERT OVERWRITE TABLE pv_users
SELECT /*+ MAPJOIN(pv) */ pv.pageid, u.age
FROM page_view pv
JOIN user u ON (pv.userid = u.userid);
```



这个查询便可以在 map 阶段全部完成 Join。此时还须设置的相关属性为：hive.join.emit.inter-1 = 1000、hive.mapjoin.size.key = 10000、hive.mapjoin.cache.numrows = 10000。hive.join.emit.inter-1 = 1000 属性定义了输出 Join 的结果前，还要判断右侧进行 Join 的操作数最多可以加载多少行到缓存中。

5. Group By 操作

进行 Group BY 操作时需要注意以下两点。

□ map 端部分聚合。事实上，并不是所有的聚合操作都需要在 reduce 部分进行，很多聚合操作都可以先在 map 端进行部分聚合，然后在 reduce 端得出最终结果。

这里需要修改的参数为：hive.map.aggr = true，用于设定是否在 map 端进行聚合，默认为 True。hive.groupby.mapaggr.checkinterval = 100000，用于设定在 map 端进行聚合操作的条目数。

□ 有数据倾斜（数据分布不均匀）时进行负载均衡。此处需要设定 hive.groupby.skewindata，当选项设定为 true 时，生成的查询计划会有两个 MapReduce 任务。在第一个 MapReduce 中，map 的输出结果集合会随机分布到 reduce 中，每个 reduce 做部分聚合操作，并输出结果。这样处理的结果是，相同的 Group By Key 有可能被分发到不同的 reduce 中，从而达到负载均衡的目的；第二个 MapReduce 任务再根据预处理的数据结果按照 Group By Key 分布到 reduce 中（这个过程可以保证相同的 Group By Key 分布到同一个 Reduce 中），最后完成最终的聚合操作。

6. 合并小文件

在第 9 章“HDFS 详解”中，我们知道文件数目过多，会给 HDFS 带来很大的压力，并且会影响处理的效率，因此，我们可以通过合并 map 和 reduce 的结果文件来消除这样的影响。需要进行的设定有以下三个：hive.merge.mapfiles = true，设定是否合并 map 输出文件，默认为 true；hive.merge.mapredfiles = false，设定是否合并 reduce 输出文件，默认为 false；hive.merge.size.per.task = 256*1000*1000，设定合并文件的大小，默认值为 256000000。

11.7 小结

本章主要对建立在 Hadoop 之上的数据仓库架构 Hive 进行了全面的讲解。

首先，介绍了 Hive 的安装和配置。由于 Hadoop 的最新版本都集成了 Hive，所以安装很简单，只需要简单地修改配置文件即可。

其次，着重介绍了 Hive 的类 SQL 语言 HQL。我们可以看到 Hive QL 有别于传统的 SQL 实现，但也有很多相似之处，HQL 既继承了传统 SQL 的优势，又结合了 Hadoop 文件系统的特性。

最后，对 Hive 的几个重要接口进行了介绍，这有助于大家更快地掌握和使用 Hive。对管理员来说，本章还给出了 Hive 的优化策略，可以为 Hive 的使用助一臂之力。



第 12 章

HBase 详解

本章内容

- ☐ HBase 简介
- ☐ HBase 的基本操作
- ☐ HBase 体系结构
- ☐ HBase 数据模型
- ☐ HBase 与 RDBMS
- ☐ HBase 与 HDFS
- ☐ HBase 客户端
- ☐ Java API
- ☐ HBase 编程实例之 MapReduce
- ☐ 模式设计
- ☐ 小结

资源分享网
PDG

12.1 HBase 简介

HBase 是 Apache Hadoop 的数据库，能够对大型数据提供随机、实时的读写访问。它目前已经是 Apache 众多开源项目中的一个顶级项目。HBase 的运行依赖于其他文件系统，它模仿并提供了基于 Google 文件系统（GFS, Google File System）中大表（Bigtable）数据库的所有功能。

HBase 的目标是存储并处理大型的数据，更具体来说是仅需使用普通的硬件配置，就能够处理由成千上万的行和列所组成的大型数据。

HBase 是一个开源的、分布式的、多版本的、面向列的存储模型。它可以直接使用本地文件系统，也可以使用 Hadoop 的 HDFS 文件存储系统。不过，为了提高数据的可靠性和系统的健壮性，并且发挥 HBase 处理大型数据的能力，还是使用 HDFS 作为文件存储系统更为稳妥。

另外，HBase 存储的是松散型数据。具体来说，HBase 存储的数据介于映射（key/value）和关系型数据之间。如图 12-1 所示，HBase 存储的数据可以理解成一种 key 和 value 的映射关系，但又不是简简单单的映射关系。除此之外它还具有许多其他的特性，我们将在本章后面详细讲述。HBase 存储的数据从逻辑上来看就像一张很大的表，并且它的数据列可以根据需要动态增加。除此之外，每个 cell（由行和列所确定的位置）中的数据又可以具有多个版本（通过时间戳来区别）。从图中可以看出，HBase 还具有这样的特点：它向下提供了存储，向上提供了运算。另外，在 HBase 之上还可以使用 Hadoop 的 MapReduce 计算模型来并行处理大规模数据，这也是它具有强大性能的核心所在。它将数据存储与并行计算完美地结合在一起。

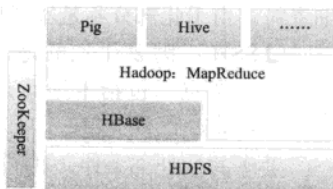


图 12-1 HBase 关系图

12.2 HBase 的基本操作

在介绍完 HBase 的基本特性之后，本节将介绍如何安装 HBase。由于它有单机、伪分布、全分布三种运行模式，因此我们将会分别进行讲解。在安装成功之后，再介绍如何对 HBase 进行详细的设置，以提高系统的可靠性和执行速度。

12.2.1 HBase 的安装

HBase 有三种运行模式，其中单机模式的配置非常简单，几乎不用对安装文件做任何修改就可以使用。从图 12-1 中不难看出，如果要运行分布式模式，Hadoop 是必不可少的。另外在对 HBase 的某些文件进行配置之前，还需要具备以下先决条件：

- Java：需要是 Java 1.6.x 以上的版本，推荐从 SUN 官网下载：<http://www.java.com/download/>。在 Ubuntu 下可以使用下面的命令来安装：

```
sudo apt-get install sun-java6-jdk
```

其他章节已经详细讲过，这里不再赘述。

□ Hadoop：由于 HBase 架构基于其他文件存储系统之上，因此在分布式模式下安装 Hadoop 是必须的。但是，如果运行在单机模式下，此条件则可以省略。

注意 在安装 Hadoop 的时候，要注意 HBase 的版本。也就是说，需要注意 Hadoop 和 HBase 之间的版本关系，如果不匹配，很可能会影响 HBase 系统的稳定性。在 HBase 的 lib 目录下可以看到对应的 Hadoop 的 jar 文件。默认情况下，HBase 的 lib 文件夹下对应的 Hadoop 版本相对稳定。如果用户想要使用其他的 Hadoop 版本，那么需要将 Hadoop 系统安装目录 `hadoop-*. *-core.jar` 文件和 `hadoop-*. *-test.jar` 文件拷贝到 HBase 的 lib 文件夹下，以替换其他版本的 Hadoop 文件。

另外，如果读者想要对 HBase 的数据存储有更好的了解，建议查看关于 HDFS 的更多详细资料。由于此部分不是本章所关注的内容，故此不再赘述。

□ SSH：需要注意的是，SSH 是必须安装的，并且要保证用户可以 SSH 到系统的其他节点（包括本地节点）。因为，我们需要使用的 Hadoop 来管理远程的 Hadoop 和 HBase 守护进程。

关于其他外部条件，我们可以在使用的过程中再具体配置，详细内容见 12.2.2 节。下面我们将具体介绍 HBase 在三种模式下的安装过程。

注意 由于本书中的 Hadoop 使用的是 0.20.2 版本，因此 HBase 使用的是与之对应的稳定版本 0.20.6。其他版本 HBase 的安装方法和过程与本版本大同小异，具体安装请参考 HBase 的官方网址：<http://hbase.apache.org/notsoquick.html>。

1. 单机模式安装

HBase 安装文件默认情况下是支持单机模式的，也就是说将 HBase 安装文件解压后就可以直接运行。在单机模式下 HBase 并不使用 HDFS。用户可以通过下面的命令将其解压：

```
tar xzf hbase-0.90.0.tar.gz
cd hbase-0.90.0
```

在运行之前，建议用户修改 `${HBase-Dir}/conf/hbase-site.xml` 文件。此文件是 HBase 的配置文件，通过它可以更改 HBase 的基本配置。另外还有一个文件为 `hbase-default.xml`，它是 HBase 的默认配置文件。我们可以通过这两个文件中的任意一个来修改 HBase 的配置参数，并且它们二者的配置方法也完全相同。但是同样一个参数如果在 `hbase-site.xml` 中配置了，那么它就会覆盖掉 `hbase-default.xml` 中的同一个配置。也就是说，同样一个配置参数，`hbase-site.xml` 中的配置将发挥作用。建议用户修改 `hbase-site.xml` 中的配置，而 `hbase-default.xml` 中的配置默认保持不变，这样当 `hbase-site.xml` 中的配置错误时，其默认配置可以保证用户能够快速地对 HBase 配置进行恢复。例如，需要修改的内容如下所示：


```

<configuration>
  <property>
    <name>hbase.rootdir</name>
    <value>file:///tmp/hbase-${user.name}/hbase</value>
  </property>
</configuration>

```

从上面可以看到，默认情况下 HBase 的数据是存储在根目录的 tmp 文件夹下的。熟悉 Linux 的用户知道，此文件夹为临时文件夹。也就是说，当系统重启的时候，此文件夹中的内容将被清空。这样用户保存在 HBase 中的数据也会丢失，这当然是用户不想看到的事情。因此，用户需要将其修改为自己希望 HBase 数据所存储的位置。

2. 伪分布模式安装

伪分布模式是一个运行在单个节点（单台机器）上的分布式模式，此种模式下 HBase 所有的守护进程将运行在同一个节点之上。由于分布式模式的运行需要依赖于分布式文件系统，因此此时必须确保 HDFS 已经成功运行。用户可以在 HDFS 系统上执行 Put 和 Get 操作来验证 HDFS 是否安装成功。关于 HDFS 集群的安装，请读者参看其他章节的介绍。

一切准备就绪后，我们开始配置 HBase 的参数（即配置 hbase-site.xml 文档）。通过设定 hbase.rootdir 参数来指定 HBase 的数据存放位置，进而让 HBase 运行在 Hadoop 之上，如图 12-1 所示。具体配置如下所示：

```

<configuration>
  ...
  <property>
    <name>hbase.rootdir</name>
    <value>hdfs://localhost:9000/hbase</value>
    <description>此参数指定了 region 服务器的位置，即数据存放位置。
  </description>
  </property>
  <property>
    <name>dfs.replication</name>
    <value>1</value>
    <description>此参数指定了 Hlog 和 Hfile 的副本个数，此参数的值不能大于 HDFS 的节点数。伪分布模式下 DataNode 只有一台，因此此参数应设置为 1。
  </description>
  </property>
  ...
</configuration>

```

注意 hbase.rootdir 指定的目录需要 Hadoop 自己创建，否则可能出现警告提示。由于目录为空，因此 HBase 在检查目录时可能会报出所需要的文件不存在这样的错误。

3. 全分布模式安装

对于全分布模式 HBase 的安装，我们需要通过 hbase-site.xml 文档来配置本机的 HBase 特性，通过 hbase-env.sh 来配置全局 HBase 集群系统的特性，也就是说每一台机器都可以

通过 `hbase-env.sh` 来了解全局 HBase 的某些特性。另外，各个 HBase 实例之间需要通过 ZooKeeper 来进行通信，因此我们还需要维护一个（一组）ZooKeeper 系统。

下面将分别介绍如何配置这三个重要方面：

（1）`conf/hbase-site.xml` 的配置

配置 `hbase.rootdir` 和 `hbase.cluster.distributed` 两个参数对于 HBase 来说是必需的。通过 `hbase.rootdir` 来指定本台机器 HBase 的存储目录，通过 `hbase.cluster.distributed` 来说明其运行模式：`true` 为全分布模式；`false` 为单机模式或伪分布模式。代码如下所示：

```
<configuration>
...
<property>
  <name>hbase.rootdir</name>
  <value>hdfs://namenode.example.org:9000/hbase</value>
  <description>区域服务器使用存储 HBase 数据库数据的目录</description>
</property>
<property>
  <name>hbase.cluster.distributed</name>
  <value>true</value>
  <description>指定 HBase 运行的模式：
false：单机模式或伪分布模式
true：全分布模式
  </description>
</property>
...
</configuration>
```

（2）`conf/regionservers` 的配置

`regionservers` 文件列出了所有运行 HBase 的机器（即 `HRegionServer`）。此文件的配置和 Hadoop 的 `slaves` 文件十分类似，每一行指定一台机器。当 HBase 启动的时候，会将此文件中列出的所有机器启动；同样，当 HBase 关闭的时候，也会同时关闭它们。

我们以四台机器为例，此四台机器的 IP 地址分别为：192.168.1.23、192.168.1.24、192.168.1.25 和 192.168.1.26。首先修改每台机器的“`/etc/hosts`”文件，代码如下所示：

```
127.0.0.1          localhost
192.168.1.23       hbase-0
192.168.1.24       hbase-1
192.168.1.25       hbase-2
192.168.1.26       hbase-3
```

需要注意的是，每台机器的 `hosts` 文件的配置都是完全相同的。

在这四台机器中，HBase Master 及 HDFS NameNode 运行在 `Hbase-0` 上，`RegionServers` 运行在 `Hbase-1`、`Hbase-2` 和 `Hbase-3` 上。完成上述操作之后，将每台机器上 HBase 安装目录下的“`conf/regionservers`”文件内容设置为：

```
hbase-1
hbase-2
hbase-3
```

这就意味着，HBase RegionServer 运行在 Hbase-1、Hbase-2、Hbase-3 三台机器上。

(3) ZooKeeper 的配置

全分布模式的 HBase 集群需要 ZooKeeper 实例运行，并且需要所有的 HBase 节点能够与 ZooKeeper 实例通信。默认情况下 HBase 自身维护着一组默认的 ZooKeeper 实例。不过，用户可以配置独立的 ZooKeeper 实例，这样能够使 HBase 系统更加健壮。

要使用独立的 ZooKeeper 实例，需要修改“conf/hbase-env.sh”配置文档，修改 HBASE_MANAGES_ZK 变量的值。HBASE_MANAGES_ZK 的默认值为 true，若将其修改为 false，这就意味着不使用默认的 ZooKeeper 实例。

关于 ZooKeeper 的安装与配置详见 ZooKeeper 相关章节。

在本书的实例中，我们将 ZooKeeper 安装到 Hbase-1、Hbase-2 和 Hbase-3 三台机器上。下面，需要修改“conf/hbase-site.xml”配置文档，代码如下所示：

```
<configuration>
...
<property>
  <name>hbase.zookeeper.quorum</name>
  <value>Hbase-1,Hbase-2,Hbase-3</value>
  <description>ZooKeeper 集群服务器的位置
</description>
</property>
...
</configuration>
```

当使用默认的 ZooKeeper 实例时，HBase 将自动地启动或停止 ZooKeeper；当使用独立的 ZooKeeper 实例时，需要用户手动启动或停止 ZooKeeper 实例。

需要注意的是，启动时需要先启动 ZooKeeper 再启动 HBase；关闭时，先关闭 HBase，然后关闭 ZooKeeper。另外，Hadoop 集群的 NameNode 即为 HBase 集群的 HMaster。

12.2.2 运行 HBase

前面说了，HBase 有三种运行模式，不同模式下启动或停止 HBase 服务的步骤稍稍有些不同，另外还有一些注意事项。下面，我们将分情况具体来讲解如何在三种模式下启动/停止 HBase 服务。

1. 单机模式

单机模式下直接运行下面命令即可：

```
start-hbase.sh
```

启动成功后用户可以看到如图 12-2 所示的界面：

```

root@ubuntu:~# start-hbase.sh
localhost: starting zookeeper, logging to /root/hadoop-0.20.2/hbase-0.20.6/bin/
./logs/hbase-root-zookeeper-ubuntu.out
starting master, logging to /root/hadoop-0.20.2/hbase-0.20.6/logs/hbase-root-mas
ter-ubuntu.out
localhost: starting regionserver, logging to /root/hadoop-0.20.2/hbase-0.20.6/bi
n/./logs/hbase-root-regionserver-ubuntu.out
root@ubuntu:~#

```

图 12-2 启动 HBase

从图中可以看到，HBase 首先启动 ZooKeeper，然后启动 HBase Master，最后启动 HBase RegionServer。

要停止 HBase 服务，直接在终端中输入下面命令即可：

```
stop-hbase.sh
```

在停止过程中用户会看到如图 12-3 所示的界面：

```

root@ubuntu:~# stop-hbase.sh
stopping master.....
localhost: stopping zookeeper.
root@ubuntu:~#

```

图 12-3 停止 HBase

从图中可以看到，在停止 HBase 的过程中，首先停止 HBase Master，然后再停止 ZooKeeper 服务。

2. 伪分布模式

由于伪分布模式的运行基于 HDFS，因此在期待运行 HBase 之前首先需要启动 HDFS。启动 HDFS 可以使用如下命令：

```
start-dfs.sh
```

详细信息参见 HDFS 章节内容。

这之后的其他步骤与单机模式相同。

3. 全分布模式

全分布模式与伪分布模式相同，在运行 HBase 之前需要保证 HDFS 已经成功启动。此时，只需要在 NameNode（即 HBase Master）上运行 start-hbase.sh 即可。

12.2.3 HBase Shell

HBase 为用户提供了一个非常方便的使用方式，我们称之为 HBase Shell。

HBase Shell 提供了大多数的 HBase 命令，通过 HBase Shell 用户可以方便地创建、删除及修改表，还可以向表中添加数据、列出表中的相关信息等。

在启动 HBase 之后，用户可以通过下面的命令进入 HBase Shell 之中，代码如下所示：

```
hbase shell
```

进入成功之后，用户将看到如图 12-4 所示的界面：

```
root@ubuntu:~# hbase shell
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Version: 0.20.6, r965666, Mon Jul 19 16:54:48 PDT 2010
hbase(main):001:0>
```

图 12-4 HBase Shell

进入 HBase Shell 之后，输入 help，可以获取 HBase Shell 所支持的命令，如表 12-1 所示：

表 12-1 HBase Shell 命令

HBase Shell 命令	描 述
alter	修改列族模式
count	统计表中行的数量
create	创建表
describe	显示表相关的详细信息
delete	删除指定对象的值（可以为表、行、列对应的值，另外也可以指定时间戳的值）
deleteall	删除指定行的所有元素值
disable	使表无效
drop	删除表
enable	使表有效
exists	测试表是否存在
exit	退出 HBase Shell
get	获取行或单元（cell）的值
incr	增加指定表、行或列的值
list	列出 HBase 中存在的所有表
put	向指定的表单元添加值
tools	列出 HBase 所支持的工具
scan	通过对表的扫描来获取对应的值
status	返回 HBase 集群的状态信息
shutdown	关闭 HBase 集群（与 exit 不同）
truncate	重新创建指定表
version	返回 HBase 版本信息

需要注意 shutdown 操作与 exit 操作之间的不同，shutdown 表示关闭 HBase 服务，必须重新启动 HBase 才可以恢复。exit 只是退出 HBase shell，退出之后完全可以重新进入。

下面我们将详细介绍常用的 HBase 命令及其使用方法。

(1) create

通过表名及用逗号分隔开的列族信息来创建表。操作如下所示：

- 1) hbase> create 't1', {NAME => 'f1', VERSIONS => 5}
- 2) hbase> create 't1', {NAME => 'f1'}, {NAME => 'f2'}, {NAME => 'f3'}

```
hbase> # 上面的命令可以简写为下面所示的格式:
hbase> create 't1', 'f1', 'f2', 'f3'
3) hbase> create 't1', {NAME=>'f1', VERSIONS=>1, TTL=>2592000, BLOCKCACHE => true}
```

以“NAME => 'f1'”举例说明，其中，列族参数的格式是箭头左侧为参数变量，右侧为参数对应的值，并用“=>”分开。

(2) list

通过此命令列出所有 HBase 中包含的表名称，操作如下所示：

```
hbase(main):011:0> list
hbase_tb
test
2 row(s) in 0.0160 seconds
hbase> list
```

(3) put

向指定的 HBase 表单元添加值，例如，向表 t1 的行 r1、列 c1:1 添加值 v1，并指定时间戳为 ts 的操作如下所示：

```
hbase> put 't1', 'r1', 'c1:1', 'value', ts1
```

(4) scan

获取指定表的相关信息，与 create 命令类似，可以通过逗号分隔的命令来指定扫描参数。例如，获取表 test 的所有值，操作如下所示：

```
hbase(main):001:0> scan 'test'
ROW          COLUMN+CELL
r1           column=c1:1, timestamp=1295692753859, value=value1-1/1
r1           column=c1:2, timestamp=1295692662360, value=value1-1/2
r1           column=c1:3, timestamp=1297476019872, value=value1-1/3
r1           column=c2:1, timestamp=1297475967537, value=value1-2/1
.....
```

获取表 test 的 c1 列的所有值，操作如下所示：

```
hbase(main):002:0> scan 'test', {COLUMNS=>'c1'}
ROW          COLUMN+CELL
r1           column=c1:1, timestamp=1295692753859, value=value1-1/1
r1           column=c1:2, timestamp=1295692662360, value=value1-1/2
r1           column=c1:3, timestamp=1297476019872, value=value1-1/3
r2           column=c1:1, timestamp=1297476064414, value=value2-1/1
2 row(s) in 0.0100 seconds
```

获取表 test 的 c1 列的前一行的所有值，操作如下所示：

```
hbase(main):012:0> scan 'test', {COLUMNS=>'c1', LIMIT=>1}
ROW          COLUMN+CELL
r1           column=c1:1, timestamp=1295692753859, value=value1-1/1
r1           column=c1:2, timestamp=1295692662360, value=value1-1/2
r1           column=c1:3, timestamp=1297476019872, value=value1-1/3
```

```
1 row(s) in 0.0120 seconds
```

(5) get

获取行或单元的值。此命令可以指定表名、行值，以及可选的列值和时间戳。

获取表 test 行 r1 的值，操作如下所示：

```
hbase(main):002:0> get 'test','r1'
COLUMN                                CELL
c1:1                                timestamp=1295692753859, value=value1-1/1
c1:2                                timestamp=1295692662360, value=value1-1/2
c1:3                                timestamp=1297476019872, value=value1-1/3
c2:1                                timestamp=1297475967537, value=value1-2/1
c2:2                                timestamp=1297476039968, value=value1-2/2
5 row(s) in 0.0450 seconds
```

获取表 test 行 r1 列 c1:1 的值，操作如下所示：

```
hbase(main):005:0> get 'test','r1',{COLUMN=>'c1:1'}
COLUMN                                CELL
c1:1                                timestamp=1295692753859, value=value1-1/1
1 row(s) in 0.0050 seconds
```

需要注意的是，COLUMN 和 COLUMNS 是不同的，scan 操作中的 COLUMNS 指定的是表的列族，get 操作中的 COLUMN 指定的是特定的列，COLUMN 的值实质上为“列族 + 列修饰符”。

另外，在 shell 中，常量不需要用引号引起来，但二进制的值需要用双引号引起来，而其他值则用单引号引起来。HBase Shell 的常量可以通过在 shell 中输入“Object.constants”命令来查看。

下面是一个使用 HBase Shell 操作的具体例子，见代码清单 12-1。

代码清单 12-1 HBase Shell 操作

```
hbase(main):004:0> create 'test','c1','c2'
0 row(s) in 1.0620 seconds
hbase(main):005:0> list
test
1 row(s) in 0.0090 seconds
hbase(main):006:0> put 'test','r1','c1:1','value1-1/1'
0 row(s) in 0.0050 seconds
hbase(main):007:0> put 'test','r1','c1:2','value1-1/2'
0 row(s) in 0.0060 seconds
hbase(main):008:0> put 'test','r1','c1:3','value1-1/3'
0 row(s) in 0.0110 seconds
hbase(main):009:0> put 'test','r1','c2:1','value1-2/1'
0 row(s) in 0.0040 seconds
hbase(main):010:0> put 'test','r1','c2:2','value1-2/2'
0 row(s) in 0.0030 seconds
hbase(main):011:0> put 'test','r2','c1:1','value2-1/1'
```



```

0 row(s) in 0.0030 seconds
hbase(main):012:0> put 'test','r2','c2:1','value2-2/1'
0 row(s) in 0.0040 seconds
hbase(main):013:0> scan 'test'
ROW          COLUMN+CELL
r1           column=c1:1, timestamp=1297513518032, value=value1-1/1
r1           column=c1:2, timestamp=1297513531036, value=value1-1/2
r1           column=c1:3, timestamp=1297513538344, value=value1-1/3
r1           column=c2:1, timestamp=1297513553055, value=value1-2/1
r1           column=c2:2, timestamp=1297513560121, value=value1-2/2
r2           column=c1:1, timestamp=1297513580833, value=value2-1/1
r2           column=c2:1, timestamp=1297513594789, value=value2-2/1
2 row(s) in 0.0260 seconds
hbase(main):014:0> get 'test','r1',{COLUMN=>'c2:2'}
COLUMN          CELL
c2:2            timestamp=1297513560121, value=value1-2/2
1 row(s) in 0.0140 seconds
hbase(main):015:0> disable 'test'
0 row(s) in 0.0930 seconds
hbase(main):016:0> drop 'test'
0 row(s) in 0.0770 seconds
hbase(main):017:0> exit

```

12.2.4 HBase 配置

关于 HBase 的所有配置参数，用户可以查看 `conf/hbase-default.xml` 文件。每个参数通过 `property` 节点来区分。其配置方式与 Hadoop 的相同：`name` 字段表示参数名，`value` 字段表示对应参数的值，`description` 字段表示参数的描述信息，相当于注释的作用。

其格式如下所示：

```

<configuration>
.....
  <property>
    <name> 配置参数 </name>
    <value> 配置参数对应取值 </value>
    <description> 描述信息 </description>
  </property>
.....
</configuration>

```

因此，如果要对 HBase 进行配置，修改 `conf/hbase-default.xml` 文件或 `conf/hbase-site.xml` 文件中的 `property` 节点即可（被 `<property></property>` 所包含的部分）。

限于篇幅，下面我们就相对比较重要的几个参数做简单的介绍。

(1) `hbase.client.write.buffer`

通过此参数设置写入缓冲区的数据大小，以字节为单位，默认写入缓冲区的数据大小为 2MB。服务器通过此缓冲区可以加快处理的速度，但是此值如果设置得过大势必加重服务器

的负担，因此一定要根据实际情况进行设置。

(2) `hbase.master.meta.thread.rescanfrequency`

HMaster 会扫描 ROOT 和 META 表的时间间隔，以毫秒为单位。默认值为 60 000 毫秒。此值不宜设置得过小，尤其是当存储数据较多的时候，否则频繁地扫描 ROOT 和 META 表将严重影响系统的性能。

(3) `hbase.regionserver.handler.count`

客户端向服务器请求服务时，服务器先将客户端的请求连接放入一个队列中。然后服务器通过轮询的方式对其进行处理。这样每一个请求就会产生一个线程。此值要根据实际情况来设置，建议设置得大一些。在服务器端由于写数据缓存所消耗的内存大小为：`hbase.client.write.buffer * hbase.regionserver.handler.count`。

(4) `hbase.hregion.max.filesize`

通过此参数可以设置 HRegion 中 HStoreFile 文件的最大值，以字节为单位。当表中的列族 (column family) 超过此值时，它将被分割。其默认大小为 256MB。

(5) `hfile.block.cache.size`

HFile/StoreFile 缓存所占 Java 虚拟机堆大小的百分比，默认值为 0.2，即 20%。将其值设置为 0 表示禁用此选项。

(6) `hbase.regionserver.global.memstore.upperLimit`

在 Region 服务器中所有的 memstore 所占用 Java 虚拟机比例的最大值。默认值为 0.4，即 40%。当 memstore 所占用的空间超过此值时，更新操作将被阻塞，并且所有的内容将被强制写出。

(7) `hbase.hregion.memstore.flush.size`

如果 memstore 缓存的内容大小超过此参数所设置的值，那么它将被写到磁盘上。默认值为 64MB。

另外，在配置文档中还有很多关于 ZooKeeper 配置的参数，诸如：`zookeeper.session.timeout`、以 `hbase.zookeeper` 开头的参数，以及以 `hbase.zookeeper.property` 开头的一些参数。这里不再赘述，关于 ZooKeeper 更详细的配置见第 15 章“ZooKeeper 详解”。

12.3 HBase 体系结构

HBase 的服务器体系结构遵从简单的主从服务器架构，它由 HRegion 服务器 (HRegion Server) 群和 HBase Master 服务器 (HBase Master Server) 构成。HBase Master 服务器负责管理所有的 HRegion 服务器，而 HBase 中所有的服务器都是通过 ZooKeeper 来进行协调，并处理 HBase 服务器运行期间可能遇到的错误的。HBase Master Server 本身并不存储 HBase 中的任何数据，HBase 逻辑上的表可能会被划分成多个 HRegion，然后存储到 HRegion Server 群中。HBase Master Server 中存储的是从数据到 HRegion Server 的映射。因此，HBase 体系结构如图 12-5 所示。

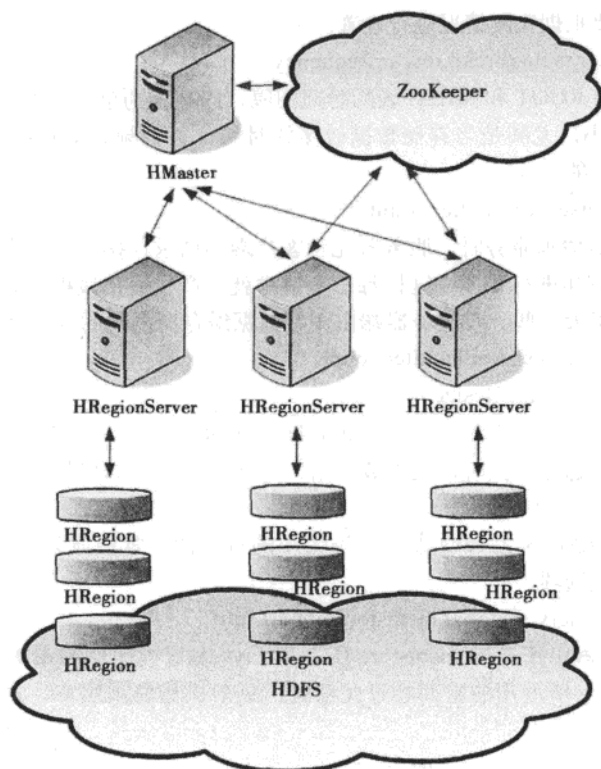


图 12-5 HBase 体系结构

1. HRegion

当表的大小超过设置值的时候，HBase 会自动地将表划分为不同的区域，每个区域包含所有行的一个子集。对用户来说，每个表是一堆数据的集合，靠主键来区分。从物理上来说，一张表被拆分成了多块，每一块就是一个 HRegion。我们用表名 + 开始 / 结束主键，来区分每一个 HRegion，一个 HRegion 会保存一个表里面某段连续的数据，从开始主键到结束主键，一张完整的表格是保存在多个 HRegion 上面的。

2. HRegion 服务器

所有的数据库数据一般是保存在 Hadoop 分布式文件系统上面的，用户通过一系列 HRegion 服务器获取这些数据，一台机器上面一般只运行一个 HRegion 服务器，且每一个区段的 HRegion 也只會被一个 HRegion 服务器维护。

当用户需要更新数据的时候，他会被分配到对应的 HRegion 服务器上提交修改，这些修改先是被写到 Hmemcache 缓存和服务器的 Hlog 文件里面（Hmemcache 是内存中的缓存，保存最近更新的数据；Hlog 是磁盘上面的记录文件，它记录着所有的更新操作），在操作写

入 Hlog 之后, commit() 调用才会将其返回给客户端。

在读取数据的时候, HRegion 服务器会先访问 Hmemcache 缓存, 如果缓存里面没有该数据, 才会回到 Hstores 磁盘上面寻找, 每一个列族都会有一个 Hstore 集合, 每个 Hstore 集合包含很多具体的 HstoreFile 文件, 这些文件都是 B 树结构的, 方便快速读取。

系统会定期调用 HRegion.flushcache() 把缓存里面的内容写到文件中, 一般这会增加一个新的 HstoreFile 文件, 而此时高速缓存就会被清空, 并且会写一个标记到 Hlog 上, 这表示上面的内容已经被写入到文件中保存。

在启动的时候, 每台 HRegion 服务器都会检查自己的 Hlog 文件, 看看最近一次执行 flushcache 之后有没有新的更新写入操作。如果没有更新, 就表示所有数据都已经更新到文件中了; 如果有更新, 服务器就会先把这些更新写入高速缓存, 然后调用 flushcache 写入到文件中。最后服务器会删除旧的 Hlog 文件, 并开始让用户访问数据。

因此, 为了节省时间可以少调用 flushcache, 但是这样会增加内存的占用量, 而且在服务器重启的时候会延长很多时间。如果定期调用 flushcache, 缓存大小将会控制在一个较低的水平上, 而且 Hlog 文件也会很快重构, 但是调用 flushcache 的时候会造成系统负载瞬间增加。

Hlog 会被定期回滚, 回滚的时候是按照时间来备份文件的, 每次回滚的时候, 系统会删除那些已经被写到文件中的更新, 回滚 Hlog 只会占用很少的时间, 建议经常回滚以减少文件尺寸。

每一次调用 flushcache 会生成一个新的 HstoreFile 文件, 从一个 Hstore 里面获取任何一个值都需要访问所有的 HstoreFile 文件, 这会十分耗时, 所以我们要定期把这些分散的文件合并到一个大文件里面, HStore.compact() 可以完成合并工作。不过这样的工作十分占用资源, 所以只有在 HstoreFile 文件的数量超过一个设定值的时候才会触发。

注意

- 1) flushcache 会建立一个新的 HstoreFile 文件, 并把缓存中所有需要更新的数据写到文件里面, 在 flushcache 之后, log 的重建次数会清零。
- 2) compact 会把所有的 HstoreFile 文件合并成一个大文件。
- 3) 和 Bigtable 不同的是, HBase 的每个更新如果被正确提交, 且 commit 没有返回错误, 那它就一定被写到记录文件里面了, 这不会造成数据丢失。

两个 HRegion 可以通过调用 HRegion.closeAndMerge() 合并成一个新的 HRegion, 在当前版本下进行此操作需要两台 HRegion 都停机。

当一个 HRegion 变得太过巨大, 超过了设定的阈值时, HRegion 服务器会调用 HRegion.closeAndSplit() 将此 HRegion 拆分为两个, 并且报告给主服务器让它决定由哪台 HRegion 服务器来存放新的 HRegion。这个拆分过程十分迅速, 因为两台新的 HRegion 最初只是保留原来 HRegionFile 文件的引用, 这个时候旧的 HRegion 会处于停止服务的状态, 当新的

HRegion 拆分完成并且把引用删除了以后, 旧的 HRegion 才会删除。

综上所述, HRegionServer 中数据的存储关系如图 12-6 所示:

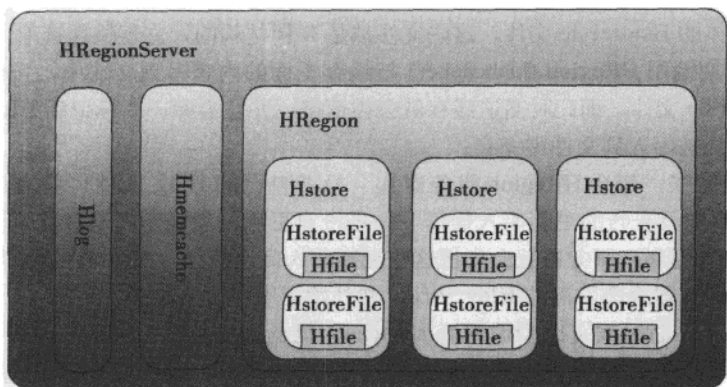


图 12-6 HRegionServer 数据存储关系图

3. HBase Master 服务器

每台 HRegion 服务器都会和 HMaster 服务器通信, HMaster 的主要任务就是要告诉每台 HRegion 服务器它要维护哪些 HRegion。

当一台新的 HRegion 服务器登录到 HMaster 服务器时, HMaster 会告诉它先等待分配数据。而当一台 HRegion 死机时, HMaster 会把它负责的 HRegion 标记为未分配, 然后再把它分配到其他 HRegion 服务器中。

若当前 HBase 已经解决了之前存在的 SPFO (单点故障), 并且 HBase 中可以启动多个 HMaster, 那么它能够通过 ZooKeeper 来保证系统中总有一个 Master 在运行。HMaster 在功能上主要负责 Table 和 Region 的管理工作, 具体包括:

- ❑ 管理用户对 Table 的增、删、改、查操作;
- ❑ 管理 HRegionServer 的负载均衡, 调整 Region 分布;
- ❑ 在 Region Split 后, 负责新 Region 的分配;
- ❑ 在 HRegionServer 停机后, 负责失效 HRegionServer 上的 Region 迁移。

4. ROOT 表和 META 表

在开始这部分内容之前, 我们先来看一下 HBase 中相关的机制是怎样的。之前我们说过 HRegion 是按照表名和主键范围来区分的, 由于主键范围是连续的, 所以一般用开始主键就可以表示相应的 HRegion 了。

不过, 因为我们有合并和分割操作, 如果正好在执行这些操作的过程中出现死机, 那么就可能存在相同的“表名和开始主键”, 这样的话只要开始主键就不够了, 这就要通过 HBase 的元数据信息来区分哪一份才是正确的数据文件, 为了区分这样的情况, 每个

HRegion 都有一个 'regionId' 来标识它的唯一性。

所以一个 HRegion 的表达符最后是：表名 + 开始主键 + 唯一 ID (tablename + startkey + regionId)。我们可以用这个识别符来区分不同的 HRegion，这些数据就是元数据 (META)，而元数据本身也是被保存在 HRegion 里面的，所以我们称这个表为元数据表 (META Table)，里面保存的就是 HRegion 标识符和实际 HRegion 服务器的映射关系。

元数据表也会增长，并且可能被分割为几个 HRegion，为了定位这些 HRegion，我们采用了一个根数据表 (ROOT table)，它保存了所有元数据表的位置，而根数据表是不能被分割的，永远只存在一个 HRegion。

在 HBase 启动的时候，主服务器先去扫描根数据表，因为这个表只会有一个 HRegion，所以这个 HRegion 的名字是被写死的。当然要把根数据表分配到一个 HRegion 服务器中需要一定的时间。

待根数据表被分配好之后，主服务器就会去扫描根数据表，获取元数据表的名字和位置，然后把元数据表分配到不同的 HRegion 服务器中。最后就是扫描元数据表，找到所有 HRegion 区域的信息，然后把它们分配给不同的 HRegion 服务器。

主服务器在内存中保存着当前活跃的 HRegion 服务器的数据，因此如果主服务器死机，整个系统也就无法访问了，这时服务器的信息也就没有必要保存到文件里面了。

元数据表和根数据表的每一行都包含一个列族 (info 列族)：

❑ info:regioninfo：包含了一个串行化的 HRegionInfo 对象。

❑ info:server：保存了一个字符串，是服务器的地址 HServerAddress.toString()。

❑ info:startcode：是一个长整型数字的字符串，它是在 HRegion 服务器启动的时候传给主服务器的，让主服务器确定这个 HRegion 服务器的信息有没有更改。

因此，一个客户端在拿到根数据表地址以后，就没有必要再连接主服务器了。主服务器的负载相对就小了很多，它只会处理超时的 HRegion 服务器，并在启动的时候扫描根数据表和元数据表，以及返回根数据表的 HRegion 服务器地址。

注意 ROOT 表包含 META 表所在的区域列表，META 表包含所有的用户空间区域列表，以及 Region 服务器地址。客户端能够缓存所有已知的 ROOT 表和 META 表，从而提高访问的效率。

12.4 HBase 数据模型

12.4.1 数据模型

HBase 是一个类似 Bigtable 的分布式数据库，它是一个稀疏的长期存储的（存在硬盘上）、多维度的、排序的映射表。这张表的索引是行关键字、列关键字和时间戳。HBase 中的数据都是字符串，没有类型。

用户在表格中存储数据，每一行都有一个可排序的主键和任意多的列。由于是稀疏存储，所以同一张表里面的每一行数据都可以有截然不同的列。

列名字的格式是 "<family>:<qualifier>"，都是由字符串组成的。每一张表有一个列族（family）集合，这个集合是固定不变的，只能通过改变表结构来改变。但是 qualifier 的值相对于每一行来说都是可以改变的。

HBase 把同一个列族里面的数据存储在同一个目录底下，并且 HBase 的写操作是锁行的，每一行都是一个原子元素，都可以加锁。

HBase 所有数据库的更新都有一个时间戳标记，每个更新都是一个新的版本，HBase 会保留一定数量的版本，这个值是可以设定的。客户端可以选择获取距离某个时间点最近的版本单元的值，或者一次获取所有版本单元的值。

12.4.2 概念视图

我们可以将一个表想象成一个大的映射关系，通过行键、行键 + 时间戳或行键 + 列（列族：列修饰符），就可以定位特定数据。由于 HBase 是稀疏存储数据的，所以某些列可以空白的，表 12-2 对应 12.2 节中创建的 test 表的数据概念视图。

表 12-2 HBase 数据的概念视图

Row Key	Time Stamp	Column Family:c1		Column Family:c2	
		列	值	列	值
r1	t7	c1:1	value1-1/1		
	t6	c1:2	value1-1/2		
	t5	c1:3	value1-1/3		
	t4			c2:1	value1-2/1
	t3			c2:2	value1-2/2
r2	t2	c1:1	value2-1/1		
	t1			c2:1	value2-1/1

从上表中可以看出，test 表有两行数据：r1 和 r2，并且有两个列族：c1 和 c2。在第一行 r1 中，列族 c1 有三条数据，列族 c2 有两条数据；在第二行 r2 中，列族 c1 有一条数据，列族 c2 有一条数据。每一条数据对应的时间戳都用数字来表示，编号越大表示数据越旧，反之表示数据越新。

12.4.3 物理视图

虽然从概念视图来看每个表格是由很多行组成的，但是在物理存储上面，它是按照列来保存的，这点在进行数据设计和程序开发的时候必须牢记。

上面的概念视图在物理存储的时候应该表现成下面的样子（如表 12-3 和表 12-4 所示）。

表 12-3 HBase 数据的物理视图 (1)

Row Key	Time Stamp	Column Family:c1	
		列	值
rl	t7	c1:1	value1-1/1
	t6	c1:2	value1-1/2
	t5	c1:3	value1-1/3

表 12-4 HBase 数据的物理视图 (2)

Row Key	Time Stamp	Column Family:c2	
		列	值
rl	t4	c2:1	value1-1/1
	t3	c2:2	value1-1/1

需要注意的是,在概念视图上面有些列是空白的,这样的列实际上并不会被存储,当请求这些空白的单元格时,会返回 null 值。

如果在查询的时候不提供时间戳,那么会返回距离现在最近的那一个版本的数据。因为在存储的时候,数据会按照时间戳来排序。

12.5 HBase 与 RDBMS

HBase 就是这样一个基于列模式的映射数据库,它只能表示很简单的键-数据的映射关系,这大大简化了传统的关系数据库。与关系数据库相比,它有如下特点:

- ❑ 数据类型: HBase 只有简单的字符串类型,所有的类型都是交由用户自己处理的,它只保存字符串。而关系数据库有丰富的类型选择和存储方式。
- ❑ 数据操作: HBase 操作只有很简单的插入、查询、删除、清空等操作,表和表之间是分离的,没有复杂的表和表之间的关系,所以不能也没有必要实现表和表之间的关联等。而传统的关系数据通常有各种各样的函数、连接操作。
- ❑ 存储模式: HBase 是基于列存储的,每个列族都由几个文件保存,不同列族的文件是分离的。传统的关系数据库是基于表格结构和行模式保存的。
- ❑ 数据维护: 确切来说, HBase 的更新操作应该不叫更新,虽然一个主键或列会对应新的版本,但它的旧版本仍然会保留,所以它实际上是插入了新数据,而不是传统关系数据库里面的替换修改。
- ❑ 可伸缩性: HBase 这类分布式数据库就是为了这个目的而开发出来的,所以它能够轻易地增加或减少(在硬件错误的时候)硬件数量,并且对错误的兼容性比较高。而传统的关系数据库通常需要增加中间层才能实现类似的功能。

当前的关系数据库基本都是从 20 世纪 70 年代发展而来的,它们都具有 ACID 特性,并且拥有丰富的 SQL 语言,除此之外它们基本都有以下的特点: 面向磁盘存储、带有索引结

构、多线程访问、基于锁的同步访问机制、基于 log 记录的恢复机制等。

而 Bigtable 和 HBase 这些基于列模式的分布式数据库，更适应海量存储和互联网应用的需求，灵活的分布式架构可以使其利用廉价的硬件设备组建一个大的数据仓库。互联网应用是以字符为基础的，而 Bigtable 和 HBase 就是针对这些应用开发出来的数据库。

由于 HBase 具有时间戳特性，所以它生来就特别适合开发 wiki、archiveorg 之类的服务，并且它原本就是作为一个搜索引擎的一部分开发出来的。

12.6 HBase 与 HDFS

伪分布模式和全分布模式下的 HBase 运行基于 HDFS 文件系统。使用 HDFS 文件系统需要设置“conf/hbase-site.xml”文件，修改 hbase.rootdir 的值，并将其指向 HDFS 文件系统的位置。此外，HBase 也可以使用其他的文件系统，不过此时需要重新设置 hbase.rootdir 参数的值。

12.7 HBase 客户端

HBase 客户端可以选择多种方式与 HBase 集群进行交互，最常用的方式为 Java，除此之外还有 Rest 和 Thrift 接口。

1. Java

HBase 是由 Java 编写的。在后面的章节中，我将详细地向读者介绍 HBase 的 Java API，用户可以通过丰富的 Java API 接口与 HBase 进行互操作，并执行各种相关操作。详细内容请见 12.8 节。

2. Rest 和 Thrift 接口

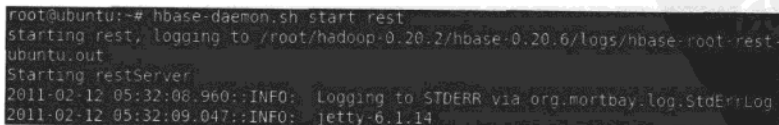
HBase 的 Rest 和 Thrift 接口支持 XML、Protobuf 和二进制数据编码等操作。

(1) Rest

用户可以通过下面命令运行 Rest：

```
hbase-daemon.sh start rest
```

运行成功后将显示如图 12-7 所示的画面：



```
root@ubuntu:~# hbase-daemon.sh start rest
starting rest, logging to /root/hadoop-0.20.2/hbase-0.20.6/logs/hbase-root-rest-ubuntu.out
Starting restServer
2011-02-12 05:32:08.960::INFO: Logging to STDERR via org.mortbay.log.StdErrLog
2011-02-12 05:32:09.047::INFO: Jetty-6.1.14
```

图 12-7 启动 HBase Rest

用户可以通过下面的命令停止 Rest 服务：


```
hbase-daemon.sh stop rest
```

停止过程如图 12-8 所示：

```
root@ubuntu:~# hbase-daemon.sh stop rest
stopping rest.
root@ubuntu:~#
```

图 12-8 停止 HBase Rest

(2) Thrift

用户可以通过下面的命令启动 Thrift 客户端，并与 HBase 进行通信：

```
hbase-daemon.sh start thrift
```

运行成功后将显示如图 12-9 所示的画面：

```
root@ubuntu:~# hbase-daemon.sh start thrift
starting thrift, logging to /root/hadoop-0.20.2/hbase-0.20.6/logs/hbase-root-thrift-ubuntu.out
root@ubuntu:~#
```

图 12-9 启动 HBase Thrift

用户可以通过下面的命令停止 Thrift 服务：

```
hbase-daemon.sh stop thrift
```

停止过程如图 12-10 所示：

```
root@ubuntu:~# hbase-daemon.sh stop thrift
stopping thrift.
root@ubuntu:~#
```

图 12-10 停止 HBase Thrift

12.8 Java API

通过前面的内容读者已经了解到，HBase 作为云环境中的数据库，与传统数据库相比拥有不同的特点。当前 HBase 的 Java API 已经比较完善了，从其所涉及的内容来讲，大体包括：关于 HBase 自身的配置管理部分、关于 Avro 部分、关于 HBase 的客户端部分、关于 MapReduce 部分、关于 Rest 部分、关于 Thrift 部分，以及关于 ZooKeeper 部分。其中关于 HBase 自身的配置管理部分又包括：HBase 配置、日志、I/O、Master、Regionserver、replication，以及安全性。

限于篇幅我们重点介绍与 HBase 数据存储管理相关的内容，其所涉及的主要类包括：HBaseAdmin、HBaseConfiguration、HTable、HTableDescriptor、Put、Get，以及 Scanner。关于 Java API 的详细内容，读者可以查看 HBase 官方网站的相关资料：<http://hbase.apache.org/apidocs/index.html>。

表 12-5 给我们描述了这几个相关类与 HBase 数据模型之间的对应关系。

表 12-5 Java API 与 HBase 数据模型之间的关系

Java 类	HBase 数据模型
HBaseAdmin	数据库 (database)
HBaseConfiguration	
HTable	表 (table)
HTableDescriptor	列族 (column family)
Put	列修饰符 (column qualifier)
Get	
Scanner	

下面我们将详细讲述这些类的功能，以及它们之间的相互关系。

1. HBaseConfiguration

关系: org.apache.hadoop.hbase.HBaseConfiguration

作用: 通过此类可以对 HBase 进行配置。

主要方法如表 12-6 所示:

表 12-6 HBaseConfiguration 类包含的方法

回传值	函 数	描 述
void	addResource (Path file)	通过给定的路径所指的文件来添加资源
void	clear()	清空所有已设置的属性
string	get (String name)	获取属性名对应的值
string	getBoolean (String name, boolean defaultValue)	获取为 boolean 类型的属性值，如果其属性值类型不为 boolean，则返回默认属性值
void	set (String name, String value)	通过设置属性名来设置值
void	setBoolean (String name, boolean value)	设置为 boolean 类型的属性值

用法示例:

```
HBaseConfiguration hconfig = new HBaseConfiguration();
hconfig.set("hbase.zookeeper.property.clientPort", "2181");
```

此方法设置了“hbase.zookeeper.property.clientPort”属性的端口号为 2181。一般情况下，HBaseConfiguration 会使用构造函数进行初始化，然后再使用其他方法（如表 12-5 介绍的方法）添加必要的配置。

2. HBaseAdmin

关系: org.apache.hadoop.hbase.client.HBaseAdmin

作用: 提供了一个接口来管理 HBase 数据库的表信息。它提供的方法包括: 创建表、删除表、列出表项、使表有效或无效，以及添加或删除表列族成员等。

主要方法如表 12-7 所示:

表 12-7 HBaseAdmin 类包含的方法

回 传 值	函 数	描 述
void	addColumn (String tableName, HColumnDescriptor column)	向一个已存在的表添加列
	checkHBaseAvailable (HBaseConfiguration conf)	静态函数, 查看 HBase 是否处于运行状态
	createTable (HTableDescriptor desc)	创建一个新表, 同步操作
	deleteTable (byte[] tableName)	删除一个已存在的表
	enableTable (byte[] tableName)	使表处于有效状态
	disableTable (String tableName)	使表处于无效状态
HTableDescriptor[]	listTables()	列出所有用户空间表项
void	modifyTable (byte[] tableName, HTableDescriptor htd)	修改表的模式, 是异步的操作, 可能需要花费一定时间
boolean	tableExists (String tableName)	检查表是否存在, 存在则返回 true

用法示例:

```
HBaseAdmin admin=new HBaseAdmin(config);
admin.disableTable("tablename");
```

上述例子通过一个 HBaseAdmin 实例 admin 调用 disableTable 方法来使表处于无效状态。

3. HTableDescriptor

关系: org.apache.hadoop.hbase.HTableDescriptor

作用: HTableDescriptor 类包含了表的名字及其对应表的列族。

主要方法如表 12-8 所示:

表 12-8 HTableDescriptor 类包含的方法

回 传 值	函 数	描 述
void	addFamily (HColumnDescriptor)	添加一个列族
HColumnDescriptor	removeFamily (byte[] column)	移除一个列族
byte[]	getName()	获取表的名字
byte[]	getValue (byte[] key)	获取属性的值
void	setValue (String key, String value)	设置属性的值

用法示例:

```
HTableDescriptor htd=new HTableDescriptor(tablename);
htd.addFamily(new HcolumnDescriptor("Family"));
```

在上述例子中, 通过一个 HColumnDescriptor 实例, 为 HTableDescriptor 添加了一个列族: Family。

4. HColumnDescriptor

关系：org.apache.hadoop.hbase.HColumnDescriptor

作用：HColumnDescriptor 维护着关于列族的信息，例如版本号、压缩设置等。它通常在创建表或为表添加列族的时候使用。列族被创建后不能直接修改，只能通过删除然后重建的方式来“修改”。并且，当列族被删除的时候，对应列族中所保存的数据也将被同时删除。

主要方法如表 12-9 所示：

表 12-9 HColumnDescriptor 类包含的方法

回传值	函 数	描 述
byte[]	getName()	获取列族的名字
byte[]	getValue (byte[] key)	获取对应属性的值
void	setValue (String key, String value)	设置对应属性的值

用法示例：

```
HtableDescriptor htd=new HtableDescriptor(tablename);
HcolumnDescriptor col=new HcolumnDescriptor("content:");
htd.addFamily(col);
```

此示例添加了一个名为 content 的列族。

5. HTable

关系：org.apache.hadoop.hbase.client.HTable

作用：此表可以用来与 HBase 表进行通信。此方法对于更新操作来说是非线程安全的，也就是说，如果有多个线程尝试与单个 HTable 实例进行通信，那么写缓冲器可能会崩溃。

主要方法如表 12-10 所示：

表 12-10 HTable 类包含的方法

回传值	函 数	描 述
void	checkAndPut (byte[] row, byte[] family, byte[] qualifier, byte[] value, Put put)	自动地检查 row/family/qualifier 是否与给定的值匹配
void	close()	释放所有的资源或挂起内部缓冲区中的更新
Boolean	exists (Get get)	检查 Get 实例所指定的值是否存在于 HTable 的列中
Result	get (Get get)	取出指定行的某些单元格所对应的值
byte[][]	getEndKeys()	获取当前已打开的表每个区域的结束键值
ResultScanner	getScanner (byte[] family)	获取当前表的给定列族的 scanner 实例
HTableDescriptor	getTableDescriptor()	获取当前表的 HTableDescriptor 实例
byte[]	getTableName()	获取表名
static boolean	isTableEnabled (HBaseConfiguration conf, String tableName)	检查表是否有效
void	put (Put put)	向表中添加值

用法示例:

```
Htable table=new Htable(conf, Bytes.toBytes(tablename));
ResultScanner scanner=table.getScanner(family);
```

6. Put

关系: org.apache.hadoop.hbase.client.Put

作用: 用来对单个行执行添加操作。

主要方法如表 12-11 所示:

表 12-11 Put 类包含的方法

回传值	函 数	描 述
Put	add (byte[] family, byte[] qualifier, byte[] value)	将指定的列和对应的值添加到 Put 实例中
Put	add (byte[] family, byte[] qualifier, long ts, byte[] value)	将指定的列和对应的值及时间戳 (版本号) 添加到 Put 实例中
byte[]	getRow()	获取 Put 实例的行
RowLock	getRowLock()	获取 Put 实例的行锁
long	getTimestamp()	获取 Put 实例的时间戳
boolean	isEmpty()	检查 familyMap 是否为空
Put	setTimestamp (long timestamp)	设置 Put 实例的时间戳

用法示例:

```
Htable table = new Htable(conf, Bytes.toBytes(tablename));
Put p = new Put(brow); // 为指定行 (brow) 创建一个 Put 操作
p.add(family, qualifier, value);
table.put(p);
```

7. Get

关系: org.apache.hadoop.hbase.client.Get

作用: 用来获取单个行的相关信息。

主要方法如表 12-12 所示:

表 12-12 Get 类包含的方法

回传值	函 数	描 述
Get	addColumn (byte[] family, byte[] qualifier)	获取指定列族和列修饰符对应的列
Get	addFamily (byte[] family)	通过指定的列族获取其对应的所有列
Get	setTimeRange (long minStamp, long maxStamp)	获取指定区间的列的版本号
Get	setFilter (Filter filter)	当执行 Get 操作时设置服务器端的过滤器

用法示例:

