



结合作者的理论积累和实际项目经验，全面介绍了开源系统Storm的系统架构、通信模型、作业及编程单元和保障机制，并精解案例，旨在为大数据流式计算提供从理论到实践的指导和参考。



# Storm： 大数据流式计算 及应用实践

□ 丁维龙 赵卓峰 韩燕波 编著



中国工信出版集团



电子工业出版社

PUBLISHING HOUSE OF ELECTRONICS INDUSTRY

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

# Storm：大数据流式 计算及应用实践

丁维龙 赵卓峰 韩燕波 编著



电子工业出版社

Publishing House of Electronics Industry

北京 · BEIJING

## 内 容 简 介

本书共分为三篇，第一篇从流式计算的原理入手，论述了大数据环境下的挑战及流式计算的基本理论与技术；第二篇详细讲解了开源工具 Storm 实现的大数据流式处理的基础，包括 Storm 的系统架构、通信模型、作业单元、数据源编程单元、数据处理编程单元、功能性保障、非功能性保障、分布式远程过程调用、事务性作业、非 Java 语言的开发等；第三篇系统性地总结了 Storm 的应用实践流程，以实际案例为例，讲解了 Storm 的系统部署、开发、调试，并分析了笔者参与的一个实际项目。

本书结合理论逐步落地实践，使读者不仅能够深入地了解当前大数据带来的挑战与机遇，还可以通过书中的案例获得更直观的感性认识，快速上手 Storm 的开发，解决个性化实践处理的需求。本书编纂严谨，非常适合业界专业人士基于 Storm 进行大数据流式处理编程与实践，也适合高等院校学生以 Storm 为参照系统，自学分布式流式数据处理技术。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有，侵权必究。

### 图书在版编目（CIP）数据

Storm：大数据流式计算及应用实践 / 丁维龙, 赵卓峰, 韩燕波编著. —北京：电子工业出版社，2015.3

ISBN 978-7-121-19568-6

I. ① S… II. ① 丁… ② 赵… ③ 韩… III. ① 数据处理软件 IV. ① TP274

中国版本图书馆 CIP 数据核字（2014）第 286291 号

责任编辑：董亚峰 特约编辑：王 纲

文字编辑：吴长莘

印 刷：北京京师印务有限公司

装 订：北京京师印务有限公司

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本：787×1092 1/16 印张：16.5 字数：400 千字

版 次：2015 年 3 月第 1 版

印 次：2015 年 3 月第 1 次印刷

定 价：48.00 元

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

质量投诉请发邮件至 [zlts@phei.com.cn](mailto:zlts@phei.com.cn)，盗版侵权举报请发邮件至 [dbqq@phei.com.cn](mailto:dbqq@phei.com.cn)。

服务热线：(010) 88258888。

# 前 言

---



信息技术行业在过去的不长时期内，经历了多种概念的出现、发展和消亡。而当今，最火热的 IT 词汇，无外乎“大数据”、“云计算”和“物联网”；而相应的支撑技术，诸如 Hadoop、HDFS 和 Cassandra 等，均得到了广泛关注和研究。特别是以 Hadoop 生态系统为代表的大规模数据处理技术，已经成为业界事实上的标准之一。Hadoop/MapReduce 这类数据处理技术和所构建的分布式系统，针对数据边界清晰的数据进行批处理操作，涉及作业的启动、中间结果写盘缓存和基于共享存储的数据传输，使得数据处理结果存在较长延迟，无法提供持续处理的能力。基于 Hadoop/MapReduce 计算框架的系统，可针对海量数据实现批处理，但在高速并发的环境下，无法满足实时、连续处理的需求。于是，针对实时数据处理，缺少编程标准化、可靠和可伸缩的计算模型和框架，成为一个巨大缺失，Twitter 公司开源的 Storm，恰好填补了这个空白。

Storm 作为开源的实时数据处理系统，针对流式数据较之 Hadoop，不仅降低了数据处理的并行编程的复杂度，也提供了数据不丢失的保证和集群节点动态部署的特性。具体来说，Storm 的关键特性如下。