```
Htable table=new Htable(conf,Bytes.toBytes(tablename));
Get g=new Get(Bytes.toBytes(row));
```

8. Result

关系: org.apache.hadoop.hbase.client.Result

作用: 存储 Get 或 Scan 操作后获取表的单行值。使用此类提供的方法能够直接方便地获取值或获取各种 Map 结构 (key-value 对)。

主要方法如表 12-13 所示:

表 12-13 Result 类包含的方法

回 传 值	函 数	描 述
boolean	containsColumn (byte[] family, byte[] qualifier)	检查指定的列是否存在
NavigableMap<byte[],byte[]>	getFamilyMap (byte[] family)	返回值格式为: Map<qualifier,value>, 获取对应列族所包含的修饰符与值的键值对
byte[]	getValue (byte[] family, byte[] qualifier)	获取对应列的最新值

用法示例:

```
Htable table = new HTable(conf, Bytes.toBytes(tablename));
Get g = new Get(Bytes.toBytes(row));
Result rowResult = table.get(g);
Bytes[] ret = rowResult.getValue( (family + ":" + column) );
```

9. ResultScanner

关系: Interface

作用: 客户端获取值的接口。

主要方法如表 12-14 所示:

表 12-14 ResultScanner 类包含的方法

回 传 值	函 数	描 述
Void	close()	关闭 scanner 并释放分配给它的所有资源
Result	next()	获取下一行的值

用法示例:

```
ResultScanner scanner = table.getScanner (Bytes.toBytes(family));
for (Result rowResult : scanner) {
    Bytes[] str = rowResult.getValue ( family , column );
}
```

下面我们使用上面所涉及的 HBase Java API 给出一个简单的例子。希望读者能够通过学习对其使用方法及特点有一个更深入的认识。见代码清单 12-2。

代码清单 12-2 HBase Java API 的简单用例

```

1. package hbase;
2.
3. import org.apache.hadoop.conf.Configuration;
4. import org.apache.hadoop.hbase.HBaseConfiguration;
5. import org.apache.hadoop.hbase.HColumnDescriptor;
6. import org.apache.hadoop.hbase.HTableDescriptor;
7. import org.apache.hadoop.hbase.KeyValue;
8. import org.apache.hadoop.hbase.client.HBaseAdmin;
9. import org.apache.hadoop.hbase.client.HTable;
10. import org.apache.hadoop.hbase.client.Result;
11. import org.apache.hadoop.hbase.client.ResultScanner;
12. import org.apache.hadoop.hbase.client.Scan;
13. import org.apache.hadoop.hbase.io.BatchUpdate;
14.
15. @SuppressWarnings("deprecation")
16. public class HBaseTestCase {
17. // 声明静态配置 HBaseConfiguration
18.     static HBaseConfiguration cfg = null;
19.     static {
20.         Configuration HBASE_CONFIG = new Configuration();
21.         cfg = new HBaseConfiguration(HBASE_CONFIG);
22.     }
23.
24. // 创建一张表, 通过 HBaseAdmin HTableDescriptor(指定表名, 添加列族) 来创建
25.     public static void creatTable(String tablename, String columnFamily) throws
        Exception {
26.         HBaseAdmin admin = new HBaseAdmin(cfg);
27.         if (admin.tableExists(tablename)) {
28.             System.out.println("table Exists!");
29.             System.exit(0);
30.         }
31.         else{
32.             HTableDescriptor tableDesc = new HTableDescriptor(tablename);
33.             tableDesc.addFamily(new HColumnDescriptor(columnFamily+""));
34.             admin.createTable(tableDesc);
35.             System.out.println("create table success!");
36.         }
37.     }
38.
39. // 添加一条数据, 通过 HTable BatchUpdate 为已经存在的表来添加数据
40.     public static void addData (String tablename, String row, String
        columnFamily, String column, String data) throws Exception {
41.         HTable table = new HTable(cfg, tablename);
42.         BatchUpdate update = new BatchUpdate(row);
43.         update.put(columnFamily+"."+column, data.getBytes());
44.         table.commit(update);
45.         System.out.println("add data success!");
46.     }

```

```

47.
48. // 显示所有数据, 通过 HTable Scan 来获取已有表的信息
49. public static void getAllData (String tablename) throws Exception{
50.     HTable table = new HTable(cfg, tablename);
51.     Scan s = new Scan();
52.     ResultScanner rs = table.getScanner(s);
53.     for(Result r:rs){
54.         for(KeyValue kv:r.raw()){
55.             System.out.print(new String(kv.getColumn()));
56.             System.out.println(new String(kv.getValue()));
57.         }
58.     }
59. }
60.
61. public static void main (String [] args) {
62.     try {
63.         String tablename="hbase_tb";
64.         HBaseTestCase.creatTable(tablename,"cl");
65.         HBaseTestCase.addData(tablename,"row1","cl","1", "this is row 1
           column cl:1");
66.         HBaseTestCase.getAllData(tablename);
67.     }
68.     catch (Exception e) {
69.         e.printStackTrace();
70.     }
71. }
72. }

```

在代码清单 12-2 中, 首先, 通过第 17 ~ 22 行创建一个默认的 HBase 配置实例: `cfg`; 然后, 通过 HBaseAdmin 接口来管理现有数据库, 见第 26 行; 第 32 ~ 34 行通过 HTableDescriptor (指定表相关信息) 和 HColumnDescriptor (指定表内列族相关信息) 来创建一个 HBase 数据库, 并设置其拥有的列族成员; 第 40 ~ 46 行通过 HTable 实例及其方法为刚刚创建的数据库添加值; 第 49 ~ 59 行通过 HTable 及其方法获取现有数据库表中的值。

如果读者想要对 HBase 的原理、运行机制, 以及编程有更深入的了解, 建议阅读 HBase 的源码。通过对 HBase 源码的深入探究, 相信读者一定能够对 HBase 有更深层次的理解。

12.9 HBase 编程实例之 MapReduce

从图 12-1 “HBase 关系图”中可以看出, 在伪分布模式和全分布模式下 HBase 是架构在 HDFS 之上的。因此完全可以将 MapReduce 编程框架和 HBase 结合起来使用。也就是说, 将 HBase 作为底层“存储结构”, MapReduce 调用 HBase 进行特殊的处理, 这样能够充分结合 HBase 分布式大型数据库和 MapReduce 并行计算的优点。

下面给出了一个 WordCount 将 MapReduce 与 HBase 结合起来使用的例子, 见代码清单 12-3。在这个例子中, 输入文件为 `hbaseinput/file01` (它包含内容: `hello world bye world`)

和文件 hbaseinput/file02（它包含内容：hello hadoop bye hadoop）。

程序首先从文件中收集数据，在 shuffle 完成之后进行统计并计算，最后将计算结果存储到 HBase 中。

代码清单 12-3 HBase 与 WordCount 的结合使用

```

1. package hbase;
2.
3. import java.io.IOException;
4.
5. import org.apache.hadoop.conf.Configuration;
6. import org.apache.hadoop.fs.Path;
7. import org.apache.hadoop.hbase.HBaseConfiguration;
8. import org.apache.hadoop.hbase.HColumnDescriptor;
9. import org.apache.hadoop.hbase.HTableDescriptor;
10. import org.apache.hadoop.hbase.client.HBaseAdmin;
11. import org.apache.hadoop.hbase.client.Put;
12. import org.apache.hadoop.hbase.mapreduce.TableOutputFormat;
13. import org.apache.hadoop.hbase.mapreduce.TableReducer;
14. import org.apache.hadoop.hbase.util.Bytes;
15. import org.apache.hadoop.io.IntWritable;
16. import org.apache.hadoop.io.LongWritable;
17. import org.apache.hadoop.io.NullWritable;
18. import org.apache.hadoop.io.Text;
19. import org.apache.hadoop.mapreduce.Job;
20. import org.apache.hadoop.mapreduce.Mapper;
21. import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
22. import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
23.
24. public class WordCountHBase
25. {
26.     public static class Map extends Mapper<LongWritable,Text,Text,IntWritable>{
27.         private IntWritable i = new IntWritable(1);
28.         public void map(LongWritable key,Text value,Context context) throws
29.             IOException, InterruptedException{
30.             String s[] =value.toString().trim().split(" ");
31.             // 将输入的每行输入以空格分开
32.             for( String m : s){
33.                 context.write(new Text(m), i);
34.             }
35.         }
36.     }
37.     public static class Reduce extends TableReducer<Text, IntWritable,
38.         NullWritable>{
39.         public void reduce(Text key, Iterable<IntWritable> values, Context
40.             context) throws IOException, InterruptedException{
41.             int sum = 0;
42.             for(IntWritable i : values){

```

```

40.             sum += i.get();
41.         }
42.         Put put = new Put(Bytes.toBytes(key.toString()));
43.         //Put 实例化, 每一个词存一行
         put.add(Bytes.toBytes("content"), Bytes.toBytes("count"), Bytes.
             toBytes(String.valueOf(sum)));
         // 列族为 content, 列修饰符为 count, 列值为数目
44.         context.write(NullWritable.get(), put);
45.     }
46. }
47.
48. public static void createHBaseTable(String tablename) throws IOException{
49.     HTableDescriptor htd = new HTableDescriptor(tablename);
50.     HColumnDescriptor col = new HColumnDescriptor("content:");
51.     htd.addFamily(col);
52.     HBaseConfiguration config = new HBaseConfiguration();
53.     HBaseAdmin admin = new HBaseAdmin(config);
54.     if(admin.tableExists(tablename)){
55.         System.out.println("table exists, trying recreate table! ");
56.         admin.disableTable(tablename);
57.         admin.deleteTable(tablename);
58.     }
59.     System.out.println("create new table: " + tablename);
60.     admin.createTable(htd);
61. }
62.
63. public static void main(String args[]) throws Exception{
64.     String tablename = "wordcount";
65.     Configuration conf = new Configuration();
66.     conf.set(TableOutputFormat.OUTPUT_TABLE, tablename);
67.     createHBaseTable(tablename);
68.     String input = args[0]; // 设置输入值
69.     Job job = new Job(conf, "WordCount table with " + input);
70.     job.setJarByClass(WordCountHBase.class);
71.     job.setNumReduceTasks(3);
72.     job.setMapperClass(Map.class);
73.     job.setReducerClass(Reduce.class);
74.     job.setMapOutputKeyClass(Text.class);
75.     job.setMapOutputValueClass(IntWritable.class);
76.     job.setInputFormatClass(TextInputFormat.class);
77.     job.setOutputFormatClass(TableOutputFormat.class);
78.     FileInputFormat.addInputPath(job, new Path(input));
79.     System.exit(job.waitForCompletion(true)?0:1);
80. }
81. }

```

在上述程序中, 第 26 ~ 34 行代码负责设置 Map 作业; 第 36 ~ 46 行代码负责设置 Reduce 作业; 第 48 ~ 61 行代码为 createHBaseTable 函数, 负责在 HBase 中创建存储 WordCount 输出结果的表。在 Reduce 作业中, 第 42 ~ 44 行代码负责将结果存储到 HBase 表中。

程序运行成功后，现在通过 HBase Shell 检查输出结果，如图 12-11 所示：

```
hbase(main):004:0> scan 'wordcount'
ROW                                COLUMN+CELL
bye                                column=content:count, timestamp=1297571391451, value=2
hadoop                             column=content:count, timestamp=1297571391452, value=2
hello                              column=content:count, timestamp=1297571391452, value=2
world                             column=content:count, timestamp=1297571391452, value=2
4 rows(s) in 0.0390 seconds
```

图 12-11 HBase WordCount 的运行结果

从输出结果中可以看出，bye、hadoop、hello、world 四个单词均出现了两次。

关于 HBase 与 MapReduce 实际应用的更多详细信息请参阅 <http://wiki.apache.org/Hadoop/HBase/MapReduce>。

12.10 模式设计

通过 HBase 与 RDBMS 的比较，我们可以了解到二者无论是在物理视图、逻辑视图还是在具体操作上都存在很大的区别。例如，HBase 中没有 Join 的概念。但是，大表的结构可以使其不需要 Join 操作就能解决 Join 操作所解决的问题。比如，在一条行记录加上一个特定的行关键字，便可以实现把所有关于 Join 的数据合并在一起。另外，Row Key 的设计也非常关键。以天气数据存储为例。假如将监测站的值作为 Row Key 的前缀，那么天气数据将以监测站聚簇存放。同时将倒序的时间作为监测站的后缀，那么同一监测站的数据将从新到旧进行排列。这样的特定存储功能可以满足用户特殊的需要。

一般来说 HBase 的使用是为了解决或优化某一问题，恰当的模式设计可以使其具有 HBase 本身所不具有的功能，并且使其执行效率得到成百上千倍的提高。

下面我们通过两个例子，让读者对 HBase 的模式设计有一个初步的认识。

12.10.1 学生表

首先，我们以一个学习数据库过程中常用的一个学生表为例。众所周知，在关系型数据库（RDBMS）中学生表的表结构如表 12-15 至表 12-17 所示。

表 12-15 学生表 (Student)

字段	S_No	S_Name	S_Sex	S_Age
描述	学号	姓名	性别	年龄

表 12-16 课程表 (Course)

字段	C_No	C_Name	C_Credit
描述	课程号	课程名	学分

表 12-17 选课表 (SC)

字 段	SC_Sno	SC_Cno	SC_Score
描 述	学号	课程号	成绩

那么在 HBase 中，数据存储的模式将如表 12-18 和表 12-19 所示。

表 12-18 HBase 中的 Student 表

Row Key	Column Family		Column Family	
	info	value	course	value
<S_No>	info:S_Name info:S_Sex info:S_Age	the name the sex the age	course:<C_No> ...	<SC_Score> ...

表 12-19 HBase 中的 Course 表

Row Key	Column Family		Column Family	
	info	value	student	value
<C_No>	info:C_Name info:C_Credit	the name the credit	student:<S_No> ...	<SC_Score> ...

从上表中可以看出，在 RDBMS 中可以完成的操作，在 HBase 中不但可以完成，还可以有更好的执行效率。在 HBase 中 Row Key 是索引，因此在 HBase 中对数据进行查询，能够比 RDBMS 有更大的速度优势。

12.10.2 事件表

首先我们给出事件表在 RDBMS 中的表结构，如表 12-20 所示：

表 12-20 事件表 (Action)

字 段	A_Id	A_UserId	A_Name	A_Time
描 述	事件 ID	用户 ID	事件名称	事件发生时间

上述事件表存储了所有用户所发生的事件信息，包括事件名称，以及事件发生的时间。由于 HBase 一般是针对某一特殊的应用存储数据的，因此我们需要首先确定用户的需求。假如用户的需求描述如下：查询某一用户最近发生的 10 个事件。那么，RDBMS 的 SQL 查询语句如下所示：

```
SELECT A_Id,A_UserId,A_Name,A_Time From Action WHERE A_UserId=*** ORDER BY A_Time
DESC LIMIT 10
```

在 HBase 中为了加快数据的查询速度，现在需要将数据以用户聚簇的方式存放，并且按照事件发生的时间倒序排列。那么在 HBase 中将有下面的存储模式，见表 12-21。

表 12-21 HBase 中 Action 表

Row Key	Column Family	
	A_Name	value
<A_UserId><Long.Max_Value-System.currentTimeMillis()><A_Id>	A_Name	the name

从上表中可以看出，数据已经按照要求聚簇存放，查询速度必然要优于 RDBMS。

12.11 小结


本章向读者介绍了 HBase 的丰富内容，包括 HBase 的特点、基本操作、体系结构、数据模型、它与其他相关产品的关系，以及如何使用 HBase 编程、设计表等内容。

通过本章的内容，读者可以了解到，HBase 是一个开源的、分布式的、多版本的、面向列的存储模型。它与传统的关系型数据库有着本质的不同，并且在某些场合中，HBase 拥有其他数据库所不具有的优势。它为大型数据的存储和某些特殊应用提供了很好的解决方案。

另外，HBase 具有三种运行模式。其中，伪分布模式和全分布模式需要以 HDFS 作为其文件存储系统。因此 HBase 可以有效地与 MapReduce 结合起来使用，充分发挥二者的优势。本章为大家介绍了几个简单的编程实例，除此之外，还向读者简单比较了 HBase 的模式与传统 RDBMS 模式在设计上的异同之处。

希望通过本章的学习，能够让读者对 HBase 有一个全面、综合的了解。限于篇幅，未能深入地讲解 HBase 相关的知识，更多的内容，读者可以到 HBase 官方网站查阅：<http://hbase.apache.org/>。另外，我们还希望读者能够阅读 HBase 的源码，这样会对 HBase 的深层机制有更深入的理解。





第 13 章

Mahout 详解

本章内容

- ☐ Mahout 简介
- ☐ Mahout 的安装和配置
- ☐ Mahout API 简介
- ☐ Mahout 中的聚类和分类
- ☐ Mahout 应用：建立一个推荐引擎
- ☐ 小结

资源分享

PDG

13.1 Mahout 简介

Apache Mahout 起源于 2008 年，当时它是 Apache Lucene 的子项目。使用 Apache Hadoop 库，可以将其功能有效地扩展到 Apache Hadoop 云平台中。Apache Lucene 是一个著名的开源搜索引擎，它实现了先进的信息检索、文本挖掘功能。在计算机科学领域中，这些概念与机器学习技术相近。正是由于这种原因，一些 Apache Lucene 的开发者最终转而开发机器学习算法了。因此，这些机器学习算法形成了最初的 Apache Mahout。不久以后，Apache Mahout 吸收了一个名为 Taste 的开源协同过滤算法的项目，经过两年的发展，2010 年 4 月 Apache Mahout 最终成为了 Apache 的顶级项目。

Apache Mahout 的主要目标是建立可伸缩的机器学习算法。这种可伸缩性是针对大规模的数据集而言的。Apache Mahout 的算法运行在 Apache Hadoop 平台下，它通过 MapReduce 模式实现。但是，Apache Mahout 并不严格要求算法的实现要基于 Hadoop 平台，单个节点或非 Hadoop 平台也可以。Apache Mahout 核心库的非分布式算法也具有有良好的性能。

Apache Mahout 是 Apache Software Foundation (ASF) 旗下的一个开源项目，提供了一些经典的机器学习算法，旨在帮助开发人员更加方便快捷地创建智能应用程序。该项目已经发展到它的第三个年头了，有了三个公共发行版本。它们分别是 Mahout-0.1、Mahout-0.2、Mahout-0.3。Apache Mahout 项目包含聚类、分类、推荐引擎、频繁项集的挖掘。Apache Mahout 虽已经实现了很多技术和算法，但是仍然还有一些算法正在开发和测试。目前 Apache Mahout 项目主要包括以下四个部分。

- 聚类：将诸如文本、文档之类的数据分成局部相关的组。
- 分类：利用已经存在的分类文档训练分类器，对未分类的文档进行分类。
- 推荐引擎（协同过滤）：获得用户的行为并从中发现用户可能喜欢的事物。
- 频繁项集的挖掘：利用一个项集（查询记录或购物目录）去识别经常一起出现的项目。

13.2 Mahout 的安装和配置

因为 Mahout 是一个开源软件，因此 Mahout 有两种安装方式：一种是下载已经编译好的二进制文件进行安装；一种是先下载源代码，然后再对源代码进行编译，最后再安装。这里以第一种方式为例详细讲解 Mahout 的安装过程和方法。

(1) 下载 Mahout

编译好的二进制文件下载地址为：<http://labs.renren.com/apache-mirror/lucene/mahout/0.3/>
下载 mahout-0.3.tar.gz。

(2) 解压下载的文件

使用下面的命令将下载的二进制文件解压到指定的文件夹中。

```
tar -zxvf mahout-0.3.tar.gz -C /root
```

参数 -C 的后面是解压目标文件夹，这里是 /root。

(3) 配置环境变量

使用下面的命令配置 Mahout 所需要的 Hadoop 环境变量:

```
export HADOOP_HOME=/root/software/hadoop-0.20.2
export HADOOP_CONF_DIR=/root/software/hadoop-0.20.2/conf
```

(4) 检查是否安装成功

我们可以用下面的命令检查一下 Mahout 是否安装成功:

```
bin/mahout -help
```

如果安装成功, 系统会自动列出 Mahout 已经实现的所有命令, 如图 13-1 所示。

```
Valid program names are:
arff.vector: : Generate Vectors from an ARFF file or directory
canopy: : Canopy clustering
cat: : Print a file or resource as the logistic regression models would see it
cleansvd: : Cleanup and verification of SVD output
clusterdump: : Dump cluster output to text
dirichlet: : Dirichlet Clustering
fkmeans: : Fuzzy K-means clustering
fpg: : Frequent Pattern Growth
itemsimilarity: : Compute the item-item-similarities for item-based collaborative filtering
kmeans: : K-means clustering
lda: : Latent Dirichlet Allocation
ldatopics: : LDA Print Topics
lucene.vector: : Generate Vectors from a Lucene index
matrixmult: : Take the product of two matrices
meanshift: : Mean Shift clustering
prepare20newsgroups: : Reformat 20 newsgroups data
recommenditembased: : Compute recommendations using item-based collaborative filtering
rowid: : Map SequenceFile<Text,VectorWritable> to {SequenceFile<IntWritable,VectorWritable>, SequenceFile<IntWritable,Text>}
rowssimilarity: : Compute the pairwise similarities of the rows of a matrix
runlogistic: : Run a logistic regression model against CSV data
seqsparse: : Sparse Vector generation from Text sequence files
seqdirectory: : Generate sequence files (of Text) from a directory
seqdumper: : Generic Sequence File dumper
seqwiki: : Wikipedia xml dump to sequence file
svd: : Lanczos Singular Value Decomposition
testclassifier: : Test Bayes Classifier
trainclassifier: : Train Bayes Classifier
trainlogistic: : Train a logistic regression using stochastic gradient descent
transpose: : Take the transpose of a matrix
vectordump: : Dump vectors from a sequence file to text
wikipediaDataSetCreator: : Splits data set of wikipedia wrt feature like country
wikipediaXMLSplitter: : Reads wikipedia data and creates ch
```

图 13-1 Mahout 实现命令图

至此 Mahout 安装完毕。

(5) 运行示例

Mahout 自带了一些示例程序, 执行下面的 Hadoop 命令, 可以运行 Canopy 聚类算法示例:

```
bin/hadoop jar /root/mahout-0.3/mahout-examples-0.3.job
org.apache.mahout.clustering.syntheticcontrol.canopy.Job
```

转到 Mahout 安装目录下, 运行以下命令可以将结果直接显示在控制台上:

```
bin/mahout vectordump --seqFile /user/root/output/data/part-00000
```

13.3 Mahout API 简介

到目前为止, Mahout 项目已经发布了三个版本的 API。这里要介绍的 Mahout 最新版本

的 API 为 Mahout Core 0.3 API^①。

Mahout Core 0.3 API 主要可以分为以下几部分：

□ 基于协同过滤的 Taste 相关的 API。

□ 聚类算法相关的 API。

□ 分类算法，主要是与 Bayes 分类相关的 API，以及部分与决策树分类相关的 API。

下面将着重介绍聚类算法和分类算法的 API，因为聚类算法是 Mahout 中最为成熟的。Apache Mahout 已经实现的聚类算法有：Canopy 聚类算法、K-Means 聚类算法、模糊 K-Means 聚类算法、Mean Shift 聚类算法、Dirichlet 过程聚类算法和 Latent Dirichlet Allocation 聚类算法。这些算法相关的 API 都可以在 org.apache.mahout.clustering 包中找到。

下面以 K-Means 算法为例进行介绍。K-Means 算法的 API 在 org.apache.mahout.clustering.kmeans 中一共包含 1 个接口和 8 个类。它们分别是 KMeansConfigKeys、KMeansCluster、KMeansClusterMapper、KMeansCombiner、KMeansDriver、KMeansInfo、KMeansMapper、KMeansReducer 和 RandomSeedGenerator。

接口 KMeansConfigKeys 一共有 4 个参数：DISTANCE_MEASURE_KEY、CLUSTER_CONVERGENCE_KEY、CLUSTER_PATH_KEY、ITERATION_NUMBER，每个参数的具体意义如表 13-1 所示。

表 13-1 接口 KMeansConfigKeys 参数表

参 数	功 能
DISTANCE_MEASURE_KEY	K-Means 聚类算法使用的距离测量方法
CLUSTER_CONVERGENCE_KEY	K-Means 聚类算法的收敛值
CLUSTER_PATH_KEY	K-Means 聚类算法的路径
ITERATION_NUMBER	K-Means 聚类算法迭代的次数

对于 API 中的 8 个类，这里将重点介绍其中的两个类：KMeansCluster 和 KMeansDriver，这两个类的方法分别如表 13-2 和表 13-3 所示。

表 13-2 类 KMeansCluster 的方法列表

方 法	描 述
KMeansCluster (DistanceMeasure measure)	初始化 K-Means 聚类算法的构造方法。 参数 DistanceMeasure 用于比较点之间的距离
emitPointToNearestCluster(Vector point, java.util.List<Cluster> clusters, org.apache.hadoop.mapred.OutputCollector<org.apache.hadoop.io.Text, KMeansInfo> output)	在所有的聚类中找到离给定点距离最近的点。距离的测算方法使用设定的 DISTANCE_MEASURE_KEY 值
outputPointWithClusterInfo(Vector point, java.util.List<Cluster> clusters, org.apache.hadoop.mapred.OutputCollector<org.apache.hadoop.io.Text, org.apache.hadoop.io.Text> output)	输出带有聚类信息的点

① <http://mahout.apache.org/javac/core/overview-summary.html>。

(续)

方 法	描 述
<code>clusterPoints(java.util.List<Vector> points, java.util.List<Cluster> clusters, DistanceMeasure measure, int maxIter, double distanceThreshold)</code>	这是一个 K-Means 聚类算法的参考实现。对给定的点执行 K-Means 聚类算法。直到聚类满足收敛条件，或者遍历达到最大次数算法才结束
<code>runKMeansIteration (java.util.List<Vector> points, java.util.List<Cluster> clusters, DistanceMeasure measure, double distanceThreshold)</code>	遍历所有的点，并将每个点分配给一个聚类

表 13-3 类 KMeansDriver 的方法列表

方 法	描 述
<code>main(java.lang.String[] args)</code> throws <code>java.lang.Exception</code>	传入的参数按照 <code>runJob</code> 方法中的参数顺序执行
<code>public static void runJob(</code> <code>java.lang.String input,</code> <code>java.lang.String clustersIn,</code> <code>java.lang.String output,</code> <code>java.lang.String measureClass,</code> <code>double convergenceDelta,</code> <code>int maxIterations,</code> <code>int numReduceTasks)</code>	参数的意义依次如下： <input type="checkbox"/> 输入点的目录路径名 <input type="checkbox"/> 初始待计算的输入点所在的路径名 <input type="checkbox"/> 输出聚类点的路径名 <input type="checkbox"/> 距离测算法的类名 <input type="checkbox"/> 收敛值 <input type="checkbox"/> 最大遍历次数 <input type="checkbox"/> 归约的次数

13.4 Mahout 中的聚类和分类

13.4.1 什么是聚类和分类

在日常生活中经常会有重复的事情发生，人们会把自己遇到的事情和记忆中的事情关联起来。例如，糖果使人们想起是甜味。因此，人们会把具有甜味的食物归类为甜食。即使人们没有甜食的概念，人们也能把甜的食物进行归类。潜意识里，人们能够自然地将甜与苦进行分类。而生活中与此类似的现象还有很多，这些现象就是分类。

下面将用一个实际的例子来介绍到底什么是分类。假设在一个两岁的宝宝面前摆放一些水果，并告诉他红色圆的是苹果，橘黄色圆的是橘子。然后，拿一个又红又大的苹果问宝宝这是不是苹果，宝宝回答是，这就是一个简单的分类过程。在这个过程中主要涉及两个阶段：第一个是建立模型阶段；第二个是使用模型阶段。建立模型就是告诉两岁的宝宝具有何种特征的水果是苹果，具有何种特征的水果是橘子；使用模型就是问宝宝又红又大的水果是不是苹果。

在日常的生活中除了前面介绍的分类型外，还有很多种不同类型的聚类。下面同样用一个实际的例子来介绍聚类。假设你是一个藏书众多的图书馆馆长，但图书馆中的书是混乱的，没有任何顺序。来到图书馆的读者不得不找遍所有书籍才能发现自己想要看的书。这个寻找

书的过程非常缓慢。对于任何一个读者来说，这都是一个很头痛的问题。如果图书按照书名的首字母来进行排列，那么在知道书名的情况下寻找一本书将会变得很容易。如果图书是按照主题进行摆放的，图书查询也会变得简单易行。将众多的图书按照主题进行排列就是一个聚类的过程。在刚刚接触这个工作的时候，你不知道这些书会有多少种主题，比如哲学、文学等，也许还会有一些你从未听说过的主题。要完成这项任务，你首先要要把它们排成一列，逐本查阅。如果发现它与之前的书主题类似，就回到前面将它们放在一起，归为一类。当读完所有的书时，一遍聚类便完成了，众多的书籍也被分成了一些类。如果你觉得第一遍聚类的结果不够精细，你可以进行第二遍聚类，直到自己满意为止。

在下面的章节中，将会详细介绍 Mahout 中的分类和聚类。

13.4.2 Mahout 中的数据表示

生活中的数据会以各种各样的形式存储，Mahout 中的数据也会以其固定的形式表示。在 Mahout 中，数据将会以向量的形式进行存储。

多数人对向量这个词并不陌生。在不同的领域，向量具有不同的实际意义。在物理中，向量用来表示力的大小和方向，或者一个移动物体的速度。在数学中，一个向量表示空间中的一个点。虽然它们代表的意义不同，但它们表示的形式是相同的。在二维空间中，所有的向量都可表示成诸如 (5, 6) 的形式，每一维中有一个数字。当计算这个二维向量时，人们常称第一个维度为 X，第二个维度为 Y。但是在现实生活中，一个向量可以是多维度的。按照顺序，向量的每一个维度依次被称为 0 维、1 维、2 维……

如上所述，向量是按照维度排列的一系列有序的值。你可能已经想到在程序设计语言中用一维数组来表示向量。使用这种方式表示向量，数组的第 i 项刚好是向量的第 i 个维度的值。这是一种很好的表示向量的方法，称为密集向量表示法。

在现实生活中，一个具有很高维度的向量经常会在很多维度上没有值。这里的没有值就是程序设计当中空的概念，在向量中它会表示为 0。在物理和数学领域，无论是高维度向量还是包含很多 0 的向量都是很少见的。但在分类算法中这种情况很常见。

使用数组表示这种向量效率太差。数组将会包含很多个 0，偶尔会有一个非 0 值。舍弃众多的 0 值，单独表示非 0 值是一种很合理的想法。当处理数百万维度带有很多 0 值的向量时，密集向量表示法的弊端变得很明显。

在这种情况下，Mahout 引入了稀疏向量，将非 0 值所在的维度与该维度的值做映射。这可以通过 Java 中的 Map 实现。当非 0 值比较少时，这种存储方式比使用基于数组的稠密存储更具优越性。但使用这种方式，程序需要更多的内存空间。

在 Mahout 中，有关向量表示的类有三个。它们分别是稠密向量 (DenseVector)、随机访问向量 (RandomAccessSparseVector) 和序列访问向量 (SequentialAccessSparseVector)。

稠密向量由一个 double 型的数组实现。当向量具有很少的非 0 值时，这种向量表示法的效率很高。它允许快速访问向量所在任何维度的值，并且能够快速按序遍历向量的所有维度。

在随机访问向量类中，向量的值存储在类似于 HashMap 的结构中，键是 int 型、值是 double 型的。只有维度上的值非 0 时，该维度值才会被存储。当一个向量的一些维度值非 0 时，用随机访问向量方式表示向量比用稠密向量表示法具有更高的内存使用效率。但是访问维度值的速度和按序遍历所有维度值的速度比较慢。

序列访问向量使用 int 和 double 的并行数组表示向量。因此，使用它按序遍历整个向量的各个维度是很快。但是随机插入和查询某一维度的值时速度要慢于随机访问向量。

这三种表示向量方式使 Mahout 的算法能够按照数据特性、数据访问方式实现。具体使用哪种表示方法是按照算法的特性来进行选择的。如果算法具有很多对向量值的随机插入和更新，就应该选择稠密向量或随机访问向量来表示向量。因为这两个向量具有快速的随机访问特性。而对于需要重复计算向量大小的 K-Means 聚类算法来说，选择序列访问向量比选择随机访问向量好。

13.4.3 将文本转化成向量

讨论完如何存储向量后，下面开始讨论如何将文本转化成向量。在信息时代，文本文件的数量呈爆炸式增长。仅 Google 搜索引擎的索引就有 200 亿的 Web 文档。文本数据是海量的，这些海量数据中蕴涵着大量的知识。公司或机构可以使用诸如聚类、分类的机器学习算法去发现这些知识。学习用向量表示文本是从海量数据发现知识的第一步。

向量空间模型 (VSM, Vector Space Model) 是最常用的相似度计算模型。Mahout 中对文本的聚类使用了这种技术。什么是向量空间模型？下面做一个简单的介绍。

假设共有十个词： w_1, w_2, \dots, w_{10} ，五篇文章： d_1, d_2, d_3, d_4 和 d_5 。统计所得的词频表如表 13-4 所示。

表 13-4 空间向量模型表

	w_1	w_2	w_3	w_4	w_5	w_6	w_7	w_8	w_9	w_{10}
d_1	1	2		5		7		9		
d_2		3		4		6	8			
d_3	10		11		12			13	14	15
d_4		5		7				4		9
d_5			2		2	3	7			

这个词频表就是空间向量模型。对于任意的两篇文档，当要计算它们的相似度时，可以选择计算两个向量的余弦值。如果余弦值为 1，则说明两篇文档完全相同；如果余弦值为 0，则说明两篇文档完全不同。总之，在 $[0,1]$ 内余弦值越大，两篇文章的相似度就越大。除了计算余弦值以外，还有其他的方法测量两篇文章的相似度，这里不作介绍。

在 Mahout 下处理的数据必然是海量数据。待处理的文本包含的所有单词就是例子中的单词 w ，待处理的文本文件就是相应的 d 。可以想象，待处理文本所包含的单词是巨大的，因此，文本向量的维度也是巨大的。示例中的具体数值是某个单词在特定文章中出现的次

数,称之为词频 (term frequency)。例如,表中的 1 代表单词 w_1 在 d_1 文档中出现一次,其词频为 1。

在一些简单的处理方法中,可以只通过词频来计算文本间的相似度,不过当某个关键词在两篇长度相差很大的文本中出现的频率相近时,会降低结果的准确性。因此通常会把词频数据正规化,以防止词频数据偏向于关键词较多,即较长的文本。如某个词在文档 d_1 中出现了 100 次,在 d_2 中出现了 100 次,仅从词频看来,这个词在这两个文档中的重要性相同,然而,再考虑另一个因素,就是 d_1 的关键词总数是 1000,而 d_2 的关键词总数是 100000,所以从总体上看,这个词在 d_1 和 d_2 中的重要性是不同的。正规化处理的方法是用词频除以所有文档的关键词总数。

当仅使用词频来区分文档时,还会遇到这样一个问题。众所周知,一篇文章会包含很多诸如一、二、你、我、他等的单词,并且这些词语会多次出现。很明显,无论使用何种距离来测算两篇文章的相似度,这些经常出现的词汇都会对结果起到很大的负面影响。但这些词语并不能区分两份文档,相似性的判断因此也变得不再准确。把文档按照相似性进行合理的聚类就更不可能了。为了解决这个问题,人们使用了 TF-IDF (Term Frequency-Inverse Document Frequency) 技术。

TF-IDF 是一种统计的方法,用于评估一个字词对于一个文件集或一个语料库中一份文件的重要程度。字词的重要性随着它在文件中出现的次数成正比增加,但同时会随着它在语料库中出现的频率成反比下降。TF-IDF 的主要思想是,如果某个词或短语在一篇文章中出现的频率 TF 高,并且在其他文章中很少出现,则认为此词或短语具有很好的类别区分能力,适合用来分类。TF-IDF 实际上是 $TF * IDF$ 。TF 代表词频 (Term Frequency),表示词条在文档中出现的频率。IDF 代表反文档频率 (Inverse Document Frequency)。IDF 的主要思想是,如果包含词语 w 的文档越少,IDF 越大,则说明词语 w 具有很好的类别区分能力。

13.4.4 Mahout 中的聚类、分类算法

Mahout 目前已经实现了 Canopy 聚类算法、K-Means 聚类算法、Fuzzy K-Means 聚类算法、Dirichlet 过程聚类算法等众多聚类算法。除此之外, Mahout 还实现了贝叶斯 (Bayes) 分类算法。这里主要介绍简单且应用广泛的 K-Means 聚类算法和贝叶斯分类算法。

K-Means 聚类算法能轻松地对几乎所有的问题进行建模。K-Means 聚类算法容易理解,并且能在并行计算机上很好的运行。学习 K-Means 聚类算法,能更容易理解聚类算法的缺点,以及其他算法对于特定数据的高效性。

K-Means 聚类算法中的 K 是聚类的数目,在算法中会强制要求用户输入。如果将新闻聚类成诸如政治、经济、文化等大类,可以选择 10 到 20 之间的数字作为 K。因为这种顶级的类别数量是很小的。如果要对这些新闻详细分类,选择 50 到 100 之间的数字也是没有问题的。假设数据库中有一百万条新闻。如果想把这一百万条新闻按照新闻谈论的内容进行聚

类, 则这个聚类数目会远远大于之前的聚类数目。因为每个聚类中的新闻数量不会太大。这就要求选择一个诸如 10 000 的聚类数值。聚类数值 K 的取值范围不定, 它既可以小至几个, 也可以大至几万个。这就对算法的伸缩性提出了很高的要求, 而 Mahout 下实现的 K-Means 聚类算法就具有很好的伸缩性。

K-Means 聚类算法主要可以分为三步。第一步是为待聚类的点寻找聚类中心; 第二步是计算每个点到聚类中心的距离, 将每个点聚类到离该点最近的聚类中去; 第三步是计算聚类中所有点的坐标平均值, 并将这个平均值作为新的聚类中心点。反复执行第二步, 直到聚类中心不再进行大范围的移动, 或者聚类次数达到要求为止。

假设有 n 个点, 需要将它们聚类成 K 个组。K-Means 聚类算法会以 K 个随机的中心点开始。算法反复执行上文中提到的第二步和第三步, 直至终止条件得到满足。接下来以 9 个点且 K 值为 3 为例, 配以相应的图示介绍 K-Means 聚类算法。

在聚类前, 首先在二维平面中随机选择 9 个点, 坐标分别为 (7, 8)、(12, 1)、(13, 6)、(13, 13)、(13, 19)、(14, 5)、(17, 16)、(19, 20)、(20, 7)。

1. 第一次聚类

1) 系统首先选取前 3 个点 (7, 8)、(12, 1)、(13, 6) 作为聚类中心, 然后计算每个点到聚类中心的距离, 该点距离哪个聚类中心的距离最小就归属于哪个聚类中心。经过计算, 点 (7, 8)、(13, 19) 为 1 个聚类, 点 (12, 1) 为 1 个聚类, 点 (13, 6)、(13, 13)、(14, 5)、(17, 16)、(19, 20)、(20, 7) 为 1 个聚类, 如图 13-2 所示。

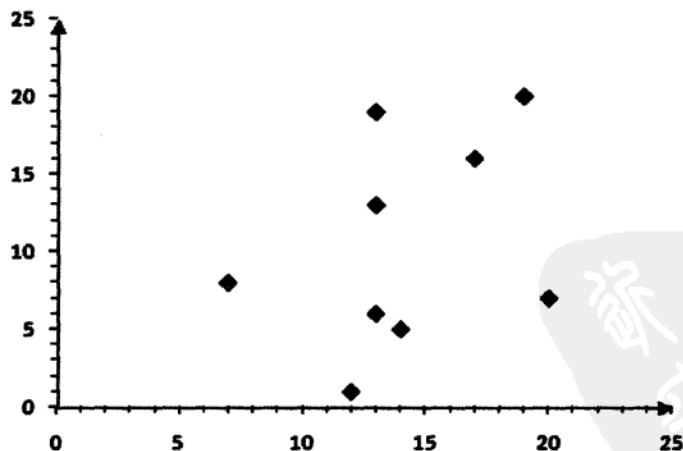


图 13-2 未聚类的 9 个点

2) 更新聚类的聚类中心, 新的聚类中心的值为聚类中所有成员的平均值。聚类 (7, 8)、(13, 19) 的新聚类中心为 (10.0, 13.5), 聚类 (12, 1) 的新聚类中心仍为 (12.0, 1.0), 聚类 (13, 6)、(13, 13)、(14, 5)、(17, 16)、(19, 20)、(20, 7) 的新聚类中心为 (16.0, 11.2),

如图 13-3 所示。

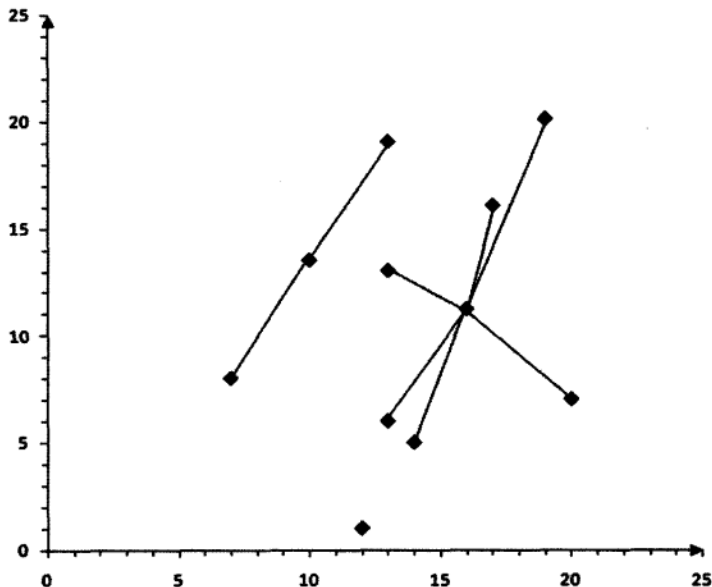


图 13-3 第一次聚类后的结果

2. 第二次聚类

1) 根据前面生成的聚类中心 (10.0, 13.5)、(12.0, 1.0)、(16.0, 11.2) 重新计算每个点到聚类中心点的距离, 根据计算出的距离对该点进行聚类。结果点 (7, 8)、(13, 13)、(13, 19) 为 1 个聚类, 点 (12, 1)、(13, 6)、(14, 5) 为 1 个聚类, 点 (17, 16)、(19, 20)、(20, 7) 为 1 个聚类。

2) 更新聚类的聚类中心, 聚类 (7, 8)、(13, 13)、(13, 19) 的新聚类中心为 (11.0, 13.3), 聚类 (12, 1)、(13, 6)、(14, 5) 的新聚类中心仍为 (13.0, 4.0), 聚类 (17, 16)、(19, 20)、(20, 7) 的新聚类中心为 (18.7, 14.3), 如图 13-4 所示。

3. 第三次聚类

根据上一步来看, 聚类结果没有发生变化, 满足收敛条件, K-Means 聚类算法结束, 如图 13-5 所示。

介绍完 K-Means 聚类算法, 下面开始介绍贝叶斯 (Bayes) 分类算法。贝叶斯 (Bayes) 分类算法是一种基于统计的分类方法, 用来预测某个样本属于某个分类的概率有多大。贝叶斯 (Bayes) 分类算法是基于贝叶斯定理的分类算法。

贝叶斯分类算法有很多变种。在这里主要介绍朴素贝叶斯分类算法。何谓朴素? 所谓朴素就是假设各属性之间是相互独立的。经过研究发现, 在大多数情况下, 朴素贝叶斯分类算法

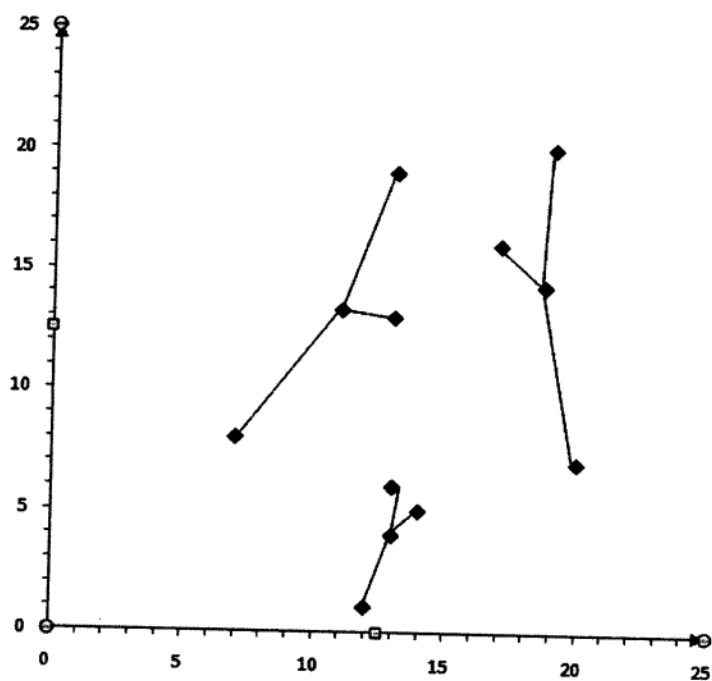


图 13-4 第二次聚类后的结果

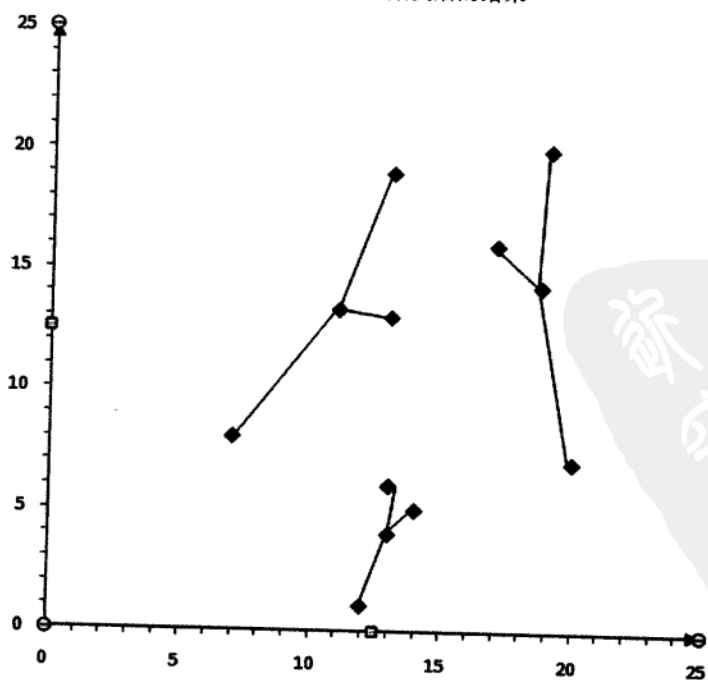


图 13-5 第三次聚类后的结果

(Naïve Bayes Classifier) 在性能上和决策树 (Decision Tree)、神经网络 (Neural Network) 相当。当针对大数据集的应用时, 贝叶斯分类算法具有方法简单、高准确率和高速度的优点。但事实上, 贝叶斯分类算法也有其缺点。缺点就是由于贝叶斯定理假设一个属性值对给定类的影响独立于其他属性的值, 而此假设在实际情况中经常是不成立的, 因此其分类准确率可能会下降。

朴素贝叶斯分类算法是一种监督学习算法, 使用朴素贝叶斯分类算法对文本进行分类, 主要有两种模型, 即多项式模型 (multinomial model) 和伯努利模型 (Bernoulli model)。Mahout 实现的贝叶斯分类算法使用的是多项式模型。对算法具体内容感兴趣的读者可以阅读 <http://people.csail.mit.edu/jrennie/papers/icml03-nb.pdf> 上的论文。本书将以一个实际的例子来简略介绍使用多项式模型的朴素贝叶斯分类算法 (Naïve Bayes Classifier)。

给定一组分类号的文本训练数据, 如表 13-5 所示。

表 13-5 文本训练数据表

文档编号	文 档	文档类别 (文档是否属于中国类)
1	中国, 北京, 中国	是
2	中国, 中国, 上海	是
3	中国, 澳门	是
4	东京, 日本, 中国	否

给定一个新的文档样本“中国、中国、中国、东京、日本”, 对该样本进行分类。该文本属性向量可以表示为 $d = (\text{中国}, \text{中国}, \text{中国}, \text{东京}, \text{日本})$, 类别集合 $Y = \{\text{是}, \text{否}\}$ 。类别“是”下共有 8 个单词, “否”类别下面共有 3 个单词。训练样本单词总数为 11。因此 $P(\text{是}) = 8/11$, $P(\text{否}) = 3/11$ 。

类条件概率计算如下:

- $P(\text{中国} | \text{是}) = (5+1)/(8+6) = 6/14 = 3/7$;
- $P(\text{日本} | \text{是}) = P(\text{东京} | \text{是}) = (0+1)/(8+6) = 1/14$;
- $P(\text{中国} | \text{否}) = (1+1)/(3+6) = 2/9$;
- $P(\text{日本} | \text{否}) = P(\text{东京} | \text{否}) = (1+1)/(3+6) = 2/9$ 。

上面 4 条语句分母中的 8, 是指“是”类别下训练样本的单词总数, 6 是指训练样本有中国、北京、上海、澳门、东京、日本共 6 个单词, 3 是指“否”类别下共有 3 个单词。有了以上的类条件概率, 开始计算后验概率:

- $P(\text{是} | d) = (3/7)^3 \times 1/14 \times 1/14 \times 8/11 = 108/184877 \approx 0.00058417$;
- $P(\text{否} | d) = (2/9)^3 \times 2/9 \times 2/9 \times 3/11 = 32/216513 \approx 0.00014780$ 。

因此, 这个文档属于类别中国。这就是 Mahout 实现的贝叶斯 (Bayes) 分类算法的主要思想。

13.4.5 算法应用实例

在 Mahout 中运行 K-Means 算法非常简单。对于不同的数据主要有以下三个步骤。

- ❑ 使用 seqdirectory 命令将待处理的文件转化成序列文件。
- ❑ 使用 seq2sparse 命令将序列文件转化成向量文件。
- ❑ 使用 kmeans 命令对数据运行 K-Means 聚类算法。

将文本文件转化成向量需要两个重要的工具。一个是 SequenceFilesFromDirectory 类，它能将一个目录结构下的文本文件转化成序列文件，这种序列文件为一种中间文本表示形式。另一个是 SparseVectorsFromSequenceFiles 类，它使用词频 (TF) 或 TF-IDF (TF-IDF weighting with n-gram generation) 将序列文件转化成向量文件。序列文件以文件编号为键、文件内容为值。下面讨论如何将文本转换成向量。

使用路透社 14578 新闻集作为示例数据。这组数据被广泛应用于机器学习的研究中，它起初是由卡内基集团有限公司和路透社共同搜集整理的，目的是发展文本分类系统。路透社 14578 新闻集分布于 22 个文档中，除最后的 reut2-0.14.sgm 包含 578 份文件外，其余的每个文件包含 1000 份文件。

路透社 14578 新闻集中的所有文件都为标准通用标记语言 SGML (Standard Generalized Markup Language) 格式，这种格式的文件与 XML 文件格式相似。可以为 SGML 文件创建一个分析器 (parser)，并将文件编号 (document ID) 和文件内容 (document text) 写到序列文件 (SequenceFiles) 中去。然后用前文提到的向量化工具将序列文件转化成向量。但是，更快捷的方式是使用 Lucene Benchmark JAR 文件提供的路透社分析器 (the Reuters Parser)。Lucene Benchmark JAR 是捆绑在 Mahout 上的，剩下的工作只是到 Mahout 目录下的 examples 文件夹下运行 org.apache.lucene.benchmark.utils.ExtractReuters 类。在这之前，需要从 <http://www.daviddlewis.com/resources/testcollections/reuters14578/reuters14578.tar.gz> 下载路透社新闻集，并将它解压到 Examples/Reuters 文件夹下。相关命令如下所示：

```
mvn -e -q exec:java
-Dexec.mainClass="org.apache.lucene.benchmark.utils.ExtractReuters"
-Dexec.args="reuters/ reuters-extracted/"
```

使用解压得到的文件夹运行 SequenceFileFromDirectory 类。使用下面的脚本命令 (Launch script) 可以实现该功能：

```
bin/mahout seqdirectory -c UTF-8
-i examples/reuters-extracted/
-o reuters-seqfiles
```

这条命令的作用是将路透社文章转化成序列文件格式，如表 13-6 所示。

表 13-6 Seqdirectory 命令参数表

参 数	描 述
--chunkSize (-chunk) chunkSize	文件按多大值进行切分，默认值是 64MB
--charset (-c) charset	输入文件使用的编码方式，多使用 UTF-8
--input (-i) input	输入文件的位置
--output (-o) output	输出文件的位置
--help (-h)	输出帮助信息

现在剩下的工作是将序列文件转化成向量文件。运行 `SparseVectorsFromSequenceFiles` 类即可实现该功能。命令如下：

```
bin/mahout seq2sparse -i reuters-seqfiles/ -o reuters-vectors -w
```

注意，在 `seq2sparse` 命令中，参数 `-w` 用来表示是否覆盖输出文件夹。Mahout 用来处理海量的数据，任何一个算法的输出都会花费很多时间。有了参数 `-w`，Mahout 就可以防止新产生的数据对未完全输出的数据进行破坏了。除此之外，`seq2sparse` 命令还有以下参数，如表 13-7 所示。

表 13-7 seq2sparse 命令参数表

参 数	描 述
-w (boolean)	该参数决定是否重写文件夹的内容。如果未设置该参数，当文件夹不存在时，则创建该文件夹；如果文件夹已经存在，则工程会抛出异常
-a (String)	所使用的分析器的名字，默认值为 <code>org.apache.lucene.analysis.standard.StandardAnalyzer</code>
-chunk (int)	分块大小以 MB 为单位。在向量化过程中，GB 或 TB 级的数据不能完全放入内存中。因此要将数据分成指定的大小，分阶段进行向量化。建议使用 Hadoop 子节点中 Java 堆大小的 80%。这里默认值为 100MB
-wt (String)	使用的加权模式。tf 为 TF；tfidf 为 TF-IDF。默认值为使用 TF-IDF
-s (int)	整个集合中单词出现最低的频率，如果一个单词出现次数小于该值，那么该单词被忽略。默认值为 2
-minDF (int)	最小的文件频率 (document frequency)，默认值是 1
-x (int)	用来滤除频率很高的词汇，该值是 0 ~ 100 之间的整数，默认值为 99
-ng (int)	最大的 N-gram 值，默认值为 2
-nr (int)	reduce 任务的数量，默认值为 1。该参数不是必须的
-seq (bool)	是否将结果以 <code>SequentialAccessVectors</code> 形式输出。如果设置，则以 <code>SequentialAccessVectors</code> 形式输出

Mahout 的 `seq2sparse` 命令的功能是从序列文件中读取数据，使用上面提到的默认参数，按照基于向量化 (vectorizer) 的字典生成向量文件。读者可以使用以下命令来检查生成文件夹：

```
ls reuters-vectors/
```

执行上述命令后，结果如下所示：

```
dictionary.file-0
tfidf/
tokenized-documents/
vectors/
wordcount/
```

输出文件夹包含一个目录文件和四个文件夹。目录文件保存着术语（term）和整数编号之间的映射。当读取算法的输出时，这个文件是非常有用的。因此，需要保留它。其他四个文件夹是向量化过程中生成的文件夹。向量化过程主要有以下几步：

第一步，标记文本文档。具体过程是使用 Lucene StandardAnalyzer 将文本文档分成个体化的单词，将结果存储在 tokenized-documents 文件夹下。

第二步，对 tokenized 文档进行迭代生成一个重要单词的集合。这个过程可能会使用单词统计、n-gram 生成，这里使用的是 unigrams 生成。

第三步，使用 TF 将标记的文档转化成向量，从而创建 TF 向量。在默认情况下，向量化是使用 TF-IDF，因此需要两步来进行，一是文档频率（document-frequency）的统计工作；二是创建 TF-IDF 向量。TF-IDF 向量在 tfidf/vectors 文件夹下。对于大多数的应用来说，仅仅需要目录文件和 tfidf/vectors 文件夹。

使用 kmeans 命令可以对数据运行 K-Means 聚类算法。命令如下：

```
bin/mahout kmeans
-i ./examples/bin/work/reuters-out-seqdir-sparse/tfidf/vectors/
-c ./examples/bin/work/clusters
-o ./examples/bin/work/reuters-kmeans
-k 20 -w
```

表 13-8 列出了 K-Means 命令参数的具体意义。

表 13-8 K-Means 聚类算法列表

参 数	描 述
--k (-k) k	K-Means 聚类算法中的 K 值，如果不指定该值……
--input (-i) input	输入文件的位置，输入文件必须是序列文件
--output (-o) output	输出文件的位置
--distance (-m) distance	使用距离测量方法，默认是欧几里得距离
--convergence (-d) convergence	收敛的阈值默认值为 0.5
--max (-x) max	最大遍历次数，默认值为 20
--numReduce (-r) numReduce	reduce 任务的数量
--vectorClass (-v) vectorClass	向量化所使用的向量名字。默认值为 RandomAccessSparseVector.class
--overwrite (-w)	决定是否对输出结果进行覆盖，如果设置，则会覆盖之前的输出结果

在已经安装好 Hadoop 和 Mahout 的前提下, 运行 Bayes 分类算法也比较简单。这里简要介绍一下 Mahout 的示例程序 20NewsGroup 的分类。实际上, 除了 20NewsGroup 示例之外, 还有 Wikimapia 数据, 但由于其数据量达到了将近 7GB, 所以对大多数初学者而言并不合适。因此这里选择数据量较小且非常经典的 20NewsGroup 示例的分类。

什么是 20NewsGroup? 20 新闻组包含 20 000 个新闻组文档, 这些文档可以被分类成 20 个新闻组。20 新闻组最初来源于 Ken Lang 的论文《Newsweeder: learning to filter netnews》。从那以后, 20 新闻组数据集合在机器学习领域越来越多地被用作实验数据。在文本聚类 and 分类方面的研究中用得最多。20 新闻组按照 20 个不同的类型进行组织, 不同的类对应不同的主题。本书用到的 20 新闻组数据可以从 <http://people.csail.mit.edu/jrennie/20Newsgroups/> 下载。在下载页面中, 一共有三种版本的 20 新闻组数据, 分别是 20news-19997.tar.gz、20news-bydate.tar.gz 和 20news-18828.tar.gz。

20news-19997.tar.gz 是最原始的版本数据。20news-bydate.tar.gz 是按照日期进行排序的, 其中的 60% 用来进行训练 Bayes 分类算法, 40% 用来测试 Bayes 分类算法。不包含重复新闻和标识新闻组的标题。20news-18828.tar.gz 不包含重复的新闻, 但是包含带有新闻来源和新闻主题的标题。这三种 20 新闻组数据都是以 tar.gz 形式存在的。读者使用 tar 命令对它们进行解压即可得到相应的数据。具体选择哪种数据对结果影响不大, 这里选择 20news-bydate.tar.gz。

介绍完 20NewsGroup 后, 下面开始介绍如何运行 Mahout 自带的 Naïve Bayes Classifier 算法示例。

首先将下载好的数据解压到文件夹下, 然后将解压好的文件上传到 HDFS 上, 使用下列命令训练 Bayes 分类器:

```
HADOOP_HOME/bin/hadoop\
  jar\
  $MAHOUT_HOME/examples/target/mahout-examples-0.3 -job.jar\
  org.apache.mahout.classifier.bayes.TrainClassifier\
  -i 20news-input\
  -o newsmodel\
  -ng 3\
  -type bayes
  -source hdfs
```

该命令将会在 Hadoop 上运行四个 MapReduce 作业。在命令执行的过程中, 可以打开浏览器在 <http://localhost:50030/jobtracker.jsp> 上监视这些作业的运行状态。

运行下面的命令测试 Bayes 分类器:

```
$HADOOP_HOME/bin/hadoop jar $MAHOUT_HOME/examples/target/mahout-examples-0.3-job.jar\
  org.apache.mahout.classifier.bayes.TestClassifier\
  -m newsmodel\
  -d 20news-input\
  -ng 3\
  -type bayes
```

这就是 Mahout 自带的 Bayes 分类算法的示例程序。如果读者想要深入了解 Mahout 的分类算法，可以自行阅读 Mahout Core API 0.3 来了解已经实现的功能。

13.5 Mahout 应用：建立一个推荐引擎

13.5.1 推荐引擎简介

每天人们都会产生各种各样的想法：喜欢一个产品、不喜欢一件事、不关心某个东西。在人们毫无察觉的情况下，这些事情在悄然发生。一个正在播放的流行歌曲可能会引起你的注意，也可能对你没有任何影响。歌曲引起你的注意可能是因为它很好听或者它很让人厌烦。同样的事情也会发生在其他的事情上。这就是人们的喜好。

每个人都有着不同的喜好，但是这些喜好会遵循着类似的规律。对于一个人来说，如果一个新的事物与他之前喜欢的事物相似，那么他很有可能会喜欢这个新事物。如果一个外国人喜欢吃中国饺子，那么他很有可能会喜欢中国的包子。因为它们都是带馅的面食。此外，如果你的朋友喜欢周国平的散文，那么你也很有可能会喜欢周国平的散文。因为朋友之间会有一些共同的喜好。

在日常生活中，预测人们的喜好是没有问题的。假设有两个人 A 和 B。对于 B 是否喜欢电影《指环王 III》的问题，大多数人只能靠猜测。但如果 A 知道 B 喜欢《指环王 I》和《指环王 II》，那么可以推测 B 喜欢《指环王 III》。如果 B 对指环王系列电影一点儿也不了解，A 基本可以断定，B 是不会喜欢《指环王 III》的。

推荐引擎就是对人们的喜好做出预测的一种技术。它会依据已经获得的各种信息，对用户的购买行为做出预测，从而达到相应目的。现实生活中，人们都经历过网站向客户推荐产品，这些推荐都是基于客户浏览信息的推荐。网站试着推断出客户的喜好，以此来向客户推荐他们可能会喜欢的产品。

卓越网使用了推荐引擎技术，在购买一本书的同时，网站会利用顾客的购买习惯和书籍之间的关系为顾客推荐他们可能会感兴趣的书籍或音像制品。例如，当某一名顾客想要购买《云计算》这本书时，在页面的下方会出现购买此商品的顾客同时购买的书籍。这样顾客就可能顺便买一本相关的书。推荐引擎技术不仅可帮助顾客更容易地发现自己想要的商品，而且可帮助商家售卖更多的商品。社交网站人人网利用推荐引擎技术，向用户推荐一些可能是用户朋友的人。对于最有可能是朋友的人，人人网会自动把这些最可能是该用户朋友的人放在最前方，以供用户选择。推荐引擎技术已经悄然地影响着人们的生活，只是人们可能并没有注意它。

13.5.2 使用 Taste 构建一个简单的推荐引擎

Taste 是 Apache Mahout 提供的一个协同过滤算法的高效实现，它是一个基于 Java 实现的可扩展的、高效的推荐引擎。Taste 既实现了最基本的基于用户的和基于内容的推荐算法，

同时也提供了扩展接口，使用户可以方便地定义和实现自己的推荐算法。同时，Taste 不仅仅适用于 Java 应用程序，它还可以作为内部服务器的一个组件以 HTTP 和 Web Service 的形式向外界提供推荐的逻辑。Taste 的设计使它能满足企业在推荐引擎上对性能、灵活性和可扩展性等方面的要求。

Taste 主要包括以下 5 个组件，具体如图 13-6 所示。

DataModel：DataModel 是用户喜好信息的抽象接口，它的具体实现支持从任意类型的数据源抽取用户喜好信息。Taste 默认提供 JDBCDataModel 和 FileDataModel，分别支持从数据库和文件中读取用户的喜好信息。

UserSimilarity 和 ItemSimilarity：UserSimilarity 用于定义两个用户间的相似度，它是基于协同过滤的推荐引擎的核心部分，可以用来计算用户的“邻居”，这里的邻居指的是与当前用户相似的用户。ItemSimilarity 是用来计算内容之间的相似度的。

UserNeighborhood：UserNeighborhood 用于基于用户相似度的推荐方法中，推荐的内容是通过找到与当前用户喜好相似的“邻居用户”而产生的。UserNeighborhood 定义了确定邻居用户的方法，具体实现一般是基于 UserSimilarity 计算得到的。

Recommender：Recommender 是推荐引擎的抽象接口，Taste 中的核心组件。在程序中，为它提供一个 DataModel，它可以计算出对不同用户的推荐内容。在实际应用中，主要使用它的实现类 GenericUserBasedRecommender 或 GenericItemBasedRecommender，分别实现基于用户相似度的推荐引擎或基于内容的推荐引擎。

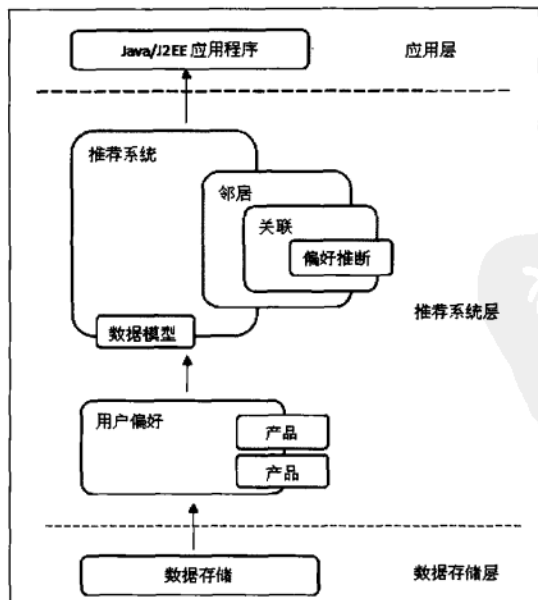


图 13-6 Taste 的主要组件图

安装 Taste 主要包括以下三部分内容：

- 如果需要 build 源代码或例子，则需要 Apache Ant 1.5+ 或 Apache Maven 2.0.10+。
- Taste 应用程序需要 Servlet 2.3+ 容器，例如 Jakarta Tomcat。
- Taste 中的 MySQLJDBCDataModel 实现需要 MySQL 4.x+ 数据库。

安装 Taste 并运行 Demo 的步骤如下：

- 1) 从 SVN 或下载压缩包中得到 Apache Mahout 的发布版本。
- 2) 从 Grouplens 网站 <http://www.grouplens.org/node/12> 下载数据源：“1 Million MovieLens Dataset”。
- 3) 解压数据源压缩包，将 movie.dat 和 ratings.dat 拷贝到 Mahout 安装目录下的 taste-web/src/main/resources/org/apache/mahout/cf/taste/example/grouplens 目录下。
- 4) 回到 core 目录下，运行“mvn install”，将 Mahout core 安装在本地库中。
- 5) 进入 taste-web，拷贝 ../examples/target/grouplens.jar 到 taste-web/lib 目录。
- 6) 编辑 taste-web/recommender.properties，将 recommender.class 设置为 org.apache.mahout.cf.taste.example.grouplens.GroupLensRecommender。
- 7) 在 Mahout 的安装目录下，运行“mvn package”。
- 8) 运行“mvn jetty:run-war”，这里需要将 Maven 的最大内存设置为 1024MB，即：MAVEN_OPTS=-Xmx1024MB。如果需要在 Tomcat 下运行，可以在执行“mvn package”后，将 taste-web/target 目录下生成的 war 包拷贝到 Tomcat 的 webapp 下，同时也需要将 Java 的最大内存设置为 1024MB，即 JAVA_OPTS=-Xmx1024MB，然后启动 Tomcat。
- 9) 访问 [http://localhost:8080/\[your_app\]/RecommenderServlet?userID=1](http://localhost:8080/[your_app]/RecommenderServlet?userID=1)，得到系统编号为 1 的用户推荐内容。Taste Demo 的运行结果界面参看图 13-7，其中每一行的第一项是推荐引擎预测的评分，第二项是电影的编号。
- 10) 同时，Taste 还提供 Web 服务访问接口，可通过以下 URL 访问：[http://localhost:8080/\[your_app\]/RecommenderService.jws](http://localhost:8080/[your_app]/RecommenderService.jws)。
- 11) WSDL 文件 ([http://localhost:8080/\[your_app\]/RecommenderService.jws?wsdl](http://localhost:8080/[your_app]/RecommenderService.jws?wsdl)) 也可以通过简单的 HTTP 请求调用这个 Web 服务：[http://localhost:8080/\[your_app\]/RecommenderService.jws?method=recommend&userID=1&howMany=10](http://localhost:8080/[your_app]/RecommenderService.jws?method=recommend&userID=1&howMany=10)。

13.5.3 简单分布式系统下基于产品的推荐系统简介

传统的推荐引擎算法多在单机上实现，它们只能处理一定量的数据。如果数据量达到一定的规模，传统的推荐引擎算法就会出现各种问题。

在传统的推荐算法中，算法会将用户喜欢的产品抽象成三个具体的数值：用户编号、产品编号和喜爱值。这里的喜爱值表示用户对产品的喜爱程度，它可以用一个具体数值来表示。例如，可以使用 1 到 5 来表示喜欢的程度。1 表示非常不喜欢；2 表示不喜欢；3 表示没有任何感觉；4 表示喜欢；5 表示非常喜欢。也可以从 1 到 5 都表示喜欢，数值越大代表越喜欢。然后通过计算产品之间的相似性来向用户推荐产品。

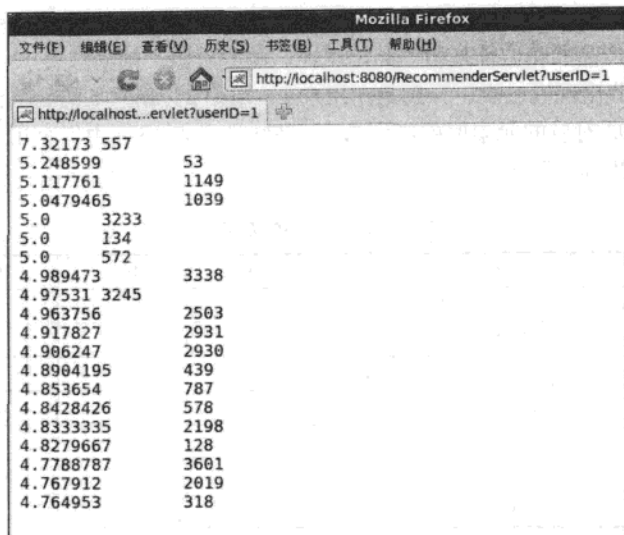


图 13-7 Taste Demo 运行结果界面

分布式系统没有使用这种方法。分布式系统下的推荐算法主要包括以下几部分：

- 计算表示产品相似性的矩阵。
- 计算表示用户喜好的向量。
- 计算矩阵与向量的乘积，为用户推荐产品。

在开始介绍推荐算法之前，首先建立一组数据，如表 13-9 所示。在这组数据中，每条记录包含三个信息：用户编号、产品编号、用户对产品的喜爱值。

表 13-9 用户购买历史表

用户编号	产品编号	喜爱值	用户编号	产品编号	喜爱值
1	101	5.0	4	101	5.0
1	102	3.0	4	103	3.0
1	103	2.5	4	104	4.5
2	101	2.0	4	106	4.0
2	102	2.5	5	101	4.0
2	103	5.0	5	102	3.0
2	104	2.0	5	103	2.0
3	101	2.5	5	104	4.0
3	104	4.0	5	105	3.5
3	105	4.5	5	106	4.0
3	107	5.0			

表 13-9 显示了 5 名顾客的购买历史。下面来介绍一种方法：使用共生矩阵来表示产品的相似性。在这里，产品的相似性是指产品出现在一起的次数。例如从表 13-9 中可以看出产品 101 和产品 102 一共出现过三次，分别是在用户 1、用户 2、用户 5 的物品清单上。那么在共生矩阵中 101 和 102 对应的元素值就应该为 3。在统计了表 13-9 中 5 个用户的购物清单后可以使用表 13-10 的矩阵来表示。

表 13-10 共生矩阵

	101	102	103	104	105	106	107
101	N/A	3	4	4	2	2	1
102	3	N/A	3	2	1	1	0
103	4	3	N/A	3	1	2	0
104	4	2	3	N/A	2	2	1
105	2	1	1	2	N/A	1	1
106	2	1	2	2	1	N/A	0
107	1	0	0	1	1	0	N/A

表 13-10 中的行和列都是产品的编号。观察可知，该矩阵是一个对称矩阵。在计算过程中可以使用一些特殊技术来对矩阵进行处理，使程序的效率更高。原因是产品 104 和产品 105 出现的次数与产品 105 和产品 104 出现的次数必然是相同的。在共生矩阵中对角线的元素是没有意义的。计算时可以使用 0 进行代替。

除了共生矩阵外，还需要一个表示用户喜好的向量。在该向量中，对于用户购买过的产品必然会有一个表示喜好的数值，对于用户没有购买的产品，选择用数字 0 来表示该用户对该产品没有任何喜好。例如对于用户 4 而言，他的向量就应该是 (5.0,0,3.0,4.5,0,4.0,0)。通过计算可以得到所有用户的喜爱值，如表 13-11 所示。

表 13-11 用户喜爱值表

	1	2	3	4	5
101	5.0	2.0	2.5	5.0	4.0
102	3.0	2.5	0	0	3.0
103	2.5	5.0	0	3.0	2.0
104	0	2.0	4.0	4.5	4.0
105	0	0	4.5	0	3.5
106	0	0	0	4.0	4.0
107	0	0	5.0	0	0

其实该表也是一个矩阵：矩阵的行值是产品编号、列值是用用户编号、行列对应的元素值为用户对产品的喜爱值。观察可以发现，矩阵中包含很多 0，可以称之为稀疏矩阵。对于稀疏矩阵，同样可以采用一些技术手段使程序效率更高。

既然已表示了产品的相似性，也表示了用户对产品的喜爱，剩下的就是如何计算推荐的产品了。其实这很简单，只要将共生矩阵与用户的列向量相乘得到一个新的列向量即可。在新的列向量中，所有可推荐产品中哪个产品对应的值最大，就是计算得到的推荐产品。

以向用户 4 推荐产品为例，如表 13-12 所示。

表 13-12 用户 4 的推荐结果

	101	102	103	104	105	106	107	4	推荐结果
101	0	3	4	4	2	2	1	5.0	38
102	3	0	3	2	1	1	0	0	37
103	4	3	0	3	1	2	0	3.0	41.5
104	4	2	3	0	2	2	1	4.5	37
105	2	1	1	2	0	1	1	0	26
106	2	1	2	2	1	0	0	4.0	25
107	1	0	0	1	1	0	0	0	9.5

从结果可以看出，用户最喜欢产品 103，但是用户已经买过 103。因此无须推荐该产品。同理 101、104、106 也都可以不推荐。在可推荐产品 102、105、107 中，选择推荐产品 102。因为 102 的计算结果是三者之中最大的。

推荐结果已经有了，下面来分析这个结果是否合理。在所有可以推荐的产品中，推荐引擎选择了计算结果最大的产品。为什么计算结果最大的产品就是最合理的推荐产品呢？

回想整个计算过程。在计算结果中处在第 2 行的计算结果 37 是矩阵第 2 行元素和用户 4 的列向量乘积。 $3 \times (5.0) + 0 \times 0 + 3 \times 3 + 2 \times (4.5) + 1 \times 0 + 1 \times (4.0) + 0 \times 0 = 37$ 。矩阵中的第 2 行表示的是所有产品和产品 102 同时出现的次数。如果用户对某个产品非常喜欢，而这个产品又和产品 102 同时出现的次数很多，那么乘积对计算结果的影响就会较大。这刚好就是推荐引擎要达到的目的，用户非常喜欢的产品和 102 很相似，推荐引擎可向用户推荐该产品。

对于大量数据，计算结果会非常大。但是没有关系，推荐引擎关注的是所有结果的大小关系，而不是具体的数值。因为最终向用户推荐的是可以推荐产品中计算结果最大的。在计算的过程中，对于不是最大的计算结果以及用户已经购买过的产品，推荐引擎无须推荐，因此也不必计算它们的结果。

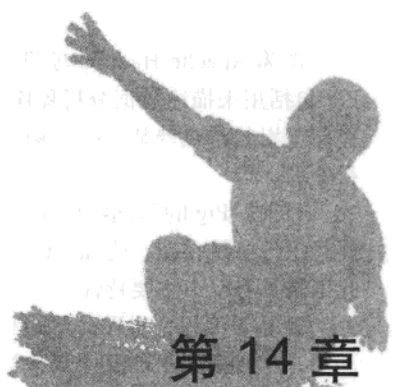
通过分析可知，推荐引擎计算出的推荐结果是合理的。但为什么它适合大规模的数据呢？下面来说明这个问题。在计算共生矩阵的时候，每次只须考虑一个向量；在计算用户向量的时候只须考虑该用户的喜好；在计算推荐结果的时候只须考虑矩阵中的一列值。这都表明，这个方法可以使用 MapReduce 编程模式。

13.6 小结

本章对 Mahout 做了简要介绍，主要有 Mahout 的详细安装过程，MahoutAPI 的介

绍，Mahout 中已经实现的聚类算法，并着重对 K-Means 聚类算法做了介绍。其中还涉及了聚类算法中的数据表示。另外，还对贝叶斯分类算法进行了简单介绍。在推荐引擎部分，着重从思想上介绍了如何在 Hadoop 云平台下实现分布式的推荐系统。Mahout 虽经过了几年的发展，但还是有很多地方值得去探索。如果读者有兴趣加入其中，可以访问 Mahout 的官网，与世界各地的开发者共同推动 Mahout 的发展。





第 14 章

Pig 详解

本章内容

- ☐ Pig 简介
- ☐ Pig 的安装和配置
- ☐ Pig Latin 语言
- ☐ 用户定义函数
- ☐ Pig 实例
- ☐ Pig 进阶
- ☐ 小结



14.1 Pig 简介

作为 Apache Hadoop 项目的子项目之一，Pig 提供了一个支持大规模数据分析的平台。Pig 包括用来描述数据分析程序的高级程序语言，以及对这些程序进行评估的基础结构。Pig 程序突出的特点就是它的结构经得起大量并行任务的检验，这使得它能够对大规模数据集进行处理。

目前，Pig 的基础结构层包括一个产生 MapReduce 程序的编译器。在编译器中，大规模并行执行已经存在（比如，Hadoop 子项目）。Pig 的语言包括一个叫做 Pig Latin 的文本语言，此语言具有以下主要特性：

- 易于编程。实现简单的和高度并行的数据分析任务非常容易。由相互关联的数据转换实例所组成的复杂任务被明确地编码为数据流，这使它们的编写更加容易，同时也更容易理解和维护。
- 自动优化。任务编码的方式允许系统自动去优化执行过程，从而使用户能够专注于语义，而非效率。
- 可扩展性。用户可以轻松编写自己的函数来进行特殊用途的处理。

14.2 Pig 的安装和配置

14.2.1 Pig 的安装条件

1. Hadoop 0.20.2

Pig 有两种运行模式：Local 模式和 MapReduce 模式。如果需要在分布式环境下运行，则需要安装 Hadoop，否则用户可以选择不安装。另外，当前 Hadoop 最新的版本为 0.20.2，当然用户也可以选择安装其他版本，不过这里建议安装最新的 Hadoop 版本。因为新的版本修正了以前版本中的一些错误，并且添加了新的特性^①。

2. Java 1.6

建议安装 Java 1.6 以上的版本。Java 环境对于 Pig 来说是必须的（推荐从 SUN 官方网站上下载）。

当下载安装完毕 Java 后，我们还需要对 Java 环境变量进行设置，将 JAVA_HOME 指向 Java 的安装位置。

如果用户使用的是 Linux 操作系统，那么以上条件就足够了。如果用户使用的是 Windows 操作系统，那么除此之外，用户还需要安装 Cygwin 和 Perl 包。本章后面的案例将以 Linux 操作系统为例进行讲解。

① 关于 Hadoop 的具体信息见其他的相关章节。

14.2.2 Pig 的下载、安装和配置

当前 Pig 最新版本为 0.7.0, 除此之外, Pig 还有 0.6.0 及 0.5.0 两个版本, 用户可以根据需要从 Apache 官方网站上下载相应的版本。本书使用最新版的 Pig 0.7.0, 安装包下载地址如下: <http://www.apache.org/dyn/closer.cgi/hadoop/pig>。

Pig 的安装包下载完成后, 需要使用 `tar -xvf pig-*.tar.gz` 命令将其解压。我们可以将 Pig 放在系统中的任意位置上, 并且只需要配置相应环境变量就可以使用 Pig 了。不过, 建议将 Pig 放在 Hadoop 目录下, 方便以后的操作。

解压完成后, 需要设置 Pig 相应环境变量, 环境变量有多种设置方法, 用户可以根据自己的需要进行选择。这里我们选择对 profile 文件进行修改, 来设置 Pig 相应环境变量。打开 “/etc/profile” 文件, 插入下面这一条语句, 保存关闭文件后需要重启系统以使环境变量设置生效:

```
export PATH=/<my-path-to-pig>/pig-n.n.n/bin:$PATH
```

当环境变量设置生效后, 我们可以通过 “\$pig -help” 命令来查看 Pig 是否安装成功。当 Pig 安装成功后会出现如图 14-1 的提示。



```
root@ubuntu:~# pig -help
Apache Pig version 0.7.0 (r941408)
compiled May 05 2010, 11:15:55

USAGE: Pig [options] [-] : Run interactively in grunt shell
      Pig [options] -e[execute] cmd [cmd ...] : Run cmd(s)
      Pig [options] [-f[file]] file : Run cmds found in file;
options include:
  -c, --log4jconf log4j configuration file, overrides log.conf
  -b, --brief brief logging (no timestamps)
  -c, --cluster clustername, kryptonite is default
  -d, --debug debug level, INFO is default
  -e, --execute commands to execute (within quotes)
  -f, --file file path to the script to execute
  -h, --help display this message
  -i, --version display version information
  -j, --jar jarfile load jarfile
  -l, --logfile logfile path to client side log file; current working directory
      is default
  -m, --param file path to the parameter file
  -p, --param key value pair of the form param=eval
  -r, --dryrun
  -t, --optimizer off optimizer rule name, turn optimizer off for this
      rule; use all to turn all rules off, optimizer is turned on by default
  -v, --verbose print all error messages to screen
  -w, --warning turn warning on; also turns warning aggregation off
  -x, --executype local|mapreduce, mapreduce is default
  -F, --stop on failure aborts execution on the first failed job; off b
      y default
  -M, --no multiquery turn multiquery optimization off; Multiquery is o
      n by default
root@ubuntu:~#
```

图 14-1 \$Pig -help

14.2.3 Pig 运行模式

Pig 有两种运行模式: Local 模式和 MapReduce 模式。当 Pig 在 Local 模式下运行时,

Pig 将只访问本地一台主机；当 Pig 在 MapReduce 模式下运行时，它将访问一个 Hadoop 集群和 HDFS 的安装位置。这时，Pig 将自动地对这个集群进行分配和回收。因为 Pig 系统可以自动对 MapReduce 程序进行优化，所以当用户使用 Pig Latin 语言进行编程的时候，不必关心程序运行的效率，Pig 系统将会自动对程序进行优化。这样能够大量节省用户编程的时间。

下面我们首先介绍 Pig 在 Local 模式下的运行方式。Pig 的 Local 模式适用于用户对测试程序进行调试，因为 Local 模式下的 Pig 将只访问本地一台主机，它可以在短时间内处理少量的数据，并且用户不必关心 Hadoop 系统对整个集群的控制，这样既能让用户使用 Pig 的功能又不至于在对集群的管理上花费太多时间。

Pig 的 Local 模式和 MapReduce 模式都有三种运行方式，分别为：Grunt Shell 方式、脚本文件方式和嵌入式程序方式。下面我们将对其进行一一介绍。

1. Local 模式

(1) Grunt Shell 方式

用户使用 Grunt Shell 方式时，首先需要使用命令开启 Pig 的 Grunt Shell，只需在 Linux 终端中输入如下命令并执行即可：

```
$pig -x local
```

这样 Pig 将进入 Grunt Shell 的 Local 模式，如果直接输入“\$pig”命令，Pig 将首先检测 Pig 的环境变量设置，然后进入相应的模式。如果没有设置 MapReduce 环境变量，Pig 将直接进入 Local 模式。图 14-2 为开启 Grunt Shell 的结果。



图 14-2 Local 模式下开启 Grunt Shell

Grunt Shell 和 Windows 中的 Dos 窗口非常类似，在这里用户可以一条一条地输入命令并对数据进行操作。

(2) 脚本文件方式

使用脚本文件作为批处理作业来运行 Pig 命令时，它实际上就是第一种运行方式中的命令集合，使用如下命令可以在本地模式下运行 Pig 脚本：

```
$pig -x local script.pig
```

其中，“script.pig”对应的是 Pig 脚本，用户在这里需要正确指定 Pig 脚本的位置，否则，系统将不能识别。例如，Pig 脚本放在“/root/pigTmp”目录下，那么这里就要写成“/root/pigTmp/script.pig”。用户在使用的时候需要注意 Pig 给出的一些提示，充分利用这些提示能够帮助用户更好地使用 Pig 进行相关的操作^①。

^① 注意：这里 script.pig 前后没有引号。

(3) 嵌入式程序方式

我们可以把 Pig 命令嵌入到主机语言中，并且运行这个嵌入式程序。和运行普通的 Java 程序相同，这里需要书写特定的 Java 程序，并且将其编译生成对应的 class 文件或 package 包，然后再调用 main 函数运行程序。

用户可以使用下面的命令对 Java 源文件进行编译：

```
$javac -cp pig-*. *-core.jar local.java
```

这里“pig-*. *-core.jar”放在 Pig 安装目录下，“local.java”为用户编写的 Java 源文件，并且“pig-*. *-core.jar”和“local.java”需要用户正确地指定相应的位置。例如，我们的“pig-*. *-core.jar”文件放在“/root/hadoop-0.20.2/”目录下，“local.java”文件放在“/root/pigTmp”目录下，所以这一条命令我们应该写成：

```
$javac -cp /root/hadoop-0.20.2/pig-0.20.2-core.jar /root/pigTmp/local.java
```

当编译完成后，Java 会生成“local.class”文件，然后用户可以通过如下命令调用执行此文件。

```
$ java -cp pig-*. *-core.jar:. local
```

2. MapReduce 模式

Pig 需要把真正的查询转换成相应的 MapReduce 作业，并提交到 Hadoop 集群去运行（集群可以是真实的分布，也可以是伪分布）。要想 Pig 能识别 Hadoop，用户需要告诉 Pig 关于 Hadoop 的版本及一些关键的信息（也就是 NameNode 和 JobTracker 的位置及端口信息）。

下面第一步首先指明 Pig 要连接的 Hadoop 版本信息，第二步详细指明 Pig 连接 Hadoop 的配置信息。

❑ 允许 Pig 连接到任何的 Hadoop.0.20.* 版本。

配置 Linux 系统的环境变量，在“/etc/profile”文件中加入如下信息：

```
export PIG_HADOOP_VERSION=20
```

❑ 指明集群中 NameNode 和 JobTracker 的位置。有以下两种方法让 Pig 识别 Hadoop 的 NameNode 和 JobTracker，采用任何一种方式均可。

方法一：

把自己的 Hadoop 的 Conf 地址添加到 Pig 的 Classpath 上：

```
export PIG_CLASSPATH=$HADOOP_INSTALL/conf/
```

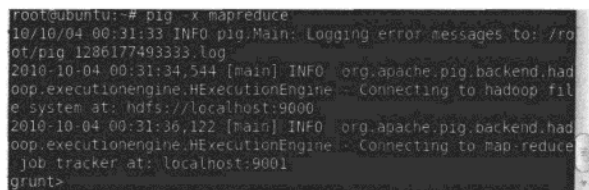
方法二：

在 pig 目录的 Conf 文件夹（可能需要自己创建）里创建一个 pig.properties 文件，然后在里面添加集群的 NameNode 和 JobTracker 信息，第二行中的 port 是用户的 Hadoop 中 JobTracker 对应的端口：

```
fs.default.name=hdfs://localhost/
mapred.job.tracker=localhost:port
```

待设置完毕并且生效之后，用户可以输入“\$pig -x mapreduce”命令进行测试，如果能够看到 Pig 连接 Hadoop 的 NameNode 和 JobTracker 的相关信息，则表明配置成功，然后用户就可以随心所欲地使用 MapReduce 模式来进行相关的 Pig 操作了。

图 14-3 为 MapReduce 配置成功后的提示信息，从图中可以看到 Pig 连接 Hadoop 的详细信息。



```

root@ubuntu:~# pig -x mapreduce
10/10/04 00:31:33 INFO pig.Main: Logging error messages to: /root/.pig/1286177493333.log
2010-10-04 00:31:34,544 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: hdfs://localhost:9000
2010-10-04 00:31:36,122 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to Map-Reduce job tracker at: localhost:9001
grunt>

```

图 14-3 MapReduce 配置成功的提示信息

配置成功之后，下面我们将针对 Pig 的 MapReduce 模式，说明如何在此模式下对 Grunt Shell 方式、脚本文件方式和嵌入式程序方式进行操作。它们和 Local 模式下的操作几乎完全相同，只不过需要将相应的参数指明为 MapReduce 模式。

(1) Grunt Shell 方式

用户在 Linux 终端下输入如下命令进入 Grunt Shell 的 MapReduce 模式：

```
$pig -x mapreduce
```

(2) 脚本文件方式

用户可以使用如下命令在 MapReduce 模式下运行 Pig 脚本文件：

```
$pig -x mapreduce script.pig
```

(3) 嵌入式程序方式

和 Local 模式相同，在 MapReduce 模式下运行嵌入式程序同样需要经过编译和执行两个步骤。用户可以使用如下两条命令，完成相应的操作。

```

javac -cp pig-0.7.0-core.jar mapreduce.java
java -cp pig-0.7.0-core.jar:. mapreduce

```

至此，Pig 系统的两个运行模式及其分别对应的三个运行方式就讲述完毕了，14.5 节和 14.6 节我们将结合实例对其做更深入的介绍，这里希望用户能够对 Pig 系统的运行模式有一个初步的印象。

14.3 Pig Latin 语言

14.3.1 Pig Latin 语言简介

Pig Latin 语言和传统的关系数据库中的数据库操作语言非常类似。但是 Pig Latin 语言更

侧重于对数据的查询和分析，而不是对数据进行修改和删除等操作。另外，由于 Pig Latin 可以在 Hadoop 的分布式云平台上运行，它的这个特点可以让其具有其他数据库所无法比拟的速度优势，它能在短时间内处理海量的数据。例如，处理系统日志文件、处理大型数据库文件、处理特定 Web 数据等。除此之外，我们在使用 Pig Latin 语言编写程序的时候，不必关心如何让程序更好地在 Hadoop 云平台上运行，因为这些任务都是由 Pig 系统自己分配的，不需要程序员参与。因此，程序员只需要专注于程序的编写即可，这样大大减轻了程序员的负担。

Pig Latin 是这样一個操作：通过对关系（relation）进行处理产生另外一组关系^①。Pig Latin 语言在书写一条语句的时候能够跨越多行，但是必须以半角的分号来结束。Pig Latin 语句通常按照下面的格式来编写：

- 1) 通过一条 LOAD 语句从文件系统中读取数据；
- 2) 通过一系列“转换”语句对数据进行处理；
- 3) 通过一条 STORE 语句把处理结果输出到文件系统中，或者使用一条 DUMP 语句把处理结果输出到屏幕上。

LOAD 和 STORE 语句有严格的语法规定，用户很容易就能掌握，关键是如何灵活使用“转换”语句对数据进行处理。

Pig Latin 语言还可以对数据进行连接操作，在 14.6 节中，我们将通过一组例子，让用户对 Pig Latin 语言的特点有更好的体会。

14.3.2 Pig Latin 的使用

这一节我们着重讲一下 Pig Latin 可以用在哪些方面。

1. 运行 Pig Latin

用户可以通过多种方式使用 Pig Latin 语句，如 14.2.3 节所述。通常，Pig 按如下方式执行 Pig Latin 语句：

- Pig 对所有语句的语法和语义进行确认；
 - 如果遇到 DUMP 或 STORE 命令，Pig 将顺序执行上面所有的语句。
- 在下面的示例中，Pig 将只确认 LOAD 和 FOREACH 语句，但不执行：

```
A = LOAD 'Student' USING PigStorage(':') AS (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray);
B = FOREACH A GENERATE Sname;
```

在下面的示例中，Pig 将确认并执行 LOAD、FOREACH 和 DUMP 语句：

```
A = LOAD 'Student' USING PigStorage(':') AS (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray);
B = FOREACH A GENERATE Sname;
DUMP B;
```

^① 这个定义适用于除LOAD和STORE之外的所有操作，LOAD和STORE分别执行从文件系统读取和写入操作。

因为 Pig 的一些命令并不会自动执行，而是需要通过其他命令来触发，也就是说如果用户连续地使用某些命令，它并不会马上执行，而是在最后一个触发操作的调用下，连续地一次性地执行完毕。

2. 查看 Pig Latin 的运行结果

Pig Latin 包括一些查看语句运行结果的操作。

❑ 使用 DUMP 操作把操作的结果显示在屏幕上，如下所示：

```
DUMP alias
```

❑ 使用 STORE 操作把操作的结果存储在文件中，如下所示：

```
STORE alias INTO 'directory' [USING function];
```

3. Pig Latin 的调试

Pig Latin 包括一些可以帮助用户进行调试的操作。

❑ 使用 DESCRIBE 操作来查看关系的模式，如下所示：

```
DESCRIBE alias;
```

❑ 使用 EXPLAIN 操作来查看对某个关系进行操作的逻辑的、物理的或 MapReduce 的执行计划，如下所示：

```
EXPLAIN [-script pignore] [-out path] [-brief] [-dot] [-param param_name =  
param_value] [-param_file file_name] alias;
```

❑ 使用 ILLUSTRATE 操作来对 Pig Latin 语句进行单步执行操作，如下所示：

```
ILLUSTRATE alias;
```

4. 注释在 Pig Latin 脚本中的使用

注释就是对代码的解释和说明。目的是为了让人和自己很容易看懂。像其他的编程语言一样，Pig Latin 脚本中也可以包含注释，下面是两种常用的注释格式：

❑ 多行注释：/*...*/

示例：

```
/*  
myscript.pig  
My script includes three simple Pig Latin Statements.  
*/
```

❑ 单行注释：--

示例：

```
A = LOAD 'Student' USING PigStorage(':'); -- 语句  
B = FOREACH A GENERATE Sname; -- foreach 语句  
DUMP B; --dump 语句
```

5. 大小写相关性

在 Pig Latin 中, 关系名、域名、函数名是区分大小写的。参数名和所有 Pig Latin 关键字是不区分大小写的。

请注意下面的示例:

- ❑ 关系名 A、B、C 等是区分大小写的;
- ❑ 域名 f1、f2、f3 等是区分大小写的;
- ❑ 函数名 PigStorage、COUNT 等是区分大小写的;
- ❑ 关键字 LOAD、USING、AS、GROUP、BY、FOREACH、GENERATE、DUMP 等是不区分大小写的, 它们也能被写成 load、using、as、group、by、foreach、generate、dump 等。

在 FOREACH 语句中, 关系 B 中的域通过位置来访问, 如下所示:

```
grunt> A = LOAD 'data' USING PigStorage() AS (f1:int, f2:int, f3:int);
grunt> B = GROUP A BY f1;
grunt> C = FOREACH B GENERATE COUNT ($0);
grunt> DUMP C;
```

14.3.3 Pig Latin 的数据类型

1. 数据模式

Pig Latin 中数据的组织形式包括: 关系 (relation)、包 (bag)、元组 (tuple) 和域 (field)。

一个关系可以按如下方式定义:

- ❑ 一个关系就是一个包 (更具体地说, 是一个外部包);
- ❑ 包是元组的集合;
- ❑ 元组是域的有序集合;
- ❑ 域是一个数据块。

一个 Pig 关系是一个由元组组成的包, Pig 中的关系和关系数据库中的表 (table) 很相似, 包中的元组相当于表中的行。但是和关系表不同的是, Pig 中不需要每一个元组包含相同数目的域, 或者相同位置的域 (同列域), 也不需要具有相同的数据类型。

另外, 关系是无序的, 这就意味着 Pig 不能保证元组按特定的顺序来执行。

2. 数据类型

表 14-1 给出了一些简单数据类型的描述及示例。限于篇幅我们不再做更详细的介绍, 具体内容用户可以在使用中慢慢体会。

表 14-1 Pig Latin 数据类型

简单数据类型	描 述	示 例
标量		
int	有符号 32 位整形	10
long	有符号 64 位整形	数据:10L 或 10l 显示: 10L
float	32 为浮点型	数据:10.5F 或 10.5f 或 10.5e2f 或 10.5E2F 显示:10.5F or 1050.0F
double	64 位浮点型	数据: 10.5 或 10.5e2 或 10.5E2 显示: 10.5 or 1050.0
数组		
chararray	字符数组使用 UTF-8 格式进行编码	hello world
bytearray	字节数组 (blob)	
复杂数据类型		
tuple	有序的字段集	(19,2)
bag	元组集合	{(19,2) (19,2), (18,1)}
map	键值对集合	[open#apache]

14.3.4 Pig Latin 关键字

Pig Latin 语言有很多关键字，但是我们不可能一一给大家介绍。在下面的第一部分内容中，给大家介绍了 Pig Latin 语言都包含哪些关键字。然后在第二部分，会就其中主要的关键字做详细介绍。

1. Pig Latin 关键字

表 14-2 给出了一些与首字母相对应的关键字。

表 14-2 Pig Latin 关键字

首字母	对应关键字
-- A	and, any, all, arrange, as, asc, AVG
-- B	bag, BinStorage, by, bytearray
-- C	cache, cat, cd, chararray, cogroup, CONCAT, copyFromLocal, copyToLocal, COUNT, cp, cross
-- D	%declare, %default, define, desc, describe, DIFF, distinct, double, du, dump
-- E	e, E, eval, exec, explain
-- F	f, F, filter, flatten, float, foreach, full
-- G	generate, group
-- H	help
-- I	if, illustrate, inner, input, int, into, is
-- J	join
-- K	kill

(续)

首字母	对应关键字
-- L	l, L, left, limit, load, long, ls
-- M	map, matches, MAX, MIN, mkdir, mv
-- N	not, null
-- O	or, order, outer, output
-- P	parallel, pig, PigDump, PigStorage, pwd
-- Q	quit
-- R	register, right, rm, rmf, run
-- S	sample, set, ship, SIZE, split, stderr, stdin, stdout, store, stream, SUM
-- T	TextLoader, TOKENIZE, through, tuple
-- U	union, using
-- V, W, X, Y, Z	
-- 符号	= != < > <= >= + - * / % ? \$. # :: () [] {}

2. 常用关键字

在 Pig Latin 常用的关键字中，我们将其分为四类：关系运算符、诊断运算符、Load/Store 函数和文件命令。下面会进行详细介绍。

(1) 关系运算符

□ Load

它的作用是从文件系统中加载数据，语法如下：

```
LOAD 'data' [USING function] [AS schema]
```

在这里“data”表示文件或目录的名字，并且要用单引号括起来。如果用户指定一个目录的名字，目录中所有的文件将被加载。中括号中的内容为可选项（如果没有特殊指明，“[]”都表示可选项），用户可只在需要的时候指明，可以省略。这里使用 schema 来指定加载数据的数据类型，如果数据类型与模式中指定的数据类型不符，那么系统将会产生一个 null，甚至会报错。

下面是几个 Load 操作的例子。

1) 不使用任何方式：

```
A = LOAD 'myfile.txt';
```

2) 使用加载函数：

```
A = LOAD 'myfile.txt' USING PigStorage('\t');
```

3) 指定模式：

```
A = LOAD 'myfile.txt' AS (f1:int, f2:int, f3:int);
```

4) 加载函数和模式均使用:

```
A = LOAD 'myfile.txt' USING PigStorage('\t') AS (f1:int, f2:int, f3:int);
```

❑ Store

它的作用是将结果保存到文件系统中, 语法如下所示:

```
STORE alias INTO 'directory' [USING function];
```

这里的“alias”是用户要存储的结果(关系)的名称, INTO 为不可省略的关键字, directory 为用户指定的存储目录的名字, 需要用单引号括起来。另外, 如果此目录已经存在, 那么 Store 操作将会失败, 输出文件将会被系统命名成 part-nnnnn 的格式。

❑ Foreach

它的作用是基于数据的列进行数据转换, 语法如下:

```
alias = FOREACH { gen_blk | nested_gen_blk } [AS schema];
```

通常我们使用“FOREACH ...GENERATE”组合来对数据列进行操作, 下面是两个简单的例子。

1) 如果一个关系 A (outer bag), FOREACH 语句可以按下面的方式使用:

```
X = FOREACH A GENERATE f1;
```

2) 如果 A 是一个 inner bag, FOREACH 语句可以按下面的方式使用:

```
X = FOREACH B {
    S = FILTER A BY 'xyz';
    GENERATE COUNT (S.$0);
}
```

对于初级用户来说, 仅需要掌握第一种操作方式, 关于 Foreach 关键字更多的内容我们将在后面进行详细的讨论。

(2) 诊断运算符

❑ Dump

它的作用是将结果显示到屏幕上, 语法如下:

```
DUMP alias
```

这里的“alias”为被操作关系的名字。

使用 DUMP 操作符来执行 Pig Latin 语句, 并且把结果输出到屏幕上。使用 DUMP 意味着使用交互式模式, 也就是说, 语句被马上执行, 但结果并没有被保存。用户可以使用 DUMP 作为一个调试设备, 用来检查用户期望的数据是否被生成了。另外用户应该有选择性地使用 DUMP, 因为它将会使多值查询优化无效, 并且可能会减慢执行。

示例:

```
A = LOAD 'student' AS (name:chararray, age:int, gpa:float);
```



```
DUMP A;
```

这里 Pig 将会把 A 中所有的数据输出到屏幕上。

❑ Describe

它的作用是返回一个名称的模式，语法如下：

```
DESCRIBE alias;
```

使用 DESCRIBE 操作符来查看指定名称的模式。

在这个例子中，使用 AS 子句来指定一个模式，如果所有的数据都符合这个模式，Pig 将使用已分配的类型。然后我们使用 DESCRIBE 操作符来查看它们的模式。

```
A = LOAD 'student' AS (name:chararray, age:int, gpa:float);
B = FILTER A BY name matches 'J.+';
C = GROUP B BY name;
D = FOREACH B GENERATE COUNT(B.age);
```

```
DESCRIBE A;
A: {group, B: {name: chararray,age: int,gpa: float}}
DESCRIBE B;
B: {group, B: {name: chararray,age: int,gpa: float}}
DESCRIBE C;
C: {group, chararray,B: {name: chararray,age: int,gpa: float}}
DESCRIBE D;
D: {long}
```

(3) Load/Store 函数：

❑ PigStorage

它的作用是加载、存储 UTF-8 格式的数据，语法如下：

```
PigStorage(field_delimiter)
```

Field_delimiter 为 PigStorage 函数的参数，用来指定函数的字段定界符。PigStorage 函数默认的字段定界符为：tab ('\t')，用户也可以指定其他字段定界符，但定界符要在单引号中指明。

PigStorage 是 LOAD 和 STORE 操作符默认的加载函数，而且能够处理简单的和复杂的数据类型。

PigStorage 对有结构的文本进行读取，并采用 UTF-8 编码进行存储。

在 Load 语句中，PigStorage 希望数据使用域定界符进行格式化。默认情况下为字符 ('\t')，用户也可以指定其他的字符。

在 Store 语句中，PigStorage 同样使用域定界符来输出数据。它的操作方法和 Load 语句相同，另外 Store 语句的记录定界符使用 ('\n')。

Load 或 Store 语句默认的域定界符均为 tab ('\t')。用户仍可以使用其他字符作为字段定界符。但是像 ^A 或 Ctrl-A 等字符应使用 UTF-16 编码格式进行编码。

在 Load 语句中, Pig 注明记录定界符为: 换行符 ('\\n')、回车返回符 ('\\r') 或 CTRL-M, 以及组合的 CR+LF 字符 ('\\r\\n')[⊖]。在 Store 语句中, Pig 使用换行符作为记录定界符。

以下提供一个示例。

在这个例子中 PigStorage 使用 tab 作为域定界符, 换行符为记录定界符, 并且下面的两条语句是等价的:

```
A = LOAD 'student' USING PigStorage('\\t') AS (name: chararray, age:int, gpa: float);
A = LOAD 'student' AS (name: chararray, age:int, gpa: float);
```

在这个例子中, PigStorage 将 X 的内容存储到文件中, 并且使用星号作为域定界符。STORE 函数将结果存储在 output 目录中。

```
STORE X INTO 'output' USING PigStorage('*');
```

(4) 文件命令

❑ cd

它的作用是将当前目录修改为其他目录, 语法如下:

```
cd [dir]
```

此处的 cd 命令和 Linux 的 cd 命令非常相似, 能够用来对文件系统进行定位。如果用户指定了一个目录, 那么这个目录将成为用户当前的工作目录, 并且用户所有其他的操作将相对于这个目录来进行。如果没有指定任何目录, 用户的根目录将成为当前的工作目录。

❑ copyFromLocal

它的作用是从本地文件系统拷贝文件或目录到 HDFS 中, 语法如下:

```
copyFromLocal src_path dst_path
```

其中, src_path 为本地系统中的文件或目录的路径, dst_path 为 HDFS 系统中的路径。

CopyFromLocal 命令让用户能够从本地文件系统中拷贝文件或目录到 Hadoop 的分布式文件系统中。

❑ ls

它的作用是显示一个目录中的内容, 语法如下:

```
ls [path]
```

此处的 ls 命令和 Linux 中的 ls 命令相似, 如果指定一个目录, 这个命令将列出被指定目录中的内容。如果不指定参数, 那么系统将列出当前工作目录中的内容。

❑ rm

它的作用是移除一个或更多的文件或目录, 语法如下:

```
rm path [path...]
```

⊖ 一定不要将这些字符用作域定界符。

此处的 rm 命令和 Linux 中的 rm 命令相似，用户可移除一个或多个文件及目录。

14.4 用户定义函数

用户可以使用 UDFs (User Defined Functions, 用户定义函数) 来编写特定的处理函数，这大大地增强了 Pig Latin 语言的功能，用户可以方便地对其功能进行扩充和完善。Pig 为用户定义函数提供了大量的支持，UDFs 几乎可以作为 Pig 所有操作符的一部分来使用。

下面我们将通过一个实例来帮助用户学习如何编写 UDFs，以及如何让 Pig 使用用户编写的 UDFs。

这里给出一个学生表 (学号, 姓名, 性别, 年龄, 所在系), 其中含有如下几条记录:

201000101: 李勇 :Boy:20: 计算机软件与理论

201000102: 王丽 :Girl:19: 计算机软件与理论

201000103: 刘花 :Girl:18: 计算机应用技术

201000104: 李肖 :Boy:19: 计算机系统结构

201000105: 吴达 :Boy:19: 计算机系统结构

201000106: 滑可 :Boy:19: 计算机系统结构

它们所对应的数据类型如下所示:

```
Student (Sno:chararray, Sname:chararray, Ssex:chararray, Sage:int, Sdept:chararray)
```

这里字段与字段之间通过冒号 (半角英文标点) 隔开, 下面我们将编写一个函数, 以将所有的小写字母转换成对应的大写字母。

14.4.1 编写用户定义函数

下面是我们编写的 UDFs 代码, 如代码清单 14-1 所示。

代码清单 14-1 UDFs 代码

```
1 package myudfs;
2 import java.io.IOException;
3 import org.apache.pig.EvalFunc;
4 import org.apache.pig.data.Tuple;
5 import org.apache.pig.impl.util.WrappedIOException;
6
7 public class UPPER extends EvalFunc <String>
8 {
9     public String exec(Tuple input) throws IOException {
10         if (input == null || input.size() == 0)
11             return null;
12         try{
13             String str = (String)input.get(0);
14             return str.toUpperCase();
15         }
```

```

15         }catch(Exception e){
16             throw WrappedIOException.wrap("Caught exception processing input row ", e);
17         }
18     }
19 }

```

代码的第 1 行表明这个函数是 myudfs 包的一部分。这个 UDF 类是 EvalFunc 类的继承，EvalFunc 是所有 eval 函数的基类。在这个例子中，这个类使用返回值类型为 Java String 的参数进行参数化。现在我们需要去实现 EvalFunc 类的 exec 函数。在这里，函数的输入是一个 tuple 集合，它们按照 Pig 脚本加载的顺序依次被调用。每当输入一个 tuple 时，UDF 将被调用一次。在我们的例子中，它是一个与学生的性别相一致的字符串域。

我们首先需要做的是处理无效的数据。这依赖于数据的格式，如果数据为字节数组那就意味着它不需要被转化为其他的数据类型；如果输入的数据为其他类型，那么就需要将数据转换成适当的数据类型；如果输入数据的格式不能被系统识别或转换，NULL 值将被返回。这就是我们例子中的第 16 行抛出一个错误的原因。在这里，WrappedIOException 是一个帮助类，帮助我们真实的异常转换为 IO 异常。

另外，注意第 10 ~ 11 行的作用，其作用是检查输入数据为 null 或空。如果为 null 或空，系统将返回 null。

很容易看出，函数的实现部分在第 13~14 行，它们使用 Java 函数将接收的输入转换为相应的大写。

如果要使用这个函数，它需要被编译并且包含在一个 jar 中。用户需要建立 pig.jar 来编译用户的 UDF。pig.jar 文件需要用户自行下载安装。用户可以使用下面的命令集从 SVN 库中检验代码并且创建 pig.jar 文件：

```

svn co http://svn.apache.org/repos/asf/hadoop/pig/trunk
cd trunk
ant

```

注意 在使用 svn 和 ant 操作之前，要确保系统已经安装了 svn 和 ant^①。

在上述操作完成之后，用户可以在自己当前的工作目录中看到 pig.jar 文件（它位于 trunk 目录下）。

pig.jar 文件创建完成之后，我们首先需要对函数进行编译，然后再创建一个包含这个函数的 jar 文件。具体操作命令如下：

```

cd myudfs
javac -cp pig.jar UPPER.java
cd ..
jar -cf myudfs.jar myudfs

```

① 这部分知识已经超出了本书的内容，具体的操作用户可以参考其他相关书籍。

14.4.2 使用用户定义函数

下面是我们所编写的 pig 脚本，它使用我们之前编写的用户定义函数对上面给出的学生表进行了相应的操作。

```
1 -- myscript.pig
2 REGISTER myudfs.jar;
3 A = LOAD 'Student' using PigStorage(',') as (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray);
4 B = FOREACH A GENERATE myudfs.UPPER(Ssex);
5 DUMP B;
```

我们使用下面的命令执行此脚本文件。其中，使用“-x mapreduce”指定函数运行的模式，如果用户只是为了对函数进行测试，建议用户在 local 模式下运行。因为对于小文件来说，MapReduce 模式的准备时间显得过长，有时候甚至让用户觉得 MapReduce 模式下文件的运行效率比 local 模式下还要低。为了验证函数的通用性，这里我们使用 MapReduce 模式。

```
java -cp pig.jar org.apache.pig.Main -x mapreduce myscript.pig
```

这个脚本的第 2 行提供了 jar 文件的位置，这个 jar 文件中包含我们刚刚编写的用户定义函数（注意：jar 文件上没有引号）。为了找到 jar 文件的位置，Pig 首先检查 classpath 环境变量，如果在 classpath 环境变量中不能找到 jar 文件，Pig 将假定地址为绝对地址或一个相对于 Pig 被调用位置的地址。如果 jar 文件仍旧不能被发现，系统将返回一个错误。

多个用户定义函数可以用在相同的脚本中。如果完全相同且合格的函数出现在多个 jar 中，那么根据 Java 语义，第一个出现的函数将被一直使用。

UDF 的名称和包名必须要完全合格，否则系统将返回一个错误：

```
java.io.IOException: Cannot instantiate:UPPER.
```

另外，函数的名称区分大小写（比如：UPPER 和 upper 是不同的），UDF 也可以包含一个或更多的参数。

待操作完成之后，我们可以在终端上看到 Pig 输出的正确结果：

```
BOY
GIRL
GIRL
BOY
BOY
BOY
```

用户定义函数还包括很多其他的内容，限于篇幅，我们只在这里做简单介绍。

14.5 Pig 实例

下面将结合第 14.2.3 节所介绍的 Pig 运行模式给出相应的例子。这里我们给出一个学生

表（学号，姓名，性别，年龄，所在系），其中含有如下几条记录：

```
201000101: 李勇 : 男 : 20 : 计算机软件与理论
201000102: 王丽 : 女 : 19 : 计算机软件与理论
201000103: 刘花 : 女 : 18 : 计算机应用技术
201000104: 李肖 : 男 : 19 : 计算机系统结构
201000105: 吴达 : 男 : 19 : 计算机系统结构
201000106: 滑可 : 男 : 19 : 计算机系统结构
```

它们所对应的数据类型如下所示：

```
Student (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray)
```

这里字段与字段之间通过冒号（半角英文标点）隔开，下面我们将在不同的运行方式下取出各个学生的姓名和年龄两个字段。

李勇	20
王丽	19
刘花	18
李肖	19
吴达	19
滑可	19

14.5.1 Local 模式

这一节我们将结合上面给出的实例具体讲解如何在 Pig 的 Local 模式下对数据进行操作。同时，我们对 Pig 在 Local 模式下的三种运行方式都进行了详细的介绍。

1. Grunt Shell

通过 14.3.3 节中对 Pig 的数据模式的介绍，我们可以了解到，记录是域的有序集合。因此，在我们对数据进行操作之前，需要按照文件中数据相应的字段和类型来加载数据。通过下面的这一条命令，可以把前面给出的例子按照对应字段和对应的数据类型进行加载：

```
grunt>>A = load '/path/Student' using PigStorage(':') as (Sno:chararray,Sname:
chararray,Ssex:chararray,Sage:int,Sdept:chararray);
```

通过 Foreach 命令，从 A 中选出 Student 相应的字段，并存储到 B 中：