1. 编程语义的简化和跨编程语言的集成。通过 Storm 的核心编程抽象，如 Topology, Spout 和 Bolt，可以便捷地接入和接收数据、更新数据和追踪数据，并可以灵活地指派数据在组件之间的分片方式。特别是由于采用了 Apache Thrift 接口，Storm 仅要求组件间遵循同一套接口描述，而不绑定两端具体的实现语言。因此，Storm 有着广泛的适用场景，可在异构环境下，集成处理数据流和更新持久化数据，并行化地连续查询指定数据流。

2. 数据处理的可靠性保障。在可用性要求高的场景下，数据是不允许丢失的，要求系统必须保证所有的数据被成功地处理。Storm 通过配置消息反馈策略，可保证到达系统的数据被最终处理，而节点的故障会被及时捕获并使得业务计算自动重新分配，保证数据处理的连续性和可靠性。同时，在数据正确性要求高的场景下，为了避免数据冗余和保证数据仅被处理一次，Storm 给出了事务型处理的概念和模型，可在



用户业务逻辑不改变的情况下，通过配置编程组件实现。这类因素也是相对于早先 Yahoo 开源同类系统的 S4 最重要的优势。

3. 计算节点的可伸缩性和集群配置简便易行。Storm 集群管理只需要通过维护有限数量的配置文件完成，保证集群中节点可管可控，其节点的动态接入和作业自动重分配也是一大亮点。若当前的计算资源成为瓶颈，可以通过水平扩展节点实现数据处理的伸缩，即增加机器和提升计算的并行度。官网的实验结果表明，10 个节点的 Storm 集群下，应用的吞吐量可达每秒 1000 000 个消息，其中还包括每秒 100 多次的数据库存取调用。

正是由于上述特点，Storm 获得了业界的广泛关注与赞誉。同时，由于活跃的社区贡献，其代码日趋成熟与稳定，在开源后的两年内，已经在数十个企业中实现了商业级应用，如 Groupon、The Weather Channel、Twitter、Yahoo、淘宝和阿里巴巴等。

本书全面介绍了 Storm 的理论基础、溯源发展、核心概念和集群配置、可靠性保障关键技术、常用的并行流模型编程范式、关键数据结构和源码解析等。本书的一大特色是，书中所有实例均来自编者所在团队的实际应用，是一个在智能交通背景下的分布式实时车牌流监控系统。希望通过理论结合实践，本书能为当前火热的大数据背景下的分布式系统开发，贡献一点微不足道的力量。

本书在撰写过程中，得到了许多同仁和朋友的帮助。云计算中心的硕士生张帅、卢帅参与了本书的校稿和整理，电子工业出版社的董亚峰和吴长莘两位编辑为本书的面世也付出大量心血。在此，对他们表示衷心的感谢！

由于时间和水平有限，书中的不妥之处在所难免，衷心希望广大读者能够批评指正，以便我们再版时修订。

编 者

2014 后 3 月

# 作者简介

---



## 丁维龙

博士，2013年1月毕业于中国科学院计算技术研究所，现任教于北方工业大学，在大规模流数据集成与分析技术北京市重点实验室从事实时数据处理与分布式系统方向的研究，已在SCI检索期刊和领域知名国际会议发表多篇学术论文，主持并参与多项科研课题。中国计算机学会（CCF）、ACM（Association for Computing Machinery）会员，目前是旗舰期刊 IEEE Transaction on Service Computing、IEEE Transactions on Industrial Informatics 和计算机学报审稿人，同时担任第七届中国传感器网络学术会议（The 7th China Conference on China Wireless Sensor Networks, CCF CWSN2013）、IEEE SDPI workshop 的程序委员会成员。

## 赵卓峰

博士，2005年1月毕业于中国科学院计算技术研究所，现任北方工业大学云计算研究中心副研究员、副主任，中国计算机学会高级会员、服务计算专委会委员，IEEE/ACM会员。作为负责人，作者承担多项国家和省部级课题，目前从事云计算、物联网等环境下新型应用系统的架构设计及开发方面的研究与工程实践，在公安应急、智能交通、科技信息服务、电子政务、先进制造等应用领域完成10余项应用实践，申报专利及软件登记等知识产权40余项，向华为、东方通科技、神州数码、万方等公司输出多项技术。