```
grunt>>B = foreach A generate Sname,Sage;
```

通过 dump 命令，将 B 中的内容输出到屏幕上：

```
grunt>>dump B;
```

下面一步将 B 的内容输出到本地文件中：

```
grunt>>store B into '/path /grunt.out';
```

现在我们可以打开 grunt.out 文件来查看我们操作的结果，如下所示：

李勇	20
王丽	19
刘花	18
李肖	19
吴达	19
滑可	19

2. 脚本文件

脚本文件实质上是 pig 命令的批处理文件。

我们给出的 script.pig 文件包含以下内容：

```
A = load '/path/Student' using PigStorage(',') as (Sno:chararray,Sname:chararray,
Ssex:chararray,Sage:int,Sdept:chararray);
B = foreach A generate Sname,Sage;
dump B;
store B into '/path/script.out';
```

用户可以看出，这个文件其实就是上面 Grunt Shell 下命令的一个集合。

我们通过下面 3 种方式来调用和执行这个脚本文件，用户可以看到，生成的 script.out 文件中的内容与 grunt.out 文件中的内容是完全相同的。

(1) Grunt Shell (本地模式)

```
$pig -x local
Grunt>exec myscript.pig;
```

或者

```
Grunt>run myscript.pig;
```

(2) 命令行 (MapReduce 模式)

```
$pig myscript.pig
```

(3) 命令行 (本地模式)

```
$pig -x local myscript.pig
```

3. 嵌入式程序

用户可以方便地使用 Java 语言来书写相应的 Pig 脚本，如代码清单 14-2 所示。

代码清单 14-2 Local 模式下用 Java 编写的 Pig 脚本

```
import java.io.IOException;
import org.apache.pig.PigServer;
public class tst_local{
    public static void main(String[] args) {
        try {
            PigServer pigServer = new PigServer("local");
            runIdQuery(pigServer, "/path/Student");// 调用函数
```

```

    }
    catch(Exception e) {}
}
public static void runIdQuery(PigServer pigServer, String inputFile) throws
    IOException {
    pigServer.registerQuery("A = load '" + inputFile + "' using PigStorage(',') as
        (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray);");
    pigServer.registerQuery("B = foreach A generate Sname,Sage; ");
    pigServer.store("B", "/path/tstJavaLocal.out");
}
}
}

```

下面我们将通过 14.2.3 节中所介绍的在嵌入式方式下运行 pig 脚本的命令来对此文件进行编译、运行。

首先，使用下面的命令对此 Java 源文件进行编译：

```
$javac -cp pig-*. *-core.jar local.java
```

待编译完成后，通过下面的命令运行 “.class” 类文件：

```
$ java -cp pig-*. *-core.jar:. local
```

然后打开生成的结果文件 “tstJavaLocal.out”，我们会发现它和前面两种方式生成的结果是完全相同的。

14.5.2 MapReduce 模式

这一节我们将结合上面给出的实例具体讲解如何在 Pig 的 MapReduce 模式下对数据进行操作。同时，我们也会对 Pig MapReduce 模式下的三种运行方式进行详细介绍。

1. Grunt Shell

在 MapReduce 模式下对 Pig 的使用其实是 Pig Local 模式和 Hadoop 操作的结合。因为运行 MapReduce 程序我们需要在 Hadoop 的 HDFS 文件系统下对文件进行操作，但是在 Linux 系统下我们又看不到 HDFS 文件系统下的文件，所以就不能使用常规的操作来“搬运”文件。这里，我们就需要使用与 HDFS 相关的命令在 HDFS 文件系统下执行 Pig 的命令了。

首先，从终端进入 Pig 的 MapReduce 模式，然后使用 copyFromLocal 命令将文件从本地复制到 HDFS 文件系统中，如下所示：

```
grunt>>copyFromLocal srcpath/Student dstpath;
```

通过 ls 命令，我们可以查看是否成功地将文件复制到相应的 HDFS 文件系统了。操作完成后，我们就可以像在 Local 模式下一样对文件进行操作了。这里，Pig 会自动地将我们的命令分散到分布式系统中去执行，然后返回给用户。

2. 脚本文件

参考 Local 模式下脚本文件的执行。

3. 嵌入式程序

参考 Local 模式下脚本文件的执行，这里我们给出 MapReduce 模式下的程序代码，用户可以看到，除了指定相应的模式之外，MapReduce 模式下的程序代码和 Local 模式没有什么不同，这是因为，所有的分布式操作将由 Pig 系统自动执行，而不需要用户在 MapReduce 的编程框架下设计程序，这就大大地减轻了用户的负担，也使用户能更容易地掌握 Pig 嵌入式程序，如代码清单 14-3 所示。

代码清单 14-3 MapReduce 模式下的 Pig 脚本

```
import java.io.IOException;
import org.apache.pig.PigServer;
public class tst_mapreduce{
public static void main(String[] args) {
    try {
        PigServer pigServer = new PigServer("mapreduce");//mapreduce 模式
        runIdQuery(pigServer, "/path/Student");// 调用函数
    }
    catch(Exception e) {}
}
public static void runIdQuery(PigServer pigServer, String inputFile) throws
    IOException {
    pigServer.registerQuery("A = load '" + inputFile + "' using PigStorage(',') as
        (Sno:chararray, Sname:chararray, Ssex:chararray, Sage:int, Sdept:chararray);");
    pigServer.registerQuery("B = foreach A generate Sname, Sage; ");
    pigServer.store("B", "/path/tstJavaMapReduce.out");
}
}
```

14.6 Pig 进阶

本节将继续介绍 Pig 在实际中的应用，为了体现 Pig 系统的特点，本节中的所有操作都将在 Hadoop MapReduce 模式下进行。另外，我们选取了一组很有特点的例子进行数据分析，相信这对用户的理解会很有帮助。

为了让用户能够更好地理解下面的操作，我们使用 Grunt Shell 方式进行数据分析，这样能够让用户更加清楚地理解 Pig 的执行过程。

14.6.1 数据实例

结合 14.5 节中的数据，我们再给出另外两个数据。

第一组数据是 14.5 节中的学生表所对应的课程表（课程号、课程名、先修课程号、学分），它包含如下几条记录：

```
01,English,,4
```

```
02,Data Structure,05,2
03,DataBase,02,2
04,DB Design,03,3
05,C Language,,3
06,Principles Of Network,07,3
07,OS,05,3
```

它们所对应的数据类型如下所示:

```
Course(Cno:chararray,Cname:chararray,Cpno:chararray,Ccredit:int)
```

另外一组数据为学生表和课程表所对应的选课表(学号、课程号、成绩),它包含如下几条记录:

```
201000101,01,92
201000101,03,84
201000102,01,90
201000102,02,94
201000102,03,82
201000103,01,72
201000103,02,90
201000104,03,75
```

它们所对应的数据类型如下所示:

```
SC(Sno:chararray,Cno:chararray,Grade:int)
```

14.6.2 Pig 数据分析

下面我们将对学生表、课程表和选课表进行数据分析操作。这一节将分为三个部分,分别计算学生的平均成绩、找出有不及格成绩的学生和找出修了先修课为“C Language”的学生。在语法上,Pig Latin 虽然没有关系数据库中的关系操作语言强大,但是因为 Pig 系统是架设在 Hadoop 的云平台之上的,所以在处理大规模数据集的时候,Pig 的效率非常高。

1. 计算每个学生的平均成绩

这里要求计算出每个学生的平均成绩,并且输出每个学生的姓名及其平均成绩。

我们先对数据进行分析。很容易看出,我们需要对学生表和选课表进行操作。首先,需要将学生表和选课表基于学号字段进行连接;然后,基于学号对学生数据进行操作,这时需要对每个学生所有的课程成绩分别求和,并除以课程总数;最后,按格式输出结果。

对于传统的关系型数据库的关系操作语言来说,为了实现这个目标,我们需要将 AVG 运算和 GROUP 运算同时使用,十分方便。下面,我们就 Pig Latin 语言给出相应的操作。

- 1) 从源数据文件学生表和选课表中读取数据。
- 2) 对学生表和选课表基于学号字段进行连接操作。
- 3) 基于学号对连接生成的表进行分组操作。
- 4) 计算每个学生的平均成绩。

5) 输出（存储）操作结果。

上面是对操作的描述，接下来就针对上述描述用 Pig Latin 语言来实现。

(1) 取数据

MapReduce 在 Hadoop 的 HDFS 文件系统中对数据进行操作，所以需要拷贝被操作的数据到 HDFS 中：

```
copyFromLocal Student Student;
copyFromLocal SC SC
```

可以使用 Hadoop 的 ls 命令查看数据是否复制成功，确认后再读取数据：

```
A = load 'Tmp/Student' using PigStorage(',') as
    (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray);
B = load 'Tmp/SC' using PigStorage(',') as
    (Sno:chararray,Cno:chararray,Grade:int);
```

(2) 连接操作

我们使用 JOIN 关键字对 A、B 两组数据基于 Sno 字段进行连接操作。JOIN 关键字的语法如下：

```
alias = JOIN alias BY {expression|('expression [, expression ...]')} (, alias
    BY {expression|('expression [, expression ...]')} ...) [USING 'replicated'
    | 'skewed' | 'merge'] [PARALLEL n];
```

下面是连接操作的命令：

```
D = Join A By Sno,B By Sno;
```

这里我们可以使用 Dump 关键字来查看 D 中存储的数据，如图 14-4 所示。



```
2010-10-04 03:48:39,769 [main] INFO org.apache.hadoop.mapreduce.lib.input
to process : 1
2010-10-04 03:48:39,770 [main] INFO org.apache.pig.backend.hadoop.execution
ut paths to process : 1
(201000101,李勇,男,20,计算机理论与,201000101,01,92)
(201000101,李勇,男,20,计算机理论与,201000101,03,50)
(201000102,王丽,女,19,计算机理论与,201000102,01,90)
(201000102,王丽,女,19,计算机理论与,201000102,02,94)
(201000102,王丽,女,19,计算机理论与,201000102,03,82)
(201000103,刘花,女,18,计算机应用技术,201000103,01,35)
(201000103,刘花,女,18,计算机应用技术,201000103,02,90)
(201000104,李莉,男,19,计算机系统结构,201000104,03,75)
grunt>
```

图 14-4 对学生表和选课表进行连接操作后的结果

(3) 分组操作

在进行分组操作之前，我们先提取必要的数据库，这不但减少了需要处理的数据量，还让我们的操作更加简单。接着，我们基于学号字段对连接操作后的数据进行分组，如下所示：

```
E = Foreach D generate A::Sno,Sname,Grade;
F = Cogroup E By (Sno,Sname);
```

我们再使用 DUMP 关键字查看一下 F 中的数据，如图 14-5 所示。接着用 DESCRIBE 分

析 F 的模式, 如图 14-6 所示。

```
((201000101,李勇),{(201000101,李勇,92),(201000101,李勇,50)})
((201000102,王丽),{(201000102,王丽,90),(201000102,王丽,94),(201000102,王丽,82)})
((201000103,刘花),{(201000103,刘花,35),(201000103,刘花,90)})
((201000104,李岗),{(201000104,李岗,75)})
grunt>
```

图 14-5 F 中的数据

```
grunt> Describe F;
F: {group: (A::Sno: chararray,A::Sname: chararray),E: {A::Sno: chararray,A::Sname: chararray,B::Grade: int}}
grunt>
```

图 14-6 F 的模式

(4) 计算学生的平均成绩

我们使用 SUM 关键字对学生成绩进行求和, 使用 COUNT 关键字来计算课程的总数, 如下所示:

```
G = Foreach F Generate group.Sname, (SUM(E.Grade)/COUNT(E));
```

下面, 我们查看一下最终的结果, 如图 14-7 所示:

```
(李勇,71L)
(王丽,88L)
(刘花,62L)
(李岗,75L)
```

图 14-7 学生平均成绩

因为 Grade 字段的数据类型为 int, 所以这里计算出的结果均为向下取整后的值。如果想要得到更为准确的数据, 用户可以将 Grade 字段的数据类型设为 Long 或 Float。

2. 找出有不及格成绩的学生

这里要求找出有不及格成绩的学生, 并且输出学生的姓名和不及格的课程和成绩。

现在对问题进行分析。我们需要使用学生表来获取学生的姓名, 使用课程表来获取学生的成绩, 并且使用课程表来获取对应成绩的课程。

首先, 我们还是需要读取源数据; 然后, 使用连接字段将数据连接在一起; 接着, 使用 FILTER 关键字过滤出我们需要的数据; 最后, 提取需要的字段将数据输出。

这里我们不再像上面那样一步步地对数据进行了分析, 下面给出 Pig Latin 操作语句:

```
A = load '/pigTmp/Student' using PigStorage(',') as
    (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray); -- 读取学生表
B = load '/pigTmp/SC' using PigStorage(',') as
    (Sno:chararray,Cno:chararray,Grade:int); -- 读取选课表
C = load '/pigTmp/Course' using PigStorage(',') as
    (Cno:chararray,Cname:chararray,Cpno:chararray,Ccredit:int); -- 读取课程表
D = Filter B By Grade < 60; -- 提前对 B 进行分析, 过滤出需要的结果, 减少操作的数据量
E = Join D By Sno,A By Sno; -- 连接操作
F = Join E By Cno,C By Cno; -- 连接操作
G = Foreach F Generate Sname,Cname,Grade; -- 输出结果
```

最后我们使用 DUMP 命令查看操作的结果，如图 14-8 所示：

```
(刘花,English,35)
(李勇,DataBase,50)
```

图 14-8 不及格成绩的学生

3. 找出修了先修课为“C Language”的学生

这里要求找出修了先修课为“C Language”的学生，并且输出学生的姓名。

现在，我们对问题进行分析，从课程表的数据结构可以看出：我们需要找出“C Language”这门课的课程号，然后找对应的“Cpno”（即此课程号的课程），最后找出修了此门课程的学生，并输出学生的姓名。

Pig Latin 语言支持嵌套的操作，所以在这一部分，我们使用嵌套语句来对数据进行操作，这样能够使 Pig Latin 语言的书写更加简便，更加有利于理解。因为嵌套的语句能够使程序的执行更加有层次感，使我们理解起来一目了然。

为了让读者便于理解，我们给出了单步的操作，如下所示：

```
A = load '/pigTmp/Student' using PigStorage(',') as
    (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray);
B = load '/pigTmp/SC' using PigStorage(',') as
    (Sno:chararray,Cno:chararray,Grade:int);
C = load '/pigTmp/Course' using PigStorage(',') as
    (Cno:chararray,Cname:chararray,Cpno:chararray,Ccredit:int);
D = load '/pigTmp/Course' using PigStorage(',') as
    (Cno:chararray,Cname:chararray,Cpno:chararray,Ccredit:int);
E = Join C By Cpno,D By Cno; -- 连接数据
F = Filter E By D::Cname == 'C Language'; -- 过滤出先修课名为 C Language 的记录
G = Foreach F Generate C::Cno; -- 找出先修课为 C Language 课程的课程号
H = Join G By Cno,B By Cno; -- 选课表和 C Language 课程的课程号做连接操作
I = Join H By Sno,A By Sno; -- 选课表与目标课程号连接结果与学生表做连接操作
J = Foreach I Generate Sname -- 输出结果
```

可以明显地看出，上面的操作十分烦琐，下面我们将上面的语句嵌套起来。

因为，等号左面和右面的操作完全是等价的，也就是说可以将模式名用对应的表达式替换。比如对于下面的句子：

```
E = Join C By Cpno,D By Cno; -- 连接数据
F = Filter E By D::Cname == 'C Language'; -- 过滤出先修课名为 C Language 的记录
```

我们可以这样写：

```
F = Filter (Join C By Cpno,D By Cno;) By D::Cname == 'C Language'; -- 过滤出先修课
    名为 C Language 的记录
```

所以，这一问题可以按下面的 Pig Latin 语句来进行操作：

```
A = load '/pigTmp/Student' using PigStorage(',') as
    (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray);
```

```

B = load '/pigTmp/SC' using PigStorage(',') as
  (Sno:chararray,Cno:chararray,Grade:int);
C = load '/pigTmp/Course' using PigStorage(',') as
  (Cno:chararray,Cname:chararray,Cpno:chararray,Ccredit:int);
D = load '/pigTmp/Course' using PigStorage(',') as
  (Cno:chararray,Cname:chararray,Cpno:chararray,Ccredit:int);
E = Foreach (Filter (Join C By Cpno,D By Cno) By D::Cname == 'C Language')
  Generate C::Cno;
F = Foreach (Join (Join B By Cno, E By Cno) By Sno,A By Sno) Generate Sname;

```

当然，如果想一步执行完也是可以的，只需要将上面操作的后两步再嵌套起来即可：

```

A = load '/pigTmp/Student' using PigStorage(',') as
  (Sno:chararray,Sname:chararray,Ssex:chararray,Sage:int,Sdept:chararray);
B = load '/pigTmp/SC' using PigStorage(',') as
  (Sno:chararray,Cno:chararray,Grade:int);
C = load '/pigTmp/Course' using PigStorage(',') as
  (Cno:chararray,Cname:chararray,Cpno:chararray,Ccredit:int);
D = load '/pigTmp/Course' using PigStorage(',') as
  (Cno:chararray,Cname:chararray,Cpno:chararray,Ccredit:int);
E = Foreach (Join (Join B By Cno,(Foreach (Filter (Join C By Cpno,D By Cno) By
  D::Cname == 'C Language') Generate C::Cno) By Cno) By Sno,A By Sno) Generate
  Sname;

```

下面，我们使用 DUMP 关键字来分别对上面三种方式查看一下运行结果，发现输出结果是完全相同的，如图 14-9 所示。



```

(王丽)
(刘花)
grunt=

```

图 14-9 修了先修课为“C Language”的学生

14.5 节通过一个简单的例子，让用户了解如何在 Local 模式和 MapReduce 模式下对数据进行操作。14.6 节则进一步通过一组复杂的例子，对如何使用 Pig Latin 语言进行复杂的操作做了更深入的介绍。

从 14.5 节和 14.6 节的实例操作中，我们可以看出 Pig Latin 语言更擅长对海量数据进行分析。另外，Pig Latin 语言还支持嵌套的操作，这样可以让 Pig Latin 语言编写的程序更加易于理解。

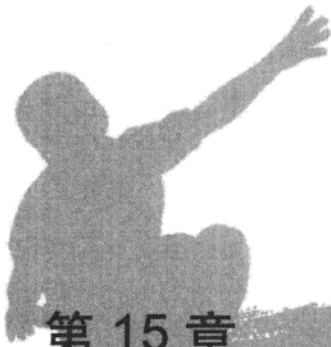
鉴于 Pig Latin 语言的如上特点，我们可以使用 Pig 对诸如日志等规则的、海量的并且需要定期维护的数据进行分析处理操作，这样可以大大地提高系统的工作效率。

14.7 小结

在本章中我们通过对 Pig 的实际操作，让用户对 Pig 有了一个新的认识。相信读完本章之后，用户可以使用 Pig 进行简单的数据处理了。Pig Latin 语言不但自身提供了很多的函数供用户使用，而且用户还可以根据实际情况结合 Java 和 Pig Latin 语言来编写具有特定功能

的函数。这体现了 Pig 的可扩展性和强大的功能。在使用 Pig 的过程中，用户还有很多技巧需要掌握，这一点可以在实际操作中慢慢地体会。另外，Pig 还处于完善阶段。从 0.5.0 版到 0.7.0 版的发展过程中，Pig 进行了很多调整，这离不开广大开发者的支持和帮助。希望用户能够通过对 Pig 的使用，向 Apache Hadoop 贡献自己的一份力量！





第 15 章

ZooKeeper 详解

本章内容

- ☐ ZooKeeper 简介
- ☐ ZooKeeper 的安装和配置
- ☐ ZooKeeper 的简单操作
- ☐ ZooKeeper 的特性
- ☐ ZooKeeper 的 Leader 选举
- ☐ ZooKeeper 锁服务
- ☐ 使用 ZooKeeper 创建应用程序
- ☐ 小结

资源如蒙
PDG

15.1 ZooKeeper 简介

ZooKeeper 是一个为分布式应用所设计的开源协调服务。它可以为用户提供同步、配置管理、分组和命名等服务。开发者意欲将 ZooKeeper 设计成一个易于编程的环境，所以它的文件系统使用了我们所熟悉的目录树结构。ZooKeeper 是使用 Java 编写的，但是它支持 Java 和 C 两种编程语言。

众所周知，协调服务非常容易出错，而且很难从故障中恢复，例如，协调服务很容易处于竞态以至于出现死锁。ZooKeeper 的设计目的是为了减轻分布式应用程序所承担的协调任务。

15.1.1 ZooKeeper 的设计目标

众所周知，分布式环境下的程序和活动为了达到协调一致的目的，通常会具有某些共同的特点，例如，简单性、有序性等。ZooKeeper 不但在这些目标的实现上有其自身的特点，并且具有其独特的优势。下面我们将简述 ZooKeeper 的设计目标。

(1) 简单化

ZooKeeper 允许分布式的进程通过共享体系的命名空间来进行协调，这个命名空间的组织与标准的文件系统非常相似，它是由一些数据寄存器组成的。用 ZooKeeper 的语法来说，这些寄存器应称为 Znode。它们和文件及目录非常相似。典型的文件系统是基于存储设备的，然而，ZooKeeper 的数据却是存放在内存当中的，这就意味着 ZooKeeper 可以达到一个高的吞吐量，并且有低的延迟。ZooKeeper 的实现非常重视高性能、高可靠性，以及严格的有序访问。

ZooKeeper 性能上的特点决定了它能够用在大型的、分布式的系统当中。从可靠性方面来说，它并不会因为一个节点的错误而崩溃。除此之外，它严格的序列访问控制意味着复杂的控制原语可以应用在客户端上。

(2) 健壮性

组成 ZooKeeper 服务的服务器必须互相知道其他服务器的存在。它们维护着一个处于内存中的状态镜像，以及一个位于存储器中的交换日志和快照。只要大部分的服务器可用，那么 ZooKeeper 服务就可用。

如果客户端连接到单个 ZooKeeper 服务器上，那么这个客户端就管理着一个 TCP 连接，并且通过这个 TCP 连接来发送请求、获得响应、获取检测事件，以及发送心跳。如果连接到服务器上的 TCP 链接断开，客户端将连接到其他的服务器上。

(3) 有序性

ZooKeeper 可以为每一次更新操作赋予一个版本号，并且此版本号是全局有序的，不存在重复的情况。ZooKeeper 所提供的很多服务也是基于此有序性的特点来完成的。

(4) 速度优势

它在读取主要负载时尤其快。ZooKeeper 应用程序在上千台机器的节点上运行。另外，

需要注意的是 ZooKeeper 有这样一个特点：当读工作比写工作更多的时候，它执行的性能会更好。

15.1.2 数据模型和层次命名空间

ZooKeeper 提供的命名空间与标准的文件系统非常相似。它的名称是由通过斜线分隔的路径名序列所组成的。ZooKeeper 中的每一个节点都是通过路径来识别的。

图 15-1 是 Zookeeper 中节点的数据模型，这种树形结构的命名空间操作方便且易于理解。

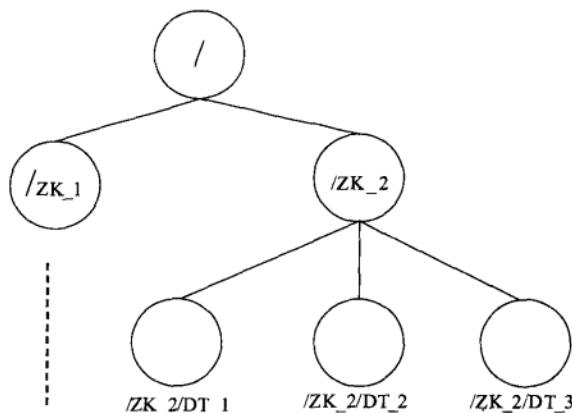


图 15-1 ZooKeeper 的层次命名空间

15.1.3 ZooKeeper 中的节点和临时节点

通过上一节的内容，读者可以了解到在 ZooKeeper 中存在着节点的概念，同时也知道了这些节点是通过像树一样的结构来进行维护的，并且每一个节点通过路径来标识及访问。除此之外，每一个节点还拥有自身的一些信息，包括：数据、数据长度、创建时间、修改时间等。从节点的这些特性（既含有数据，又通过路径来表示）可以看出，它既可以被看作是一个文件，又可以被看作是一个目录，因为它同时具有二者的特点。为了便于表达，后面我们将使用 Znode 来表示所讨论的 ZooKeeper 节点。

具体地说，Znode 维护着数据、ACL（Access Control List，访问控制列表）、时间戳等包含交换版本号信息的数据结构，它通过对这些数据的管理来让缓存中的数据生效，并且执行协调更新操作。每当 Znode 中的数据更新时它所维护的版本号就会增加，这非常类似于数据库中计数器时间戳的操作方式。

另外 Znode 还具有原子性操作的特点：在命名空间中，每一个 Znode 的数据将被原子地读写。读操作将读取与 Znode 相关的所有数据，写操作将替换掉所有的数据。除此之外，每一个节点都有一个访问控制列表，这个访问控制列表规定了用户操作的权限。

ZooKeeper 中同样存在临时节点。这些节点与 session 同时存在，当 session 生命周期结

束时，这些临时节点也将被删除。临时节点在某些场合下也发挥着非常重要的作用，例如 Leader 选举、锁服务等。

15.1.4 ZooKeeper 的应用

ZooKeeper 成功地应用于大量的工业程序中。它在 Yahoo! 被用于雅虎消息代理 (Yahoo! Message Broker) 的协调和故障恢复服务。雅虎消息代理是一个高度可扩展的发布 - 订阅系统，它管理着上千的总联机程序和信息系统 (TOPICS: Total On-line Program and Information Control System)，另外它还用于为 Yahoo! crawler 获取服务并恢复错误故障。除此之外，一些 Yahoo! 广告系统也同样使用 ZooKeeper 来实现可靠的服务。

15.2 ZooKeeper 的安装和配置

在这一节中，我们将首先向读者介绍如何在不同的环境下安装并配置 ZooKeeper 服务，然后具体介绍如何通过 ZooKeeper 配置文件对 ZooKeeper 进行配置管理，最后向读者介绍如何在不同环境下启动 ZooKeeper 服务。

15.2.1 在集群上安装 ZooKeeper

ZooKeeper 有不同的运行环境，包括：单机环境、集群环境和集群伪分布环境。这里，我们将分别介绍不同环境下如何安装 ZooKeeper 服务并简单介绍它们的区别与联系。

1. 单机下安装 ZooKeeper

(1) ZooKeeper 的下载

如果读者是第一次使用 ZooKeeper，那么我们建议首先尝试在单机模式下配置 ZooKeeper 服务器。因为，在单机模式下配置和使用相对来说要简单得多，并且易于帮助我们理解 ZooKeeper 的工作原理。这对读者进一步学习使用 ZooKeeper 会有很大的帮助。

从 Apache 官方网站下载一个 ZooKeeper 的最新稳定版本，网址如下：

<http://hadoop.apache.org/zookeeper/releases.html>

作为国内用户来说，选择最近的源文件服务器所在地，能够节省不少的时间，比如：

<http://labs.renren.com/apache-mirror//hadoop/zookeeper/>

(2) ZooKeeper 的安装

ZooKeeper 要求有 Java 的环境才能运行，并且需要是 Java 6 以上的版本，读者可以从 SUN 官网上下载，然后对 Java 的环境变量进行设置，这部分内容前面已经详细介绍过了，故此不再赘述。除此之外，为了今后操作方便，我们需要对 ZooKeeper 的环境变量进行配置，方法如下，在 /etc/profile 文件中加入如下的内容：

```
#Set ZooKeeper Enviroment
```

```
export ZOOKEEPER_HOME=/root/hadoop-0.20.2/zookeeper-3.3.1
export PATH=$PATH:$ZOOKEEPER_HOME/bin:$ZOOKEEPER_HOME/conf
```

ZooKeeper 服务器包含在单个 JAR 文件中，安装此服务需要用户创建一个配置文档，并对其进行设置。我们在 ZooKeeper-*.*. 目录（本书以当前 ZooKeeper 的最新版 3.3.1 为例，故在下文中此“ZooKeeper-*.*. ”都将写为“ZooKeeper-3.3.1”）的 conf 文件夹下创建一个 zoo.cfg 文件，它包含如下的内容：

```
tickTime=2000
dataDir=/var/zookeeper
clientPort=2181
```

在这个文件中，我们需要指定 dataDir 的值，它指向了一个目录，这个目录在开始的时候应为空。下面是每个参数的含义：

- ❑ tickTime：基本事件单元，以毫秒为单位。它用来指示心跳，最小的 session 过期时间为两倍的 tickTime。
- ❑ dataDir：存储内存中数据库快照的位置，如果不设置参数，更新事务的日志将被存储到默认位置。
- ❑ clientPort：监听客户端连接的端口。

使用单机模式时用户需要注意：在这种配置方式下没有 ZooKeeper 副本，所以如果 ZooKeeper 服务器出现故障，ZooKeeper 服务将会停止。

代码清单 15-1 是我们根据自身情况所设置的 ZooKeeper 配置文档：zoo.cfg。

代码清单 15-1 zoo.cfg

```
# The number of milliseconds of each tick
tickTime=2000

# the directory where the snapshot is stored.
dataDir=/root/hadoop-0.20.2/zookeeper-3.3.1/snapshot/data

# the port at which the clients will connect
clientPort=2181
```

2. 在集群下安装 ZooKeeper

(1) 支持的平台

ZooKeeper 可以在不同的系统上运行，表 15-1 是关于这方面的一个简单说明：

表 15-1 ZooKeeper 支持的平台

系 统	支持的平台	是否服务器、客户端
GNU/Linux	开发和生产平台	服务器和客户端
Sun Solaris	开发和生产平台	服务器和客户端

(续)

系 统	支持的平台	是否服务器、客户端
FreeBSD	开发和生产平台	仅用作客户端
Win32	开发平台	服务器和客户端
MacOSX	开发平台	服务器和客户端

(2) 集群部署

为了获得可靠的 ZooKeeper 服务，用户应该在一个集群上部署 ZooKeeper。只要集群上大多数的 ZooKeeper 服务启动了，那么总的 ZooKeeper 服务将是可用的。另外，最好使用奇数台机器。例如，拥有 4 台机器的 ZooKeeper 只能处理一台机器的故障，如果两台机器发生故障，余下的两台机器并不能组成一个可用的 ZooKeeper ensemble。然而，如果 ZooKeeper 拥有 5 台机器，那么它就能处理 2 台机器的故障了。

这之后的操作和单机模式的安装类似，我们同样需要对 Java 环境进行设置，下载最新的 ZooKeeper 稳定版本并配置相应环境变量。不同之处在于每台机器上 `conf/zoo.cfg` 配置文件的参数设置不同，可参考下面的配置：

```
tickTime=2000
dataDir=/var/zookeeper/
clientPort=2181
initLimit=5
syncLimit=2
server.1=zoo1:2888:3888
server.2=zoo2:2888:3888
server.3=zoo3:2888:3888
```

更多关于 ZooKeeper 参数的设置请参看 15.2.2 节。“`server.id=host:port:port.`”标识了不同的 ZooKeeper 服务器，这些服务器作为集群的一部分应该知道 ensemble 中的其他机器，用户可以从“`server.id=host:port:port.`”中读取相关的信息。参数中 `host` 和 `port` 非常直接。在服务器的 `data`（`dataDir` 参数所指定的目录）目录下创建一个文件名为 `myid` 的文件，这个文件中仅含一行的内容，它所指定的是自身的 `id` 值。比如，服务器“1”应该在 `myid` 文件中写入“1”。而且这个 `id` 值必须是 ensemble 中唯一的，且大小在 1 到 255 之间。在这一行配置中，第一个端口（`port`）是从（`follower`）机器连接到主（`leader`）机器的端口，第二个是用来进行 `leader` 选举的端口。在这个例子中，每台机器使用三个端口，分别是：`clientPort`，2181；`port`，2888；`port`，3888。

我在拥有三台机器的 Hadoop 集群上测试使用了 ZooKeeper 服务，下面的代码清单 15-2 就是根据自身情况所设置的 ZooKeeper 配置文档。

代码清单 15-2 zoo.cfg

```
# The number of milliseconds of each tick
tickTime=2000
```

```
# The number of ticks that the initial
# synchronization phase can take
initLimit=10

# The number of ticks that can pass between
# sending a request and getting an acknowledgement
syncLimit=5

# the directory where the snapshot is stored.
dataDir=/root/hadoop-0.20.2/zookeeper-3.3.1/snapshot/d1

# the port at which the clients will connect
clientPort=2181

server.1=IP1:2887:3887
server.2=IP2:2888:3888
server.3=IP3:2889:3889
```

清单中的 IP 分别对应的是配置分布式 ZooKeeper 的 IP 地址。当然，也可以通过机器名访问 ZooKeeper，但是需要在 Ubuntu 的 host 环境中进行设置，这部分内容不是本书的重点，不再赘述。读者可以查阅 Ubuntu 及 Linux 的相关资料。

3. 在集群伪分布模式下安装 ZooKeeper

通过前面的章节，读者了解到 Hadoop 可以在伪分布模式下模拟分布式 Hadoop 的运行。与它不同的是，ZooKeeper 不但可以在伪分布模式下运行单机的 ZooKeeper（即只有一个 Zookeeper 服务），而且可以在伪分布模式下模拟集群模式 ZooKeeper 的运行，我们索性将其称之为“集群伪分布模式”，以区别“单机伪分布模式”。我们知道，伪分布模式下 Hadoop 的操作和分布式模式下有着很大的不同，但是在集群伪分布模式下对 ZooKeeper 的操作却和集群模式下没有本质的区别。显然，集群伪分布模式为我们体验 ZooKeeper 和做一些尝试性的实验，提供了很大的便利。比如，我们在实验的时候，可以先使用少量数据在集群伪分布模式下进行测试。当测试可行的时候，再将其移植到集群模式下进行真实的数据实验。这样，不但保证了它的可行性，同时大大提高了实验的效率。

那么，如何配置 ZooKeeper 的集群伪分布模式呢？其实很简单。用心的读者可以发现，在 ZooKeeper 的配置文档中，clientPort 参数是用来设置客户端连接 ZooKeeper 的端口。在 server.1=IP1:2887:3887 中，IP1 指示的是组成 ZooKeeper 服务的机器 IP 地址，2887 为进行 leader 选举的端口，3887 是组成 ZooKeeper 服务的机器之间的通信端口。在集群伪分布模式下我们使用每个配置文档模拟一台机器，也就是说，需要在单台机器上运行多个 ZooKeeper 实例。但是，必须要保证各个配置文档的 clientPort 不能冲突。

下面是我们所配置的集群伪分布模式，分别通过 zoo1.cfg、zoo2.cfg、zoo3.cfg 来模拟有三台机器的 ZooKeeper 集群。详见代码清单 15-3 至代码清单 15-5。

代码清单 15-3 zoo1.cfg

```
# The number of milliseconds of each tick
tickTime=2000

# The number of ticks that the initial
# synchronization phase can take
initLimit=10

# The number of ticks that can pass between
# sending a request and getting an acknowledgement
syncLimit=5

# the directory where the snapshot is stored.
dataDir=/root/hadoop-0.20.2/zookeeper-3.3.1/d_1

# the port at which the clients will connect
clientPort=2181

server.1=localhost:2887:3887
server.2=localhost:2888:3888
server.3=localhost:2889:3889
```

代码清单 15-4 zoo2.cfg

```
# The number of milliseconds of each tick
tickTime=2000

# The number of ticks that the initial
# synchronization phase can take
initLimit=10

# The number of ticks that can pass between
# sending a request and getting an acknowledgement
syncLimit=5

# the directory where the snapshot is stored.
dataDir=/root/hadoop-0.20.2/zookeeper-3.3.1/d_2

# the port at which the clients will connect
clientPort=2182

#the location of the log file
dataLogDir=/root/hadoop-0.20.2/zookeeper-3.3.1/logs

server.1=localhost:2887:3887
server.2=localhost:2888:3888
server.3=localhost:2889:3889
```

代码清单 15-5 zoo3.cfg

```
# The number of milliseconds of each tick
tickTime=2000

# The number of ticks that the initial
# synchronization phase can take
initLimit=10

# The number of ticks that can pass between
# sending a request and getting an acknowledgement
syncLimit=5

# the directory where the snapshot is stored.
dataDir=/root/hadoop-0.20.2/zookeeper-3.3.1/d_2

# the port at which the clients will connect
clientPort=2183

#the location of the log file
dataLogDir=/root/hadoop-0.20.2/zookeeper-3.3.1/logs

server.1=localhost:2887:3887
server.2=localhost:2888:3888
server.3=localhost:2889:3889
```

从上述三个代码清单可以看到，它们除了 clientPort 不同之外，dataDir 也不同。另外，不要忘记在 dataDir 所对应的目录中创建 myid 文件来指定对应的 ZooKeeper 服务器实例。

注意 ZooKeeper 的集群伪分布模式兼有 ZooKeeper 单机模式和集群模式的优点，下述内容如果没有特殊声明，均为在 ZooKeeper 集群伪分布模式下的操作。

15.2.2 配置 ZooKeeper

ZooKeeper 的功能特性是通过 ZooKeeper 配置文件来进行控制管理（zoo.cfg 配置文件）的。这样的设计其实有它自身的原因。通过前面对 ZooKeeper 的配置可以看出，在对 ZooKeeper 集群进行配置的时候，它的配置文档是完全相同的（对于集群伪分布模式来说，只有很少的部分是不同的）。这样的配置方式使得在部署 ZooKeeper 服务的时候非常方便。如果服务器使用不同的配置文件，必须要确保不同配置文件中的服务器列表相匹配。

在设置 ZooKeeper 配置文档的时候，某些参数是可选的，但是某些参数是必须的。这些必须的参数就构成了 ZooKeeper 配置文档的最低配置要求。另外，如果需要对 ZooKeeper 进行更详细的配置，读者可以参考我们下面将要讲述的内容。

1. 最低配置

下面是在最低配置要求中必须配置的参数：

- clientPort：监听客户端连接的端口。
- dataDir：存储内存中数据库快照的位置。

注意：应该谨慎地选择日志存放的位置，使用专用的日志存储设备能够大大提高系统的性能，如果将日志存储在比较繁忙的存储设备上，那么将会在很大程度上影响系统的性能。

- tickTime：基本事件单元，以毫秒为单位。它用来控制心跳和超时，默认情况下最小的会话超时时间为两倍的 tickTime。

2. 高级配置

下面是高级配置要求中可选的配置参数，用户可以使用下面的参数来更好地规定 ZooKeeper 的行为：

- dataLogDir：这个操作让管理机器把事务日志写入“dataLogDir”所指定的目录中，而不是“dataDir”所指定的目录。这将允许使用一个专用的日志设备，并且帮助我们避免日志和快照之间的竞争。配置如下：

```
#the location of the log file
dataLogDir=/root/hadoop-0.20.2/zookeeper-3.3.1/log/data_log
```

- maxClientCnxns：这个操作将限制连接到 ZooKeeper 的客户端的数量，并且限制并发连接的数量，它通过 IP 来区分不同的客户端。此配置选项可以用来阻止某些类别的 Dos 攻击。将它设置为 0 或忽略而不进行设置将会取消对并发连接的限制。

例如，我们将 maxClientCnxns 的值设置为 1，如下所示：

```
#set maxClientCnxns
maxClientCnxns=1
```

启动 ZooKeeper 之后，首先用一个客户端连接到 ZooKeeper 服务器之上。之后若有第二个客户端尝试对 ZooKeeper 进行连接，或者有某些隐式的对客户端的连接操作，将会触发 ZooKeeper 的上述配置。系统会提示相关信息，如图 15-2 所示：

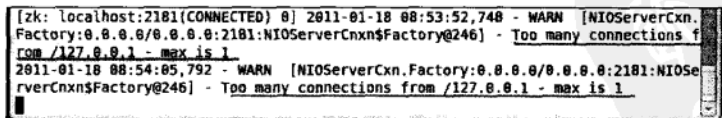


图 15-2 ZooKeeper maxClientCnxns 异常

- minSessionTimeout 和 maxSessionTimeout：最小的会话超时时间和最大的会话超时时间。在默认情况下，最小的会话超时时间为 2 倍的 tickTime 时间，最大的会话超时时间为 20 倍的会话超时时间。系统启动时，会显示相应的信息，如图 15-3 所示：

从下图中可以看出，minSessionTimeout 及 maxSessionTimeout 的值均为 -1，现在我们来设置系统的最小和最大的会话超时时间，如下所示：

```
INFO [main:Environment@97] - Server environment:user.name=root
INFO [main:Environment@97] - Server environment:user.home=/root
INFO [main:Environment@97] - Server environment:user.dir=/root
INFO [main:ZooKeeperServer@666] - tickTime set to 2000
INFO [main:ZooKeeperServer@669] - minSessionTimeout set to 1
INFO [main:ZooKeeperServer@678] - maxSessionTimeout set to 1
INFO [main:NIOServerCnxnFactory@143] - binding to port 0.0.0.0/0.0.0.0:21
```

图 15-3 默认会话超时时间

```
#set minSessionTimeout
minSessionTimeout=1000

#set maxSessionTimeout
maxSessionTimeout=10000
```

在配置 `minSessionTimeout` 及 `maxSessionTimeout` 的值时需要注意，如果将此值设置得太小的话，会话很可能刚刚建立便由于超时而不得不退出。一般情况下，不能将此值设置得比 `tickTime` 的值还小。

3. 集群配置

- `initLimit`：此配置表示，允许 follower（相对于 leader 而言的“客户端”）连接并同步到 leader 的初始化连接时间，它是以 `tickTime` 的倍数来表示的。当初始化连接时间超过设置倍数的 `tickTime` 时间时，则连接失败。
- `syncLimit`：此配置表示，leader 与 follower 之间发送消息时，请求和应答的时间长度。如果 follower 在设置的时间内不能与 leader 通信，那么此 follower 将被丢弃。

15.2.3 运行 ZooKeeper

1. 单机模式下运行 ZooKeeper

如果读者已经按照 15.2.1 节中的第 1 点正确地配置了 ZooKeeper 的环境变量，那么我们现在可以直接在终端运行 ZooKeeper 的 `sh` 脚本了，从而启动 ZooKeeper 的服务。

读者可以通过下面的命令来启动 ZooKeeper 服务：

```
zkServer.sh start
```

这个命令默认情况下执行 ZooKeeper 的 `conf` 文件夹下的 `zoo.cfg` 配置文件。当运行成功时用户会看到类似如下的提示界面：

```
root@ubuntu:~# zkServer.sh start
JMX enabled by default
Using config: /root/hadoop-0.20.2/zookeeper-3.3.1/bin/../conf/zoo.cfg
Starting zookeeper ...
STARTED
... ..
2011-01-19 10:04:42,300 - WARN [main:QuorumPeerMain@105] - Either no config or no
quorum defined in config, running in standalone mode
... ..
```

```

2011-01-19 10:04:42,419 - INFO [main:ZooKeeperServer@660] - tickTime set to 2000
2011-01-19 10:04:42,419 - INFO [main:ZooKeeperServer@669] - minSessionTimeout
set to -1
2011-01-19 10:04:42,419 - INFO [main:ZooKeeperServer@678] - maxSessionTimeout
set to -1
2011-01-19 10:04:42,560 - INFO [main:NIOServerCnxn$Factory@143] - binding to port
0.0.0.0/0.0.0.0:2181
2011-01-19 10:04:42,806 - INFO [main:FileSnap@82] - Reading snapshot /root/
hadoop-0.20.2/zookeeper-3.3.1/data/version-2/snapshot.200000036
2011-01-19 10:04:42,927 - INFO [main:FileSnap@82] - Reading snapshot /root/
hadoop-0.20.2/zookeeper-3.3.1/data/version-2/snapshot.200000036
2011-01-19 10:04:42,950 - INFO [main:FileTxnSnapLog@208] - Snapshotting:
400000058

```

从上面可以看出，运行成功后，系统会列出 ZooKeeper 运行的相关环境配置信息。

2. 集群模式下运行 ZooKeeper

在集群模式下需要用户在每台 ZooKeeper 机器上运行第一部分中的命令，这里不再赘述。

3. 集群伪分布模式下运行 ZooKeeper

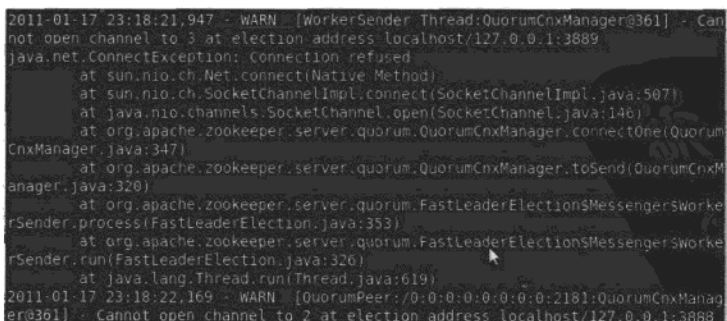
在集群伪分布模式下，我们只有一台机器，但是要运行三个 ZooKeeper 服务实例。此时，如果再使用上述命令式肯定是行不通的。其实，通过下面三条命令就能运行前面所配置的 ZooKeeper 服务了。如下所示：

```

zkServer.sh start zoo1.cfg
zkServer.sh start zoo2.cfg
zkServer.sh start zoo3.cfg

```

在运行完第一条命令之后，读者将会发现一些系统错误提示，如图 15-4 所示。



```

2011-01-17 23:18:21,947 - WARN [WorkerSender Thread:QuorumCnxManager@361] - Can
not open channel to 3 at election address localhost/127.0.0.1:3889
java.net.ConnectException: Connection refused
    at sun.nio.ch.Net.connect(Native Method)
    at sun.nio.ch.SocketChannelImpl.connect(SocketChannelImpl.java:507)
    at java.nio.channels.SocketChannel.open(SocketChannel.java:146)
    at org.apache.zookeeper.server.quorum.QuorumCnxManager.connectOne(Quorum
CnxManager.java:347)
    at org.apache.zookeeper.server.quorum.QuorumCnxManagerToSend(QuorumCnxM
anager.java:326)
    at org.apache.zookeeper.server.quorum.FastLeaderElection$Messenger$Worke
rSender.process(FastLeaderElection.java:353)
    at org.apache.zookeeper.server.quorum.FastLeaderElection$Messenger$Worke
rSender.run(FastLeaderElection.java:326)
    at java.lang.Thread.run(Thread.java:619)
2011-01-17 23:18:22,169 - WARN [QuorumPeer:0:0:0:0:0:2181:QuorumCnxManag
er@361] - Cannot open channel to 2 at election address localhost/127.0.0.1:3888

```

图 15-4 集群伪分布异常提示

产生如图 15-4 所示的异常信息是由于 ZooKeeper 服务的每个实例都拥有全局的配置信息，它们在启动的时候会随时地进行 Leader 选举操作（此部分内容下面将会详细讲述）。此

时第一个启动的 ZooKeeper 需要和另外两个 ZooKeeper 实例进行通信。但是，另外两个 ZooKeeper 实例还没有启动起来，因此也就产生了这样的异常信息。

我们直接将其忽略即可，待把图示中的“2号”和“3号”ZooKeeper 实例启动起来之后，相应的异常信息就会自然而然的消失了。

4. ZooKeeper 四字命令

ZooKeeper 支持某些特定的四字命令字母与其交互。它们大多是查询命令，用来获取 ZooKeeper 服务的当前状态及相关信息。用户在客户端可以通过 telnet 或 nc 向 ZooKeeper 提交相应的命令。ZooKeeper 常用的四字命令见表 15-2。

表 15-2 ZooKeeper 四字命令

ZooKeeper 四字命令	功能描述
conf	输出相关服务配置的详细信息
cons	列出连接到服务器的所有客户端的详细连接/会话信息。包括“接受/发送”的包数量、会话 id、操作延迟、最后的操作执行等信息
dump	列出未经处理的会话和临时节点
envi	输出关于服务环境的详细信息（区别于 conf 命令）
reqs	列出未经处理的请求
ruok	测试服务是否处于正确状态。如果确实如此，那么服务返回“imok”，否则不做任何响应
stat	输出关于性能和连接的客户端列表
wchs	列出服务器 watch 的详细信息
wchc	通过 session 列出服务器 watch 的详细信息，它的输出是一个与 watch 相关的会话列表
wchp	通过路径列出服务器 watch 的详细信息。它输出一个与 session 相关的路径

图 15-5 是 ZooKeeper 四字命令的一个简单用例。

```

root@ubuntu-laptop:~# echo ruok | nc 10.77.20.23 2181
imokroot@ubuntu-laptop:~# echo conf | nc 10.77.20.23 2181
clientPort=2181
dataDir=/root/hadoop-0.20.2/zookeeper-3.3.1/d_1/version-2
dataLogDir=/root/hadoop-0.20.2/zookeeper-3.3.1/d_1/version-2
tickTime=2000
maxClientCnxns=10
minSessionTimeout=4000
maxSessionTimeout=40000
serverId=1
initLimit=10
syncLimit=5
electionAlg=3
electionPort=3887
quorumPort=2887
peerType=0
root@ubuntu-laptop:~#

```

图 15-5 ZooKeeper 四字命令用例

5. ZooKeeper 命令行工具

在成功启动 ZooKeeper 服务之后，输入下述命令，连接到 ZooKeeper 服务：

```
zkCli.sh -server 10.77.20.23:2181
```

连接成功后，系统会输出 ZooKeeper 的相关环境及配置信息，并在屏幕上输出 “Welcome to ZooKeeper” 等信息。

输入 help 之后，屏幕会输出可用的 ZooKeeper 命令，如图 15-6 所示：

```
[zk: 10.77.20.23:2181(CONNECTED) 1] help
ZooKeeper -server host:port cmd args
connect host:port
get path [watch]
ls path [watch]
set path data [version]
delquota [-n|-b] path
quit
printwatches on|off
create [-s] [-e] path data acl
stat path [watch]
close
ls2 path [watch]
history
listquota path
setAcl path acl
getAcl path
sync path
redo cmdno
addauth scheme auth
delete path [version]
setquota -n|-b val path
```

图 15-6 ZooKeeper 命令

15.3 ZooKeeper 的简单操作

15.3.1 使用 ZooKeeper 命令的简单操作步骤

1) 使用 ls 命令来查看当前 ZooKeeper 中所包含的内容：

```
[zk: 10.77.20.23:2181(CONNECTED) 1] ls /
[zookeeper]
```

2) 创建一个新的 Znode，使用 create /zk myData。这个命令创建了一个新的 Znode 节点 “zk”，以及与它关联的字符串：

```
[zk: 10.77.20.23:2181(CONNECTED) 2] create /zk myData
Created /zk
```

3) 再次使用 ls 命令来查看现在 ZooKeeper 中所包含的内容：

```
[zk: 10.77.20.23:2181(CONNECTED) 3] ls /
[zk, zookeeper]
```

此时看到，zk 节点已经被创建。

4) 下面我们运行 get 命令来确认第二步中所创建的 Znode 是否包含我们所创建的字符串：

```
[zk: 10.77.20.23:2181(CONNECTED) 4] get /zk
myData
Zxid = 0x40000000c
time = Tue Jan 18 18:48:39 CST 2011
Zxid = 0x40000000c
mtime = Tue Jan 18 18:48:39 CST 2011
pZxid = 0x40000000c
cversion = 0
dataVersion = 0
aclVersion = 0
ephemeralOwner = 0x0
dataLength = 6
numChildren = 0
```

5) 下面我们通过 set 命令来对 zk 所关联的字符串进行设置:

```
[zk: 10.77.20.23:2181(CONNECTED) 5] set /zk shenlan211314
cZxid = 0x40000000c
ctime = Tue Jan 18 18:48:39 CST 2011
mZxid = 0x40000000d
mtime = Tue Jan 18 18:52:11 CST 2011
pZxid = 0x40000000c
cversion = 0
dataVersion = 1
aclVersion = 0
ephemeralOwner = 0x0
dataLength = 13
numChildren = 0
```

6) 下面我们将刚才创建的 Znode 删除:

```
[zk: 10.77.20.23:2181(CONNECTED) 6] delete /zk
```

7) 最后再次使用 ls 命令查看 ZooKeeper 所包含的内容:

```
[zk: 10.77.20.23:2181(CONNECTED) 7] ls /
[zookeeper]
```

经过验证, zk 节点已经被删除。

15.3.2 ZooKeeper API 的简单使用

1. ZooKeeper API 简介

ZooKeeper API 共包含 5 个包, 分别为: org.apache.zookeeper、org.apache.zookeeper.data、org.apache.zookeeper.server、org.apache.zookeeper.server.quorum 和 org.apache.zookeeper.server.upgrade。其中 org.apache.zookeeper 包含 ZooKeeper 类, 它是我们编程时最常用的类文件。这个类是 ZooKeeper 客户端库的主要类文件。如果要使用 ZooKeeper 服务, 应用程序首先必须创建一个 Zookeeper 实例, 这时就需要使用此类。一旦客户端和

ZooKeeper 服务建立起了连接, ZooKeeper 系统将会给此连接会话分配一个 ID 值, 并且客户端将会周期性地向服务器发送心跳来维持会话的连接。只要连接有效, 客户端就可以调用 ZooKeeper API 来做相应的处理。

ZooKeeper 类提供了表 15-3 所示的几类主要方法。

表 15-3 ZooKeeper 类方法描述

功 能	描 述
create	在本地目录树中创建一个节点
delete	删除一个节点
exists	测试本地是否存在目标节点
get/set data	从目标节点上读取 / 写数据
get/set ACL	获取 / 设置目标节点访问控制列表信息
get children	检索一个子节点上的列表
sync	等待要被传送的数据

2. ZooKeeper API 的使用

这里通过一个例子来简单介绍如何使用 ZooKeeper API 编写自己的应用程序, 见代码清单 15-6。

代码清单 15-6 ZooKeeper API 的使用

```

1. import java.io.IOException;
2.
3. import org.apache.zookeeper.CreateMode;
4. import org.apache.zookeeper.KeeperException;
5. import org.apache.zookeeper.Watcher;
6. import org.apache.zookeeper.ZooDefs.Ids;
7. import org.apache.zookeeper.ZooKeeper;
8.
9. public class demo {
10. // 会话超时时间, 设置为与系统默认时间一致
11. private static final int SESSION_TIMEOUT=30000;
12.
13. // 创建 ZooKeeper 实例
14. ZooKeeper zk;
15.
16. // 创建 Watcher 实例
17. Watcher wh=new Watcher(){
18.     public void process(org.apache.zookeeper.WatchedEvent event)
19.     {
20.         System.out.println(event.toString());
21.     }
22. };
23.

```

```

24. // 初始化 ZooKeeper 实例
25. private void createZKInstance() throws IOException
26. {
27.     zk=new ZooKeeper("localhost:2181",demo.SESSIION_TIMEOUT,this.wh);
28.
29. }
30.
31. private void ZKOperations() throws IOException,InterruptedException,
    KeeperException
32. {
33.     System.out.println("\n1. 创建 ZooKeeper 节点 (znode: zoo2, 数据: myData2, 权
        限: OPEN_ACL_UNSAFE, 节点类型: Persistent");
34.     zk.create("/zoo2", "myData2".getBytes(), Ids.OPEN_ACL_UNSAFE,
        CreateMode.PERSISTENT);
35.
36.     System.out.println("\n2. 查看是否创建成功:");
37.     System.out.println(new String(zk.getData("/zoo2", false, null)));
38.
39.     System.out.println("\n3. 修改节点数据");
40.     zk.setData("/zoo2", "shenlan211314".getBytes(), -1);
41.
42.     System.out.println("\n4. 查看是否修改成功:");
43.     System.out.println(new String(zk.getData("/zoo2", false, null)));
44.
45.     System.out.println("\n5. 删除节点");
46.     zk.delete("/zoo2", -1);
47.
48.     System.out.println("\n6. 查看节点是否被删除:");
49.     System.out.println("节点状态: ["+zk.exists("/zoo2", false)+""]");
50. }
51.
52. private void ZKClose() throws InterruptedException
53. {
54.     zk.close();
55. }
56.
57. public static void main(String[] args) throws IOException,InterruptedException
    ,KeeperException {
58.     demo dm=new demo();
59.     dm.createZKInstance();
60.     dm.ZKOperations();
61.     dm.ZKClose();
62. }
63. }

```

此类包含两个主要的 ZooKeeper 函数，分别为 createZKInstance() 和 ZKOperations()。其中 createZKInstance() 函数负责对 ZooKeeper 实例 zk 进行初始化。ZooKeeper 类有两个构造函数，这里使用 “ZooKeeper (String connectString, int sessionTimeout, Watcher

watcher)” 对其进行初始化。因此，我们需要提供初始化所需的连接字符串信息、会话超时时间，以及一个 watcher 实例。第 17 行到第 23 行的代码是程序所构造的一个 watcher 实例，它能够输出所发生的事件。

ZKOperations() 函数是我们所定义的对节点的一系列操作。它包括：创建 ZooKeeper 节点（第 33 行到第 34 行代码）、查看节点（第 36 行到第 37 行代码）、修改节点数据（第 39 行到第 40 行代码）、查看修改后的节点数据（第 42 行到第 43 行代码）、删除节点（第 45 行到第 46 行代码）、查看节点是否存在（第 48 行到第 49 行代码）。另外，需要注意的是：在创建节点的时候，需要提供节点的名称、数据、权限，以及节点类型。此外，使用 exists 函数时，如果节点不存在将返回一个 null 值。关于 ZooKeeper API 的更多详细信息，读者可以查看 ZooKeeper 的 API 文档，地址如下所示：

<http://hadoop.apache.org/zookeeper/docs/r3.3.1/api/index.html>

代码清单 15-6 中程序运行的结果如下所示：

```
1. 创建 ZooKeeper 节点 (znode: zoo2, 数据: myData2, 权限: OPEN_ACL_UNSAFE, 节点类型: Persistent
11/01/18 05:07:16 INFO zookeeper.ClientCnxn: Socket connection established to
localhost/127.0.0.1:2181, initiating session
11/01/18 05:07:16 INFO zookeeper.ClientCnxn: Session establishment complete on
server localhost/127.0.0.1:2181, sessionId = 0x12d97fd5d39000a, negotiated timeout =
30000
```

```
WatchedEvent state:SyncConnected type:None path:null
```

2. 查看是否创建成功：

```
myData2
```

3. 修改节点数据

4. 查看是否修改成功：

```
shenlan211314
```

5. 删除节点

6. 查看节点是否被删除：

```
节点状态: [null]
```

15.4 ZooKeeper 的特性

15.4.1 ZooKeeper 的数据模型

ZooKeeper 拥有一个层次的命名空间，这和分布式的文件系统非常相似。唯一不同的地方是命名空间中的每个节点可以有和它自身或它的子节点相关联的数据。这就好像是一个文件系统，只不过文件系统中的文件还可以具有目录的功能。另外，指向节点的路径必须使用规范的绝对路径来表示，并且以斜线 “/” 来分隔。需要注意的是，在 ZooKeeper 中不允许使

用相对路径。

1. Znode

ZooKeeper 目录树中的每一个节点对应着一个 Znode。每个 Znode 维护着一个属性结构，它包含数据的版本号（dataVersion）、时间戳（ctime、mtime）等状态信息。ZooKeeper 正是使用节点的这些特性来实现它的某些特定功能的。每当 Znode 的数据改变时，它相应的版本号将会增加。每当客户端检索数据时，它将同时检索数据的版本号。并且若一个客户端执行了某个节点的更新或删除操作，它也必须提供要被操作的数据的版本号。如果所提供的数据版本号与实际的不匹配，那么这个操作将会失败。

Znode 是客户端要访问的 Zookeeper 的主要实体，它包含以下几个主要特征：

- ❑ **Watches**：客户端可以在节点上设置 watch（我们称其为监视器）。当节点的状态发生改变时（数据的增、删、改等操作）将会触发 watch 对应的操作。当 watch 被触发时，Zookeeper 将会向客户端发送仅且发送一个通知，因为 watch 只能被触发一次。
- ❑ **数据访问**：ZooKeeper 中的每个节点上存储的数据需要被原子性的操作。也就是说，读操作将获取与节点相关的所有数据，写操作也将替换掉节点的所有数据。另外，每一个节点都拥有自己的 ACL（访问控制列表），这个列表规定了用户的权限，即限定了特定用户对目标节点可以执行的操作。
- ❑ **临时节点**：ZooKeeper 中的节点有两种，分别为：临时节点和永久节点。节点的类型在创建时即被确定，并且不能改变。ZooKeeper 临时节点的生命周期依赖于创建它们的会话。一旦会话结束，临时节点将被自动删除，当然也可以手动删除。另外，需要注意的是，ZooKeeper 的临时节点不允许拥有孩子节点。相反，永久节点的生命周期不依赖于会话，并且只有在客户端显示执行删除操作的时候，它们才被删除。
- ❑ **顺序节点（唯一性保证）**：当创建 Znode 节点的时候，用户可以请求在 ZooKeeper 的路径结尾添加一个递增的计数。这个计数对于此节点的父节点来说是唯一的，它的格式为“%010d”（10 位数字，没有数值的数据位用 0 填充，例如 0000000001）。当计数值大于 $2^{32}-1$ 时，计数器将会溢出。

2. ZooKeeper 中的时间

ZooKeeper 中有多种记录时间的形式，其中包括如下几个主要属性：

- ❑ **Zxid**：致使 ZooKeeper 节点状态改变的每一个操作都将使节点接收到一个 zxid 格式的时间戳，并且这个时间戳是全局有序的。也就是说，每一个节点改变都将产生一个唯一的 zxid。如果 zxid1 的值小于 zxid2 的值，那么 zxid1 所对应的事件发生在 zxid2 所对应的事件之前。实际上，ZooKeeper 的每个节点维护着三个 zxid 值，分别为：cZxid、mZxid 和 pZxid。cZxid 是节点的创建时间所对应的 Zxid 格式时间戳，mZxid 是节点的修改时间所对应的 Zxid 格式时间戳。
- ❑ **版本号**：对节点的每一个操作都将致使这个节点版本号增加。每个节点维护着三个版本号，它们分别为：dataVersion（节点数据版本号）、cversion（子节点版本号）、

aclVersion (节点所拥有的 ACL 的版本号)。

3. 节点属性结构

通过上面的介绍，我们可以了解到，一个节点自身拥有表示其状态的许多重要属性，表 15-4 给出了详细的介绍：

表 15-4 ZooKeeper 节点属性

属 性	描 述
cZxid	节点被创建的 Zxid 值
mZxid	节点被修改的 Zxid 值
ctime	节点被创建的时间
mtime	节点最后一次的修改时间
dataVersion	节点被修改的版本号
cversion	节点所拥有的子节点被修改的版本号
aclVersion	节点的 ACL 被修改的版本号
ephemeralOwner	如果此节点为临时节点，那么它的值为这个节点拥有者的会话 ID；否则，它的值为 0
dataLength	节点数据域的长度
numChildren	节点拥有的孩子节点个数

15.4.2 ZooKeeper 会话及状态

ZooKeeper 客户端通过句柄为 ZooKeeper 服务建立一个会话。这个会话一旦被创建，句柄将以 CONNECTING 状态开始启动。客户端将尝试连接到其中一个 ZooKeeper 服务器，如果连接成功，它的状态变为 CONNECTED。在一般情况下只有上述这两种状态。如果一个可恢复的错误发生，比如会话终结或认证失败，或者如果应用程序明确地关闭了句柄，句柄将会转入关闭状态。

ZooKeeper 的状态转换如图 15-7 所示：

为了创建一个客户端会话，应用程序必须提供一个由主机 (IP 或主机名) 和端口所组成的连接字符串，这个字符串标识了要连接的目标主机及主机端口。ZooKeeper 客户端将选择服务器列表中的任意一个服务器并尝试连接。如果连接失败，那么客户端将自动地尝试连接服务列表中的其他服务器，直到连接成功。

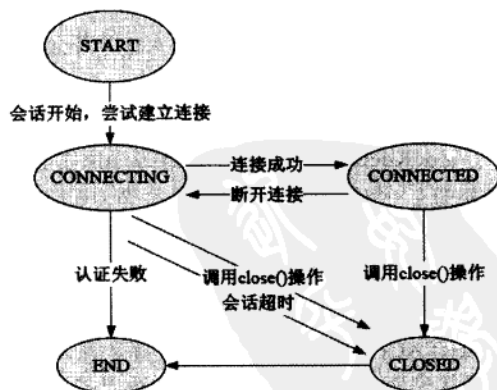


图 15-7 ZooKeeper 状态转换图

15.4.3 ZooKeeper Watches

ZooKeeper 可以为所有的读操作设置 watch，这些读操作包括：exists()、getChildren() 及 getData()。watch 事件是一次性的触发器，当 watch 的对象状态发生改变时，将会触发此对象上 watch 所对应的事件。

在使用 watch 时需要注意，watch 是一次性触发器，并且只有在数据发生改变时，watch 事件才会被发送给客户端。例如：如果一个客户端进行了 getData("/znode1,true) 操作，并且之后 "/znode1" 的数据被改变或删除了，那么客户端将获得一个关于 "/znode1" 的事件。如果 /znode1 再次改变，那么将不再有 watch 事件发送给客户端，除非客户端为另一个读操作重新设置了一个 watch。

watch 事件将被异步地发送给客户端，并且 ZooKeeper 为 watch 机制提供了有序的一致性保证。理论上，客户端接收 watch 事件的时间要快于其看到 watch 对象状态变化的时间。

ZooKeeper 所管理的 watch 可以分为两类：一类是数据 watch (data watches)；另一类是孩子 watch (child watches)。getData() 和 exists() 负责设置数据 watch，getChildren() 负责设置孩子 watch。我们可以通过操作返回的数据来设置不同的 watch。getData() 和 exists() 返回关于节点数据的信息，getChildren() 返回孩子列表。因此，setData() 将触发设置了数据 watch 的对应事件。一个成功的 create() 操作将触发 Znode 的数据 watch，以及孩子 watch。一个成功的 delete() 操作将触发数据 watch 和孩子 watch，因为 Znode 被删除的时候，它的 child watch 也将被删除。

Watch 由客户端所连接的 ZooKeeper 服务器在本地维护，因此 watch 可以非常容易地设置、管理和分派。当客户端连接到一个新的服务器时，任何的会话事件都将可能触发 watch。另外，当从服务器断开连接的时候，watch 将不会被接收。但是，当一个客户端重新建立连接的时候，任何先前注册过的 watch 都会被重新注册。

15.4.4 ZooKeeper ACL

ZooKeeper 使用 ACL 来对 Znode 进行访问控制。ACL 的实现和 UNIX 文件访问许可非常相似：它使用许可位来对一个节点的不同操作进行允许或禁止的权限控制。但是，和标准的 UNIX 许可不同的是，ZooKeeper 节点有 user（文件的拥有者）、group 和 world 三种标准模式，并且没有节点所有者的概念。

需要注意的是，一个 ACL 和一个 ZooKeeper 节点相对应。并且，父节点的 ACL 与孩子节点的 ACL 是相互独立的。也就是说，ACL 不能被孩子节点所继承，父节点所拥有的权限与孩子节点所拥有的权限没有任何关系。

表 15-5 为访问控制列表所规定的权限。

ZooKeeper ACL 的使用依赖于验证，它支持如下几种验证模式：

- ❑ world：代表某一特定的用户（客户端）。
- ❑ auth：代表任何已经通过验证的用户（客户端）。

□ digest：通过用户名密码进行验证。

□ ip：通过客户端 IP 地址进行验证。

表 15-5 ACL 权限

权 限	权 限 描 述
CREATE (创建)	创建孩子节点
READ (读)	从节点获取数据或列出节点的所有孩子节点
WRITE (写)	设置节点的数据
DELETE (删除)	删除孩子节点
ADMIN (管理员)	可以设置权限

当会话建立的时候，客户端将会进行自我验证。

另外，ZooKeeper Java API 支持三种标准用户权限，它们分别为：

```
ZOO_OPEN_ACL_UNSAFE;
ZOO_READ_ACL_UNSAFE;
ZOO_CREATOR_ALL_ACL;
```

ZOO_OPEN_ACL_UNSAFE 对于所有的 ACL 来说都是完全开放的：任何应用程序都可以在节点上执行任何操作，比如可以创建、列出并删除孩子。ZOO_READ_ACL_UNSAFE 对于任意的应用程序来说，仅仅具有读权限。ZOO_CREATOR_ALL_ACL 将授予节点创建者所有的权限。需要注意的是，在设置此权限之前，创建者必须已经通过了服务器的认证。

15.4.5 ZooKeeper 的一致性保证

ZooKeeper 是一种高性能、可扩展的服务。ZooKeeper 的读写速度非常快，并且读的速度要比写更快。另外，在进行读操作的时候，ZooKeeper 依然能够为旧的数据提供服务。这些都是由 ZooKeeper 所提供的一致性保证的，它具有如下特点：

(1) 顺序一致性

客户端的更新顺序与它们被发送的顺序相一致。

(2) 原子性

更新操作要么成功要么失败，没有第三种结果。

(3) 单系统镜像

无论客户端连接到哪一个服务器，他将看到相同的 ZooKeeper 视图。

(4) 可靠性

一旦一个更新操作被应用，那么在客户端再次更新它之前，其值将不会改变。这会保证产生下面两种结果：

1) 如果客户端成功地获得了正确的返回代码，那么说明更新已经成功。如果不能够获得返回代码（由于通信错误、超时等原因），那么客户端将不知道更新操作是否生效。

2) 当故障恢复的时候, 任何客户端都能看到的执行成功的更新操作将不会回滚。

(5) 实时性

在特定的一段时间内, 客户端看到的系统需要被保证是实时的(在十几秒的时间里)。在此时间段内, 任何系统的改变将被客户端看到, 或者被客户端侦测到。

这些一致性得到保证后, ZooKeeper 更高级功能的设计与实现将会变得非常容易, 例如: leader 选举、队列, 以及可撤销锁等机制的实现。

15.5 ZooKeeper 的 Leader 选举

ZooKeeper 需要在所有的服务(可以理解为服务器)中选举出一个 Leader, 然后让这个 Leader 来负责管理集群。此时, 集群中的其他服务器则成为此 Leader 的 Follower。并且, 当 Leader 出现故障的时候, ZooKeeper 要能够快速地在 Follower 中选举出下一个 Leader。这就是 ZooKeeper 的 Leader 机制, 下面我们将简单介绍在 ZooKeeper 中, Leader 选举(Leader Election)是如何实现的。

此操作实现的核心思想是: 首先创建一个 EPHEMERAL 目录节点, 例如 “/election”。然后每一个 ZooKeeper 服务器在此目录下创建一个 SEQUENCE|EPHEMERAL 类型的节点, 例如 “/election/n_”。在 SEQUENCE 标志下, ZooKeeper 将自动地为每一个 ZooKeeper 服务器分配一个比前面所分配的序号要大的序号。此时创建节点的 ZooKeeper 服务器中拥有最小编号的服务器将成为 Leader。

在实际的操作中, 还需要保障: 当 Leader 服务器发生故障的时候, 系统能够快速地选出下一个 ZooKeeper 服务器作为 Leader。一个简单的解决方案是, 让所有的 Follower 监视 Leader 所对应的节点。当 Leader 发生故障时, Leader 所对应的临时节点会被自动删除, 此操作将会触发所有监视 Leader 服务器的 watch。这样这些服务器将会收到 Leader 出现故障的消息, 并进而进行下一次的 Leader 选举操作。但是, 这种操作将会导致“从众效应”的发生, 尤其是当集群中服务器众多并且带宽延迟比较大的时候更为明显。

在 ZooKeeper 中, 为了避免从众效应的发生, 它是这样来实现的: 每一个 Follower 为 Follower 集群中对应着比自己节点序号小的节点中序号最大的节点设置一个 watch。只有当 Follower 所设置的 watch 被触发时, 它才进行 Leader 选举操作, 一般情况下它将成为集群中的下一个 Leader。很明显, 此 Leader 选举操作的速度是很快的。因为每一次 Leader 选举几乎只涉及单个 Follower 的操作。

15.6 ZooKeeper 锁服务

在 ZooKeeper 中, 完全分布的锁是全局同步的。这也就是说, 在同一时刻, 不会有两个不同的客户端认为他们持有了相同的锁。这一节将向大家介绍在 ZooKeeper 中的各种锁机制是如何实现的。

15.6.1 ZooKeeper 中的锁机制

ZooKeeper 将按照如下方式实现加锁的操作：

1) ZooKeeper 调用 create() 方法来创建一个路径格式为 “_locknode_/lock-” 的节点，此节点的类型为 sequence（连续）和 ephemeral（临时）。也就是说，创建的节点为临时节点，并且所有的节点连续编号，即为 “lock-i” 的格式。

2) 在创建的锁节点上调用 getChildren() 方法，以获取锁目录下的最小编号节点，并且不设置 watch。

3) 步骤 2 中获取的节点恰好是步骤 1 中客户端创建的节点，那么此客户端会获得此种类型的锁，然后退出操作。

4) 客户端在锁目录上调用 exists() 方法，并且设置 watch 来监视锁目录下序号相对自己次小的连续临时节点的状态。

5) 如果监视节点的状态发生变化，则跳转到第 2 步，继续进行后续操作，直到退出锁竞争。

ZooKeeper 的解锁操作非常简单，客户端只需要将加锁操作步骤 1 中创建的临时节点删除即可。

注意 1) 一个客户端解锁之后，将只可能有一个客户端获得锁，因此每一个临时的连续节点对应着一个客户端，并且节点之间没有重叠；2) 在 ZooKeeper 的锁机制中没有轮询和超时。

ZooKeeper 中锁机制流程图如图 15-8 所示。

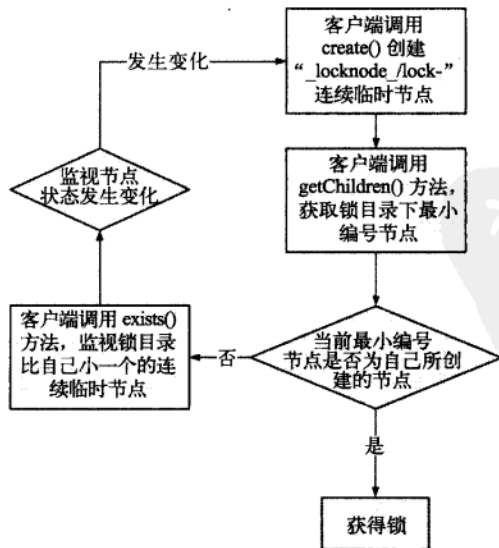


图 15-8 ZooKeeper 锁机制流程图

15.6.2 ZooKeeper 提供的一个写锁的实现

在 ZooKeeper 安装目录的 recipes 目录下有一个 ZooKeeper 分布式写锁的实现方式 (ZooKeeper_Dir/src/recipes/lock 目录)。

其中, 加锁的实现如代码清单 15-7 所示。

代码清单 15-7 lock

```

1. public synchronized boolean lock() throws KeeperException, InterruptedException {
2.     if (isClosed()) {
3.         return false;
4.     }
5.     ensurePathExists(dir);
6.
7.     return (Boolean) retryOperation(zop);
8. }
```

在加锁操作的实现中, 首先调用 isClosed() 方法来检查锁的状态, 如果没有获得锁, 则调用 ensurePathExists() 方法来设置一个监视器。这正如我们在 15.6.1 节中所描述的步骤。

解锁的实现如代码清单 15-8 所示。

代码清单 15-8 unlock

```

1. public synchronized void unlock() throws RuntimeException {
2.
3.     if (!isClosed() && id != null) {
4.         try {
5.
6.             ZooKeeperOperation zopdel = new ZooKeeperOperation(){
7.                 public boolean execute() throws KeeperException,
8.                     InterruptedException
9.                 {
10.                    zookeeper.delete(id, -1);
11.                    return Boolean.TRUE;
12.                }
13.            };
14.            zopdel.execute();
15.        } catch (InterruptedException e) {
16.            LOG.warn("Caught: " + e, e);
17.            //set that we have been interrupted.
18.            Thread.currentThread().interrupt();
19.        } catch (KeeperException.NoNodeException e) {
20.
21.        } catch (KeeperException e) {
22.            LOG.warn("Caught: " + e, e);
23.            throw (RuntimeException) new RuntimeException(e.getMessage()).
24.                initCause(e);
25.        }
26.    }
27. }
```

```

25.         finally {
26.             if (callback != null) {
27.                 callback.lockReleased();
28.             }
29.             id = null;
30.         }
31.     }
32. }

```

解锁的操作主要是通过代码中的第 6 ~ 12 行来实现的，只需要删除锁对应的临时节点即可。

注意 当此操作出现故障的时候，我们不需要重复这个解锁操作。另外，在不能重新连接的时候，我们也不需要做任何处理，因为 ZooKeeper 会自动地删除临时节点，并且在服务器出现故障的时候，此临时节点也会随着服务的结束而自动删除。

15.7 使用 ZooKeeper 创建应用程序

本节将通过一组简单的 ZooKeeper 应用程序实例来向读者展示 ZooKeeper 的某些功能。这一节所实现的主要功能包括：创建组、加入组、列出组成员，以及删除组。

为了避免某些重复性的操作，我们创建了一个本应用程序的基类：ZooKeeperInstance。它主要实现了 Zookeeper 对象的实例化操作。详见代码清单 15-9：

代码清单 15-9 ZooKeeperInstance

```

1. package app;
2.
3. import java.io.IOException;
4.
5. import org.apache.zookeeper.WatchedEvent;
6. import org.apache.zookeeper.Watcher;
7. import org.apache.zookeeper.ZooKeeper;
8.
9. public class ZooKeeperInstance {
10.     // 会话超时时间，设置为与系统默认时间一致
11.     public static final int SESSION_TIMEOUT=30000;
12.
13. // 创建 ZooKeeper 实例
14. ZooKeeper zk;
15.
16. // 创建 Watcher 实例
17. Watcher wh=new Watcher(){
18.     public void process(WatchedEvent event){
19.         System.out.println(event.toString());
20.     }

```



```

21.     };
22.
23. // 初始化 Zookeeper 实例
24. public void createZKInstance() throws IOException{
25.     zk=new ZooKeeper("localhost:2181",ZooKeeperInstance.SESSION_
        TIMEOUT,this.wh);
26. }
27.
28. // 关闭 ZK 实例
29. public void ZKclose() throws InterruptedException{
30.     zk.close();
31. }
32. }

```

在 ZooKeeper 中的组机制也同样是通过 ZooKeeper 节点来实现的。一个 Znode 为一个目录，即代表着一个组。这里我们创建了一个名为“/ZKGroup”的组。详见代码清单 15-10：

代码清单 15-10 CreateGroup

```

1. package app;
2.
3. import java.io.IOException;
4.
5. import org.apache.zookeeper.CreateMode;
6. import org.apache.zookeeper.KeeperException;
7. import org.apache.zookeeper.ZooDefs.Ids;
8.
9. public class CreateGroup extends ZooKeeperInstance {
10.
11. // 创建组
12. // 参数: groupPath
13. public void createPNode(String groupPath) throws KeeperException,
        InterruptedException{
14.     // 创建组
15.     String cGroupPath=zk.create(groupPath, "group".getBytes(), Ids.OPEN_
        ACL_UNSAFE, CreateMode.PERSISTENT);
16.     // 输出组路径
17.     System.out.println(" 创建的组路径为: "+cGroupPath);
18. }
19.
20. public static void main(String[] args) throws IOException, KeeperException,
        InterruptedException{
21.     CreateGroup cg=new CreateGroup();
22.     cg.createZKInstance();
23.     cg.createPNode("/ZKGroup");
24.     cg.ZKclose();
25. }
26. }

```

在创建组操作完成之后，我们需要将成员加入到组当中，也就是将节点加入到组节点的目录下，成为其孩子节点。在创建节点之前，首先需要调用 exists() 函数来判断组目录是否存在。此程序中创建了一个 MultiJoin() 函数，它通过一个计数器为组节点创建了 10 个孩子节点。详见代码清单 15-11：

代码清单 15-11 JoinGroup

```

1. package app;
2.
3. import java.io.IOException;
4.
5. import org.apache.zookeeper.CreateMode;
6. import org.apache.zookeeper.KeeperException;
7. import org.apache.zookeeper.ZooDefs.Ids;
8.
9. public class JoinGroup extends ZooKeeperInstance {
10. // 加入组操作
11. public int Join(String groupPath,int k) throws KeeperException,
    InterruptedException{
12.     String child=k+"";
13.     child="child_"+child;
14.
15.     // 创建的路径
16.     String path=groupPath+"/"+child;
17.     // 检查组是否存在
18.     if(zk.exists(groupPath,true) != null){
19.         // 如果存在，加入组
20.         zk.create(path,child.getBytes(), Ids.OPEN_ACL_UNSAFE,
            CreateMode.PERSISTENT);
21.         return 1;
22.     }
23.     else{
24.         System.out.println(" 组不存在! ");
25.         return 0;
26.     }
27. }
28.
29. // 加入组操作
30. public void MultiJoin() throws KeeperException, InterruptedException{
31.     for(int i=0;i<10;i++){
32.         int k=Join("/ZKGroup",i);
33.         // 如果组不存在则退出
34.         if(0==k)
35.             System.exit(1);
36.     }
37. }
38. public static void main(String[] args) throws IOException, KeeperException,
    InterruptedException{

```

```

39.      JoinGroup jg=new JoinGroup();
40.      jg.createZKInstance();
41.      jg.MultiJoin();
42.      jg.ZKclose();
43.  }
44.}

```

在加入组操作完成之后，我们通过 `getChildren()` 函数来列出所有组的成员（即获取组目录下的所有孩子节点）。详见代码清单 15-12：

代码清单 15-12 ListMembers

```

1. package app;
2.
3. import java.io.IOException;
4. import java.util.List;
5.
6. import org.apache.zookeeper.KeeperException;
7.
8. public class ListMembers extends ZooKeeperInstance {
9.     public void list(String groupPath) throws KeeperException, InterruptedException{
10.         // 获取所有子节点
11.         List<String> children=zk.getChildren(groupPath, false);
12.         if(children.isEmpty()){
13.             System.out.println("组 "+groupPath+" 中没有组成员存在！");
14.             System.exit(1);
15.         }
16.         for(String child:children)
17.             System.out.println(child);
18.     }
19.
20. public static void main(String[] args) throws IOException, KeeperException,
    InterruptedException{
21.     ListMembers lm=new ListMembers();
22.     lm.createZKInstance();
23.     lm.list("/ZKGroup");
24. }
25.}

```

在执行删除组操作时，我们首先需要删除组目录下的所有成员，当组目录为空时，再将组目录删除。那么，这过程中首先就需要调用 `getChildren()` 函数，获取组目录的所有成员，然后调用 `delete()` 函数将其一一删除，最后删除组目录。详见代码清单 15-13：

代码清单 15-13 DelGroup

```

1. package app;
2.
3. import java.io.IOException;
4. import java.util.List;

```

```

5.
6. import org.apache.zookeeper.KeeperException;
7.
8. public class DelGroup extends ZooKeeperInstance {
9. public void delete(String groupPath) throws KeeperException, InterruptedException{
10.     List<String> children=zk.getChildren(groupPath, false);
11.     // 如果不空, 则进行删除操作
12.     if(!children.isEmpty()){
13.         // 删除所有孩子节点
14.         for(String child:children)
15.             zk.delete(groupPath+"/"+child, -1);
16.     }
17.     // 删除组目录节点
18.     zk.delete(groupPath, -1);
19. }
20.
21. public static void main(String args[]) throws IOException, KeeperException,
    InterruptedException{
22.     DelGroup dg=new DelGroup();
23.     dg.createZKInstance();
24.     dg.delete("/ZKGroup");
25.     dg.ZKclose();
26. }
27.}

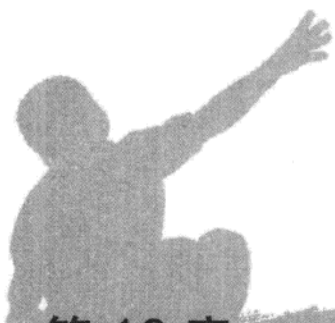
```

限于篇幅, 本章只介绍了关于 ZooKeeper 的一些基本知识, 希望读者通过本章的学习能够对 ZooKeeper 的机制有一个全面的了解。另外, 希望读者能够亲自动手编写 ZooKeeper 程序, 这样可以促进对 ZooKeeper 更深入的了解。

15.8 小结

ZooKeeper 作为 Hadoop 项目的一个子项目, 是 Hadoop 集群管理中一个必不可少的模块。它主要用来控制集群中的数据, 如管理 Hadoop 集群中的 NameNode, 以及 HBase 中的 Master Election、Server 之间的状态同步等。除此之外, 它还在其他多种场合中发挥着重要的作用。

本章介绍了 ZooKeeper 的基本知识, 以及 ZooKeeper 的配置、使用和管理等内容。另外还深入挖掘了 ZooKeeper 重要功能的实现机制, 并介绍了它的某些应用场景。ZooKeeper 作为一个用于协调分布式程序的服务, 必将在更多的场合发挥越来越重要的作用。



第 16 章

Avro 详解

本章内容

- ☐ Avro 简介
- ☐ Avro 的 C/C++ 实现
- ☐ Avro 的 Java 实现
- ☐ GenAvro (Avro IDL) 语言
- ☐ Avro SASL 概述
- ☐ 小结

资源分享

PDG

16.1 Avro 简介

Avro 作为 Hadoop 下相对独立的子项目，是一个数据序列化的系统。类似于其他序列化系统，Avro 可以将数据结构或对象转化成便于存储或传输的格式，特别是设计之初它可以用来支持数据密集型应用，适合于大规模数据的存储和交换。总之，Avro 可以提供以下一些特性和功能：

- ❑ 丰富的数据结构类型；
- ❑ 快速可压缩的二进制数据形式；
- ❑ 存储持久数据的文件容器；
- ❑ 远程过程调用（RPC）；
- ❑ 简单的动态语言结合功能。

Avro 和动态语言结合后，读写数据文件和使用 RPC 协议都不需要生成代码了，代码作为一种可选的优化只需要在静态类型语言中实现。

Avro 依赖于模式（Schema）。Avro 数据的读 / 写操作很频繁，而这些操作都需要使用模式，这样可减少写入每个数据资料的开销，使得序列化快速而又轻巧。这种数据及其模式的自我描述方便了动态脚本语言的使用。

当 Avro 数据存储到文件中时，它的模式也会随之存储，这样任何程序就都可以对文件进行了。如果读取数据时使用的模式与写入数据时使用的模式不同，那也很容易解决，因为读取和写入的模式都是已知的。图 16-1 表示的是 Avro 的主要作用，它将用户定义的模式和具体的数据编码成二进制序列存储在对象容器文件中，假设用户定义了包含学号、姓名、院系和电话的学生模式，而 Avro 对其进行编码后存储在 student.db 文件中，其中存储数据的模式放在文件头的元数据中，这样即使读取的模式与写入的模式不同，也可以迅速地读出数据，如果另一个程序需要获取学生的姓名和电话，只需定义包含姓名和电话的学生模式，然后用此模式去读取容器文件中的数据即可。

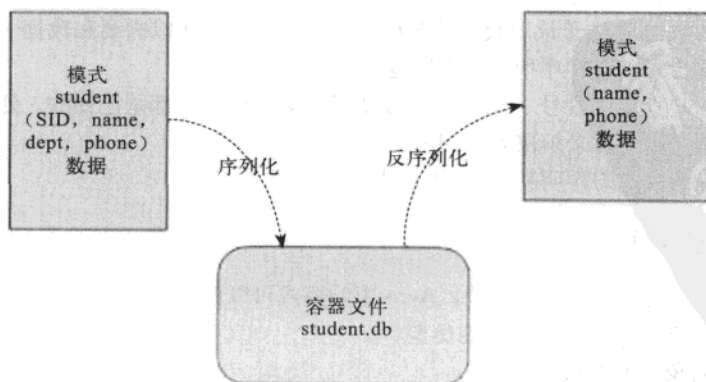


图 16-1 Avro 的主要作用

当在 RPC 中使用 Avro 时,服务器和客户端可以在握手连接时交换模式。服务器和客户端有彼此全部的模式,因此含有相同名字段、缺失字段和多余字段等的信息之间通信时,需要处理的一致性问题的就可以容易地解决。如图 16-2 所示,协议中定义了用于传输的消息,消息使用框架后放入缓冲区中进行传输,由于传输初始时就交换了各自的协议定义,因此即使传输双方使用的协议不同,所传输的数据也能够正确解析,具体过程将在后面介绍。

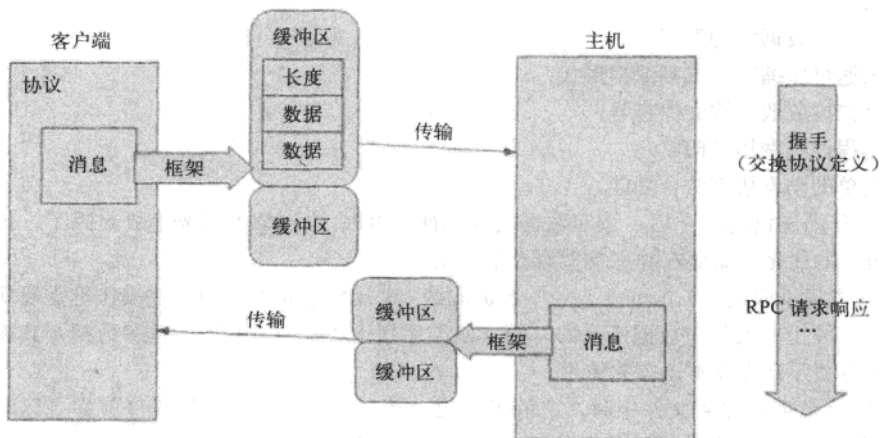


图 16-2 RPC 使用 Avro

Avro 模式是用 JSON (一种轻量级的数据交换模式) 定义的,这对于已经拥有 JSON 库的语言来说就可以容易地实现。

Avro 提供与诸如 Thrift 和 Protocol Buffers 等系统相似的功能,但是在一些基础方面还是有区别的,主要有:

- ❑ 动态类型: Avro 并不需要与生成代码、模式和数据存放在一起,而整个数据的处理过程并不生成代码、静态数据类型,等等。这方便了数据处理系统和语言的构造。
- ❑ 未标记的数据: 由于读取数据的时候模式是已知的,所以需要和数据一起编码的类型信息就很少了,这样序列化的规模也就小了。
- ❑ 不需要用户指定字段号: 即使模式改变了,新旧模式也都是已知的,处理数据时可以通过使用字段名称来解决差异问题。

下面详细介绍模式的声明和 Avro 的具体使用。

16.1.1 模式声明

模式声明主要是定义数据的类型,Avro 中的模式可以使用 JSON 通过以下方式来表示。

- 1) JSON 字符串,指定已定义的类型。
- 2) JSON 对象,其形式为:

```
{ "type": "typeName" ...attributes... }
```


其中, `typeName` 可以是原生的或衍生的类型名称, 本章没有定义的属性可以视为元数据, 但是其不能影响序列化数据的格式。

3) JSON 数组, 表示嵌入类型的联合。

声明的类型必须是 Avro 所支持的数据类型, 其中包括原始类型 (Primitive Types) 和复杂类型 (Complex Types), 下面分别介绍之。

原始类型的名称包括以下几部分。

- ❑ `null`: 没有值;
- ❑ `boolean`: 二进制值;
- ❑ `int`: 32 位有符号整数;
- ❑ `long`: 64 位有符号整数;
- ❑ `float`: 单精度 (32 位) IEEE 754 浮点数;
- ❑ `double`: 双精度 (64 位) IEEE 754 浮点数;
- ❑ `bytes`: 8 位无符号字节序列;
- ❑ `string`: unicode 字符序列。

原始类型没有特定的属性, 其名称可以通过类型来定义, 如模式 "string" 相当于:

```
{ "type": "string" }
```

Avro 支持 6 种复杂类型: 记录 (records)、枚举 (enums)、数组 (arrays)、映射 (maps)、联合 (unions) 和固定型 (fixed), 下面一一介绍。

(1) 记录

记录使用类型名称 "record" 并且支持以下属性。

- ❑ `name`: 提供记录名称的 JSON 字符串 (必须)。
- ❑ `namespace`: 限定名称的 JSON 字符串。
- ❑ `doc`: 向模式使用者提供说明的 JSON 字符串 (可选)。
- ❑ `aliases`: 字符串的 JSON 数组, 为记录提供代替名称 (可选)。
- ❑ `field`: 一个 JSON 数组, 用来列出字段 (必须)。每个字段就是一个 JSON 对象且拥有以下属性。
 - `name`: 提供记录名称的 JSON 字符串 (必须)。
 - `doc`: 为使用者提供字段说明的 JSON 字符串 (可选)。
 - `type`: 定义模式的 JSON 对象, 或者记录定义的 JSON 字符串 (必须)。
 - `default`: 该字段的默认值用于读取缺少该字段的实例 (可选)。如表 16-1 所示, 允许的值依赖于字段的模式类型。联合字段的默认值对应于联合中的第一个模式。字节和固定字段的默认值是 JSON 字符, 这里 0~255 的 Unicode 映射到 0~255 的 8 位无符号字节。
 - `order`: 指定该字段如何影响记录的排序 (可选)。有效的值有 "ascending" (默认)、"descending" 或 "ignore"。

❑ **aliases**：字符串的 JSON 数组，为该字段提供可选的名称（可选）。

表 16-1 字段默认值

Avro 类型	JSON 类型	例 子
null	null	null
boolean	boolean	true
int,long	integer	1
float,double	number	1.1
bytes	string	"\u00FF"
string	string	"foo"
record	object	{"a":1}
enum	string	"FOO"
array	array	[1]
map	object	{"a":1}
fixed	string	"\u00ff"

例如，一个 64 位的链表可以定义为：

```
{
  "type": "record",
  "name": "LongList",
  "aliases": ["LinkedLongs"],           // 别名
  "fields" : [
    { "name": "value", "type": "long" }, // 每个元素都含有长整型
    { "name": "next", "type": ["LongList", "null"] } // 下一元素
  ]
}
```

(2) 枚举

枚举使用类型名称“enum”并且支持以下几种类型。

- ❑ **name**：提供实例名称的 JSON 字符串（必须）。
- ❑ **namespace**：限定名称的 JSON 字符串。
- ❑ **aliases**：字符串的 JSON 数组，为枚举提供替代名称（可选）。
- ❑ **doc**：对模式使用者提供说明的 JSON 字符串（可选）。
- ❑ **symbols**：列出标记的 JSON 数组（必须）。枚举中的所有标记必须是唯一的，不允许有重复的标记。

例如，纸牌游戏可以定义为：

```
{ "type": "enum",
  "name": "Suit",
  "symbols" : ["SPADES", "HEARTS", "DIAMONDS", "CLUBS"]
}
```

(3) 数组

数组使用类型名称“array”并且支持一个属性。

□ items：数组项目的模式。

例如，字符串数组可以定义为：

```
{ "type": "array", "items": "string" }
```

(4) 映射

映射使用类型名称“map”并且支持一个属性。

□ values：映射值的模式。

映射值默认为字符串，例如，从字符串到长整型的映射可以声明为：

```
{ "type": "map", "values": "long" }
```

(5) 联合

联合主要使用 JSON 数组表示，例如可以用 ["string", "null"] 声明一个是字符串或 null 的模式，联合使用的类型名称为“unim”。除了指定的记录、固定型（fixed）和枚举外，对于相同的类型，联合只能包含一个模式。例如，联合中不允许包含两个数组类型或两个映射类型，但是允许包含不同名称的两种类型。联合中不能直接包含其他的联合。

(6) 固定型

固定型使用类型名称“fixed”并且支持以下属性。

□ name：固定型名称的字符串（必须）。

□ namespace：限定名称的字符串。

□ aliases：提供替代名称的字符串的 JSON 数组（可选）。

□ size：说明每个值的字节数的整型（可选）。

例如，16 字节大小的固定型可以声明为：

```
{ "type": "fixed", "size": 16, "name": "md5" }
```

记录、枚举和固定型都是指定的类型，其全名由两部分组成：名称和命名空间。全名为由点分开的名字序列，其名称部分和记录字段名字必须：

□ 以字母 A～Z 或 a～z 或 _ 开头；

□ 后面应只含有 A～Z、a～z、0～9 或 _。

在记录、枚举和固定型的定义中，全名可以通过以下几种方式定义：

□ 指定名称和命名空间。如使用名称 "name": "X" 和命名空间 "namespace": "org.foo" 来表示全名 org.foo.X。

□ 指定全名。如果指定的名称中包含点，则可以使用名称作为全名，并且任何指定的命名空间将被忽略。比如使用名称 "name": "org.foo.X" 来表示全名 org.foo.X。

□ 只指定名称，且名称中不包含点。这种情况下命名空间取自外层的模式或协议，比如指定名称 "name": "X"，其所在记录定义的字段则为 org.foo.Y，即全名为 org.foo.X。

总结以上后两种情况可得：如果名称中包含点，则是全名；如果不包含点，则命名空间默认为外层定义的命名空间。原始类型没有命名空间并且它们的名称也不能定义为任何命名空间。

命名的类型和字段可以拥有别名。为方便模式的发展和处理不同的数据集，在实现中可以选择使用别名将写者的模式（writer's schema）映射成读者的模式（reader's schema）。使用别名可以改变写者的模式，例如，如果写者的模式命名为 "Foo"，而读者的模式命名为 "bar" 且别名为 "Foo"，那么在读取时即使将 "Foo" 称作 "Bar" 也能实现。同理，如果数据曾经写成字段名为 "x" 的记录，那么即使是字段名为 "y" 别名为 "x" 的记录也能读取，尽管 "x" 写成了 "y"。

16.1.2 数据序列化

模式声明后就可以根据模式写入数据了，当数据存储或传输时需要对其序列化，需要注意的是 Avro 数据和其模式会一起被序列化。基于 Avro 的 RPC 系统必须保证远程的数据接收器拥有写入数据的模式副本，因为读取数据时写入数据的模式是知道的，所以 Avro 数据本身不需要标记类型信息。

通常，序列化和还原序列化过程（见图 16-1）可以看成是对模式深度优先、从左到右的遍历过程，并在遍历过程中序列化或还原序列化遇到的原始类型。

Avro 指定两种序列化编码：二进制和 JSON。在这两种序列化编码中，因为二进制编码速度快且生成的数据量小，所以大多数的应用程序使用二进制编码。但是基于调试和网络的应用程序，有时使用 JSON 编码比较合适。下面先介绍一下各种类型的二进制编码。

原始类型的二进制编码有如下几种。

□ null 编码成零字节。

□ boolean 编码成单字节，值为 0（false）或 1（true）。

int 和 long 使用可变长度的 ZigZag 编码^①，如表 16-2 所示。

表 16-2 ZigZag 编码举例

value (值)	hex (十六进制)
0	00
-1	01
1	02
-2	03
2	04
⋮	⋮
-64	7f
64	80 01
⋮	⋮

① ZigZag 编码原使用于 Protocol Buffers，是一种将有符号数映射成无符号数的编码方式。

- float 占 4 字节，使用类似于 Java 中 floatToIntBits 的方法可以将浮点数转化成 32 位的整型，然后编码成低字节序的格式。
- double 占 8 字节，使用类似于 Java 中 doubleToLongBits 的方法可以将双精度数转化为 64 位的整型，然后编码成低字节序的格式。
- bytes 根据数据的字节编码成长整型。
- string 根据 UTF-8 字符集编码成长整型。

如果 UTF-8 字符集中 'f'、'o'、'o' 的十六进制分别为 66 6f 6f，并且字符串 "foo" 含有 3（编码成十六进制 06）个字符，那么 "foo" 编码为 06 66 6f 6f。

复杂类型的二进制编码方法有如下几种。

（1）记录（records）

通过对声明的每个字段值按顺序编码来对记录进行编码。换句话说，记录的编码由每个字段的编码串联而成。例如，记录模式的代码如下：

```
{
  "type": "record",
  "name": "test",
  "fields": [
    { "name": "a", "type": "long" },
    { "name": "b", "type": "string" }
  ]
}
```

以上代码中，假设 a 字段的值为 27（十六进制为 36），b 字段的值为 "foo"（十六进制为 06 66 6f 6f），那么记录的编码仅仅是这些编码的串联，即十六进制序列 36 06 66 6f 6f。

（2）枚举（enums）

枚举按整型编码，其中整型代表每个标志在模式中的位置（从 0 开始）。例如，枚举模式的代码如下：

```
{ "type": "enum", "name": "Foo", "symbols": ["A", "B", "C", "D"] }
```

上面的例子序列化时将被编码成整型 0~3，其中 0 代表 "A"，3 代表 "D"。

（3）数组（arrays）

数组编码成一系列块，每个块包含一个长整型的数值，长整型的数值组成数组项，其中最后一块为 0 表示是数组的结尾，每个数组项的模式都会编码。如果块中的值为负数，则取绝对值，紧跟数值后面的块的大小为长整型，表示块中的字节数。如果只映射记录部分字段，则利用块大小可以跳过部分数据。例如，数组模式为：

```
{ "type": "array", "items": "long" }
```

对于包含项 3 和 27 的数组，其数组包含两个长整型值，其中数组个数“2”使用 ZigZag 编码成十六进制为 04，而 3 和 27 编码成十六进制分别为 06 和 36，最后以 0 结尾，其数组编码为：04 06 36 00。

这种块表示方法允许读 / 写其大小超过内存缓冲的数组，因为在不知道数组长度的情况下就可以开始写入。

(4) 映射 (maps)

映射的编码和数组相似，也是编码成一系列块，每个块包含一个长整型值，然后是键值对，值为 0 的块表示映射的结尾。如果块的值为负数，则取其绝对值。紧跟数值后面的块的大小为长整型，表示块中的字节数。如果只映射记录的部分字段，则利用块大小可以跳过部分数据。

同数组一样，在不知道映射长度的情况下就可以写入，因此这种块的表示方法也允许读 / 写大小超过内存缓冲的映射。

(5) 联合 (unions)

联合编码时先写入一个长整型值表示联合中每个模式值的位置（从 0 开始），再对联合中的值编码。例如，联合模式 ["string", "null"] 应如此编码：null 为整数 1（联合中 null 的索引，使用 ZigZag 编码成十六进制 02）；字符串 "a" 为整数 0（联合中 "string" 的索引）；随后为序列化的字符串 61，所以最后这个联合编码应为 00 02 61。

(6) 固定型 (fixed)

固定型实例编码可使用模式中声明的字节数。对于编码成 JSON 类型，除了联合以外，还可以参照表 16-1 中 JSON 类型与 Avro 类型的对应关系进行编码。联合的 JSON 编码如下所示：

- ❑ 如果值为 null，那么按照 JSON null 来编码。
 - ❑ 否则，按照带有名称 / 值的 JSON 对象进行编码，其中名称为类型名称，值为递归编码值。对于 Avro 的命名类型（记录、固定性和枚举）将使用用户指定的名称，对于其他类型将使用类型名。
- 例如，对于联合模式 ["null", "string", "Foo"]（其中 Foo 是记录名称），应如此编码：
- ❑ null 作为 null 编码；
 - ❑ 字符串 "a" 按照 {"string": "a"} 编码；
 - ❑ Foo 实例按照 {"Foo": {...}} 编码，这里 {...} 表示一个 Foo 实例的 JSON 编码。

需要注意的是，正确处理 JSON 编码数据仍需要模式，因为 JSON 编码并不区分 int 和 long、float 和 double、记录 (records) 和映射 (maps)、枚举 (enums) 和字符串 (strings) 等。

16.1.3 数据排列顺序

对象化前最常使用的操作就是排序，在 Avro 确定了数据标准排列顺序以后，就允许系统写入的数据被另外的系统高效地排序了，这是个很重要的优化。即使 Avro 二进制数据还没有反序列化成对象，也可以对其进行高效排序。

要对拥有相同模式的数据项进行比较，可以采用对模式的深度优先、从左到右递归遍历的方式。遇到不能匹配的项即按原来顺序，比如 boolean 类型的数据和 int 类型的数据不能匹配，那就不用进行排序。具体来说，相同模式的两个项进行比较时需遵从下面

的规则。

- ❑ null 数据总是相等的。
- ❑ boolean 类型中 false 排在 true 的前面。
- ❑ int、long、float 和 double 数据按照数值升序排列。
- ❑ bytes 和 fixed 数据根据 8 位无符号值按照字典序进行比较。
- ❑ string 数据根据 Unicode 按字典序进行比较，要注意的是，对字符串而言，既然 UTF-8 作为二进制编码使用，那么按字节排序和按字符串二进制数据排序是相同的。
- ❑ array 数据根据元素按字典序进行比较。
- ❑ enum 数据根据枚举模式中符号的位置进行排序。例如，枚举的符号位 ["z","a"] 把 "z" 排在 "a" 前面。
- ❑ union 数据先按照联合的分支进行排序，然后按照分支的类型排序。例如，联合 ["int","string"] 中，所有整型将排在所有字符型值前，而整型和字符型各自按照上面的规则排序。
- ❑ record 数据根据字段按字典序排序。如果字段指定其顺序为：
 - "ascending"：其值排序的顺序不变；
 - "descending"：其值排序的顺序反转；
 - "ignore"：排序时其值将忽略。
- ❑ map 数据不进行比较。试图比较包含映射的数据是非法的，除非映射是“有序”的，否则“忽略”记录字段。

16.1.4 对象容器文件

序列化后的数据需要存入文件中，Avro 包含一个简单的对象容器文件，一个文件拥有一个模式，文件中所有存储的对象必须根据模式使用二进制编码写入。对象存储在可以压缩的块中，块之间使用同步机制为 MapReduce 处理提供高效的文件分离。文件中可能包含用户随意指定的元数据。那么一个文件包含：

- ❑ 文件头。
 - ❑ 一个或多个文件数据块。
- 其中文件头包含：
- ❑ 4 个字节，分别是 ASCII 码的 o、b、j、l。
 - ❑ 包含模式的文件元数据。
 - ❑ 为此文件随机生成的 16 字节同步器。
- 文件的元数据包含：
- ❑ 指示元数据的一个键 / 值对的长整型。
 - ❑ 每个对的字符串键和字节值。

所有以 "Avro" 开头的元数据属性是保留的，以下文件元属性主要用于：

- ❑ avro.schema：包含存储在文件中的对象的模式，如 JSON 数据（必须）。



□ `avro.codec`：编解码器名称，其编码器用来压缩诸如字符串的数据块。需要实现支持 "null" 和 "deflate" 编解码器，如果没有编解码器，那假设为 "null"。

"null" 编解码器不需要对数据解压缩，而 "deflate" 编解码器使用文档 RFC1951 中指定的 deflate 算法写入数据块并使用 zlib 库实现，要注意的是这个格式（不像 RFC1950 中的 zlib 格式）没有校验和。

一个文件头需要按照如下模式进行描述：

```
{
  "type": "record", "name": "org.apache.avro.file.Header",
  "fields" : [
    { "name": "magic", "type": { "type": "fixed", "name": "Magic", "size": 4 } },
    { "name": "meta", "type": { "type": "map", "values": "bytes" } },
    { "name": "sync", "type": { "type": "fixed", "name": "Sync", "size": 16 } },
  ]
}
```

文件数据块包括：

- 一个长整型，用于指示块中对象的数目。
- 一个长整型，用于表示使用编解码器后，所在块中序列化对象的字节大小。
- 序列化对象，如果编解码器是指定的，则用它进行压缩。
- 16 字节的文件同步器。

这样，即使不用反序列化，每个块的二进制数据也可以高效获得或跳过。这种块的大小、对象数目和同步器的结合可以检测出坏的块并且帮助保持数据的完整性。

图 16-3 表示了对象容器文件的具体格式。

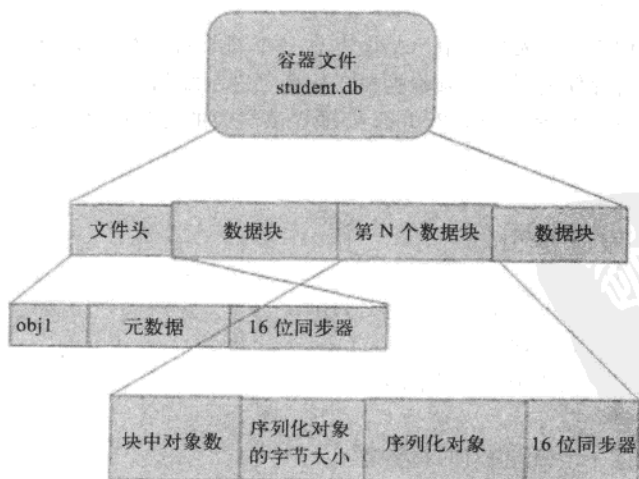


图 16-3 对象容器文件的具体格式

16.1.5 协议声明

当 Avro 用于 RPC 时, Avro 使用协议描述远程过程调用 RPC 接口。和模式一样, 它们是用 JSON 文本来定义的。协议是带有以下属性的 JSON 对象:

- protocol: 协议名称的字符串 (必须)。
- namespace: 限定名称的可选字符串。
- doc: 描述协议的可选字符串。
- types: 指定类型 (记录、枚举、固定型和错误) 定义的可选列表。错误的定义和记录一样, 只不过错误使用 "error" 而记录使用 "record", 要注意不允许对指定类型的向前引用。
- messages: 一个可选的 JSON 对象, 其键是消息名称, 值是对象, 任意两个消息不能拥有相同的名称。

模式中定义的名称和命名空间规则也同样适用于协议。下面介绍的协议消息可以拥有以下属性:

- doc: 消息的可选描述。
- request: 指定的类型化的参数模式列表 (这和记录声明中的字段有相同的形式)。
- response: 响应模式。
- error union: 所声明的错误模式的联合 (可选)。有效的联合会在声明的联合前面加上 "string", 允许传递未声明的“系统”错误。例如, 如果声明的错误联合是 ["AccessError"], 那么有效的联合是 ["string", "AccessError"]。如果没有错误声明, 那么有效的错误联合是 ["string"]。使用有效联合错误可以序列化, 且协议的 JSON 声明只能包含声明过的联合。
- one-way: 布尔参数 (可选)。

处理请求参数列表相当于处理没有名称的记录。既然读取的记录字段列表和写入的记录字段列表可以不同, 那么调用者和响应者的请求参数也可以不同, 这种区别的解决方法与记录字段间差异的解决方式相同。只有当回应的类型是“null”并且没有错误列出的时候, one-way 参数才为真。

下面来举一个简单的 HelloWorld 协议的例子, 它可以定义为:

```
{
  "namespace": "com.acme", // 名称的限定
  "protocol": "HelloWorld", // 协议名称
  "doc": "Protocol Greetings", // 协议的说明
  "types": [
    { "name": "Greeting", "type": "record", "fields": [
      { "name": "message", "type": "string" } ] },
    { "name": "Curse", "type": "error", "fields": [
      { "name": "message", "type": "string" } ] } ],
  "messages": { // 消息
    "hello": {
```

```

"doc": "Say hello.",
"request": [{"name": "greeting", "type": "Greeting" }],
"response": "Greeting",
"errors": ["Curse"]
}
}
}

```

16.1.6 协议传输格式

消息可以通过不同的传输机制进行传输，而传输中的消息则是一些字节序列，那么传输机制需要支持：

- 请求信息的传送。
- 响应信息的接收。

服务器会对客户机的请求信息发送响应信息，这种响应机制就是特定传输，例如在 HTTP 中，由于 HTTP 直接支持请求和响应，所以这种传输是透明的，但是利用同一套接字传输多种不同客户线程的时候需要用特定的标识来区分不同客户的信息。

传输可能是无状态的也可能是有状态的。在无状态传输中，是假定消息发送没有建立连接状态。而有状态传输则建立了连接，这个连接可以用来传输不同的消息。下面我们会在握手（handshake）部分深入分析。

当用 HTTP 协议进行传输时，每个 Avro 消息交换都是一对 HTTP 请求/响应。一个 Avro 协议的所有消息共享一个 HTTP 服务器上的 URL，正常的和错误的 Avro 消息都应该使用 200（OK）响应代码。尽管 Avro 请求和响应是 HTTP 请求和响应的整个内容，但也可能使用大量的编码。HTTP 请求和响应的内容类型应该指定为“avro/binary”而且请求应该使用 POST 方法生成。Avro 使用 HTTP 作为无状态传输。

Avro 消息经过框架处理后由一系列缓冲区组成，消息框架是消息和传输之间的一层，用来优化某些操作。经过框架处理后的消息数据格式如下（见图 16-4）。

- 由一系列缓冲区组成，其中缓冲区包括：
 - 4 个字节，用大端字节（big-endian）方法[⊖]表示的缓冲区长度。
 - 缓冲区数据。
- 最后以空字节（zero-length）的缓冲区结束。

对于请求和响应消息格式，框架是透明的，任何消息可以表示为一个或多个缓冲。框架使得消息接收者更高效地从不同的渠道获取不同的缓冲，也使得开发者更高效地向不同的目的地存储不同的缓冲。特别是当复制大量二进制对象时，它可以减少读/写的次数。例如，如果 RPC 参数中包含一个 MB 大小的文件数据，那么一方面，数据可以从文件描述符直接复制到套接字上；另一方面，数据可以直接写入文件描述符而不需要进入用户空间。

⊖ 存放字节顺序的方法，大端方式将高位存放在低地址，小端方式将高位存放在高地址。

一个简单且值得推荐的框架策略是相对于那些大于正常输出缓冲区的单个二进制对象建立的新段的。小的对象可以附加在缓冲区中，而较大的对象可以写入自己的缓冲区中。当读者需要读取大的对象时，可以直接处理整个缓冲区而不用复制。

使用握手的目的是确保客户机和服务器均有对方的协议定义，这样客户机可以正确地对响应反序列化，而服务器可以正确地对请求反序列化。客户机和服务器都应在高速缓冲区中保留最近的协议，这样在大多数情况下，可以不用多余往返网络交换或全部传输协议就能完成握手。

在完成握手过程后执行 RPC 请求和响应，对于无状态的传输，在所有请求和响应之前都要进行握手，而对于有状态的传输，在成功响应之前，握手过程应该附加在请求和响应上，之后就不需要握手了。

握手过程使用以下记录模式，代码如下：