## 韩燕波

博士，毕业于柏林工业大学，现任北方工业大学教授、北方工业大学云计算研究中心主任。中国计算机学会服务计算专业委员会副主任、中国计算机学会大数据专家委员会委员、中国电子学会云计算专家委员会委员、计算机学报编委。曾就职于德国国家计算机研究中心、德国弗郎霍夫软件技术研究所和美国大规模分布系统实验室等机构，归国后 2000 年被聘为中科院计算技术研究所研究员，入选中科院海外杰出人才计划（百人计划），任网络重点实验室研究员、博士生导师、中科院研究生院教授。主要研究领域包括分布式系统、互联网服务、业务流程管理和协同等，在多个领域主持完成了多项 863、973 和自然基金重点项目，发表论文 140 余篇，出版专著 4 部。申报或合作申报发明专利和软件登记 50 项，其中 5 项已向工业界转化。

# 目 录

---



## 第一篇 基础篇 流式数据处理概论

第 1 章 大数据环境下的云计算与物联网 .....	3
1.1 云计算与物联网 .....	3
1.1.1 云计算 .....	3
1.1.2 物联网 .....	6
1.2 大数据下的新挑战 .....	8
1.2.1 大数据及其特征 .....	8
1.2.2 大数据处理的技术挑战 .....	11
1.3 本章小结 .....	14
第 2 章 流式计算的理论与技术 .....	15
2.1 流式数据与流式实时计算 .....	15
2.1.1 流式数据 .....	15
2.1.2 流式实时计算 .....	18
2.2 流式数据处理的系统与应用 .....	20
2.2.1 发展与挑战 .....	20
2.2.2 Hadoop 2.0 生态圈 .....	22
2.3 Storm .....	27
2.3.1 起源与发展: Twitter 的开源与影响 .....	27
2.3.2 功能 .....	29
2.3.3 特色: 可扩展、可靠的分布式流式数据处理 .....	30
2.4 其他开源流式数据处理系统 .....	34
2.4.1 Yahoo S4 .....	34
2.4.2 Spark Streaming .....	37
2.4.3 Facebook Puma .....	41
2.5 本章小结 .....	42

第 3 章 实际案例：城市道路车辆数据的实时监控分析系统.....	43
3.1 背景与需求分析.....	43
3.1.1 背景 .....	43
3.1.2 数据处理的业务需求 .....	45
3.2 数据处理系统的架构设计与技术选型.....	46
3.2.1 架构设计 .....	46
3.2.2 技术选型 .....	48
3.3 本章小结 .....	49
<b>第二篇 系统篇 流式数据处理系统 Storm 的基础原理</b>	
第 4 章 Storm 的系统架构 .....	53
4.1 系统架构与部署模式 .....	53
4.1.1 系统架构 .....	53
4.1.2 单机/分布式部署.....	56
4.1.3 本地模式 .....	58
4.2 系统节点 .....	59
4.2.1 Zookeeper: 协调节点.....	59
4.2.2 nimbus: 主控节点 .....	63
4.2.3 supervisor: 工作节点 .....	65
4.2.4 UI: 控制台节点 .....	68
4.3 本章小结 .....	70
第 5 章 Storm 的通信模型 .....	71
5.1 Thrift: 可扩展、跨语言的通信软件框架.....	71
5.1.1 Thrift 的基础概念 .....	71
5.1.2 基于 Thrift 的数据通信 .....	74
5.2 Thrift 在 Storm 中的应用：系统节点间的通信 .....	75
5.2.1 接口的定义与实现 .....	75
5.2.2 客户端与 Storm 系统的通信 .....	82
5.3 ZeroMQ 在 Storm 中的应用：作业任务间的通信 .....	83
5.3.1 ZeroMQ: 面向分布式并发应用的高性能异步消息处理库 .....	83
5.3.2 Tuple 与 declareOutputFields( ): 数据项结构及声明 .....	86
5.4 Storm 可配置的通信机制 .....	89
5.5 本章小结 .....	90
第 6 章 Storm 的作业单元：Topology.....	91
6.1 Topology 的构成 .....	91
6.2 Stream: 组件间的数据传递 .....	93