```
{
  "type": "record",
  "name": "HandshakeRequest", "namespace": "org.apache.avro.ipc",
  "fields": [
    { "name": "clientHash",
      "type": { "type": "fixed", "name": "MD5", "size": 16 } },
    { "name": "clientProtocol", "type": [ "null", "string" ] },
    { "name": "serverHash", "type": "MD5" },
    { "name": "meta", "type": [ "null", { "type": "map", "values": "bytes" } ] }
  ]
}
{
  "type": "record",
  "name": "HandshakeResponse", "namespace": "org.apache.avro.ipc",
  "fields": [
    { "name": "match",
      "type": { "type": "enum", "name": "HandshakeMatch",
        "symbols": [ "BOTH", "CLIENT", "NONE" ] } },
    { "name": "serverProtocol",
      "type": [ "null", "string" ] },
    { "name": "serverHash",
      "type": [ "null", { "type": "fixed", "name": "MD5", "size": 16 } ] },
    { "name": "meta", "type": [ "null", { "type": "map", "values": "bytes" } ] }
  ]
}
```

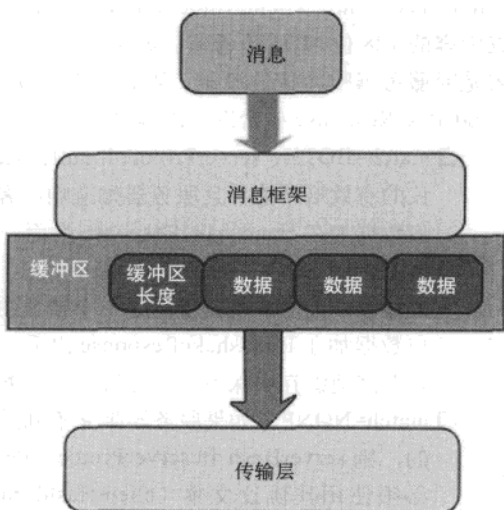


图 16-4 消息的封装

客户机在每个请求前面加上 HandshakeRequest, 表示包含客户机和服务器协议 (clientHash!=null, clientProtocol=null, serverHash!=null) 的哈希值, 这里哈希值是 JSON 协议内容的 128 位 MD5 哈希值。如果客户机没有连接到给定的服务器, 那么它发送的哈希值就是对服务器哈希值的猜测, 否则它会发送之前从服务器中获得的哈希值。服务器响应的 HandshakeResponse 包含以下内容之一。

- ❑ match=BOTH, serverProtocol=null, serverHash=null。如果客户机发送的是服务器协议的有效哈希值并且服务器知道响应客户机哈希值的协议, 那么请求是完整的并且响应数据加在 HandshakeResponse 后面。
- ❑ match=CLIENT, serverProtocol!=null, serverHash!=null。如果服务器事先知道客户机的协议, 而客户机却发送了一个错误的服务器协议哈希值, 那么请求是完整的并且响应数据加在 HandshakeResponse 之后。之后客户机必须使用返回的协议来处理响应并且在高速缓存中保留这个协议和与服务器通信的哈希值。
- ❑ match=NONE。如果服务器事先不知道客户机的协议, 且服务器的协议哈希值是错误的, 则 serverHash 和 serverProtocol 的值可能也为 non-null。在这种情况下, 客户机必须使用其协议文本 (clientHash!=null, clientProtocol!=null, serverHash!=null) 重新提交它的请求并且服务器应该以正确的方式响应 (match=BOTH, serverProtocol=null, serverHash=null)。还有 meta 字段是保留字段, 用于以后增加握手的功能。

一次调用包括请求消息和与之对应的结果响应或错误消息。请求和响应包含可扩展的元数据, 两种消息都会如上进行框架处理。调用请求的格式包括以下几种:

- ❑ 请求元数据, 即类型值的映射。
- ❑ 消息名称, 即一个 Avro 字符串。
- ❑ 消息参数, 根据消息请求声明对参数进行序列化。

当消息声明为单向的并通过成功握手响应建立有状态连接的, 就不需要发送响应数据了。否则需要发送, 发送的调用请求的格式如下:

- ❑ 响应元数据, 即类型字节的映射。
- ❑ 单字节的错误标志布尔值, 然后:
 - 如果错误标志为假, 消息响应, 序列化每个消息响应模式。
 - 如果错误标志为真, 即为错误, 序列化每个消息有效错误联合模式。

16.1.7 模式解析

无论是从 RPC 还是文件中获得 Avro 数据, 由于模式已知, 读者都可以解析数据, 但是那个模式可能并不完全就是所期望的模式。例如, 如果数据写入的软件版本与读者不同, 那么记录中的一些字段可能会增加或减少。下面将详述如何解决这种模式区别。

我们称用来写数据的模式为写者的模式, 应用程序期望的模式为读者的模式。两个模式之间是否匹配可按照下面的规则进行判断。

1) 如果两个模式符合以下情况之一则为匹配, 否则为不匹配并产生错误:

- ☐ 模式都是数组且项类型匹配。
- ☐ 模式都是映射且值类型匹配。
- ☐ 模式都是枚举且名称匹配。
- ☐ 模式都是固定型且大小和名称匹配。
- ☐ 模式都是记录且名称相同。
- ☐ 模式是其中之一为联合。
- ☐ 两个模式拥有相同的原始类型。
- ☐ 写者的模式可以提升为读者的模式, 如下所示:
 - int 可以转化为 long、float 或 double ;
 - long 可以转化为 float 或 double ;
 - float 可以转化为 double。

2) 如果两个都记录, 则

- ☐ 字段的顺序可以不同, 因为字段是通过名称来匹配的。
- ☐ 有相同名称字段的模式记录是递归解析的。
- ☐ 如果写者的记录中包含读者记录中没有的字段, 那么写者字段的值将被忽略。
- ☐ 如果读者记录模式中有一个为默认值的字段, 并且写者的模式中没有相同名称的字段, 那么读者的这个字段应该使用默认值。
- ☐ 如果读者记录模式中有一个没有默认值的字段, 并且写者的模式中没有相同名称的字段, 那么将发出错误信号。

3) 如果两个都是枚举, 且写者的符号并不在读者的枚举中, 那么产生错误。

4) 如果两个都是数组, 解析算法递归应用于读者和写者数组项的模式。

5) 如果两个都是映射, 解析算法递归应用于读者和写者映射值的模式。

6) 如果两个都是联合, 对读者联合中匹配写者联合模式的第一个模式进行递归解析, 如果没有匹配的, 将产生错误。

7) 如果读者为联合, 而写者的不是, 对读者联合中匹配所选写者模式的第一个模式进行递归解析, 如果没有匹配的, 将产生错误。

8) 如果写者的是联合, 读者的不是, 且如果读者的模式匹配所选写者的模式, 那么对它进行递归解析, 如果它们不匹配, 那么将产生错误。

模式解析时将忽略模式中协议说明的“doc”字段, 因此, 序列化时模式中的“doc”部分将被抛弃。

16.2 Avro 的 C/C++ 实现

本节主要介绍 Avro 的 C/C++ 实现, 其中在 Avro C 库中已经嵌入 Jansson (Jansson 为编译和操控 JSON 数据的 C 语言库), 这样可以将 JSON 解析成模式结构。目前 C/C++ 实现支

持：所有原始和复杂数据类型的二进制编码和解码；向 Avro 对象容器文件进行存储；模式解析、提升和映射；写入 Avro 数据的有效方式和无效方式，但 C 语言接口暂不支持远程过程调用 RPC。

Avro C 为所有模式和数据对象进行引用计数，当引用数降为零时便释放内存。例如，创建和释放一个字符串：

```
avro_datum_t string = avro_string("This is my string");
...
avro_datum_decref(string);
```

当考虑创建更加详细的模式和数据结构时就会有一点复杂，例如，创建带有字符串字段的记录：

```
avro_datum_t example = avro_record("Example");
avro_datum_t solo_field = avro_string("Example field value");
avro_record_set(example, "solo", solo_field);
...
avro_datum_decref(example);
```

在这个例子中，solo_field 数据没有被释放，因为它有两个引用：原来的引用和隐藏在记录 Example 中的引用。调用 avro_datum_decref(example) 只能将引用数减少为一。如果想结束 solo_field 模式，则需要用 avro_datum_decref(solo_field) 完全删除 solo_field 数据并释放。

一些数据类型是可以“包装”和“给予”的，这可以让 C 程序员自由地决定谁负责内存的分配回收，以字符串为例，建立一个字符串数据有三种方式：

```
avro_datum_t avro_string(const char *str);
avro_datum_t avro_wrapstring(const char *str);
avro_datum_t avro_givestring(const char *str);
```

如果使用 avro_string，那么 Avro C 会复制字符串并且当不再引用时再释放它。在有些情况下，特别是当处理大量数据时要避免这种内存复制，这时需要使用 avro_wrapstring 和 avro_givestring。如果使用 avro_wrapstring，那么 Avro C 不做任何内存处理，它只保存指向数据的指针，这时需要自己来释放字符串。需要注意的是，当使用 avro_wrapstring 时，在用 avro_datum_decref() 取消引用数据前不要释放字符串。如果使用 avro_givestring，那么 Avro C 在数据取消引用之后会释放字符串，从某种程度上说，avro_givestring 将释放字符串的“责任”给了 Avro C。另外，如果没有使用如 malloc 或 strdup 分配堆给字符串，则不要把“责任”给 Avro C。例如，不能这样做：

```
avro_datum_t bad_idea = avro_givestring("This isn't allocated on the heap");
```

写入数据时可以使用下面的函数：

```
int avro_write_data(avro_writer_t writer,
avro_schema_t writers_schema, avro_datum_t datum);
```

如果省略 `writers_schema` 值, 那么数据在发送给写数据的函数前必须检验数据格式的正确性。如果已经确定数据是正确的, 那么可以设置 `writers_schema` 为 `NULL`, 这时 Avro C 不会检查格式。需要注意的是, 写入 Avro 文件对象容器的数据总是要进行验证。

下面介绍一个简单例子, 例子中建立了学生信息的数据库, 并向数据库中读写记录:

```
/*student.c*/
#include <avro.h>
#include <inttypes.h>
#include <stdio.h>
#include <stdlib.h>
#include <unistd.h>

avro_schema_t student_schema;
/*id 用于添加记录时为学生建立学号*/
int64_t id = 0;

/* 定义学生模式, 拥有字段学号、姓名、学院、电话和年龄 */
#define STUDENT_SCHEMA \
"{\"type\": \"record\", \\\n\  \"name\": \"Student\", \\\n\  \"fields\": [\\\n\    {\"name\": \"SID\", \"type\": \"long\"}, \\\n\    {\"name\": \"Name\", \"type\": \"string\"}, \\\n\    {\"name\": \"Dept\", \"type\": \"string\"}, \\\n\    {\"name\": \"Phone\", \"type\": \"string\"}, \\\n\    {\"name\": \"Age\", \"type\": \"int\"}]}"

/* 把 JSON 定义的模式解析成模式的数据结构 */
void init(void)
{
    avro_schema_error_t error;
    if(avro_schema_from_json(STUDENT_SCHEMA,
                             sizeof(STUDENT_SCHEMA),
                             &student_schema, &error)){
        fprintf(stderr, "Failed to parse student schema\n");
        exit(EXIT_FAILURE);
    }
}

/* 添加学生记录 */
void add_student(avro_file_writer_t db, const char *name, const char *dept, const
char *phone, int32_t age)
{
    avro_datum_t student = avro_record("Student", NULL);

    avro_datum_t sid_datum = avro_int64(++id);
    avro_datum_t name_datum = avro_string(name);
    avro_datum_t dept_datum = avro_string(dept);
    avro_datum_t age_datum = avro_int32(age);
```

```

        avro_datum_t phone_datum = avro_string(phone);

        /* 创建学生记录 */
        if (avro_record_set(student, "SID", sid_datum)
            || avro_record_set(student, "Name", name_datum)
            || avro_record_set(student, "Dept", dept_datum)
            || avro_record_set(student, "Age", age_datum)
            || avro_record_set(student, "Phone", phone_datum)) {
            fprintf(stderr, "Failed to create student datum structure");
            exit(EXIT_FAILURE);
        }

        /* 将记录添加到数据库文件中 */
        if (avro_file_writer_append(db, student)) {
            fprintf(stderr, "Failed to add student datum to database");
            exit(EXIT_FAILURE);
        }

        /* 解除引用, 释放内存空间 */
        avro_datum_decref(sid_datum);
        avro_datum_decref(name_datum);
        avro_datum_decref(dept_datum);
        avro_datum_decref(age_datum);
        avro_datum_decref(phone_datum);
        avro_datum_decref(student);

        fprintf(stdout, "Successfully added %s\n", name);
    }

    /* 输出数据库中的学生信息 */
    int show_student(avro_file_reader_t db,
                    avro_schema_t reader_schema)
    {
        int rval;
        avro_datum_t student;

        rval = avro_file_reader_read(db, reader_schema, &student);

        if (rval == 0) {
            int64_t i64;
            int32_t i32;
            char *p;
            avro_datum_t sid_datum, name_datum, dept_datum,
                phone_datum, age_datum;

            if (avro_record_get(student, "SID", &sid_datum) == 0) {
                avro_int64_get(sid_datum, &i64);
                fprintf(stdout, "%PRIu64", i64);
            }

            if (avro_record_get(student, "Name", &name_datum) == 0) {

```



```

        avro_string_get(name_datum, &p);
        fprintf(stdout, "%12s ", p);
    }
    if (avro_record_get(student, "Dept", &dept_datum) == 0) {
        avro_string_get(dept_datum, &p);
        fprintf(stdout, "%12s ", p);
    }
    if (avro_record_get(student, "Phone", &phone_datum) == 0) {
        avro_string_get(phone_datum, &p);
        fprintf(stdout, "%12s ", p);
    }
    if (avro_record_get(student, "Age", &age_datum) == 0) {
        avro_int32_get(age_datum, &i32);
        fprintf(stdout, "%d", i32);
    }
    fprintf(stdout, "\n");

    /* 释放记录 */
    avro_datum_decref(student);
}
return rval;
}

int main(void)
{
    int rval;
    avro_file_reader_t dbreader;
    avro_file_writer_t db;
    avro_schema_t extraction_schema, name_schema,
    phone_schema;
    int64_t i;
    const char *dbname = "student.db";

    init();

    /* 如果 student.db 存在, 则删除 */
    unlink(dbname);
    /* 创建数据库文件 */
    rval = avro_file_writer_create(dbname, student_schema, &db);
    if (rval) {
        fprintf(stderr, "Failed to create %s\n", dbname);
        exit(EXIT_FAILURE);
    }

    /* 向数据库文件中添加学生信息 */
    add_student(db, "Zhanghua", "Law", "15201161111", 25);
    add_student(db, "Lili", "Economy", "15201162222", 24);
    add_student(db, "Wangyu", "Information", "15201163333", 25);
    add_student(db, "Zhaoxin", "Art", "15201164444", 23);
    add_student(db, "Sunqin", "Physics", "15201165555", 25);

```

```

add_student(db, "Zhouping", "Math", "15201166666", 23);
avro_file_writer_close(db);

fprintf(stdout, "\nPrint all the records from database\n");

/* 读取并输出所有的学生信息 */
avro_file_reader(dbname, &dbreader);
for (i = 0; i < id; i++) {
    if (show_student(dbreader, NULL)) {
        fprintf(stderr, "Error printing student\n");
        exit(EXIT_FAILURE);
    }
}
avro_file_reader_close(dbreader);

/* 输出学生的姓名和电话信息 */
extraction_schema = avro_schema_record("Student", NULL);
name_schema = avro_schema_string();
phone_schema = avro_schema_string();
avro_schema_record_field_append(extraction_schema,
    "Name", name_schema);
avro_schema_record_field_append(extraction_schema, "Phone", phone_schema);

/* 只读取每个学生的姓名和电话 */
fprintf(stdout,
    "\n\nExtract Name & Phone of the records from database\n");
avro_file_reader(dbname, &dbreader);
for (i = 0; i < id; i++) {
    if (show_student(dbreader, extraction_schema)) {
        fprintf(stderr, "Error printing student\n");
        exit(EXIT_FAILURE);
    }
}
avro_file_reader_close(dbreader);
avro_schema_decref(name_schema);
avro_schema_decref(phone_schema);
avro_schema_decref(extraction_schema);

/* 最后释放学生模式 */
avro_schema_decref(student_schema);
return 0;
}

```

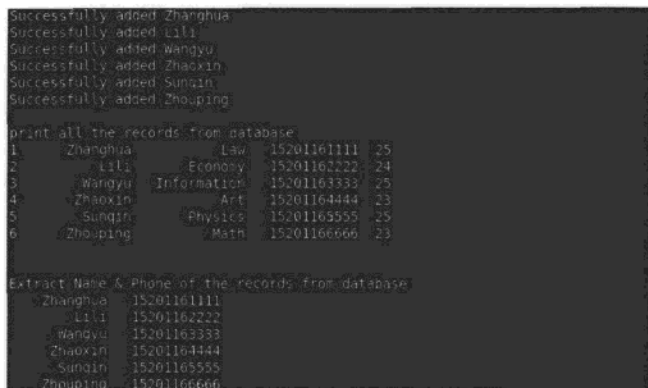
如果要编译上面的 C 文件，则需要安装 Avro C。首先可以从网站 <http://www.apache.org/dyn/closer.cgi/avro/> 选择镜像下载 avro-c-1.4.0.tar.gz 文件，使用命令 `tar -zxvf avro-c-1.4.0.tar.gz` 解压后进入其目录并使用命令 `./configure` 和 `make`、`make install` 进行编译安装。注意，需要在 root 的权限下进行安装。安装成功后，在编译 C 语言前需要将 libavro 加入动态链接库中，使用命令：

```
export LD_LIBRARY_PATH=/usr/local/lib:$LD_LIBRARY_PATH
```

然后对程序进行编译：

```
gcc -o student -lavro student.c
```

运行生成的执行文件可得到如图 16-5 所示的结果。运行时在当前目录下生成 student.db 对象容器文件，可以使用命令 cat 查看文件中的内容——先存储学生的模式，然后存储学生的记录信息，具体内容可参见 16.1.4 节对象容器文件和图 16-3 中的对象容器文件格式。



```
Successfully added Zhanghua
Successfully added Lili
Successfully added Wangyu
Successfully added Zhaoxin
Successfully added Sunqin
Successfully added Zhouping

print all the records from database
1 Zhanghua Law 15201161111 25
2 Lili Economy 15201162222 24
3 Wangyu Information 15201163333 25
4 Zhaoxin Art 15201164444 23
5 Sunqin Physics 15201165555 25
6 Zhouping Math 15201166666 23

Extract Name & Phone of the records from database
Zhanghua 15201161111
Lili 15201162222
Wangyu 15201163333
Zhaoxin 15201164444
Sunqin 15201165555
Zhouping 15201166666
```

图 16-5 运行结果

下面介绍 Avro 的 C++ 应用程序接口。虽然 Avro 并不需要使用代码生成器，但是使用代码生成工具可以更简单地使用 Avro C++ 库。代码生成器既可以读取模式并输出模式数据的 C++ 对象，还可以产生代码来完成序列化或反序列化对象等所有复杂的译码工作。如果使用 C++ 核心库来编写自己的序列化器或解析器，其产生的代码也可以说明如何使用这些库。下面举一个使用模式的简单例子，此例用来表示一个虚数：

```
{
  "type": "record",
  "name": "complex",
  "fields" : [
    { "name": "real", "type": "double" },
    { "name": "imaginary", "type": "double" }
  ]
}
```

假设 JSON 可用来表示存储在名为 imaginary 文件中的模式，那么产生代码分成两步。第一步：

```
precompile < imaginary > imaginary.flat
```

预编译会将模式转化为代码生成器所使用的中间格式，中间文件是模式的文本形式，它

是通过对模式类型树深度优先遍历得到的。

第二步：

```
python scripts/gen-cppcode.py --input=example.flat --output=example.hh
--namespace=Math
```

上面的命令告诉代码生成器去读取模式作为输入，并且在 example.hh 中生成 C++ 头文件。可选参数指定对象放置的命名空间，如果没有指定命名空间，仍可得到默认的命名空间。下面是所产生代码的开始部分：

```
namespace Math {
struct complex {
    complex () :
        real(),
        imaginary()
    { }
    double real;
    double imaginary;
};
```

以上代码是用 C++ 表示的模式，它创建记录、默认构造函数并为记录的每个字段建立成员。下面是序列化数据的例子：

```
void serializeMyData()
{
    Math::complex c;
    c.real = 10.0;
    c.imaginary = 20.0;

    // Writer 是实际 I/O 和缓冲结果的对象
    avro::Writer writer;

    // 在对象上调用 writer
    avro::serialize(writer, c);

    // 这时，writer 将序列化后的数据存储在缓冲区中
    InputBuffer buffer = writer.buffer();
}
```

使用生成的代码，调用对象的 avro::serialize() 函数可以序列化数据，通过调用 avro::InputBuffer 对象可以获取数据，通过网络可以发送文件。下面读取序列化的数据到对象中：

```
void parseMyData(const avro::InputBuffer &myData)
{
    Math::complex c;

    // Reader 为实际 I/O 读取的对象
    avro::Reader reader(myData);
```

```

// 在对象上调用 reader
avro::parse(reader, c);

// 此时, C 中存放的是反序列化后的数据
}

```

在下面的代码中 `avro::serialize()` 函数和 `avro::parse()` 函数可用于处理用户数据类型, 具体实现如下:

```

template <typename Serializer>
inline void serialize(Serializer &s, const complex &val, const boost::true_type &) {
    s.writeRecord();
    serialize(s, val.real);
    serialize(s, val.imaginary);
    s.writeRecordEnd();
}

template <typename Parser>
inline void parse(Parser &p, complex &val, const boost::true_type &) {
    p.readRecord();
    parse(p, val.real);
    parse(p, val.imaginary);
    p.readRecordEnd();
}

```

以下内容也可加入 `avro` 命名空间中:

```

template <> struct is_serializable<Math::complex> : public boost::true_type{};

```

这样为复杂结构建立了类型特征, 告诉 Avro 对象的序列化和解析功能可用。

除了上面介绍的使用 Avro C++ 代码生成器来读写对象外, Avro C++ 也可以读入 JSON 模式。库函数提供了一些工具来读取存储在 JSON 文件或字符串中的模式, 如:

```

void readSchema()
{
    // My schema is stored in a file called "example"
    std::ifstream in("example");

    avro::ValidSchema mySchema;
    avro::compileJsonSchema(in, mySchema);
}

```

上面代码读取文件并将 JSON 模式解析成 `avro::ValidSchema` 类型的对象。如果模式是无效的, 将无法建立有效模式 (`ValidSchema`) 对象并抛出异常。那么如何从 JSON 存储的模式中建立有效模式对象呢?

有效模式 (`ValidSchema`) 可以保证开发者实际写入的类型匹配模式所期望的类型。现在重写序列化函数并检查模式:

```

void serializeMyData(const ValidSchema &mySchema)
{
    Math::complex c;
    c.real = 10.0;
    c.imaginary = 20.0;

    // ValidatingWriter 保证序列化写入正确类型的数据
    avro::ValidatingWriter writer(mySchema);

    try {
        avro::serialize(writer, c);
        // 这时, ostream "os" 存储序列化后的数据
    }
    catch (avro::Exception &e) {
        std::cerr << "ValidatingWriter encountered an error: " << e.what();
    }
}

```

这段代码和前面的区别就是用 ValidatingWriter 代替了 Writer object。如果序列化函数错误地写入不匹配模式的类型, 那么 ValidatingWriter 将抛出异常。ValidatingWriter 会在写入数据的时候增加很多处理过程。对于产生的代码则没有必要进行验证, 因为自动生成的代码是匹配模式的。然而, 在写入和测试自己序列化的代码时加上安全验证还是需要的。解析数据时也可以使用有效模式, 它不仅可以确保解析器读取的类型匹配模式有效, 还提供了接口, 通过该接口可以查询下一个期望的类型和记录成员字段的名称。下面的例子介绍了如何使用 API:

```

void parseMyData(const avro::InputBuffer &myData, const avro::ValidSchema
&mySchema)
{
    // 手动解析数据, 解析对象将数据绑定到模式上
    avro::Parser<ValidatingReader> parser(mySchema, myData);

    assert( nextType(parser) == avro::AVRO_RECORD);

    // 开始解析
    parser.readRecord();

    Math::complex c;

    std::string recordName;
    assert( currentRecordName(parser, recordName) == true);
    assert( recordName == "complex");
    std::string fieldName;
    for(int i=0; i < 2; ++i) {
        assert( nextType(parser) == avro::AVRO_DOUBLE);
        assert( nextFieldName(parser, fieldName) == true);
        if(fieldName == "real") {
            c.real = parser.readDouble();

```

```

    }
    else if (fieldName == "imaginary") {
        c.imaginary = parser.readDouble();
    }
    else {
        std::cout << "I did not expect that!\n";
    }
}

parser.readRecordEnd();
}

```

上面的代码表明，如果编译时不知道模式，也可以通过写出解析数据的代码在运行时读取模式，并且查询 `ValidatingReader` 来了解序列化数据的内容。

在自己的代码中使用对象来建立模式是允许的，每个原始类型和复合类型都有模式对象并且它们拥有共同的 `Schema` 基类。下面是一个为复数记录数组建立模式的例子：

```

void createMySchema()
{
    // 首先建立复数类型：
    avro::RecordSchema myRecord("complex");

    // 在记录中加入字段（每个字段又是一个模式）：
    myRecord.addField("real", avro::DoubleSchema());
    myRecord.addField("imaginary", avro::DoubleSchema());

    // 这个复数记录和之前使用的一样，下面为这些记录的数组建立模式
    avro::ArraySchema complexArray(myRecord);
    // 如果模式是无效的将抛出
    avro::ValidSchema validComplexArray(complexArray);
    // 这样建立好了模式
    // 输出到屏幕上
    validComplexArray.toJson(std::cout);
}

```

以上代码建立的模式可能是无效的，因此为了使用模式，需要将它转化为 `ValidSchema` 对象。执行上述代码可以得到：

```

{
  "type": "array",
  "items": {
    "type": "record",
    "name": "complex",
    "fields": [
      {
        "name": "real",
        "type": "double"
      },
      {

```

```

        "name": "imaginary",
        "type": "double"
    }
}
}

```

随着时间的变化，程序模式期望的数据可能与之前存储的数据不同，为了把一个模式转化为另一个模式，Avro 提供了不完全一样的模式规则。在这种情况下，代码生成工具就有用了，对于每个生成的结构都会建立一个用来读取数据的特别索引结构，即使数据是用不同的模式写的。在 `example.hh` 中的索引结构如下：

```

class complex_Layout : public avro::CompoundOffset {
public:
    complex_Layout(size_t offset) :
        CompoundOffset(offset)
    {
        add(new avro::Offset(offset + offsetof(complex, real)));
        add(new avro::Offset(offset + offsetof(complex, imaginary)));
    }
};

```

数据前若是 `float` 类型而不是 `double` 类型，根据模式解决规则，`floats` 可以升级为 `doubles`，只要新旧模式都有用，就会建立一个动态的解析器来读取代码生成结构的数据。如下所示：

```

void dynamicParse(const avro::ValidSchema &writerSchema,
                  const avro::ValidSchema &readerSchema) {

    // 实例化布局对象
    Math::complex_Layout layout;

    // 创建知道类型布局和模式的模式解析器
    resolverSchema(writerSchema, readerSchema, layout);

    // 设置 reader
    avro::ResolvingReader reader(resolverSchema, data);

    Math::complex c;

    // 执行解析
    avro::parse(reader, c);

    // 这时，C 中存放的是反序列化后的数据
}

```

16.3 Avro 的 Java 实现

本节主要介绍 Avro 在 Java 中的实现，Java API 现在的版本是 1.5.1，其中主要的包有如

下几个。

- ❑ org.apache.avro : Avro 内核类。
- ❑ org.apache.avro.file : 存放 Avro 数据的文件容器相关类。
- ❑ org.apache.avro.generic : Avro 数据的一般表示类。
- ❑ org.apache.avro.io : Avro 输入 / 输出工具类。
- ❑ org.apache.avro.io.parsing : Avro 格式的 LL(1) 语法实现。
- ❑ org.apache.avro.ipc : 进程间调用支持类。
- ❑ org.apache.avro.ipc.stats : 收集和显示 IPC 统计数据工具类。
- ❑ org.apache.avro.ipc.trace : 追踪 RPC 递归调用的相关类。
- ❑ org.apache.avro.mapred : 使用 Avro 数据运行 Hadoop MapReduce, 其 Map 和 Reduce 功能用 Java 实现。
- ❑ org.apache.avro.mapred.tether : 使用 Avro 数据运行 Hadoop MapReduce, 其 Map 和 Reduce 功能在子进程运行。
- ❑ org.apache.avro.reflect : 使用 Java 映射为存在的类生成格式和协议。
- ❑ org.apache.avro.specific : 为格式和协议生成特定的 Java 类。
- ❑ org.apache.avro.tool : Avro 命令行工具类。
- ❑ org.apache.avro.util : 普通工具类。

关于上面各包中包含的类的具体使用可参见 Java API, 下面通过简单的例子来介绍各类的用法。下面是用 Java 实现学生信息的存储和读取:

```
/*student.java*/
import java.io.File;
import java.io.IOException;

import org.apache.avro.Schema;
import org.apache.avro.file.DataFileReader;
import org.apache.avro.file.DataFileWriter;
import org.apache.avro.generic.GenericData;
import org.apache.avro.generic.GenericDatumReader;
import org.apache.avro.generic.GenericDatumWriter;
import org.apache.avro.generic.GenericData.Record;
import org.apache.avro.util.Utf8;

public class student {
    String fileName = "student.db";
    String prefix = "{ \"type\": \"record\", \"name\": \"Student\", \"fields\": [";
    String suffix = "]}";
    String fieldSID= "{ \"name\": \"SID\", \"type\": \"int\"}";
    String fieldName = "{ \"name\": \"Name\", \"type\": \"string\"}";
    String fieldDept = "{ \"name\": \"Dept\", \"type\": \"string\"}";
    String fieldPhone="{ \"name\": \"Phone\", \"type\": \"string\"}";
    String fieldAge = "{ \"name\": \"Age\", \"type\": \"int\"}";
    Schema studentSchema = Schema.parse(prefix + fieldSID + "," + fieldName + "," +
```

```

        fieldDept + "," + fieldPhone + "," + fieldAge + suffix);
    Schema extractSchema = Schema.parse(prefix + fieldName + "," + fieldPhone +
        suffix);
    int SID=0;

    public static void main(String[] args) throws IOException {
        student st = new student();
        st.init();
        st.print();
        st.printExtraction();
    }

    /**
     * 初始化添加学生记录
     */
    public void init() throws IOException {
        DataFileWriter<Record> writer = new DataFileWriter<Record>(
            new GenericDatumWriter<Record>(studentSchema)).create(
                studentSchema, new File(fileName));

        try {
            writer.append(createStudent("Zhanghua", "Law", "15201161111", 25));
            writer.append(createStudent("Lili", "Economy", "15201162222", 24));
            writer.append(createStudent("Wangyu", "Information", "15201163333", 25));
            writer.append(createStudent("Zhaoxin", "Art", "15201164444", 23));
            writer.append(createStudent("Sunqin", "Physics", "15201165555", 25));
            writer.append(createStudent("Zhouping", "Math", "15201166666", 23));

        } finally {
            writer.close();
        }
    }

    /**
     * 将学生信息添加到记录中
     */
    private Record createStudent(String name, String dept, String phone, int age) {
        Record student = new GenericData.Record(studentSchema);
        student.put("SID", (++SID));
        student.put("Name", new Utf8(name));
        student.put("Dept", new Utf8(dept));
        student.put("Phone", new Utf8(phone));
        student.put("Age", age);
        System.out.println("Successfully added "+name);
        return student;
    }

    /**
     * 输出学生信息
     */
    public void print() throws IOException {
        GenericDatumReader<Record> dr = new GenericDatumReader<Record>();
    }

```

```

        dr.setExpected(studentSchema);
        DataFileReader<Record> reader = new DataFileReader<Record>(new
            File(fileName), dr);
        System.out.println("\nprint all the records from database");
        try {
            while (reader.hasNext()) {
                Record student = reader.next();
                System.out.println(student.get("SID").toString()+" "+student.
                    get("Name")+" "+student.get("Dept")+" "+student.get("Phone")+"
                    "+student.get("Age").toString());
            }
        } finally {
            reader.close();
        }
    }
}
/**
 * 输出学生姓名和电话
 */
public void printExtraction() throws IOException {
    GenericDatumReader<Record> dr = new GenericDatumReader<Record>();
    dr.setExpected(extractSchema);
    DataFileReader<Record> reader = new DataFileReader<Record>(new
        File(fileName), dr);
    System.out.println("\nExtract Name & Phone of the records from
        database");
    try {
        while (reader.hasNext()) {
            Record student = reader.next();
            System.out.println(student.get("Name").toString() + " " + student.
                get("Phone").toString() + "\t");
        }
    } finally {
        reader.close();
    }
}
}
}

```

编译 student.java 不仅需要从网站 <http://www.apache.org/dyn/closer.cgi/avro/> 下载 avro-1.4.0.jar 等相关类, 还需要从网站 <http://wiki.fasterxml.com/JacksonDownload> 下载 jackson-core-asl-1.6.1.jar 和 jackson-mapper-asl-1.6.1.jar 这些 Java 中 JSON 生成的解析相关类。编译后运行文件的结果如图 16-6 所示, 同时生成 student.db 文件, 可以通过查看该文件中的内容来了解对象容器文件的格式。

16.4 GenAvro (Avro IDL) 语言

为了让开发者在声明模式时使用一种与诸如 Java、C++、Python 等普通编程语言相似的方法, Avro 提供了 GenAvro 语言。GenAvro 是声明 Avro 模式的高级语言 (最新版本中称

为 Avro IDL), 虽然它目前还没有完全确定下来, 但不会有重大的变化。之前在其他构架如 Thrift、Protocol、CORBA 中使用过接口描述语言 (IDL) 的开发者可能会对 Avro IDL 语言有亲切感。

```
Successfully added Zhanghua
Successfully added Lili
Successfully added Wangyu
Successfully added Zhaoxin
Successfully added Sunqin
Successfully added Zhouping

print all the records from database
1 Zhanghua Law 15201161111 25
2 Lili Economy 15201162222 24
3 Wangyu Information 15201163333 25
4 Zhaoxin Art 15201164444 23
5 Sunqin Physics 15201165555 25
6 Zhouping Math 15201166666 23

Extract Name & Phone of the records from database
Zhanghua 15201161111
Lili 15201162222
Wangyu 15201163333
Zhaoxin 15201164444
Sunqin 15201165555
Zhouping 15201166666
```

图 16-6 编译运行 student 文件

每个 Avro IDL 文件定义了单一的 Avro 协议, 它产生一个 JSON 格式的 Avro 协议文件, 其扩展名为 .avpr。为了使 Avro IDL (新版本中为 .avdl) 文件转化为 .avpr 文件, 必须使用 idl 工具进行处理, 例如:

```
$ java -jar avroj-tools-1.4.0.jar idl src/test/idl/input/namespaces.avdl /tmp/namespaces.avpr
$ head /tmp/namespaces.avpr
{
  "protocol" : "TestNamespace",
  "namespace" : "avro.test.protocol",
  ...
}
```

这个 idl 工具也可以处理从 stdin 输入的数据或输出到 stdout 的数据, 更多的信息可以用 idl --help 命令查询。一个 Avro IDL 文件只包含一个协议定义, 较小的协议可由以下代码定义:

```
protocol MyProtocol {
}
```

这相当于以下的 JSON 协议定义:

```
{
  "protocol" : "MyProtocol",
  "types" : [ ],
  "messages" : {
  }
}
```

使用 @namespace 注解后, 协议的命名空间可能会改变, 代码如下:

```
@namespace("mynamespace")
protocol MyProtocol {
}
```

在 Avro IDL 中可以通过使用 @namespace 为所注解的元素指定属性。Avro IDL 中的协议包含以下项目：

□ 指定模式的定义，包括记录、错误、枚举和固定型。

□ RPC 消息的定义。

□ 外部协议和模式文件的引用。

引入文件可以用以下三种方式之一：

□ 引入 IDL 文件使用语句 `import idl "foo.avdl"`。

□ 引入 JSON 协议文件使用语句 `import protocol "foo.avpr"`。

□ 引入 JSON 模式文件使用语句 `import schema "foo.avsc"`。

下面介绍各种类型的定义方法。

1) 定义枚举。在 Avro IDL 中使用类似于 C 或 Java 的语法来定义枚举，代码如下：

```
enum Suit {
    SPADES, DIAMONDS, CLUBS, HEARTS
}
```

需要注意的是，不像 JSON 格式，匿名的枚举是无法定义的。

2) 定义固定长度的字段。定义一个固定长度的字段使用以下语法：

```
fixed MD5(16);
```

该例子定义了一个包含 16 字节称为 MD5 的固定长度类型。

3) 定义记录和错误。在 Avro IDL 中定义记录的语法类似于 C 中的结构体定义，代码如下：

```
record Employee {
    string name;
    boolean active;
    long salary;
}
```

以上例子定义了一个带有三个字段称为“Employee”的记录，错误类型的定义只需要将 record 改为 error 就可以了，代码如下：

```
error Kaboom {
    string explanation;
    int result_code;
}
```

记录和错误中的字段包括类型和名称，也可以有属性注解。Avro IDL 语言中引用的类型必须为以下之一：

□ 原始类型。

□ 已命名的模式，该模式在同一协议中且使用前已经定义。

□ 复杂类型（数据、映射或联合）。

下面分别介绍它们。

1) Avro IDL 支持的原始类型与 Avro 的 JSON 格式支持的类型一样，包括 int、long、string、boolean、float、double、null 和 bytes。

2) 如果相同的 Avro IDL 文件中已经定义了指定的模式且为原始类型，那么可以通过名称直接引用，代码如下：

```
record Card {
    Suit suit; // 引用之前定义的枚举类型 Card
    int number;
}
```

3) 复杂类型。数组类型的定义方法与 C++ 或 Java 中的定义方式类似。任何类型 t 的数组都会写为 array<t>。例如，字符串的数组写为 array<string>，记录 Foo 的多维数组写为 array<array<Foo>>。映射类型和数组类型相似，包含类型 t 的数组写为 map<t>，和 JSON 模式格式一样，所有的映射包含 string 类型的键。联合类型写为 union { typeA, typeB, typeC, ... }，例如，下面这个记录包含可选的字符串字段：

```
record RecordWithUnion {
    union { null, string } optionalString;
}
```

需要注意的是，Avro IDL 中联合的限制与 JSON 格式的一样，即记录不能包含相同类型的多种元素。

使用 Avro IDL 协议定义 RPC 消息的语法与 C 语言头文件或 Java 接口的方法声明相似。比如带有参数 foo 和 bar 且返回 int 值的 RPC 消息定义为：

```
int add(int foo, int bar);
```

定义一个没有返回值的消息可以使用别名 void，相当于 Avro 的 null 类型，如下所示：

```
void logMessage(string message);
```

如果在相同的协议之前已经定义了一个错误类型，那么可以使用下面语法声明消息抛出这个错误：

```
void goKaboom() throws Kaboom;
```

如果定义一个 one-way 的消息，只需在参数后面使用关键字 oneway，代码如下：

```
void fireAndForget(string message) oneway;
```

最后介绍其他的 Avro IDL 语言特征。

(1) 注释

Avro IDL 语言支持所有的 Java 类型注释。每行 // 后面的内容将被忽略，用 /* 和 */ 可以

注释多行内容。

(2) 区别标识

当语言保留字需要用来作为标识时，需要用符号 “`” 来区别标识。例如，定义一个带有名称 `error` 的消息：

```
void `error`();
```

这个语法可以使用在任何有标识的地方。

(3) 排序和命名空间的注释

在 Avro IDL 中 Java 风格的注释可以用来给类型增加额外的属性。例如，指定记录中字段的排序顺序可以使用 `@order`，如下所示：

```
record MyRecord {
  @order("ascending") myAscendingSortField;
  @order("descending") myDescendingField;
  @order("ignore") myIgnoredField;
}
```

当然注释也可以放在字段类型的前面，如：

```
record MyRecord {
  @java-class("java.util.ArrayList") array string myStrings;
}
```

类似地，当定义一个指定模式时，使用 `@namespace` 可以修改命名空间，如：

```
@namespace("org.apache.avro.firstNamespace")
protocol MyProto {
  @namespace("org.apache.avro.someOtherNamespace")
  record Foo {}

  record Bar {}
}
```

这里在 `firstNamespace` 命名空间中定义了一个协议，记录 `Foo` 定义在 `someOtherNamespace` 中，`Bar` 定义在 `firstNamespace` 中且从容器中继承了默认值。

对于类型和字段的别名可以用注释 `@aliases` 来指定，如：

```
@aliases(["org.old.OldRecord", "org.ancient.AncientRecord"])
record MyRecord {
  string @aliases(["oldField", "ancientField"]) myNewField;
}
```

下面是 Avro IDL 文件的完整例子：

```
/**
 * An example protocol in Avro IDL
 */
```

```

@namespace("org.apache.avro.test")
protocol Simple {

    @aliases(["org.foo.KindOf"])
    enum Kind {
        FOO,
        BAR, // the bar enum value
        BAZ
    }

    fixed MD5(16);

    record TestRecord {
        @order("ignore")
        string name;

        @order("descending")
        Kind kind;

        MD5 hash;

        union { MD5, null } @aliases(["hash"]) nullableHash;

        array<long> arrayOfLongs;
    }

    error TestError {
        string message;
    }

    string hello(string greeting);
    TestRecord echo(TestRecord `record`);
    int add(int arg1, int arg2);
    bytes echoBytes(bytes data);
    void `error`() throws TestError;
    void ping() oneway;
}

```

16.5 Avro SASL 概述

SASL (Simple Authentication and Security Layer, 简单验证安全层) 是网络协议中提供验证和安全的框架, 它将验证机制从用户程序协议中分离出来, 使得采用 SASL 的程序可以使用任何 SASL 所支持的验证机制, 同样也支持代理验证。SASL 提供的数据安全层能够提供数据完整性和数据加密服务, 支持 SASL 的用户协议也支持 SASL 服务所需的安全传输层协议, 其中安全传输层协议是为因特网上通信提供安全性的加密协议。开发者可通过 SASL 对通用 API 进行编码, 此方法避免了对特定机制的依赖。采用 SASL 的协议需要定义 SASL profile, 即如何使用 SASL 进行验证协商, 下面对 Avro RPC 采用的 SASL 进行介绍。

SASL 协商过程可以看成是客户端和服务端使用特定的 SASL 机制在连接的基础上进行一系列消息的交互。客户端通过发送带有初始消息（可能为空）的机制名称（这里是 SASL）来协商过程。协商过程一直伴随着消息的交换直到某一方表明协商成功或失败。消息的内容由具体的机制决定，如果协商成功就可以通过连接进行会话，否则将被抛弃。一些机制在协商之后会继续处理会话的数据（如对数据进行加密），而一些机制会指定会话数据传输不需修改。

Avro SASL 协商使用 4 个单字节命令，分别是：

- ❑ 0：START（开始），使用于客户端初始消息中。
- ❑ 1：CONTINUE（继续），使用于协商进行中。
- ❑ 2：FAIL（失败），协商失败。
- ❑ 3：COMPLETE（完成），成功完成协商。

开始消息的格式是：

| 0 | 4 字节的机制名称的长度 | 机制名称 | 4 字节的有效负载的长度 | 有效负载数据 |

继续消息的格式是：

| 1 | 4 字节的有效负载的长度 | 有效负载数据 |

失败消息的格式是：

| 2 | 4 字节的消息长度 | UTF-8 的消息 |

完成消息的格式是：

| 3 | 4 字节的有效负载的长度 | 有效负载数据 |

协商以客户端发送 START 命令开始，START 命令中包含客户端选定的机制名称和指定机制的有效负载数据。然后，服务器和客户端交换一些 CONTINUE 消息，每个消息包含由安全机制生成的下个消息的有效负载数据。一旦客户端或服务器发送 FAIL 消息，协商就会失败，失败消息中包含 UTF-8 编码的文本。只要接收到或发送了 FAIL 消息，或者在协商过程中发生任何错误，基于此次连接的通信就必须结束。如果客户端或服务器发送 COMPLETE 消息，那么协商将成功完成，会话数据可以通过此次连接进行传输直到一方关闭。

如果 SASL QOP（Quality Of Protection，品质保证）没有进行协商，则基于此次连接的读 / 写无须修改，特别是传输的消息使用了 Avro 框架并采用了下面的形式：

| 4 字节的框架长度 | 框架数据 | ... | 4 个零字节 |

如果 SASL QOP 协商成功，则此次连接后的消息传输使用 QOP。写数据时使用安全机制对非空的框架进行封装，读取数据时需要解开。完整的框架必须传送到安全机制进行解封装，之后传送到应用程序中。如果在封装、解封装或框架处理时发生错误，那么此次连接的通信必须结束。

SASL 的匿名机制很容易实现，特别之处是，一个初始的匿名请求可以用以下静态序列作为前缀：

| 0 | 009 | ANONYMOUS | 0000 |

如果服务器使用匿名机制，则它应检查所接受到的请求前缀，即开始消息的机制名称是

否为“ANONYMOUS”，然后对负载为 COMPLETE 消息的初始响应的前面加上前缀：

```
| 3 | 0000 |
```

如果匿名服务器接收到带有其他机制名称的请求，那么它将发送 FAIL 消息：

```
| 2 | 0000 |
```

注意匿名机制不会在客户端和服务端之间增加多余的往返，START 消息附加在初始请求中而 COMPLETE 和 FAIL 消息则附加在初始响应中。

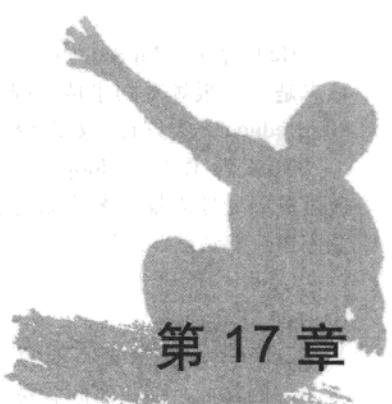
16.6 小结

本章内容主要包括：16.1 节首先说明如何声明 Avro 模式，以及如何对数据进行序列化；然后介绍对象容器文件的具体格式和 RPC 中 Avro 的使用方法，包括协议的声明、协议传输的格式等；最后介绍如何解析获取的数据，重点说明如何处理写入模式和读取模式的不同。16.2 节主要介绍在 C 和 C++ 中如何使用 Avro，主要叙述函数的使用，其中引用关于学生模式的具体例子来详细介绍。16.3 节首先介绍 Java 中使用 Avro 所需要的一些包，后面给出了上节中学生模式例子的 Java 实现程序。16.4 节主要介绍了 GenAvro 语言，说明如何用类似高级语言的方法来声明一个 Avro 模式。16.5 节简单介绍了 Avro 的简单验证安全层，具体说明了通信双方如何进行协商。

Avro 作为一个数据序列化系统，为数据密集型动态应用程序提供了数据存储和交换的平台，它的最大特点就是模式和数据在一起，也就是在反序列化时写入的模式和读出的模式都是已知的，这为 Avro 带来了很多好处，比如生成的数据文件很小等。

今后 Avro 可能会替换 Hadoop 现有的 RPC，作为 Hadoop 的子项目，Avro 的很多特性是为 Hadoop 准备的：容器文件中的同步器可以使 MapReduce 快速地分离文件；不需要生成代码有利于 Avro 用于 Hive 和 Pig；大规模存储较小的数据文件有利于减少数据量，等等。由于 Avro 的数据结构特性，并且拥有多语言支持的优势，于是可以帮助 Hadoop 在跨版本、多语言等方面提高性能。





第 17 章

Chukwa 详解

本章内容

- ☐ Chukwa 简介
- ☐ Chukwa 架构
- ☐ Chukwa 的可靠性
- ☐ Chukwa 集群搭建
- ☐ Chukwa 数据流的处理
- ☐ Chukwa 与其他监控系统比较
- ☐ 小结



17.1 Chukwa 简介

Hadoop 中 MapReduce 最初的主要作用就是用于日志处理，但是使用 MapReduce 处理日志是一件很烦琐的事情，因为集群中机器的日志在不断地增加，会生成大量小文件，而 MapReduce 其实只有在处理少量的大文件数据时才能产生最好的效用。

Chukwa 作为 Hadoop 的子项目弥补了这一缺陷。同时它也是一个高可靠性的应用，能通过扩展处理大量的客户端请求，并且能汇聚多路客户端的数据流。Chukwa 也非常适合商业应用，特别是在云环境上，并且它已经成功地使用于多个场景中。

Chukwa 的开发主要面向四类群体：Hadoop 使用者、集群运营人员、集群的管理者、Hadoop 开发者。

- ❑ Hadoop 使用者：他们一般会想了解作业运行的状态，以及还有多少资源可以用于新的作业，因此他们需要得到的是作业日志和作业输出。
- ❑ 集群运营人员：他们需要了解硬件故障，异常状态、资源的消耗情况。
- ❑ 集群的管理者：他们需要了解在什么样的成本下能够提供什么样的服务，这就意味着他们需要一个工具去分析集群系统或单个用户过去的使用状况，并利用分析出的信息预测将来的需求。他们也要了解系统的一些特征值，如一个任务的平均等待时间。
- ❑ Hadoop 开发者：Hadoop 的开发人员通常需要了解系统的运行情况，Hadoop 的运行瓶颈、失效模式等。

Chukwa 作为 Hadoop 软件家族中的一员，依赖于其他 Hadoop 的子项目使用，比如：以 HDFS 作为存储层，以 MapReduce 作为计算模型，以 Pig 作为高层的数据处理语言。Chukwa 系统的最大开销被限制在整个集群系统可用资源的 5% 以内。

Chukwa 是一个分布式系统，它采用的是流水式数据处理方式和模块化结构的收集系统，在每一个模块中有一个简单规范的接口，这将有用于将来更新，而不需要打破现行的编码结构。流水式模式就是利用其分布在各个节点客户端（Agent）的采集器（Adaptors）收集各个节点被监控的信息，然后以块的形式通过 Http post 汇集到收集器（Collector）上，由它处理后转储到 HDFS 中，之后这些数据由 Archiving 处理（去除重复数据和合并数据）提纯，再由分离解析器（Demux）利用 MapReduce 将这些数据转换成结构化记录，并存储到数据库中，HICC（Hadoop Infrastructure Care Center）通过调用数据库里数据，向用户展示可视化后的系统状态。

图 17-1 展示了 Chukwa 流水式数据处理结构。

下面的章节将从 Chukwa 架构出发，介绍系统中的各个模块，并且讲解 Chukwa 如何实现系统的可靠性。在对 Chukwa 整个系统框架及原理有所了解后，读者可以根据 Chukwa 集群搭建章节的介绍，搭建一个自己的 Chukwa 系统，来监控 Hadoop 集群，这样就可以与其他监控系统有一个比较，希望读者可以结合自己的实际使用感受，进一步了解 Chukwa 监控系统的特点。

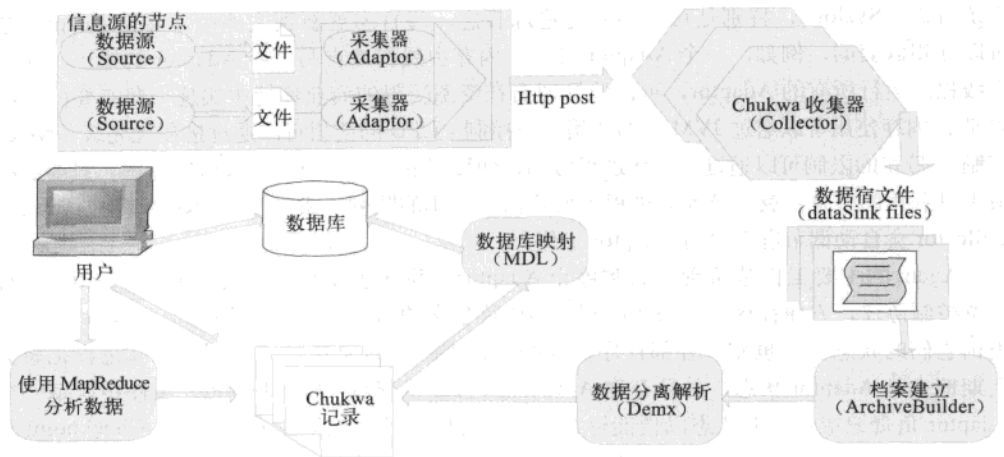


图 17-1 Chukwa 流水线数据处理结构

17.2 Chukwa 架构

Chukwa 有三个主要组成部分：

- ❑ 客户端（Agent），它运行在每一个被监控的机器上，并且传送源数据到收集器（Collector）中。
- ❑ 收集器（Collector）和分离解析器（Demux）：收集器接收从 Agent 传来的数据，并且不断地将其写到 HDFS 中，而分离解析器则进行数据抽取并解析变换成有用的记录。
- ❑ HICC（Hadoop Infrastructure Care Center）：是一个门户样式的网页界面，用于数据的可视化。

17.2.1 客户端（Agent）及其数据模型

在 Chukwa 中，Agent 的主要目的是：使内部进程通信协议能够兼容处理本地的日志文件。

随着分布式计算处理的开始或结束，分布存放的文件和套接字将会不断增加或减少，这种变化是需要被监控的，因此要在每一台机器上配置 Agent。现在绝大多数的监控系统都要求通过特殊的协议传送数据，Chukwa 也不例外，所以在 Chukwa 中，Agent 不直接负责接收数据，取而代之的是提供一个可执行环境：提供可配置的承载数据模块（采集器（Adaptor））。这些 Adaptor 在文件系统或被监控的应用中的功能是读取数据，Adaptor 的输出是一个逻辑上的比特流，单个数据流对应单个文件，或者在相应套接字上接收对应的数据包或一系列重复调用的 UNIX 程序。数据流被存储成序列块，每一个数据块由一些流级别的元数据（Stream-level metadata）加上一个数组比特构成。启动 Adaptor 可以通过 UNIX 命令来完成。Adaptor 能够扫描目录，追踪新创建的文件。这样它便能够接收 UDP 消息，包括

系统日志 (Syslog)，特别是可以不断地追踪日志，写日志更新到文件中，并且 Adaptor 是可以互相嵌套的，例如，一个 Adaptor 可以在内存中缓存来自另一个 Adaptor 的输出。在单个线程内运行所有的 Adaptor，可以让管理员在资源受限的商业环境中实施一些必要的资源限制：内存使用可以通过 JVM 堆的使用进行控制；CPU 的使用可以通过进程优先级 (Nice) 控制；带宽的限制可以通过 Agent 进程协调，即设置它在网络中的最大传输速率，只要超过最大可利用的带宽，就在 Agent 进程中设置固定大小的队列，或者当 Collector 响应缓慢时，Collector 会自动调节进程中的 Adaptor 工作。

Agent 的主要工作是负责开始和停止 Adaptor，并且通过网络传输数据。Agent 支持行定位控制协议，方便程序对 Agent 控制。该协议包含的命令有：启动和停止 Adaptor，以及查询它们的状态，这也允许外部程序在开始读日志时重新配置 Chukwa。Agent 进程也将会定期地查询 Adaptor 状态，并且存储 Adaptor 状态在检查点 (Checkpoint) 文件中，每一个 Adaptor 负责记录足够的状态以便能够在需要的时候完整地恢复原先的状态，Checkpoint 只是包含状态，因此 Checkpoint 文件是很小的，一般每一个 Adaptor 的 Checkpoint 文件只有几百比特。

Agent 和 Adaptor 会自动设置一些元数据，但是其中有两个元数据是必须要用户自己定义的：集群名字和数据类型。集群名字被设置在 conf/Chukwa-env.sh 中，是在每一个进程当中的全局变量。数据类型描述了由 Adaptor 实例收集的数据类型，启动实例时，它必须已经指定。下面的表 17-1 列举了块的元数据字段。

表 17-1 块的元数据字段

字 段	含 义
数据源 (Source)	产生块的节点
集群 (Cluster)	标明集群节点
数据类型 (Datatype)	数据输出格式
序列号 (Sequence ID)	流数据中块的偏移量
名字 (name)	数据源的名字

Adaptor 需要以序列号 (Sequence ID) 作为参数，以便于在崩溃后能重新恢复之前的状态，启动 Adaptor，通常会把序列号置为 0。但是有时候也会为了其他的目的置为其他值，例如，只想追踪文件的下半部分。

17.2.2 收集器 (Collector) 和分离解析器 (Demux)

现在介绍 Chukwa 架构中 Collector 的模型。如果每一个 Agent 都直接往 HDFS 中写入数据，这将会产生许多小文件，所以 Chukwa 使用 Collector 技术，由单个 Collector 线程处理多个来自于 Agent 的数据，每一个 Collector 将它接受到的数据写到单个输出文件中，这个文件放在数据宿 (dataSink) 目录下，这就减少了单个机器或单位时间内 Adaptor 产生的文件数，同时也减少了整个集群产生的文件数。在某种意义上讲，Collector 的存在减轻了大

量的低速率数据源和优化过的少量高速率文件系统间写入的匹配问题，Collector 会定期地关闭它们的输出文件，同时重新命名该文件来标记其可以被进一步处理，并且又开始写另一个新的文件。这个过程被称为文件轮转。一个 MapReduce 作业定期压缩收集到的日志文件并且将它们合并成一个文件。

Chukwa 不同于其他监控系统的地方，就是它利用了 Collector 技术，在 Collector 中没有实施任何可靠性策略，Chukwa 的可靠性是依赖于系统端到端的协议。在 Chukwa 的可靠路径中，Collectors 以标准的 Hadoop 序列文件（sequencefile）格式写数据。这种格式使 MapReduce 的多路处理更加容易。

Chukwa Agent 在分配 Collector 时，也没有实施动态负载均衡方法，而是由 Agent 随机选择 Collector 轮询，直到有一个可以工作为止，而后 Agent 将独占该 Collector，直到 Agent 接收到报错信息，这时它们就会转移到一个新的 Collector 上。如图 17-2 所示：

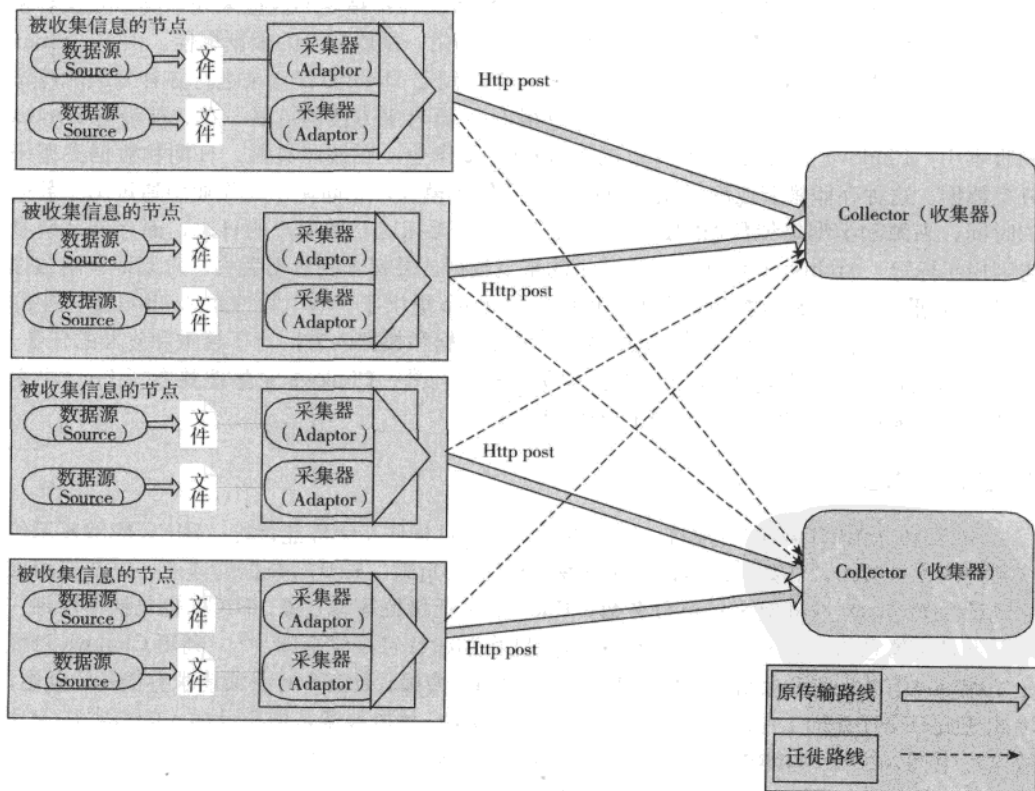


图 17-2 Agent 的可靠性实施

该方法的好处是在文件系统写数据之前，限定了由于 Collector 故障而影响的 Agent 的数量，这也避免了故障扩散，否则会发生每一个 Agent 都被迫对任意一个 Collector 所发生的

故障做出响应的情况。而这种情况下的缺点是 Collector 间的负载不均衡。但在实际的应用中该问题造成的影响不大, Collector 不会太饱和。

为了处理过载的情况, Agent 会重新询问 Collector 是否有特定的方法, 如果 Agent 往一个 Collector 中写入数据时失败, 则该 Collector 被标记为“坏的”(bad), 并且该 Agent 将在再次写入之前等待一个系统设置的时间。因此, 如果所有 Collector 过载, 一个 Agent 将会询问每一个 Collector, 其结果都将会是访问失败, 这样的话, Agent 会等待几分钟后, 再次访问 Collector。

在 Collector 端筛选数据有许多优点, 如果 Collector 是 IO- 约束型, 不是 CPU 约束型, 这就意味着 CPU 资源可以根据作业的状态进行分配, 进一步说, 也就是 Collector 是无状态的, 只须在机器间简单增加更多的 Collector 即可。

Demux 的功能是抽取记录并解析, 使之变换成可以利用的记录, 为了减少文件数目和降低分析难度。Demux 的实现是通过在数据类型和配置文件中指定的数据来处理类, 并执行相应的数据分析工作的。一般是把非结构化的数据结构化, 抽取其中的数据属性。由于 Demux 的本质是一个 MapReduce 作业, 所以可以根据需求制定 Demux 作业来进行各种复杂的逻辑分析。Chukwa 提供的 Demux interface 可以用 Java 语言来方便的扩展。在之前没有 Demux 的版本中, Chukwa 引入了 Archiving 的 MapReduce 作业, 它按照集群、日期和数据类型来分类数据。这种存储模型匹配了使用数据的传统作业模式, 它简化了写作业中通过基于数据的时间、来源和类型提纯数据的过程, 例如, 用 14 天标记存储用户的日志, 而存储系统日志则用年标记。Archiving 作业能去除其中的重复数据, 并探查数据丢失, 重复地调用该作业可让数据随时间不断压缩到一个大文件中。Chukwa 提供了一些工具搜索 Archiving 产生的文件。Chukwa 支持用正则表达式来查询文件的元数据和数据内容, 对于繁重和复杂的任务, 用户可以运行特定的 MapReduce 作业去收集数据。此外, Chukwa 完整地整合了 Pig (见第 14 章 Pig 详解), 用来提供更加强大的搜索功能。

17.2.3 HICC

HICC 作为 Chukwa 的子项目, 它的重要功能是可视化系统性能指标。HICC 能够显示传统系统的度量数据, 例如系统资源空闲比率、CPU 的负载、磁盘写数据的速度, 以及应用层的统计数据 (如本地机架内 map 任务数、Hadoop 块迁移数量等) 等。HICC 也包括使用每一个节点日志信息的 SALSA 作业执行模型状态机和 Mochi 可视化框架^[1,2]。利用 Chukwa 可视化功能可以清楚地看到集群中的作业是否在被均匀的传播。HDFS 对于读请求有很长的延迟, 因此在执行交互查询工作时, 反应会比较慢, 而 HICC 抽取数据是使用批插入的方式往 SQL 数据库中插入经由 MapReduce 处理收集的数据的。MapReduce 作业默认为每 5 分钟执行一次, 因此显示数据至少比实时慢 5 分钟。HICC 也可以支持集群性能的调试和 Hadoop 作业执行的可视化等应用。在这些应用中, 延迟并不是问题。目前, HICC 不需要 Chukwa 的可靠性传输, 但是它依赖于 Chukwa 收集数据和 MapReduce 处理数据。

可以从图 17-3 中清晰地看到 HICC 在 Chukwa 中的作用。

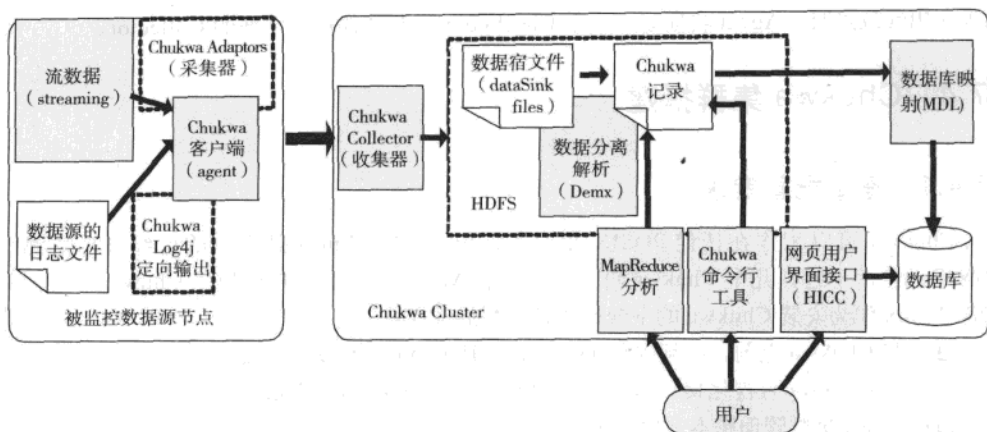


图 17-3 HICC 可视化流程演示图

17.3 Chukwa 的可靠性

容错能力是 Chukwa 设计的一个重要指标。即使在系统崩溃、网络连接中断的情况下，也不能丢失数据。Chukwa 方案与其他分布式系统本质的不同是其分布式存储日志的方式。它会将数据源的相应状态写入数据节点，由 Agent 管理节点崩溃的情况，Agent 通常会为自己的状态设置检查点（Checkpoint）。该 Checkpoint 描述了每一个当前被监控的数据流，并清查有多少来自流中的数据已经被提交到 dataSink 上。在节点崩溃后，Chukwa 使用后台管理工具去重启 Agent。

在 Agent 进程恢复后，每一个 Adaptor 将从最近的 Checkpoint 状态重启。这意味着 Agent 将重新发送没有提交的数据，或者重新发送由最后的 Checkpoint 记录之后所提交的数据。在恢复过程中所产生的重复块将通过 archiving 作业滤除掉。跟踪文件状态的 Adaptor 通过文件的定位固定偏移量来恢复文件内容，并且 Adaptor 也能够监控临时数据源，如网络的套接字。在这种情况下，Adaptor 通过重新发送数据就能很容易地恢复丢失的数据，因此丢失数据将不会是一个大的麻烦，例如：丢失一分钟的系统度量信息，因为在默认提供封装好的库的 Adaptor 中，已经缓存了来自不稳定数据源的数据所以就可以建立不带容错机制的 Collector。Agent 将检查 Collector 对于文件系统的状态，这个状态会起到侦测系统故障并且从故障中恢复的作用，恢复则完全由 Agent 处理，并不需要从失效的 Collector 中获取信息。然后 Agent 发送数据到 Collector，Collector 将写数据存储到 HDFS 文件中，并且也定位了数据在文件中的位置。这个位置很容易就能确定，因为每一个文件仅是通过一个 Collector 写的，唯一需要满足的要求就是排列数据和增加其长度。

Collector 将不监控已写入的文件，也不存储每一个 Agent 状态，轮询 Collector 而不是直接对文件系统访问是为了减少文件系统主节点的负载，把 Agent 从存储系统的烦琐中解脱

出来。出现故障时, Agent 将恢复上一个 Checkpoint, 并且选择一个新的 Collector。

17.4 Chukwa 集群搭建

17.4.1 基本配置要求

Chukwa 可以工作在任何 POSIX 平台上, 但是 GNU/Linux 是唯一的已经被广泛测试的商用平台, 不过, 几个 Chukwa 研发团队也在 Mac OS X 上成功使用了 Chukwa。目前将 GNU/Linux 作为安装 Chukwa 的平台是比较理想的选择。

- ❑ 安装 Chukwa 之前, 必须安装 Java 1.6 和 Hadoop 0.18 或以上版本。
- ❑ 安装 Chukwa 可视化接口 HICC, 需要安装 MySQL 5.1.30 或以上版本。
- ❑ Chukwa 集群管理脚本需要安装 SSH。

17.4.2 安装 Chukwa

最简单的 Chukwa 部署包括三个组件:

- ❑ 一个 Hadoop 集群, Chukwa 需要利用其存储的数据。
- ❑ 一个 Collector 进程, 利用它写数据到 HDFS。
- ❑ 一个或多个 Agent 进程, 它发送监控数据到 Collector, 并且拥有 Agent 进程的节点作为被监控点。

1. 安装 Chukwa

- ❑ 首先需要在官网 (<http://incubator.apache.org/chukwa/>) 上下载 Chukwa 0.4.0, 并且解压。
- ❑ 设置环境变量 `export CHUKWA_HOME=/root/chukwa-0.4.0`, `/root/chukwa-0.4.0` 是 chukwa-0.4.0 的安装目录, 读者可以根据自己的需要新建安装目录。
- ❑ 拷贝 `$CHUKWA_HOME/chukwa-hadoop-0.4.0-client.jar` 和 `$CHUKWA_HOME/lib/json.jar` 到 `$HADOOP_HOME/lib` 中。
- ❑ 在 `conf/chukwa-env.sh` 中设置 `JAVA_HOME`: `export JAVA_HOME = /usr/lib/jvm/java-6-sun-1.6.0.20`。
- ❑ 在 `conf/chukwa-env.sh` 中设置 `HADOOP_HOME` 地址: `export HADOOP_HOME = "/home/user/hadoop-0.20.2"`, 并设置 `HADOOP_CONF_DIR` 地址: `export HADOOP_CONF_DIR = "/home/user/hadoop-0.20.2/conf"`。
- ❑ 拷贝 Chukwa 文件到集群所有需要监控的节点中, 每一个节点都将运行 Collector (这一步需要在配置好 Chukwa 后执行)。

2. Agent 配置

建立 Agent 进程, 首先必须配置的是 `$CHUKWA_HOME/conf/collectors`, 这个文件包括一系列运行的 Collector 节点。Agent 将从列表中随机挑 Collector 接口发送数据, 如果发送失

败，将转换到列表中的另一个 Collector 上。

配置文件格式如下：

```
http://<collector1HostName>:<collector1Port>/
http://<collector2HostName>:<collector2Port>/
http://<collector3HostName>:<collector3Port>/
```

首先，把 CHUKWA_HOME/conf/initial_Adaptors.template 文件改 CHUKWA_HOME/conf/initial_Adaptors。这个文件主要用于设置 Chukwa 监控哪些日志，并且是以什么方式、什么频率来监控等。

其次，把 CHUKWA_HOME/conf/Agent.template 文件改为 CHUKWA_HOME/conf/Agent，并且修改该文件，在其中输入运行 Agent 进程的节点名，并且这些节点必须已经在操作系统中注册过的。

修改后的文件如下：

```
Master
Slave1
Slave2
Slave3
```

Master 为被监控的 NameNode 节点，Slave1、Slave2、Slave3 是运行 Agent 进程被监控的子节点。

还有一些可选的配置在 CHUKWA_HOME/conf/chukwa-Agent-conf.xml 中设置，其中最重要的属性是集群名，用于标示被监控的节点，这个值被存储在每一个被收集到的块中，以区分不同的集群，比如设置 cluster 名称：cluster="demo"，如下所示：

```
<property>
<name>chukwaAgent.tags</name>
<value>cluster="demo"</value>
<description>The cluster's name for this Agent</description>
</property>
```

另一个可选的节点是 chukwaAgent.checkpoint.dir，这个目录是 Chukwa 运行 Adaptor 的定期检查点，它是不可共享的目录，并且只能是本地目录，不能是网络文件系统目录。

3. 为监控配置 Hadoop

为了从 Hadoop 中收集日志，需要安装 Hadoop 0.20.0 或以上的版本，因为只有该版本才能让 Chukwa 从中获取 MapReduce 作业日志，而且也需要修改一些 Hadoop 的配置文件。

- ❑ 拷贝 CHUKWA_HOME/conf/ 路径下的 hadoop-log4j.properties 文件到 HADOOP_HOME/conf/ 路径下，并把文件名变更为 log4j.properties。
- ❑ 拷贝 CHUKWA_HOME/conf/hadoop-metrics.properties 文件到 HADOOP_HOME/conf/hadoop-metrics.properties 中。
- ❑ 在 HADOOP_HOME/conf/ 路径下，把 hadoop-metrics.properties 文件中的 @CHUKWA_

LOG_DIR @ 修改成 CHUKWA 日志存放的地址, 如: CHUKWA_HOME/logs。

4. Collector 配置

修改 \$CHUKWA_HOME/conf/chukwa-collector-conf.xml 中的参数, 如下所示:

```
<property>
  <name>writer.hdfs.filesystem</name>
  <value>hdfs://Master:9000/</value>
  <description>HDFS to dump to</description>
</property>
```

Writer.hdfs.filesystem 中的 hdfs://Master:9000/ 是 HADOOP 分布式文件系统的地址, Chukwa 将利用它来存储数据, 它可以根据实际地址修改。

Chukwa Collectors 负责接收来自 Agent 的数据, 并且存储数据, Collector 将写所有收到的数据到 HDFS 中, 在 collector-conf.xml 中节点 writer.hdfs.filesystem 会设置分布式文件地址, 该节点是 Collector 成功运行的必要条件。

下面的属性设置用于指定 sink data 地址 (见代码内容), /chukwa/logs/ 就是它在 HDFS 中的地址, 在默认情况下, Collector 监听 8080 端口 (代码如下所示), 不过这是可以修改的, 各个 Agent 将会向该端口发消息。

```
<property>
  <name>chukwaCollector.outputDir</name>
  <value>/chukwa/logs/</value>
  <description>Chukwa data sink directory</description>
</property>
<property>
  <name>chukwaCollector.http.port</name>
  <value>8080</value>
  <description>The HTTP port number the collector will listen on</description>
</property>
```

5. 开始启动, 中止 Collector

- ❑ 在单个节点上运行 Collector 进程可以使用 bin/chukwa collector 命令。
- ❑ Collector 可以作为守护进程运行, 其脚本命令是 bin/start-collectors.sh, 它将远程登录 (ssh) 到在 conf/collectors 配置中列出的 Collector 地址, 并且启动一个 Collector 进程在后台运行。
- ❑ 脚本命令 bin/stop-collectors.sh 则是用来关闭 Collector 进程的, 检查 Collector 是否工作正常。
- ❑ 可以在浏览器中输入 http://collectorhost:collectorport/chukwa?ping=true, 其中, collectorhost 是 Collector 的节点, collectorport 是相应的端口, 如果 Collector 运行正常, 一些统计数据将在页面中显示, 如:

Date:1287908152100

Now:1287908171872

```
numberHTTPConnection in time window:0
numberchunks in time window:0
lifetimechunks:0
```

6. Demux 和 HICC 创建

(1) 启动 Chukwa 进程

```
CHUKWA_HOME/start-all.sh
```

(2) 创建定期调度的采样作业配置

```
CHUKWA_HOME/bin/downSampling.sh --config <path to chukwa conf> -n add
```

(3) 安装数据库

这里介绍的是在 Ubuntu 9.04 中安装数据库。

建立和配置数据库：在 Ubuntu 菜单栏中依次选择 system → Administration → Synaptic Package Manager，选择安装 5.1 版本的 mysql-server 和 mysql-client。如图 17-4 所示。

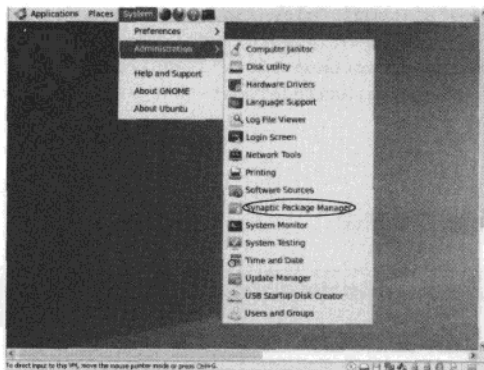


图 17-4 选择 Synaptic Package Manager

以管理员身份进入数据库管理系统：`mysql -u root -p`。

❑ 创建 Chukwa 数据库：`create database chukwa`。

❑ 在数据库 Chukwa 中生成记录：`mysql -u root -p<password> source chukwa < $CHUKWA_HOME/conf/database_create_tables.sql`。这需要在 ubuntu 命令中输入。

❑ 修改 `CHUKWA_HOME/conf/jdbc.conf` 配置文件，如下所示：

```
<clustername>=jdbc:mysql://localhost:3306/<clustername>?user=root&password=<password>
```

第一个 Clustername 是集群名字。

第一个 Clustername 是建立的数据库的名字。

password 是输入数据库 root 用户的密码。

❑ 在 `conf/chukwa-env.sh` 中设置 Database driver: `export JDBC_DRIVE=com.mysql.jdbc.Driver` 和 Database URL prefix: `JDBC_URL_PREFIX=jdbc:mysql://`

(4) 创建 HICC

安装 HICC 需要做以下工作:

❑ 为了使用 HICC, 需要把 Chukwa 0.30 版本中的 `chukwa-0.3.0/bin/dbAdmin.sh` 拷贝到 `$CHUKWA_HOME/bin/` 目录下。

❑ 开始运行 HICC, 输入命令 `$CHUKWA_HOME/bin/chukwa hicc`。

在 Web 地址栏输入 `http://<server>:8080/hicc`。即可看到 Chukwa 的可视化界面。

其中, Server 是主机名, 8080 端口是 jetty 端口, 可以根据需要将 `$CHUKWA/webapps/hicc.war` 文件里 `/WEB-INF/` 目录下的 `jetty.xml` 文件修改, 如下所示:

```
<Call name="addConnector">
  <Arg>
    <New class="org.mortbay.jetty.nio.SelectChannelConnector">
      <Set name="host"><SystemProperty name="jetty.host" /></Set>
      <Set name="port"><SystemProperty name="jetty.port" default="8081"/></Set>
      <Set name="maxIdleTime">30000</Set>
      <Set name="Acceptors">2</Set>
      <Set name="statsOn">false</Set>
      <Set name="confidentialPort">8443</Set>
      <Set name="lowResourcesConnections">5000</Set>
      <Set name="lowResourcesMaxIdleTime">5000</Set>
    </New>
  </Arg>
</Call>
```

Chukwa HICC 界面如图 17-5 所示。

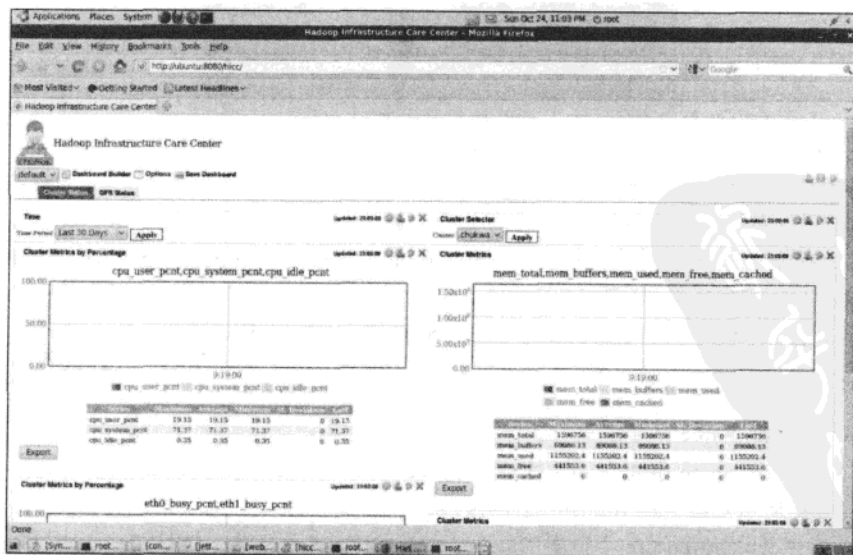


图 17-5 HICC 界面

在 Cluster Status 表单可以看到监控的集群之运行情况, 如图 17-6 所示。

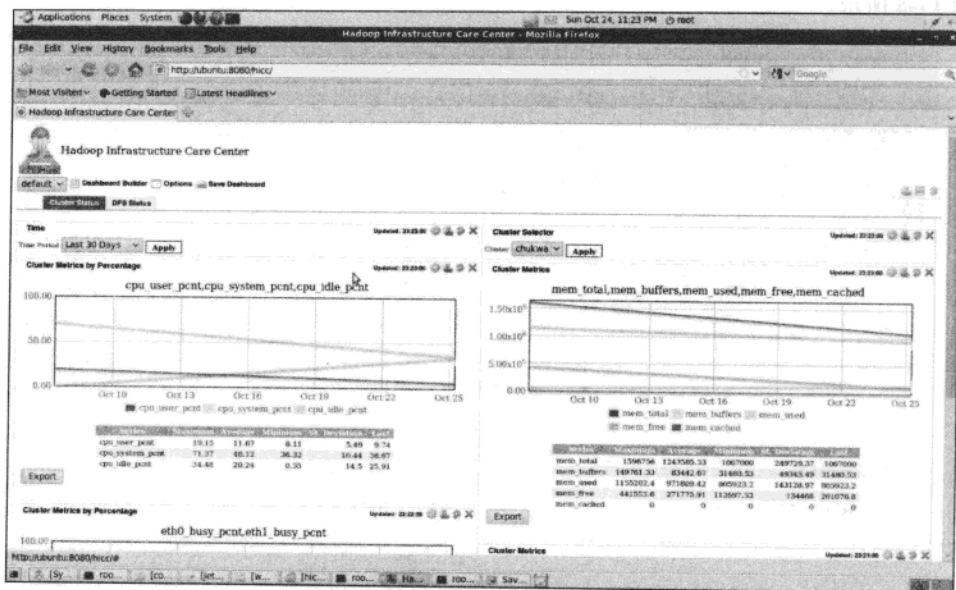


图 17-6 HICC 监控的集群运行

在 DFS Status 表单可以看到分布式文件系统的状态, 如图 17-7 所示。

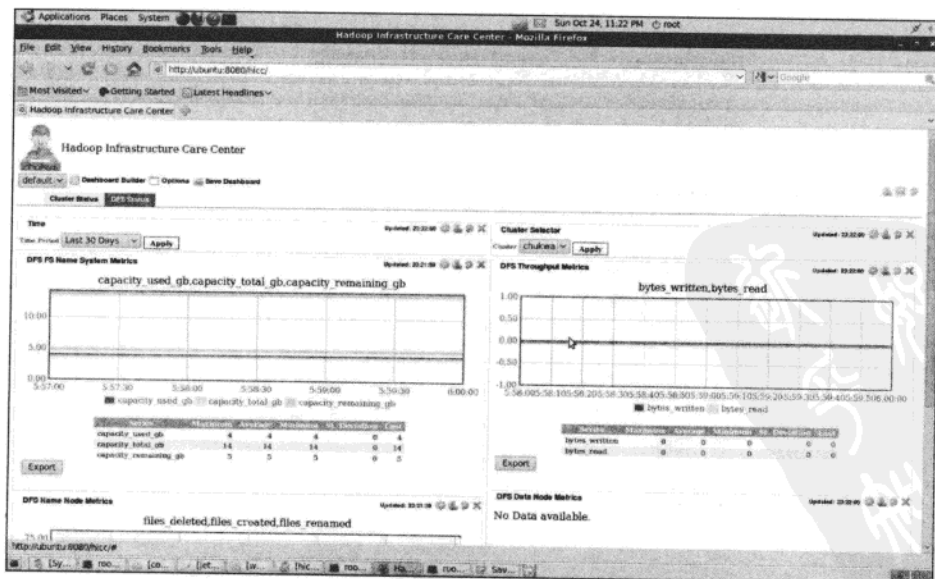


图 17-7 HICC 监控的分布式文件系统运行

也可以点击菜单栏中 Options 选项的 add widget (窗件), 向网页中添加需要监控的窗件, 如图 17-8 所示。

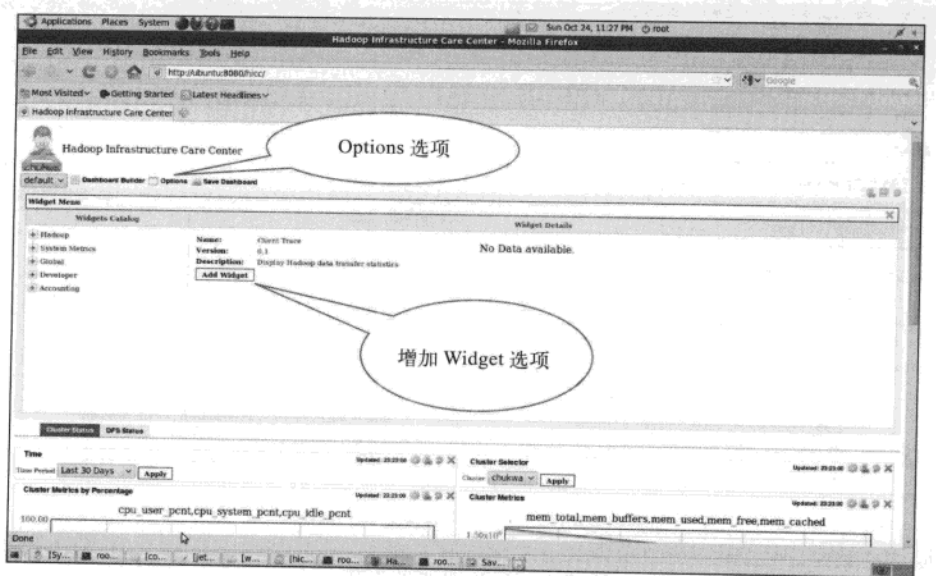


图 17-8 向 HICC 添加信息窗

点击信息窗右上角的齿轮，可以打开并选择需要显示的度量指标，如图 17-9 所示：

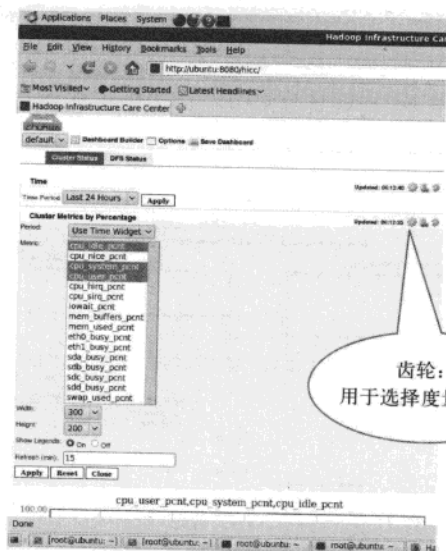


图 17-9 在 HICC 的窗件中选择需要显示的度量指标

7. 启动 CHUKWA 的过程:

- ❑ 启动 Hadoop。
- ❑ 启动 Chukwa : \$CHUKWA_HOME/bin/start-all.sh (在 Hadoop 启动时会有一段时间处于 safemode (安全模式), 在离开 safemode 后再执行该命令)。
- ❑ 启动 dbAdmin : \$CHUKWA_HOME/bin/dbAdmin.sh。
- ❑ 启动 downSampling: \$CHUKWA_HOME/bin/downSampling.sh --config <path to chukwa conf> -n add。
- ❑ 启动 HICC : \$CHUKWA_HOME/bin/chukwa hicc。

17.5 Chukwa 数据流的处理

原始日志收集和聚集的流程, 就是基于 Chukwa 分布式文件系统 (DFS) 的。Chukwa 文件在 HDFS 中的存储结构如图 17-10 所示:

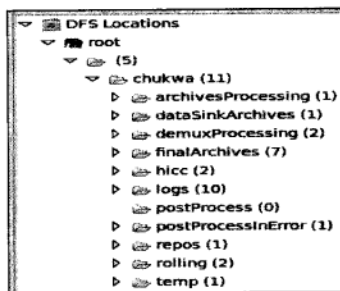


图 17-10 Chukwa 分布式文件系统 (DFS) 的结构

存储流程和结构如下所示:

1) Collector 写块到 logs/ 目录下的 *.chukwa 文件中, 直到达到块的大小 (64MB) 或超时了, Collector 关闭块, 并且将 logs/*.chukwa 改为 logs/*.done 后缀的文件。

2) DemuxManager 每 20 秒检查一次 *.done 文件。

如果这些文件存在, 则移动它们到 demuxProcessing/mrInput, 之后 Demux 将在 demuxProcessing/mrInput 目录下执行 MapReduce 作业。

如果 Demux 在三次以内成功整理完成 MapReduce 文件, 则将 demuxProcessing/mrOutput 中的文件移动到 dataSinkArchives/[yyyyMMdd]/*/*.done, 否则移动执行完 MapReduce 的文件 demuxProcessing/mrOutputdataSinkArchives/ InError/[yyyyMMdd]/*/*.done。

3) 每隔几分钟 PostProcessManager 将执行聚集、排序和去除重复文件作业, 并且将 postProcess/ demuxOutputDir_*/[clusterName]/[data Type]/[data Type]_[yyyyMMdd]_[HH].Re 移动到 repos/[clusterName]/[data Type]/[yyyyMMdd]/[HH]/ [mm]/ [data Type]_ [yyyyMMdd] _ [HH]_[N].[N].evt。

4) HourlyChukwaRecordRolling 将会每小时运行一次 MapReduce 作业, 然后将每小时的日志数据划分为以 5 分钟为日期单位的日志, 并且移动 Repos/[clusterName]/[dataType]/[yyyyMMdd]/[HH]/[mm]/[dataType]_[yyyyMMdd]_[mm].[N].evt 文件到 temp/hourlyRolling/[clusterName]/[dataType]/[yyyyMMdd] 和 repos/[clusterName]/[dataType]/[yyyyMMdd]/[HH]/[dataType]_HourlyDone_[yyyyMMdd]_[HH].[N].evt 上, 同时将文件保留到 Repos/[clusterName]/[dataType]/[yyyyMMdd]/[HH]/rotateDone/ 路径下。

5) DailyChukwaRecordRolling 在凌晨 1:30 运行 MapReduce 作业, 归类以小时为单位的日志到以日为单位的日志中, 同时保留在 repos/[clusterName]/[dataType]/[yyyyMMdd]/rotateDone/ 路径下。

6) ChukwaArchiveManager 大约每半个小时使用 MapReduce 作业聚集和移除 dataSinkArchives 中的数据, 移动 dataSinkArchives/[yyyyMMdd]/*/*.done 到 archivesProcessing/mrInput 和 archivesProcessing/mrOutput, 以及 finalArchives/[yyyyMMdd]/*/*chukwaArchive-part-* 中。

7) 以下目录下的文件将随时间的增长而增加, 因此需要定期清理。

☐ finalArchives/[yyyyMMdd]/*

☐ repos/[clusterName]/[dataType]/[yyyyMMdd]/*/*.evt

17.6 Chukwa 与其他监控系统比较

在了解了 Chukwa 的特点和如何使用之后, 读者或许会问 Chukwa 监控系统与其他监控系统相比有什么特点, 下面我将通过介绍其他监控系统特点, 来帮助读者了解 Chukwa 所具有的特点。

Splunk^[3-6] 是一个日志收集和索引分析的商业化系统, 它依赖于集中的存储和收集架构, 它不考虑传输日志的可靠性。然而高级日志分析领域又有这样的需求, 为了满足这样的需求, 许多大型互联网公司都已经建了大集群监控和分析的高级工具。

有一些专门的日志收集系统, 在这些系统当中, Scribe^[7] 是一个开源的监控系统, 它的元数据模型比 Chukwa 简单, 消息是 key-value 对, 其优点是灵活, 但是它的缺点是要求用户设计自己的元数据标准, 这使得用户之间很难分享源代码。Scribe 的部署由多个服务器组成, 它们被安排在有向非循环图中, 其中的每一个节点对是否提交和存储接收信息是有具体规定的, 相比 Chukwa 而言, Scribe 没有被设计成兼容传统的应用, 被监控的系统必须通过 Thrift RPC 服务发消息给 Scribe, 其优点在于避免通常情况下的本地写开销, 使消息可以无误地传输。它的缺点是在对不适用于 Scribe 的数据源进行收集时, 需要额外的处理。相对而言, Chukwa 处理这样的问题就平滑得多。Scribe 传送上的可靠性也弱于 Chukwa。一旦数据被提交到 Scribe 服务器, 服务器将负责该数据, 而这个服务器为了后续传送会长时间地缓存数据, 这意味着 Scribe 服务器故障可能会造成数据丢失, 同时也没有一个端到端的传输保证, 这是因为原始发送者没有保留一个副本。客户端在向多个服务器发送消息时, 如果其发送在失败之前没有找到正常工作的 Scribe server, 数据将丢失。

另一个相关的系统是 Artemis，它是由 Microsoft 研究设计的，用来调试大规模 Dryad^[8] 集群，Artemis 被设计为专门针对上下文：它只是在本地处理日志，使用 DryadLINQ^[9] 作为它的处理引擎。该架构的优点是避免了网络中多个副本的冗余，也使系统资源可以重复利用于正在分析和已经分析了的结果，缺点是如果一个节点坏掉或暂时不能用，查询会给出错误的结果。Artemis 没有被设计成为使用长期可靠的存储，因为那需要除本地以外的副本，另外，在本地分析对于监控商业服务来说也是不理想的，因为其分析数据的功能可能会干扰正在被监控的系统。

还有一些其他监控系统工具，像 Astrolabe、Pier 和 Ganglia^[10-12] 被设计成能帮助用户查询分布式系统监控信息的系统。在所有的情况下，每个被监控机器上的客户端存储了一定量的数据，用于答复查询，客户端不收集和存储大数据集的结构化数据，所以它们不适合一般目的的编程模型。它们通过在系统中应用一个特殊的数据聚集策略来实现可扩展性，但是会耗费较大的系统性能，相比而言，由于 Chukwa 从收集过程中剥离了分析过程，所以每一个部分部署都可以独立的扩展。

Ganglia 擅长实时故障侦测，相对而言，Chukwa 则会有分钟级的延迟，但考虑到系统能在分钟级处理海量数据，并且能够敏锐侦测到运行变化，同时会对故障诊断有所帮助，而工程师一般不能对秒级别的事件有所反应，所以有分钟级的延时是被允许的。

17.7 小结

Chukwa 作为 Hadoop 的子项目，既能帮助 Hadoop 处理其日志，也能利用 MapReduce 对日志进行分析处理。在 Chukwa 的帮助下，Hadoop 用户能够清晰了解系统运行的状态，分析作业运行的状态及 HDFS 的文件存储状态，从而让我们对整个分布式系统状态有形象直观的了解。

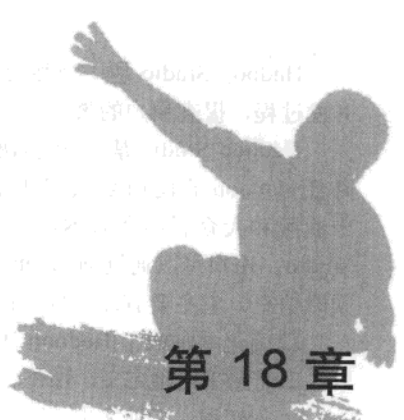
和 Hadoop 一样，Chukwa 也是一个分布式系统，它虽然构建在 Hadoop 之上，但是本身也有自己的特点。它利用分布在各个节点上的 Agent 的 Adaptor 收集各个节点被监控的信息，然后以块的形式通过 Http post 汇集到 Collector，在由它处理后转储到 HDFS 中。之后这些数据由 Archiving 处理（去除重复数据和合并数据）提纯，再由 Demux 利用 MapReduce 将这些数据转换成结构化记录，并存储到数据库中，HICC 通过调用数据库里数据，向用户展示可视化后的系统状态。

要想利用好 Chukwa 这个工具，就必须对 Hadoop 的各个配置项都有清晰的认识。同时 Chukwa 这个项目自身也在不断完善中，感兴趣的读者可以持续跟进。以下是其官网地址：<http://incubator.apache.org/chukwa/>。

参考文献

[1] Tan J, Pan X, Kavulya S, Gandhi R, and Narasimhan P. *SALSA: Analyzing Logs as StAte*

- Machines*. In First USENIX Workshop on Analysis of System Logs (WASL '08), San Diego, CA, December 2008.
- [2] Tan J, Pan X, Kavulya S, Gandhi R, and Narasimhan P. *Mochi: Visual Log-Analysis Based Tools for Debugging Hadoop*. In Workshop on Hot Topics in Cloud Computing (HotCloud '09), San Diego, CA, June 2009.
- [3] Techies get 'the Wikipedia' of glitches , <http://management.silicon.com/itpro/0,39024675,39157789,00.html>
- [4] BryanB, Dave K, Nicolas B, Eric M, Julien S, Michael L, Eric M, Chris I, Philippe B, Jennifer S G, Steve M, Paul G. *Security Power Tools*. Sebastopol : O'Reilly Media, Inc.. 2007-9.
- [5] Max S, Derrick B, Jonathan G, Andrew H, John S. *Nagios 3 Enterprise Network Monitoring: Including Plug-Ins and Hardware Devices*. Syngress. 2008-2.
- [6] Splunk Inc. IT Search for Log Management, Operations, Security and Compliance. <http://www.splunk.com/>, 2009.
- [7] <https://github.com/facebook/scribe>.
- [8] Cretu-Ciocarlie G F, Budiu M, and Goldszmidt M. *Hunting for problems with Artemis*. In First USENIX Workshop on Analysis of System Logs (WASL '08), San Diego, CA, December 2008.
- [9] Yu Y, Isard M, Fetterly D, Budiu M, Erlingsson U, Gunda P, and Currey J. *DryadLINQ: A system for general-purpose distributed data-parallel computing using a high-level language*. In 8th USENIX Symposium on Operating Systems Design and Implementation (OSDI '08), San Diego, CA, December 2008.
- [10] Van Renesse R, Birman K P, and Vogels W. *Astrolabe: A Robust and Scalable Technology for Distributed System Monitoring, Management, and Data Mining*. ACM TOCS, 21(2):164-206, 2003.
- [11] Huebsch R, J. Hellerstein M, Lanham N, Loo B T, Shenker S, and Stoica I. *Querying the Internet with PIER*. Proceedings of 19th International Conference on Very Large Databases (VLDB), pages 321-332, 2003.
- [12] Massie M L, Chun B N, Culler D E. *The Ganglia Distributed Monitoring System: Design, Implementation, and Experience*. Parallel Computing [M], 30(7):817-840, 2004.



第 18 章

Hadoop 的常用插件与开发

本章内容

- ☐ Hadoop Studio 简介和使用
- ☐ Hadoop Eclipse 简介和使用
- ☐ Hadoop Streaming 简介和使用
- ☐ Hadoop Libhdfs 简介和使用
- ☐ 小结

资源库
PDG

18.1 Hadoop Studio 简介和使用

Hadoop Studio 是一个强大的 Hadoop 插件，它具有众多优点，能够简化用户的 Hadoop 开发过程，提高用户的效率。

Hadoop Studio 是一个加快 Hadoop 开发进程的可视化开发环境。Hadoop Studio 通过降低 Hadoop 的使用复杂度让用户在更少的步骤内完成更多的事情以提高效率。Studio 有专业版和大众版两个版本，大众版仅需要注册就可以获得，本章介绍的 Studio 都指大众版 Studio。用户可以通过 Hadoop Studio 强大的 GUI 部署 Hadoop 任务，并监控 Hadoop 任务的实时信息。它主要有以下优点：

- ❑ 简化并加快了 Hadoop 任务模型建立、开发和调试的进程。
- ❑ 能够实时地定义、管理、可视化和监视作业、集群和文件系统等；能够查看任务的实时工作情况；能够让用户通过观察输入输出和中间结果的工作流程图来管理任务的执行时间。
- ❑ 具有很强的移植性，能够被部署在任何操作系统和任何版本的私有或公有 Hadoop 云系统上，且其服务能通过代理服务器和防火墙而不受影响。

Hadoop Studio 的优点决定了无论用户是只有极少 MapReduce 或 Hadoop 开发经验的 Java 程序员，还是熟练的并行程序开发者，它都能简化用户的工作，提高其工作效率。而这主要是从设计、部署、调试和可视化四个方面来实现的。

- ❑ 设计：由于 Studio 能够仿真 Hadoop 系统，所以用户初期建立 MapReduce 任务模型时就不需要一个集群，这可以帮助用户迅速上手。
- ❑ 部署：无论用户使用的是私有网络内的集群还是公共网络上的集群，Studio 都能简化用户任务的部署而且不受服务器和防火墙的影响。在 Hadoop Studio 环境下，用户只需要简单几步便可以启动计算任务：首先在 Hadoop Jobs 中添加生成好的 JAR 包，然后选择要执行的主类，添加依赖项，并选择执行任务的目标 Cluster 节点和目标 Filesystems 即可启动。
- ❑ 调试：MapReduce 编程中最具挑战性的领域之一就是在集群上调试 MapReduce 任务。Studio 提供了可视化工具和任务实时监控，并支持图表化 Hadoop 任务执行状态（包括作业类型、完成情况、执行状态、起止时间、报错信息、输出结果等）和查看任务计数器，这都使得调试 MapReduce 变得容易起来。
- ❑ 可视化：强大的图形用户界面能够使用户不用关注分布式平台的细节就可以编写程序、调试程序、管理集群和文件系统、配置任务信息和日志文件等，这都为用户节省了时间。同时图形界面还能让用户通过实时查看输入输出和中间结果的流程图等其他任务信息来管理任务的执行情况。

18.1.1 Hadoop Studio 的安装和配置

Hadoop Studio 专注于简化数据处理。为了满足其广泛的可用性，Hadoop Studio 开发和部署环境的要求都设计得很简单。表 18-1 是它的开发和部署环境要求：

表 18-1 Hadoop Studio 开发环境

	环境要求
Hadoop 版本要求	Apache Hadoop (0.19, 0.20)、Amazon EMR 或 S3、Cloudera、IBM InfoSphere BigInsights 或 Yahoo!
客户端配置	Mac、PC、Linux
操作系统	Mac、Microsoft Windows、Linux
开发环境	Eclipse (版本 3.5 或更高) Netbeans (版本 6.7 或更高)

从这个表中可以看出 Studio 的开发可以基于 Eclipse 和 Netbeans 进行。下面，我们以基于 Eclipse（安装在 Windows XP 系统上）的大众版 Hadoop Studio 为例介绍其安装和使用方法。

基于 Eclipse 安装 Hadoop Studio 需要一个集成开发环境（IDE）、Java 平台和 Java SE，并且首先需要有以下软件的支持，如表 18-2 所示。

表 18-2 Hadoop Studio 的软件支持

Eclipse IDE 版本	3.5（或更高）
Java 版本	1.6（或更高）
开发环境	Java 环境或 Java SE

安装 JDK 和 Eclipse 的过程不再赘述，重点介绍安装好 JDK 和 Eclipse 之后如何安装基于 Eclipse 的大众版 Hadoop Studio。

Hadoop Studio 是 Eclipse 的一个插件，在启动 Eclipse 之后依次点击 Help 菜单下的 Install New Software → 弹出的 Install 窗口 → Add，然后在弹出的 Add Repository 窗口中填入以下信息：

Name : Karmasphere Studio Plugin

Location : http://updates.karmasphere.com/dist/<<serial_key>>/Eclipse/site.xml

填完之后点击 OK。接下来在 Install 窗口下会出现可能需要安装的插件，选择 Karmasphere Studio Community Edition 或 Karmasphere Studio Professional Edition，之后一直点击 Next 并选择 I accept the terms of the license agreements。接下来 Hadoop Studio 插件将会自动下载并安装。中途如果出现 Security Warning 窗口，选择 OK 安装就会继续。安装结束后重启 Eclipse，安装便完成了，并且 Hadoop Studio 能够正常使用了。

18.1.2 Hadoop Studio 的使用举例

下面以本地 MapReduce 任务的开发、调试和部署，以及远程部署为例介绍 Hadoop Studio 的使用情况。

1. 本地开发、调试和部署

(1) 本地的开发和调试

Hadoop Studio 的任务开发工具允许用户开发并调试 MapReduce 任务。Hadoop Studio 可以降低 MapReduce 编程的入门门槛，因为有了它用户可以在不需要集群支持的情况下，不

断开和调试自己的任务以避免延误整个工程的开发周期。接下来介绍两个工作流程并说明如何在本地部署它们，让读者熟悉 Hadoop Studio 开发工具的使用方法，这两个工作流程中一个使用 MapReduce 预定义类（WordCount Workflow），另一个使用 MapReduce 自定义类（Pi Project Workflow）。

❑ WordCount workflow

- 1) 创建一个名为 WordCountProject 的 Java 工程（详细过程略）；
- 2) 创建一个 MapReduce 工程。

为了使用 Hadoop Studio，我们需要引用 Karmasphere 和 Hadoop libraries，右键点击步骤 1) 中创建的 Java 工程，选择 Build Path>Add Libraries，接着选择弹出窗口中的 Hadoop Libraries from Karmasphere 并点击 Next，在弹出的窗口中选择 Karmasphere Client for Hadoop 并点击 finish。然后右键点击 WordCountProject 工程，选择 New>Other，接着选择 Hadoop Jobs 下的 Hadoop Map-Reduce Job (Karmasphere API) 并点击 finish。再然后展开 Eclipse 工作区面板中的 WordCountProject，双击 src 下的 HadoopJob.workflow，就会出现下面的界面（如图 18-1 所示）：

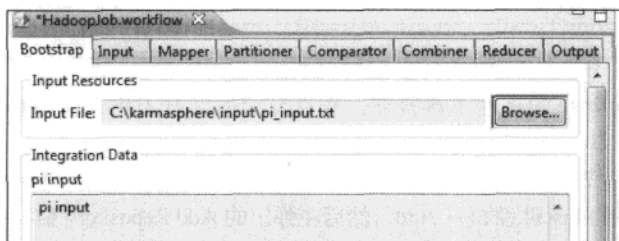


图 18-1 Hadoop Studio 工程配置图

点击窗口界面中第一行的 Bootstrap 按钮，再点击相应页面中的 Browse 按钮，打开文件系统上的一个文件，然后保存工程。这样 Studio 就会为你的工程生成所有的代码并且进行编译。在图 18-1 的窗口中有很多选项卡，点击各个选项卡用户可以查看自己工程所处的状态和输入数据在各个时间点对工程的影响。用户也可以点击选项卡设置对应的工程配置参数。

现在先点击 Input 选项卡，在 Class name 一栏中输入 org.apache.Hadoop.mapred.TextInputFormat，将输入文件的格式设定为 TextInputFormat。再点击 Mapper 选项卡，在 Class name 一栏输入 org.apache.Hadoop.mapred.lib.TOKENCountMapper，为 Mapper 的计数令牌设定类的格式。接着点击 Partitioner 选项卡，在 Class name 一栏输入 org.apache.Hadoop.mapred.lib.HashPartitioner，以设定 Partitioner 的类。跟着点击 Comparator 选项卡，在 Class name 一栏输入 org.apache.Hadoop.io.Text.Comparator 选定 Text Comparator。然后点击 Combiner，在 Classname 一栏输入 org.apache.Hadoop.mapred.lib.IdentityReducer，选定 Identity Reduce。再然后点击 Reducer，在 Class name 一栏输入 org.apache.Hadoop.mapred.

lib.LongSumReducer, 选定 Long Sum Reducer。最后点击 Output 选项卡, 在 Class name 中输入 org.apache.Hadoop.mapred.TextOutputFormat, 以设定 output 的数据类型。

❑ Pi Project 工作流

1) 创建新的 Java 工程, 按照上面 WordCount-Project 中添加工作流的步骤添加一个工作流。如图 18-2 所示。

2) 右键点击 PiProject, 选择 New>Other。选择 New 窗口中 Hadoop Types 下的 Hadoop Mapper, 然后点击 finish。Package Explorer 中的 PiProject 工程下会出现 HadoopMapper.java, 按照同样的步骤添加 HadoopReducer.java。双击 HadoopMapper.java 打开界面, 输入下面的代码:

```
/*
 * HadoopMapper.java
 *
 */

import java.util.Random;
import java.io.IOException;

import org.apache.Hadoop.mapred.MapReduceBase;
import org.apache.Hadoop.mapred.Mapper;
import org.apache.Hadoop.mapred.OutputCollector;
import org.apache.Hadoop.mapred.Reporter;

import org.apache.Hadoop.io.Text;
import org.apache.Hadoop.io.LongWritable;
/**
 *
 */
public class HadoopMapper extends MapReduceBase implements Mapper<Text,Text,Text,
    LongWritable> {

    public void map(Text key, Text value, OutputCollector<Text, LongWritable> output,
        Reporter reporter)
        throws IOException {

        Random generator = new Random();
        int i;

        final int iter = 100000;

        for (i = 0; i < iter; i++)
        {
```

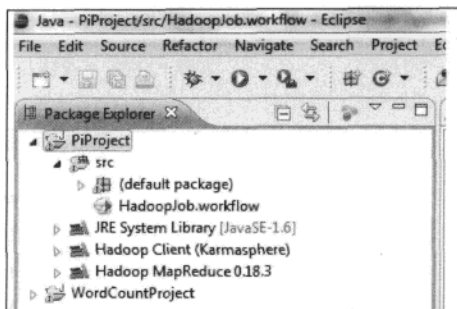


图 18-2 Hadoop Studio 新创建的工程图

```

        double x = generator.nextDouble();
        double y = generator.nextDouble();

        double z;

        z = x*x + y*y;

        if (z <= 1)
            output.collect(new Text("VALUE"), new LongWritable(1));
        else
            output.collect(new Text ("VALUE"), new LongWritable(0));
    }
}
}

```

再双击 HadoopReducer.java 打开界面，输入下面代码：

```

/*
 * HadoopReducer.java
 *
 */

import java.io.IOException;
import java.util.Iterator;

import org.apache.Hadoop.mapred.MapReduceBase;
import org.apache.Hadoop.mapred.OutputCollector;
import org.apache.Hadoop.mapred.Reducer;
import org.apache.Hadoop.mapred.Reporter;

import org.apache.Hadoop.io.Text;
import org.apache.Hadoop.io.LongWritable;
import org.apache.Hadoop.io.DoubleWritable;

public class HadoopReducer extends MapReduceBase implements Reducer<Text,
    LongWritable,Text,DoubleWritable> {
    public void reduce(Text key, Iterator<LongWritable> value, OutputCollector<Text,
        DoubleWritable> output, Reporter reporter)
        throws IOException {

        double pi = 0;
        double inside = 0;
        double outside = 0;

        while (value.hasNext())
        {
            if (value.next().get() == (long)1)
                inside++;
            else

```