6.2.1 概述 .....	93
6.2.2 Stream Grouping: 流组模式 .....	94
6.2.3 自定义流组 .....	101
6.3 构建 Topology .....	104
6.3.1 TopologyBuilder 与 Config .....	104
6.3.2 Topology 构建示例 .....	106
6.3.3 Topology 常见的编程模式 .....	107
6.4 本章小结 .....	109
<b>第 7 章 Storm 的数据源编程单元: Spout .....</b>	<b>110</b>
7.1 Spout 的接口与实现 .....	110
7.1.1 Spout 与接口层次 .....	110
7.1.2 ISpout 和 IComponent 接口 .....	111
7.1.3 接口的实现类及实例 .....	113
7.2 Spout 的使用模式 .....	115
7.2.1 直接连接 .....	115
7.2.2 队列连接 .....	119
7.3 Spout 与数据的可靠性 .....	121
7.3.1 可靠的 Spout 与不可靠的 Spout .....	121
7.3.2 可靠的 Spout 的数据项管理 .....	122
7.4 本章小结 .....	125
<b>第 8 章 Storm 的数据处理编程单元: Bolt .....</b>	<b>126</b>
8.1 Bolt 的接口与实现 .....	126
8.1.1 Bolt 与接口层次 .....	126
8.1.2 IBolt 和 IComponent 接口 .....	127
8.1.3 接口的实现类及实例 .....	131
8.2 Bolt 与数据的可靠性 .....	133
8.2.1 可靠的 Bolt 与不可靠的 Bolt .....	133
8.2.2 可靠的 Bolt 的数据项管理 .....	133
8.2.3 IBasicBolt 和 BaseBasicBolt .....	136
8.3 * 本章小结 .....	137
<b>第 9 章 Storm 的保障机制 .....</b>	<b>138</b>
9.1 Storm 的功能性保障: 多粒度的并行化 .....	138
9.1.1 并发模型 .....	138
9.1.2 并行度配置 .....	139
9.1.3 可插拔的自定义调度器 .....	144



9.2 Storm 的非功能性保障：多级别的可靠性 .....	149
9.2.1 不同级别的容错机制 .....	149
9.2.2 记录级容错：保障数据项不丢失 .....	151
9.2.3 记录级容错的原理：acker 任务与追踪算法 .....	157
9.3 本章小结 .....	164
<b>第 10 章 Storm 的高层使用模式 .....</b>	<b>165</b>
10.1 分布式远程过程调用 .....	165
10.1.1 概述 .....	165
10.1.2 DRPC 的构建与使用 .....	166
10.1.3 Storm 的 DRPC 原理 .....	171
10.2 事务型作业 .....	173
10.2.1 概述 .....	173
10.2.2 Transactional Topology 的构建与使用 .....	175
10.2.3 Transactional Topology 的编程接口与事务型作业的实现 .....	179
10.2.4 CoordinatedBolt 的原理 .....	181
10.3 非 Java 语言的开发 .....	182
10.3.1 支持多语言的协议 .....	182
10.3.2 Shell 组件 .....	187
10.4 本章小结 .....	189
<b>第三篇 应用篇 基于流式数据处理系统 Storm 的开发</b>	
<b>第 11 章 Storm 的系统部署 .....</b>	<b>193</b>
11.1 系统环境 .....	193
11.2 依赖程序的安装 .....	194
11.2.1 libuuid, libuuid-devel, gcc-c++, libtool .....	194
11.2.2 ZeroMQ 和 JZMQ .....	196
11.3 Storm 的安装与配置 .....	198
11.3.1 Zookeeper 的安装与配置 .....	198
11.3.2 单机模式和集群模式下 Storm 的安装、配置和启动 .....	200
11.3.3 Storm 各节点的服务启动 .....	203
11.4 Storm 集群水平扩展工作节点 .....	206
11.5 本章小结 .....	207
<b>第 12 章 Storm 应用的开发与调试 .....</b>	<b>208</b>
12.1 Eclipse 环境下的 Storm 工程 .....	208
12.1.1 Eclipse 开发环境 .....	208
12.1.2 将 Storm-starter 组织为 Eclipse 工程 .....	210