```

        outside++;
    }

    pi = (4*inside)/(inside + outside);

    output.collect(new Text ("pi"), new DoubleWritable(pi));
}
}

```

右键点击 Eclipse 菜单栏中的 Project 选项卡，查看 Build Automatically 项是否选中，如果没有选中，就点击 Project 下的 Build Project。需要注意的是，Build Automatically 对于 Studio 生成的 Hadoop Job 默认是选中的。之后点击 HadoopJob.workflow 下的 Bootstrap，接着点击 Browse 选择输入文件，并点击 input 选项卡设定输入格式为 org.apache.Hadoop.mapred.KeyValueTextInputFormat。然后点击 mapper 选项卡输入 HadoopMapper 选定 HadoopMapper。点击 Partitioner 选项卡输入 org.apache.Hadoop.mapred.lib.HashPartitioner 选定 Hash Partitioner。点击 Comparator 选项卡，输入 org.apache.Hadoop.io.TextComparator 选定 Text Comparator。点击 Reducer 选项卡，输入 HadoopReducer 选定 Hadoop Reducer。最后点击 Output 选项卡输入 org.apache.Hadoop.mapred.TextOutputFormat 选定输出数据格式。

(2) 本地任务部署

Hadoop Studio 使用户能够将自己的本地任务部署成线程模式。这里将介绍将本地工作流任务和 JAR 包任务部署成线程模式的详细步骤，包括工作流和 JAR 文件。需要注意的是，如果读者使用的是 Windows 系统，则需要先安装 Cygwin 模拟 Linux 环境。

□ 部署工作流

打开上面已经创建的 PiProject 工程中的工作流，点击 Eclipse 工具栏中最后一个 Deploy 按钮，设定 Deployment 窗口中 Target Cluster 和 Data Filesystem 的参数值，分别为 In-Process Thread(0.20.2) 和 Local Filesystem C:\，然后点击 OK。当工作流在本地部署完成时，在 Output 窗口下就可以看到实时执行状态了。

□ 部署 JAR 包

首先还是打开上面已经创建的 PiProject 工程中的工作流，然后依次选择 Eclipse → Window → Open Perspective → Other → Hadoop，点击 OK 之后会打开 Hadoop 视图。在 Jobs 上右键点击选择 New Job，输入 Job Name，选择 Job Type 为 Hadoop Job from pre-existing JAR file，点击 Next，然后选择你要部署的 Jar 文件并点击 Next。接着选择 Default Cluster 为 In-Process Thread (0.19.3)，设定 Default Arguments 为 pi 10 10000。最后右键点击新建的 Job，选择 Execute Job。到此 Jar 文件的部署已经完成。同样在 Jar 文件部署完成后，就可以在 Output 窗口中查看 Job 的实时执行状态了。

2. 集群部署

(1) 新建 Hadoop HDFS

为了使用 HDFS，我们首先需要在 Hadoop 视图下创建一个文件系统。Hadoop Studio 允

许用户通过 Socket 或 SSH 连接、浏览、读写一个 HDFS。它有一个内置的用来展示本地文件系统的选项。

首先，让我们打开 Hadoop 视图创建一个文件系统选项。右键点击 Filesystem 选择 New Filesystem，在打开的窗口中输入文件系统的名字并设定 Filesystem Type 为 Hadoop HDFS Filesystem。接下来配置运行 HDFS 的 NameNode。如果计划通过 SSH 连接，那么需要将 NameNode Host 配置成 localhost，然后再配置 NameNode Port、Hadoop Version、Username、Group 并点击 finish，接下来将连接类型配置成 DIRECT，之后点击 finish 完成文件系统选项的创建。右键点击创建的文件系统选项选择 Open Filesystem 可以浏览文件系统项目，Studio 将会创建同文件系统的连接，并打开 Filesystem Browser 窗口以便于用户查看管理文件系统。

(2) 监控 HDFS

Hadoop Studio 可以图形化地描述 HDFS 文件系统的状态，在需要查看的文件系统上点击右键选择 Monitor status 就可以查看。当然前提是用户已经创建了文件系统。

(3) 创建 Hadoop 集群

要在分布式 Hadoop 集群上部署、调试、监控，需要先在 Hadoop 视图下创建集群选项。Hadoop Studio 允许用户在集群上运行自己的任务并通过图表监控集群的状态。

Hadoop Studio 有一个内置模拟集群的选项。用户可以用它来运行任务，在测试小数据量上任务的运行情况时显得尤其有用。在这里将创建一个 Hadoop 集群。首先需要添加一个新的 JobTracker 集群。打开 Hadoop 视图右键点击 Hadoop Clusters 并选择 New Cluster，在出现的窗口输入 Cluster Name，选择 Cluster Type 为 Hadoop Cluster (JobTracker)，再设定正确的 Hadoop Version 和 Default Filesystem。点击 Next 之后再配置集群，输入 JobTracker Host、JobTracker Port 和 Username，之后点击 Next。接下来配置 Hadoop 集群的通信机制，有直接通信、Socket 和可选 SSH。我们这里设置为直接通信。然后点击 finish 完成 Hadoop 集群的创建。

(4) 监控正在运行的任务

Hadoop Studio 可以解释并显示用户 Hadoop 集群在运行任务时保存的日志文件和错误诊断信息。这个功能使用户可以监控自己任务的执行情况，并且分析任务的执行结果。当 MapReduce job 运行时，Hadoop Studio 会切换到 Job Monitor 视图，在这个视图上用户可以看到任务的执行信息。Job Monitor 视图列出了集群上所有任务的信息，选择一个任务并点击 Task Monitor 按钮就可以看到这个任务的 Summary、Timeline、Logs、Tasks 和 Config 等信息。如果需要监控集群的状态，可以右键点击对应集群并选择 Monitor Status。Monitor Status 窗口，就会列出 Map Attempts、Reduce Attempts、Task Trackers 和 User Accounting 的统计表格。

(5) 部署运行任务

Hadoop Studio 允许用户部署三种类型的任务：工作流、JAR 文件和流任务。这里我们将介绍部署这三种类型任务的步骤。需要注意的是如果集群通过 SSH 连接，那么部署的 Job 必须是通过 Hadoop Client 创建的（具体过程参考本章相关内容）。另一个需要注意的是保证 Hadoop 的版本与集群的 Hadoop 版本一致。

□ 工作流

创建一个工作流（过程略），然后点击 deploy，在弹出的 Deployment 窗口中输入 Job Name，选择 Target Cluster 和 Data Filesystem，键入输入和输出的参数，然后点击 OK。需要注意的是，输入参数的时候每个参数都要占一行或同行参数之间需要空格隔开，并且输出目录应为空或不存在。接下来点击 OK，之后工作流就会部署运行了，在 Output 窗口里可以查看运行情况。

□ 流任务

打开 Hadoop 视图，右键点击 Jobs 选择 New Job。输入 Job Name，选择 Job Type 为 Hadoop Streaming Job，点击 Next。再输入 Input Location 和 Output Location，再点击 Next。接着选择 Mapper 和 Reducer 的 types 为 Raw Command，在 Mapper 和 Reducer 中输入 /bin/cat，然后点击 finish，如果是自己编写的代码那么就需要在设置 Mapper 和 Reducer 时选择 Upload。接下来右键点击新建的 Job 选择 Execute Job，在确认各项参数无误之后点击 OK，这样，任务就会部署运行，同样可以在 Output 中查看状态。

□ JAR 文件

在集群上部署 Jar 文件需要用到 Hadoop Services。具体步骤是打开 Hadoop 视图，右键点击 Jobs，选择 New Job，输入 Job Name，选择 Job Type 为 Hadoop Job from pre-existing JAR file，点击 Next。在弹出的窗口中浏览文件系统并选择 Primary Jar file，然后输入 Main Class，点击 Next。接下来需要选择默认集群和默认参数，配置完成之后点击 finish。参数输入格式的要求和 Job Workflow 中相同，然后右键点击新创建的 Job，选择 Execute Job，确认参数无误之后点击 OK，这样，任务就会部署运行，同样可以在 Output 中查看状态。

到这里 Hadoop Studio 的使用方法已经介绍完毕，我们从本机和集群两个角度分别介绍了不同任务的部署和运行，同时还介绍了如何使用 Hadoop Studio 监控用户任务，以及利用其用户界面简化 Hadoop 任务的创建、调试、监控和执行。Hadoop Studio 的可视化设计和全面的功能大大降低了基于 Hadoop 项目的开发难度，值得所有 Hadoop 使用者和开发者使用。

18.2 Hadoop Eclipse 简介和使用

Hadoop 是一个强大的并行框架，它允许任务在其分布式集群上并行处理。但是编写、调试 Hadoop 程序都有很大的难度。正因为如此，Hadoop 的开发者开发出了 Hadoop Eclipse 插件，它在 Hadoop 的开发环境中嵌入了 Eclipse，从而实现了开发环境的图形化，降低了编程难度。在安装插件，配置 Hadoop 的相关信息之后，如果用户创建 Hadoop 程序，插件会自动导入 Hadoop 编程接口的 JAR 文件，这样用户就可以在 Eclipse 的图形化界面中编写、调试、运行 Hadoop 程序（包括单机程序和分布式程序），也可以在其中查看自己程序的实时状态、错误信息和运行结果，还可以查看、管理 HDFS 及其文件。总的来说，Hadoop Eclipse 插件安装简单，使用方便，功能强大，尤其是在 Hadoop 编程方面，是 Hadoop 入门和 Hadoop 编程必不可少的工具。

18.2.1 Hadoop Eclipse 安装和配置

Hadoop Eclipse 插件有很多版本, 比如 Hadoop 中官方下载包中的版本、IBM 的版本等。下面将以 Hadoop 官方下载包中的插件为例介绍安装和使用方法。安装插件之前要先安装 Hadoop 和 Eclipse (这部分内容略去, 直接介绍插件的安装)。

1. 安装环境要求

操作系统: 无要求

软件: ☐ Eclipse 3.2.2

☐ Java 1.5 或更高版本

☐ Hadoop

2. 安装步骤

1) 将 Hadoop Eclipse plugin 移动到 Eclipse 的插件文件夹中。找到 Hadoop Eclipse plugin, 它在 ...\Hadoop-0.20.2\contrib\Eclipse-plugin 目录下, 将这个文件夹目录下的文件拷贝到 ...\Eclipse\plugins 目录下。

2) 在 Eclipse 中打开 Hadoop 视图。依次选择: Eclipse → Window → perspective → Other, 然后选择 Map/Reduced 并点击 OK。Eclipse 界面中会出现 Hadoop 视图。左边的 Project Explorer 中会出现 DFS Locations, 下方的选项卡中会出现 Map/Reduce Locations 选项卡。

3) 在下方选项卡中选中 Map Reduce Locations, 然后在出现的空白处右键点击选择 New Hadoop location..., 这时会弹出配置 Hadoop location 的窗口。按照下面的提示正确配置 Hadoop。

Location Name - localhost

Map/Reduce Master :

Host - localhost

Port - 9001

DFS Master :

Host - localhost

Port - 9000

User name - 系统用户名

配置完成之后点击 finish。MapReduce Locations 下就会出现新配置的 MapReduce location。Eclipse 界面左边的 DFS location 下面也出现新配置的 DFS, 点击 “+” 可以查看其结构。

到此, Hadoop Eclipse 插件已经安装完成, 下面举例介绍如何使用。

18.2.2 Hadoop Eclipse 的使用举例

首先打开 Hadoop 视图 (图略), 然后右键点击 Project Explorer 空白处选择 New → Project。

在创建工程向导中选择创建 MapReduce 工程，然后输入工程名，点击 finish。此时 Project Explorer 中会出现新创建的工程。接下来就是编写具体的 MapReduce 代码了。这个时候有两种做法。一种是右键点击新建工程然后新建一个 class，并输入自己完整的 MapReduce 的代码以新建 class 代码区。注意，代码中的类名要和创建类时输入的类名相同，代码编写完之后直接选择 Run as Java Application 即可。另外一种方法是分别建立 MapReduce Driver、Mapper、Reducer。下面详细介绍第二种做法。

首先在刚才新创建的 Hadoop 工程上右键点击选择 New → Other → Map/Reduce → MapReduceDriver，然后点击 Next，输入类名 TestDriver 之后点击 finish。如果生成的代码中有下面两行内容：

```
conf.setInputPath(new Path("src"));
conf.setOutputPath(new Path("out"));
```

说明 Hadoop Eclipse 插件和 Hadoop 的 API 不匹配，需要将这两段代码改成：

```
conf.setInputFormat(TextInputFormat.class);
conf.setOutputFormat(TextOutputFormat.class);

FileInputFormat.setInputPaths(conf, new Path("In"));
FileOutputFormat.setOutputPath(conf, new Path("Out"));
```

同时还需要确认 MapReduce 工程下已经创建了输入文件夹 In 且没有输出文件夹 Out。然后在 TestDriver 类名上点击右键选择 Run As → Run on Hadoop，再选择之前已经配置的 Hadoop server，点击 finish，接下来就可以看到 Eclipse 开始运行 TestDriver 了。这里需要注意的问题有三个：

1) 若点击 Run on Hadoop 无反应，这主要是因为 Hadoop Eclipse 和 Eclipse 的版本不匹配，这个时候有两种办法，一种是更换二者的版本直到匹配为止；另一种办法是选择 Run as Java Application。

2) 若任务执行失败，出错提示为 Java space heap。这主要是因为 Eclipse 执行任务时内存不够，导致任务失败，解决的办法是选中工程并点击 Run → Run Configurations，再点击出现窗口中间的 Arguments 选项卡，在 VM arguments 中写入：-Xms512m -Xmx512m，然后点击 Apply，接下来就可以正常执行程序了。这句话的主要作用是配置这个工程可以使用的内存最小值与最大值都是 512MB。

3) 如何调试 MapReduce 程序。安装有 Hadoop Eclipse 插件的 Eclipse 可以调试 MapReduce 程序，调试的办法就是正常 Java 程序在 Eclipse 中的调试办法，即设置断点，启动 Debug，按步调试。

18.2.3 Hadoop Eclipse 插件开发

可能有时候会因为 Eclipse 版本问题或操作系统版本问题使得 Hadoop 提供的 Eclipse plugin 不太好用。其实这是可以自己生成的，步骤如下：

1) 修改 \$HADOOP_HOME/src/contrib/build-contrib.xml。

将文件中 Eclipse 的路径改成自己电脑中的 Eclipse 路径。

2) 修改 \$HADOOP_HOME/src/contrib/Eclipse-plugin/src/java/org/apache/Hadoop/Eclipse-launch/HadoopApplicationLaunchShortcut.java。

将 import org.Eclipse.jdt.internal.debug.ui.launcher.JavaApplicationLaunchShortcut; 改为 import org.Eclipse.jdt.debug.ui.launchConfigurations.JavaApplicationLaunchShortcut;

3) 执行:

```
cd $HADOOP_HOME
ant compile
ln -sf $HADOOP_HOME/docs $HADOOP_HOME/build/docs
ant package
```

需要注意的是,在执行最后一句的时候可能遇到这个问题: build.xml:908: 'java5.home' is not defined. Forrest requires Java 5. Please pass -Djava5.home=<base 5="" distribution> to Ant on the command-line.。出现这个问题意味着程序没有找到 define java5 的 home 目录,这是因为系统安装的 JDK 是 Java 6。解决办法是下载安装一个 Java 5。然后设定正确的 JDK 路径,接下来修改 \$HADOOP_HOME/build/contrib/Eclipse-plugin/Hadoop-0.20.3-dev-Eclipse-plugin.jar 名字为 Hadoop-0.20.2-Eclipse-plugin.jar,这样就可以生成自己的 Hadoop Eclipse plugin 了。

4) 把这个 jar 包放到 Eclipse 的 plugins 目录下,然后重启 Eclipse 就可以使用了。

18.3 Hadoop Streaming 简介和使用

Hadoop Streaming 是 Hadoop 的一个工具,它帮助用户创建和运行一类特殊的 MapReduce 作业,这些特殊的 MapReduce 作业由一些可执行文件或脚本文件充当 mapper 或 reducer。也就是说 Hadoop Streaming 允许用户用非 Java 的编程语言编写 MapReduce 程序,然后 Streaming 用 STDIN (标准输入) 和 STDOUT (标准输出) 来和我们编写的 Map 和 Reduce 进行数据交换,并提交给 Hadoop。命令格式如下:

```
$HADOOP_HOME/bin/hadoop jar $HADOOP_HOME/hadoop-streaming.jar \
-input myInputDirs \
-output myOutputDir \
-mapper /bin/cat \
-reducer /bin/wc
```

1. Streaming 的工作原理

在上面的命令里, mapper 和 reducer 都是可执行文件,它们从标准输入中按行读入数据,并把计算结果发送给标准输出。Streaming 工具会创建一个 MapReduce 作业,并把它发送给合适的集群,同时监视这个作业的整个执行过程。

如果一个可执行文件被用于 mapper,则在其初始化时,每一个 mapper 任务会把这个可执行文件作为一个单独的进程启动。mapper 任务运行时,它把输入切分成行,并把结果提供

给可执行文件对应进程的标准输入。同时，它会收集可执行文件进程标准输出的内容，并把收到的每一行内容转化成 key/value 对，作为输出。在默认情况下，一行中第一个 tab 之前的部分被当作 key，之后的（不包括 tab）被当作 value。如果没有 tab，则整行内容被当作 key 值，value 值为 null。具体的转化策略会在下文中讨论。

如果一个可执行文件被用于 reducer，每个 reducer 任务同样会把这个可执行文件作为一个单独的进程启动。reducer 任务运行时，它把输入切分成行，并把结果提供给可执行文件对应进程的标准输入。同时，它会收集可执行文件进程标准输出的内容，并把每一行内容转化成 key/value 对，作为输出。默认情况下，一行中第一个 tab 之前的部分被当作 key，之后的（不包括 tab）被当作 value。

用户也可以使用 java 类作为 mapper 或 reducer。本节最初给出的命令与这里的命令等价：

```
$HADOOP_HOME/bin/Hadoop jar $HADOOP_HOME/Hadoop-streaming.jar \
    -input myInputDirs \
    -output myOutputDir \
    -mapper org.apache.hadoop.mapred.lib.IdentityMapper \
    -reducer /bin/wc
```

用户可以设定 stream.non.zero.exit.is.failure 的值为 true 或 false，从而表明 streaming task 的返回值非零时是 Failure 还是 Success。默认情况下，streaming task 返回非零时表示失败。

2. 将文件打包到提交的作业中

利用 Streaming 用户可以将任何可执行文件指定为 mapper/reducer。这些可执行文件可以事先存放在集群上，也可以用 -file 选项让可执行文件成为作业的一部分，并且会一起打包提交。例如：

```
$HADOOP_HOME/bin/hadoop jar $HADOOP_HOME/hadoop-streaming.jar \
    -input myInputDirs \
    -output myOutputDir \
    -mapper myPythonScript.py \
    -reducer /bin/wc \
    -file myPythonScript.py
```

上面的例子描述了一个用户把可执行 Python 文件指定为 mapper。其中的选项“-file myPythonScript.py”使可执行 Python 文件作为作业的一部分被上传到集群的机器上。

除了可执行文件外，其他 mapper 或 reducer 需要用到的辅助文件（比如字典、配置文件等）也可以用这种方式打包上传。例如：

```
$HADOOP_HOME/bin/hadoop jar $HADOOP_HOME/hadoop-streaming.jar \
    -input myInputDirs \
    -output myOutputDir \
    -mapper myPythonScript.py \
    -reducer /bin/wc \
    -file myPythonScript.py \
    -file myDictionary.txt
```

3. Streaming 选项与用法

(1) 只使用 Mapper 的作业

有时候只需要使用 map 函数处理输入数据。这时只须把 `mapred.reduce.tasks` 设置为零，MapReduce 框架则不会创建 reducer 任务，mapper 任务的输出就是整个作业的最终输出。

为了做到向下兼容，Hadoop Streaming 也支持 “-reduce None” 选项，它与 “-jobconf mapred.reduce.tasks=0” 等价。

(2) 为作业指定其他插件

和其他普通的 MapReduce 作业一样，用户可以为 Streaming 作业指定数据格式，命令如下：

```
-inputformat JavaClassName
-outputformat JavaClassName
-partitioner JavaClassName
-combiner JavaClassName
```

如果不指定输入格式，程序会默认使用 `TextInputFormat`。因为 `TextInputFormat` 得到的 key 值是 `LongWritable` 类型的（key 值并不是输入文件中的内容，而是 value 偏移量），所以 key 会被丢弃，只有 value 被用管道方式发给 mapper。

另外，用户提供的定义输出格式的类需要能够处理 Text 类型的 key/value 对。如果不指定输出格式，则默认会使用 `TextOutputFormat` 类。

(3) Hadoop Streaming 中的大文件和档案

任务依据 `-cacheFile` 和 `-cacheArchive` 选项在集群中分发文件和档案，选项的参数是用户已上传至 HDFS 的文件或档案的 URI。这些文件和档案在不同的作业间缓存。用户可以通过 `fs.default.name.config` 配置参数的值得到文件所在的 host 和 `fs_port`。

下面是使用 `-cacheFile` 选项的例子：

```
-cacheFile hdfs://host:fs_port/user/testfile.txt#testlink
```

在上面的例子里，URL 中 # 后面的内容是建立在任务当前工作目录下的符号链接的名字。这个任务的当前工作目录下有一个 “testlink” 符号链接，它指向 `testfile.txt` 文件在本地的拷贝位置。如果有多个文件，选项可以写成：

```
-cacheFile hdfs://host:fs_port/user/testfile1.txt#testlink1 -cacheFile hdfs://
host:fs_port/user/testfile2.txt#testlink2
```

`-cacheArchive` 选项用于把 jar 文件拷贝到任务当前的工作目录中，并自动把 jar 文件解压缩。例如：

```
-cacheArchive hdfs://host:fs_port/user/testfile.jar#testlink3
```

在上面的例子中，`testlink3` 是当前工作目录下的符号链接，它指向 `testfile.jar` 解压后的目录。

下面是使用 `-cacheArchive` 选项的另一个例子。其中，`input.txt` 文件有两行内容，分别是

两个文件的名字：testlink/cache.txt 和 testlink/cache2.txt。“testlink”是指向档案目录（jar 文件解压后的目录）的符号链接，这个目录下有“cache.txt”和“cache2.txt”两个文件。代码如下所示：

```
$HADOOP_HOME/bin/Hadoop jar $HADOOP_HOME/Hadoop-streaming.jar \
    -input "/user/me/samples/cachefile/input.txt" \
    -mapper "xargs cat" \
    -reducer "cat" \
    -output "/user/me/samples/cachefile/out" \
    -cacheArchive 'hdfs://Hadoop-nn1.example.com/user/me/samples/
        cachefile/cachedir.jar#testlink' \
    -jobconf mapred.map.tasks=1 \
    -jobconf mapred.reduce.tasks=1 \
    -jobconf mapred.job.name="Experiment"

$ ls test_jar/
cache.txt  cache2.txt

$ jar cvf cachedir.jar -C test_jar/ .
added manifest
adding: cache.txt(in = 30) (out= 29)(deflated 3%)
adding: cache2.txt(in = 37) (out= 35)(deflated 5%)

$ Hadoop dfs -put cachedir.jar samples/cachefile

$ Hadoop dfs -cat /user/me/samples/cachefile/input.txt
testlink/cache.txt
testlink/cache2.txt

$ cat test_jar/cache.txt
This is just the cache string

$ cat test_jar/cache2.txt
This is just the second cache string

$ Hadoop dfs -ls /user/me/samples/cachefile/out
Found 1 items
/user/me/samples/cachefile/out/part-00000 <r 3> 69

$ Hadoop dfs -cat /user/me/samples/cachefile/out/part-00000
This is just the cache string
This is just the second cache string
```

4. 为作业指定附加配置参数

用户可以使用“-jobconf <n>=<v>”增加一些配置变量。例如：

```
$HADOOP_HOME/bin/Hadoop jar $HADOOP_HOME/Hadoop-streaming.jar \
    -input myInputDirs \
    -output myOutputDir \
```

```
-mapper org.apache.Hadoop.mapred.lib.IdentityMapper\
-reducer /bin/wc \
-jobconf mapred.reduce.tasks=2
```

在上面的例子中，`-jobconf mapred.reduce.tasks=2` 表明用两个 reducer 完成作业。

关于 `jobconf` 参数的更多细节可以参考 Hadoop 安装包中的 `Hadoop-default.html` 文件。

5. 其他选项

Streaming 命令的其他选项如表 18-3 所示：

表 18-3 Streaming 命令选项表

选 项	可选 / 必须	描 述
<code>-cluster name</code>	可选	在本地 Hadoop 集群与一个或多个远程集群之间进行切换
<code>-dfs host:port or local</code>	可选	覆盖作业的 HDFS 配置
<code>-jt host:port or local</code>	可选	覆盖作业的 JobTracker 配置
<code>-additionalconfspec specfile</code>	可选	用一个类似于 <code>Hadoop-site.xml</code> 的 XML 文件保存所有配置，从而不需要用多个 <code>"-jobconf name=value"</code> 类型的选项单独为每个配置的变量赋值了
<code>-cmdenv name=value</code>	可选	传递环境变量给 streaming 命令
<code>-cacheFile fileNameURI</code>	可选	指定一个上传到 HDFS 的文件
<code>-cacheArchive fileNameURI</code>	可选	指定一个上传到 HDFS 的 jar 文件，这个 jar 文件会被自动解压缩到当前工作目录下
<code>-inputreader JavaClassName</code>	可选	为了向下兼容：指定一个 record reader 类（而不是 input format 类）
<code>-verbose</code>	可选	详细输出

使用 `-cluster <name>` 实现“本地”Hadoop 和一个或多个远程 Hadoop 集群间的切换。默认情况下，使用 `Hadoop-default.xml` 和 `Hadoop-site.xml`。当使用 `-cluster <name>` 选项时，会使用 `$HADOOP_HOME/conf/Hadoop-<name>.xml`。

下面的选项可改变 temp 目录：

```
-jobconf dfs.data.dir=/tmp
```

下面的选项指定其他本地 temp 目录：

```
-jobconf mapred.local.dir=/tmp/local
-jobconf mapred.system.dir=/tmp/system
-jobconf mapred.temp.dir=/tmp/temp
```

在 streaming 命令中设置环境变量：

```
-cmdenv EXAMPLE_DIR=/home/example/dictionaries/
```

18.3.1 Hadoop Streaming 的使用举例

Hadoop Streaming 插件是 Hadoop 安装包其中的一个 Jar 文件，具体位置在 `...\\Hadoop-`

0.20.2\contrib\streaming 目录下，所以 Hadoop Streaming 插件是直接使用的，只需要在执行 Hadoop 程序时输入命令 Hadoop Streaming 就可以了，无须安装，在编写 MapReduce 程序时，只要按照整个框架要求并根据自己的需要编写出符合对应语言格式的程序，然后用下面的命令格式将程序提交给 Hadoop 就可以了：

```
$HADOOP_HOME/bin/hadoop jar $HADOOP_HOME/hadoop-streaming.jar \
-input myInputDirs \
-output myOutputDir \
-mapper /bin/cat \
-reducer /bin/wc
```

需要注意的是，程序执行所需要的支持文件也要在提交程序的同时提交到 Hadoop 集群上，这在上文已有说明，不再赘述。下面以一个用 PHP 语言编写的 WordCount 使用 Hadoop Streaming 提交的程序为例，来说明此插件的使用方法（Linux 系统下需要安装 PHP 环境，命令为 `sudo apt-get install php5-client`）。

程序代码举例如下所示。

(1) Mapper.php

```
#!/usr/bin/php
<?php

$word2count = array();

// 标准输入 STDIN (standard input)
while (($line = fgets(STDIN)) !== false) {
    // 移除小写与空格
    $line = strtolower(trim($line));
    // 切词
    $words = preg_split('/\W/', $line, 0, PREG_SPLIT_NO_EMPTY);
    // 将字+1
    foreach ($words as $word) {
        $word2count[$word] += 1;
    }
}

// 结果写到 STDOUT (standard output)
foreach ($word2count as $word => $count) {
    echo $word, chr(9), $count, PHP_EOL;
}
?>
```

(2) Reduce.php

```
#!/usr/bin/php
<?php

$word2count = array();
```



```
// 输入为 STDIN
while (($line = fgets(STDIN)) != false) {
    // 移除多余的空白
    $line = trim($line);
    // 每一行的格式为 ( 字 "tab" 数字 ), 记录到 ($word, $count)
    list($word, $count) = explode(chr(9), $line);
    // 转换格式 string -> int
    $count = intval($count);
    // 求总的频数
    if ($count > 0) $word2count[$word] += $count;
}

// 此行非必要内容, 但可让 output 排列更完整
ksort($word2count);

// 将结果写到 STDOUT (standard output)
foreach ($word2count as $word => $count) {
    echo $word, chr(9), $count, PHP_EOL;
}

?>
```

执行情况如下:

```
$ bin/Hadoop jar contrib/streaming/Hadoop-0.20.2-streaming.jar \
-mapper /opt/Hadoop/mapper.php -reducer /opt/Hadoop/reducer.php -input lab4_input
-output stream_out2
```

下面来看一下结果:

```
$ bin/Hadoop dfs -cat stream_out2/part-00000
```

18.3.2 使用 Hadoop Streaming 时常见的问题

1. 如何处理多个文件, 每个文件一个 map ?

需要处理多个文件时, 用户可以采用多种途径, 这里以在集群上压缩 (zipping) 多个文件为例, 用户可以使用以下几种方法:

(1) 使用 Hadoop Streaming 和用户编写的 mapper 脚本程序

先生成一个文件, 文件中包含所有要压缩的文件在 HDFS 上的完整路径。每个 map 任务获得一个路径名作为输入。

然后创建一个 mapper 脚本程序, 实现如下功能: 获得文件名, 把该文件拷贝到本地, 压缩该文件并把它发到期望的输出目录中。

(2) 使用现有的 Hadoop 框架

在 main 函数中添加如下命令:

```
FileOutputFormat.setCompressOutput(conf, true);
```

```
FileOutputFormat.setOutputCompressorClass(conf, org.apache.Hadoop.
    io.compress.GzipCodec.class);
conf.setOutputFormat(NonSplitableTextInputFormat.class);
conf.setNumReduceTasks(0);
```

编写 map 函数:

```
public void map(WritableComparable key, Writable value,
    OutputCollector output,
    Reporter reporter) throws IOException {
    output.collect((Text)value, null);
}
```

注意输出的文件名和原文件名不同。

2. 如果在 Shell 脚本里设置一个别名, 并放在 `-mapper` 之后, Streaming 会正常运行吗? 例如, `alias cl='cut -f1'`, `-mapper "cl"` 会运行正常吗?

脚本里是无法使用别名的, 但是允许变量替换, 例如:

```
$ Hadoop dfs -cat samples/student_marks
alice 50
bruce 70
charlie 80
dan 75

$ c2='cut -f2'; $HADOOP_HOME/bin/Hadoop jar $HADOOP_HOME/Hadoop-streaming.jar \
    -input /user/me/samples/student_marks
    -mapper \"$c2\" -reducer 'cat'
    -output /user/me/samples/student_out
    -jobconf mapred.job.name='Experiment'

$ Hadoop dfs -ls samples/student_out
Found 1 items/user/me/samples/student_out/part-00000    <r 3>    16

$ Hadoop dfs -cat samples/student_out/part-00000
50
70
75
80
```

3. 在 Streaming 作业中用 `-file` 选项运行一个分布式的超大可执行文件 (例如, 3.6GB) 时, 如果得到错误信息 “No space left on device” 如何解决?

由于配置变量 `stream.tmpdir` 指定了一个目录, 所以会在这个目录下进行打 jar 包的操作。stream.tmpdir 的默认值是 `/tmp`, 用户需要将这个值设置为一个有更大空间的目录:

```
-jobconf stream.tmpdir=/export/bigspace/...
```

4. 如何设置多个输入目录?

可以使用多个 `-input` 选项设置多个输入目录:

```
Hadoop jar Hadoop-streaming.jar -input '/user/foo/dir1' -input '/user/foo/dir2'
```

5. 如何生成 gzip 格式的输出文件?

除了纯文本格式的输出, 用户还可以让程序生成 gzip 文件格式的输出, 只需将 Streaming 作业中的选项设置为“-jobconf mapred.output.compress=true -jobconf mapred.output.compression.codec=org.apache.Hadoop.io.compress.GzipCode”。

6. 在 Streaming 中如何自定义 input/output format ?

在 Hadoop 0.14 版本以前, 不支持多个 jar 文件。所以当指定自定义的类时, 用户需要把它们和原有的 streaming jar 打包在一起, 并用这个自定义的 jar 包替换默认的 Hadoop streaming jar 包。在 0.14 版本以后, 就无须打包在一起了, 只需要正常的编译运行。

7. Streaming 如何解析 XML 文档?

用户可以使用 StreamXmlRecordReader 来解析 XML 文档, 如下所示:

```
Hadoop jar Hadoop-streaming.jar -inputreader "StreamXmlRecord,begin=BEGIN_STRING,end=END_STRING" ...
```

Map 任务会把 BEGIN_STRING 和 END_STRING 之间的部分看作一条记录。

8. 在 Streaming 应用程序中如何更新计数器?

Streaming 进程能够使用 stderr 发出计数器信息。应该把 reporter:counter:<group>,<counter>,<amount> 发送到 stderr 来更新计数器。

9. 如何更新 Streaming 应用程序的状态?

Streaming 进程能够使用 stderr 发出状态信息。可把 reporter:status:<message> 发送到 stderr 来设置状态。

18.4 Hadoop Libhdfs 简介和使用

Libhdfs 是一个基于 C 编程接口的为 Hadoop 分布式文件系统开发的 JNI (Java Native Interface), 它提供了一个 C 语言接口以结合管理 DFS 文件和文件系统, 并且它会在 \${HADOOP_HOME}/libhdfs/libhdfs.so 中预编译, 它是 Hadoop 分布式结构中的一部分。

18.4.1 Hadoop Libhdfs 安装和配置

在安装 Libhdfs 之前首先需要安装 Hadoop 的分布式文件系统 HDFS。当用户有一个正在运行的工作集时, 进入 src/c++/libhdfs 目录, 使用 makefile 文件安装 Libhdfs。一旦安装 Libhdfs 成功, 用户可以通过它连接到自己的程序。

18.4.2 Hadoop Libhdfs API 简介

这节将介绍 Libhdfs 提供的各种用来管理 DFS 的 API。我们按照管理对象(单个文件、文件系统)对 API 进行了分类, 如下所示:

1. FileSystem API

Libhdfs 不仅提供了通用文件系统管理 API, 比如创建文件夹、复制文件、移动文件等,

还提供了一些特殊功能的 API，比如获取备份文件信息等。

在启动时应该在任何文件或文件系统操作之前先用 HDFSconnect API 将 Libhdfs 和 DFS 连接起来。还有一个类似的 hdfsdisconnect API 负责连接的清除。

通用操作如下所示：

- 1) hdfsCopy (也适用文件系统)
- 2) hdfsMove (也适用文件系统)
- 3) hdfsRename
- 4) hdfsDelete

Libhdfs 还提供了在 DFS 上管理目录的 API，如下所示。

- 1) hdfsCreateDirectory
- 2) hdfsSetWorkingDirectory
- 3) hdfsGetWorkingDirectory
- 4) hdfsListDirectory / hdfsGetPathInfo / hdfsFreeFileInfo

查询 filesystem 各种属性的 API 如下所示。

- 1) hdfsGetHosts
- 2) hdfsGetDefaultBlockSize
- 3) hdfsGetUsed / hdfsGetCapacity

2. File APIs

Libhdfs 提供了一些类似于 POSIX (Portable Operating System Interface) 的接口来实现单个文件上的操作，比如 create、read/write、query 等，如下所示：

- 1) hdfsOpenFile / hdfsCloseFile
- 2) hdfsRead / hdfsWrite
- 3) hdfsTell / hdfsSeek
- 4) hdfsFlush
- 5) hdfsAvailable

18.4.3 Hadoop Libhdfs 的使用举例

Libhdfs 可以用在 POSIX 线程编写的线程应用程序中。无论是与 JNI 的全局还是局部引用交互，使用者都必须显式地调用 hdfsConvertToGlobalRef / hdfsDeleteGlobalRef API。

下面是一个 Libhdfs 应用的程序：

```
#include "hdfs.h"

int main(int argc, char **argv) {

    hdfsFS fs = hdfsConnect("default", 0);
    const char* writePath = "/tmp/testfile.txt";
    hdfsFile writeFile = hdfsOpenFile(fs, writePath, O_WRONLY|O_CREAT, 0, 0, 0);
```

```

    if(!writeFile) {
        fprintf(stderr, "Failed to open %s for writing!\n", writePath);
        exit(-1);
    }
    char* buffer = "Hello, World!";
    tSize num_written_bytes = hdfsWrite(fs, writeFile, (void*)buffer,
        strlen(buffer)+1);
    if (hdfsFlush(fs, writeFile)) {
        fprintf(stderr, "Failed to 'flush' %s\n", writePath);
        exit(-1);
    }
    hdfsCloseFile(fs, writeFile);
}

```

接下来再介绍一些具体问题的解决办法。

(1) 如何连接到库

使用 Libhdfs 源文件目录 (\${HADOOP_HOME}/src/c++/libhdfs/Makefile) 下的 makefile 或使用下面的命令：gcc above_sample.c -I\${HADOOP_HOME}/src/c++/libhdfs -L\${HADOOP_HOME}/libhdfs -lhdfs -o above_sample 来连接库。

(2) CLASSPATH 配置问题

使用 Libhdfs 最常见的问题就是在运行一个使用了 Libhdfs 的程序时 CLASSPATH 没有正确配置 Libhdfs。请确保在每个运行 Hadoop 所必需的 Hadoop jar 包中对其进行了正确配置。另外，目前还没有使用程序自动生成 CLASSPATH 的方法，但有一个很好的解决办法就是引用 \${HADOOP_HOME} 和 \${HADOOP_HOME}/lib 下的所有 JAR 包，并且正确配置 hdfs-site.xml 中的目录。

18.5 小结

本章介绍了使用 Hadoop 开发的 4 种常用插件，分别是 Hadoop Studio、Hadoop Eclipse、Hadoop Streaming 和 Hadoop LibHdfs。Hadoop Studio 是一个加快 Hadoop 开发进程的可视化开发环境。Hadoop Studio 通过降低 Hadoop 的使用复杂度让用户在更少的步骤内完成更多的事情来提高生产率。用户可以通过 Hadoop Studio 强大的 GUI，部署 Hadoop 任务，并监控 Hadoop 任务的实时信息。无论用户的开发经验有多少，Studio 都能从设计、部署、调试和可视化四个方面来简化用户的工作，提高工作效率。Hadoop Studio 全面强大的功能使其使用范围甚广。

Hadoop Eclipse 插件将 Hadoop 的开发环境图形化。在安装插件并配置 Hadoop 的相关信息之后，若用户创建 Hadoop 程序，插件会自动导入 Hadoop 编程接口的 jar 文件，这样，用户就可以在 Eclipse 的图形化界面中编写、调试、运行 Hadoop 程序（单机程序和分布式程序都可以）了，也可以在其中查看自己程序的实时状态、错误信息和运行结果了，还可以查看、管理 HDFS 和其他文件。总的来说，Hadoop Eclipse 插件安装简单，使用方便，功能强

大，尤其是在 Hadoop 编程方面，是 Hadoop 入门和 Hadoop 编程必不可少的工具。

Hadoop Streaming 是 Hadoop 的一个工具，它帮助用户创建和运行一类特殊的 MapReduce 作业，这些特殊的 MapReduce 作业由一些可执行文件或脚本文件充当 mapper 或 reducer。本章也举例说明了它的使用方法。

Libhdfs 是一个基于 C 编程接口，为 Hadoop 分布式文件系统开发的 JNI (Java Native Interface)，它提供了一个 C 语言接口以结合管理 DFS 文件和文件系统。它在 `${HADOOP_HOME}/libhdfs/libhdfs.so` 中预编译，是 Hadoop 分布式结构中的一部分。其丰富的 API 方便了用户对于 HDFS 和 HDFS 文件的管理。在这部分内容的最后给出了 Libhdfs 使用的具体例子，并说明了一些常见问题的解决办法。

本章详细介绍了 Hadoop 开发常用的四种插件，从安装步骤到使用方法，再到常见问题的解决方法，希望能帮忙读者提高使用和开发 Hadoop 的效率。





附录 A

云计算在线检测平台

本章内容

- ☐ 平台介绍
- ☐ 结构和功能
- ☐ 检测流程
- ☐ 使用
- ☐ 小结

资源分享

PDG

A.1 平台介绍

MapReduce 日趋流行,带起了普通程序员相关知识的热潮,它的学习资料也日趋丰富起来。但是 MapReduce 运行所需的并行环境却成为了入门者最大的障碍,主要原因是并行环境的硬件要求高,配置复杂,同时现有的学习资料中鲜有编程实战方面的指导,更多专注在 MapReduce 的理论知识上。综合这些情况,我们开发了云计算在线检测平台 (<http://cloudcomputing.ruc.edu.cn/>),为读者提供理论知识测试和利用理论知识进行实战的机会,此平台提供了运行程序的并行环境,避免入门者将精力都浪费在环境配置上,让其配合本书学习实践。

云计算在线检测平台是一个 MapReduce 程序检测平台。此平台基于 Hadoop 集群提供了 MapReduce 并行程序运行的分布式环境,它旨在为 MapReduce 的入门者提供简单具体的编程练习,使其初步掌握 MapReduce 框架的编程思想,并拥有使用 MapReduce 并行化解决实际问题的能力。用户可以根据平台提供的问题背景,开发自己的并行程序,并提交运行,平台会根据运行结果反馈给用户一定的信息,以便进行修改或进一步优化。用户也可以在网站上进行分布式系统理论知识的测试,以提高理论水平。同时,此平台结合分布式系统架构 Hadoop、MySQL 技术和 Tomcat 技术,提供了在线的分布式并行运行环境,为用户运行他所提交的并行程序。平台主要由前台用户接口和后台程序运行两部分组成,前台用户接口负责同用户的交互,包括保存用户提交的代码、返回程序的检测结果等,后台程序运行负责将前台收集的用户代码运行并检测结果,同时将检测的结果交给前台并维护网站用户的信息,提供整个网站的网络服务。

云计算在线检测平台兼顾实战和理论,能让用户在进行理论测试和了解其缺陷所在的过程中掌握开源分布式系统架构 Hadoop 的相关知识和 MapReduce 的理论知识,能让用户在编程提交和修改再提交的过程中切身体验到如何利用分布式系统 Hadoop、MapReduce 编程,以及如何用 MapReduce 并行程序来解决实际问题。总体来说,此平台能够提高用户的理论水平和实战能力,是 MapReduce 入门者不错的入门指导。

A.2 结构和功能

正如前一节中所介绍的,云计算在线检测平台由两大部分组成,分别是前台用户接口和后台程序运行,下面分别对它们进行介绍。

A.2.1 前台用户接口的结构和功能

前台用户接口的功能结构如图 A-1 所示。它主要包括四部分内容:用户完全服务、实例编程练习、分布式系统理论知识测试、帮助功能(指网站的使用帮助、Hadoop 介绍文档,以及网站的中文页面)。下面分别详细介绍这四个功能块。

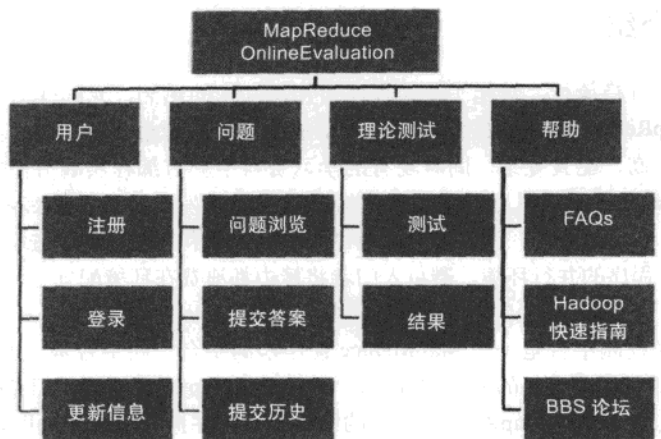


图 A-1 前台界面结构图

用户完全服务主要包括用户注册、登录和更改个人信息等。用户注册是指用户在 Register 页面完成新用户的注册，平台只对注册用户提供代码检测服务。用户注册页面需要提供用户名、注册码（选填）、密码、单位、邮箱等信息，注册成功之后用户就可以使用注册的用户名和密码登录了，同时邮箱会收到一封注册邮件，以防止用户忘记用户名和密码以致无法登录。在注册时如果用户名已注册、密码重复错误、验证码输入错误等都会导致注册失败。可以在首页的右上角直接使用用户名和密码登录，也可以在 login 页面完成登录操作。登录成功的用户可以选择 login out。更改个人信息是指更改个人密码等信息，如果用户期望做出更改，可以在 update your info 页面完成。

MapReduce 实例编程练习主要包括题目浏览、提交题目、查看提交记录、查看提交源码、查看检测结果。在云计算在线检测的平台上，开发小组设计了很多基于 MapReduce 并行框架能够解决的实际问题。用户可以在 problem 页面详细浏览各个问题的背景，以及输入输出要求和注意事项。然后利用自己的 MapReduce 理论知识，针对具体实例问题来编写自己的实例解决代码，并可以在 submit solution 页面提交代码。网站会运行用户提交的代码，然后在网站上反馈相应的结果。用户可以在 My submission 页面中查看自己的提交记录，也可以点击每一条记录中的 source 连接查看自己提交的代码，同时也可以点击 result 栏下面的连接查看检测结果。

分布式系统理论知识测试主要指用户点击首页的 theory test 之后会出现一份限时半小时的试卷，共 20 道选择题。这些选择题都在平台中随机从题库里选出来的，题目是关于分布式系统的理论知识，主要集中在 Hadoop 及其子项目上。用户答题完毕点击提交或在页面上停留的时间超过半小时，所有答案就会提交。平台通过对比之后会将每道题目的正确答案及你的回答一起返回，并计算出此次测试的分数。

帮助功能主要指网站的使用帮助、网站对应的中文页面，以及 Hadoop 的介绍文档和用

于讨论的 BBS 板块。网站的使用帮助在首页的 FAQs 页面下，主要是网站在实际使用中要注意的事项。网站对应的中文页面可以点击 Chinese 链接进入。中文页面也提供了与英文页面完全相同的服务。Hadoop 介绍文档指网站上的 Hadoop Quick Start 链接和它所提供的在线文档，主要为用户提供一些 Hadoop 分布式系统的初步认识和安装说明。BBS 板块允许用户在平台上交流 MapReduce 的学习经验，以及对平台上题目的交流，同时还可以留下自己关于平台使用的疑问。

A.2.2 后台程序运行的结构和功能

后台程序运行的功能结构如图 A-2 所示。后台中的主要模块也是四部分：Tomcat 服务器、MySQL 数据库、Hadoop 分布式环境、shell 文档。下面详细介绍这四个功能块。

Tomcat 服务器：Tomcat 服务器担当网站的 Web 服务器角色，保证用户能够从网络上访问到平台，并将开发小组基于 JSP 技术开发的网页呈现在用户的电脑上。

MySQL 数据库：MySQL 数据库里主要是网站的信息，包括用户个人信息、用户提交记录、网站题库等。基于 JSP 技术开发的网页通过调用 MySQL 的接口，获取用户请求的信息，并将其呈现在网页上或将网页上提交的信息保存到数据库中。

Hadoop 分布式环境：Hadoop 分布式环境是整个后台的核心所在，因为它是云计算在线检测平台提供特色服务的核心。开发小组首先在多台计算机上安装好 Hadoop 分布式系统，形成一个分布式环境，然后再在集群上配置网站提供服务，这就可以保证为用户提交的代码提供并行程序运行所需的真实分布式环境。Hadoop 集群的主要功能就是运行用户提交的代码，给出结果。

Shell 文档：Shell 文档在检测平台的系统中扮演着人体中血液的角色。它首先将网页保存下来的用户代码进行预处理，比如检测是否是正确的 Java 程序等，然后再对预处理之后的结果进行预编译，成功之后再将代码提交到 Hadoop 上，接着再收集 Hadoop 的运行结果，然后与标准结果进行比对，最后将比对的结果分类返回给前台网页，呈现到用户面前。综合来说，Shell 文档将前台功能块和后台功能块串联了起来，以便为用户提供连贯的服务。

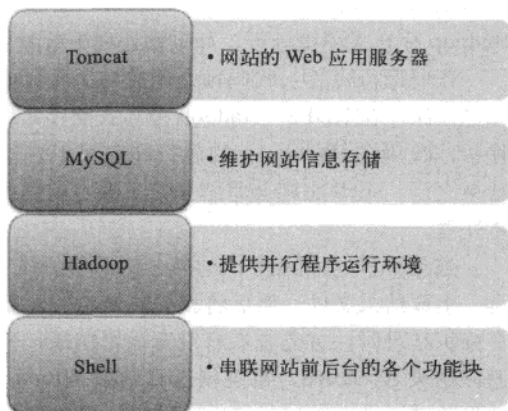


图 A-2 后台结构图

A.3 检测流程

经过前面两节的介绍，读者对整个平台已经有了一个直观整体的认识，那么这个平台是如何运行的呢？它的运行流程是什么？本节将详细介绍云计算在线检测平台检测用户代码的

流程。

总体来说，平台对用户代码的检测流程主要包括代码保存、代码预处理、代码运行和结果分析返回、结果显示五个阶段，下面将分别介绍这五个阶段。

代码保存阶段：用户在网页上粘贴自己的代码，点击 submit 提交之后，网站会把用户的代码保存在服务器上一个唯一的文件中，并在后台数据库中保存这一次提交的信息和代码路径。

代码预处理阶段：用户在提交代码之后，网站在进行代码保存的同时还会调用 Shell 文档来进行代码的预处理。shell 文档被调用运行之后就会开始用户代码的预处理。首先 shell 文档会按照调用的路径参数从本地找到用户代码，然后检测用户代码，比如程序是否是可运行的 Java 代码、是否符合 MapReduce 框架要求等。如果预处理成功就会将代码提交给 Hadoop 分布式环境运行，如果预处理失败就会直接返回并将错误原因呈现到网页界面上。

代码运行阶段：代码预处理成功之后会被提交到 Hadoop 分布式环境上，Hadoop 调用事先已经保存在 HDFS（Hadoop 分布式文件系统）上的输入数据来运行代码。在平台的处理过程中，代码在 Hadoop 上的运行和在线下自己提交代码到 Hadoop 上的运行相同。代码运行结束之后，shell 文档会将结果信息重定向到代码文件里唯一的结果信息文件中，以交给下一步处理。

结果分析返回阶段：结果分析返回阶段主要是分析 Hadoop 运行的结果信息，对结果分类，生成结果文件，然后将相关的信息写入数据库，供平台显示代码运行结果时调用。Shell 在分析结果时，首先查看有没有输出结果，如果有输出结果就和标准输出进行对比，正确就返回结果 Accepted，错误就返回结果 Wrong Answer。如果没有结果，再将输出信息同一些结果关键词进行匹配，然后返回匹配成功的那一类错误信息。

结果显示：用户在 My Submission 界面点击 result 一栏的结果链接之后，页面会调用数据库接口，搜索此次提交记录在数据库中的对应记录。找到之后，页面直接获取结果信息文件的路径，然后将其内容显示在页面上，如果代码有误，会告知用户错误所在，用户可进行调整之后再重新提交。

A.4 使用

前面介绍了云计算在线检测平台的理论内容，本节将从功能使用、题目介绍、返回结果说明、使用注意事项四个方面详细介绍平台的使用方法。

A.4.1 功能使用

本附录第 2 节介绍了平台中前台用户接口和后台程序运行的结果和功能块。而与用户直接相关的就是前台功能的使用。下面用三个使用实例来说明如何使用前台功能。

1. 如何注册用户，如何修改信息？

注册功能的使用流程如下：

- 1) 在首页点击 Register 链接, 进入注册界面;
- 2) 填写个人信息, 包括用户名、注册码(选填)、密码、邮箱、单位、国家、验证码等;
- 3) 根据提示进行调整, 比如如果提示用户名已存在, 就需要换一个用户名, 如果提示密码重复错误, 就需要重新输入密码等;
- 4) 注册成功, 如果注册完之后可以进入注册成功界面, 就表示注册成功了, 界面上显示的是自己除密码以外的所有注册信息, 同时用户所注册的邮箱会收到一封包含用户名和密码的注册邮件, 以防止忘记用户名或密码。

使用修改信息功能的流程如下:

- 1) 登录之后在首页点击 Update your info 链接, 进入信息修改页面;
- 2) 填写要修改的个人信息;
- 3) 点击提交之后就会进入修改成功界面, 界面显示修改的信息。

2. 如何提交自己的代码并查看结果?

1) 登录之后点击具体题目下的 submit 按钮, 进入代码提交页面, 或者点击 submit solution 链接直接进入提交页面, 再或者在首页的 problem 一栏下输入 problem ID 直接进入 problem, 然后点击提交进入代码提交页面;

2) 在代码提交页面的空白处粘贴自己的代码, 点击提交;

3) 提交之后页面自动跳入仅包含此次提交信息的页面, 在这个页面中用户能够查看自己提交的代码, 同时页面还能够在代码运行结束之后自动更新 result 一栏的状态, 并显示运行结果(此处采用了 AJAX 技术, 由于存在技术兼容问题, 所以只有 firefox 支持), 更新之后用户可以点击查看运行结果;

4) 用户想查看结果和自己的代码, 也可以点击 My Submission 链接, 进入自己的提交记录页面, 点击特定记录后的 source 就可以查看提交的代码了, 点击 result 一栏可以查看具体的结果信息。

3. 如何进行理论测试?

- 1) 登录后点击 Theory test 进入理论测试界面;
- 2) 根据具体的题目选择正确答案, 然后提交(理论测试每份试卷限时 30 分钟, 如果在页面上停留的时间超过 30 分钟, 平台也会自己提交页面现有答案);
- 3) 提交之后页面自动跳入结果页面, 显示每到题目的回答是否正确。

A.4.2 返回结果介绍

在平台上提交代码之后在提交历史中的 result 一栏就可以看到结果。那么都有什么结果? 都代表什么意思? 针对具体的错误用户应该如何应对? 下面将进行详细介绍。

Accepted: 表示用户提交的代码已经被接受, 而用户代码被接受的前提是代码能够正确运行, 并且以平台的测试数据作为输入数据执行后的输出结果和平台的标准输出结果完全相同。但需要提醒的是, 由于 MapReduce 编程框架的原因, 平台上的这些题目完全可以在

MapReduce 结构中的 map 或 reduce 阶段独立完成，只是这种做法没有完全发挥 MapReduce 并行运行的效率，不是最优的办法。所以如果用户的代码被 Accepted 了，用户还需要审视自己的代码，检查它是否最大程度利用了并行运行来提高效率。

Compile Error：表示用户代码编译错误，出现这种情况说明在用户的代码存在语法问题，在进行普通的 Java 程序编译时出错了。用户可以点击 result 栏的错误结果链接去查看具体的语法错误位置，并进行修改。用户也可以在本地进行普通的 Java 编译，待通过之后再提交到平台上。

MapReduce Error：表示代码在 Hadoop 上运行时出现错误并没有输出结果。这种情况出现的可能原因比较多，主要包括：常见的 Java 程序逻辑错误、MapReduce 逻辑错误等。Java 程序逻辑错误又主要包括数组越界、未初始化等，MapReduce 逻辑错误则主要包括输入输出类型不匹配等。在遇到 MapReduce Error 时相对比较麻烦，需要用户仔细核对代码，找出逻辑错误的地方进行修改，然后再尝试提交。

Wrong Answer：表示代码能够在 Hadoop 上正常运行并有输出结果，只是用户的输出结果和标准结果并不匹配。出现这种情况时，用户首先要检查自己代码的输出格式是否正确，比如顺序是否和实例输出相同。然后再检查结果是否完整，是不是漏掉了某些结果等，最后检查是不是程序逻辑错误导致的结果错误。

Runtime Error：表示代码执行的时间太长，也就是说用户代码在 Hadoop 上执行的时间超出了正常的执行时间。出现这种情况的原因主要是用户程序存在死循环或平台同时提交的程序太多，使运行效率降低了。用户只需要查看是否存在死循环代码并在平台空闲的时间提交就可以了。

Memory Exceed：表示程序运行时内存溢出，即用户代码中过多使用内存或无限申请内存的代码（这主要针对主函数中的代码，如果在 MapReduce 中出现类似的代码会返回 MapReduce Error）。出现这种情况，就需要用户在自己的代码中仔细查找是否有过多使用内存或无限申请内存的代码。

Evil Code：表示提交的程序中存在恶意代码，也就是说用户代码中存在系统调用代码或意图更改平台服务器配置的代码等。这就需要用户清除代码中根本用不到的代码和一些恶意代码了。

以上介绍了平台运行用户提交的代码之后所返回的各种结果及其出现的原因和应对策略。错误的根本原因是代码问题，所以用户遇到问题需要耐心审视自己代码，修改其中不正确的代码和逻辑，删除无用代码。

A.4.3 使用注意事项

这一节主要向读者介绍平台使用的一些注意事项，这部分内容也可以参考平台 FAQs 中的内容。

- ❑ Java 程序主类的名字必须为 MyMapre（否则编译错误）。存在这个限制的原因是需要统一所有提交的代码，然后由 Shell 文档再将其提交到 Hadoop 上运行，所以不能为

每个用户的代码写专门的 Shell 文档。

- ❑ 在配置 MapReduce 程序的输入输出时必须使用下面这两个语句（原因和前一个注意事项相同）：

旧 API

```
FileInputFormat.setInputPaths(conf, new Path(args[0]));
FileOutputFormat.setOutputPath(conf, new Path(args[1]));
```

新 API

```
FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
```

- ❑ MapReduce 程序必须处于一个 Java 源文件内，它不支持引用其他文件的类。也就是说必须把 Map、Reduce、Combine 等类写到一个文件内。
- ❑ 平台对同时运行的 MapReduce 程序数量有限制。因为系统资源有限，但 Hadoop 平台及 MapReduce 程序在处理少量数据时的表现并不是很好（即使运行少量数据，WordCount 程序也需要花费 20 多秒的时间），所以需要用户耐心等待提交程序的检测结果，而且不要同时提交多个程序，以免占用过多的平台资源。

A.5 小结

本附录主要介绍了云计算在线检测平台。平台以 Hadoop 集群作为并程序的运行环境，为 MapReduce 的入门者提供了兼顾实战和理论的训练，使其初步掌握 MapReduce 框架和 Hadoop 系统的理论知识，同时具有使用 MapReduce 并行化解决实际问题的能力。

在附录的第 2 节中介绍了平台的各个组成部分及其功能。平台主要包括前台用户接口和后台程序运行。前台主要包括用户完全服务、实例编程练习、分布式系统理论知识测试、帮助功能。前台主要完成与用户的交互和用户服务的功能。后台主要包括 Tomcat 服务器、MySQL 数据库、Hadoop 分布式环境、shell 文档，它为前台功能提供支持。接着又介绍了用户代码的检测流程，主要是用户提交之后网页保存用户代码、启动 Shell 调用用户提交的代码进行代码预处理、预处理成功后代码会提交到 Hadoop 上运行，然后分析并返回用户程序执行的结果、最后将用户的结果信息显示在前台界面上。最后一节对网站的使用进行了介绍，主要是一些功能使用的举例，比如注册更新信息、提交代码、理论测试等。同时本节还介绍了用户代码运行之后返回的各个结果所表示的意思，以及原因和如何应对。

云计算在线检测平台能够帮助用户掌握 MapReduce 编程框架和 Hadoop 分布式系统的理论知识，并且在实践中让其拥有利用 MapReduce 框架解决实际问题的能力，是 MapReduce 入门者不错的选择。

[General Information]

书名=HADOOP实战

作者=陆嘉恒著

页数=441

出版社=北京市：机械工业出版社

出版日期=2011.10

SS号=12865849

DX号=000008178770

URL=<http://book1.duxiu.com/bookDetail.jsp?dxNumber=000008178770&d=B5B27685B4BC0BA810D985281BEC847A>