---

12.2 Storm 应用的开发、调试与部署.....	212
12.2.1 本地开发与调试.....	212
12.2.2 远程部署 .....	213
12.3 常见问题与应对技巧 .....	215
12.3.1 ZeroMQ 版本.....	215
12.3.2 Zookeeper 日志清理.....	216
12.3.3 Topology 作业的打包与远程部署 .....	216
12.4 本章小结 .....	217
<b>第 13 章 项目案例分析.....</b>	<b>218</b>
13.1 业务计算的设计.....	218
13.1.1 需求分析 .....	218
13.1.2 概要设计 .....	219
13.2 业务计算的实现.....	220
13.2.1 Topology 的构建 .....	220
13.2.2 JmsSpout 的实现 .....	222
13.2.3 三个 Bolt 的实现.....	224
13.3 本章小结 .....	229
<b>附录 .....</b>	<b>230</b>
<b>参考文献 .....</b>	<b>244</b>
<b>后记 .....</b>	<b>249</b>

## 第一篇 基础篇

---

.....

# 流式数据处理概论



# 第1章

## 大数据环境下的云计算与物联网



根据维基百科的定义，大数据是指无法在一定时间内用常规软件工具对其进行抓取、管理和处理的数据集合。数据已成为与自然资源、人力资源一样重要的战略资源，隐含巨大的价值，已引起学术界和工业界的重视。大数据的有效组织和使用，将对社会发展产生巨大推动作用，孕育着前所未有的机遇。O'Reilly 公司断言：“数据是下一个 Intel inside 未来属于将数据转换成产品的公司和人们。”

大数据的繁荣，离不开云计算和物联网技术的发展。一方面，以云计算为代表的资源虚拟化为数据处理提供了高效、健壮的基础设施，以物联网为代表的网络形态产生了规模庞大的原始数据，逐步积累了大数据；另一方面，大数据的相关研究和技术，必将影响现有的云计算和物联网的发展趋势。云计算、物联网等信息技术的发展使得物理世界、信息世界和人类社会已融合成一个三元世界（the ternary human-cyber-physical universe），大数据是形成统一的三元世界的纽带。

本章主要与大数据相关的云计算、物联网的发展和兴盛，分析大数据来源、特征和数据处理的需求。

### 1.1 云计算与物联网

#### 1.1.1 云计算

云计算（Cloud Computing）是一种基于互联网的计算方式，通过这种方式，共享的软硬件资源和信息可以按需提供给计算机和其他设备。

云计算为数据的处理提供了便利的基础设施条件，包括更便宜的分布式存储、计算设备和网络。一方面，技术和业务需求的双重推动会让越来越多的政府机构、公司企业和个人意识到数据是巨大的经济资产，将带来全新的创业方向、商业模式和投资机会；另一方

面，大数据处理的兴起也将改变云计算的发展方向，云计算正在进入以 AaaS (Analysis as a Service，分析即服务) 为主要标志的时代。

云计算是大规模计算机系统自客户-服务器的转变之后的又一种巨变：用户不再需要了解基础设施的细节，甚至不必具有相应的专业知识，而是采用一种基于互联网的新的 IT 服务和交付模式，通常涉及通过互联网来提供动态易扩展而且经常是虚拟化的资源。在软件即服务 (SaaS, Software as a Service) 的服务模式当中，用户能够访问服务软件及数据。服务提供者则维护基础设施及平台，以维持服务正常运作。这使得企业能够通过外包硬件、软件维护及支持服务给服务提供者来降低 IT 运营费用。另外，由于应用程序是集中供应的，更新可以实时发布，无须用户手动更新或安装新的软件，使得企业能够更迅速地部署应用程序，并降低管理的复杂度及维护成本，同时允许 IT 资源的迅速重新分配以适应企业需求的快速改变。

云计算通过大规模分布式资源的共享形成规模效应，集成大量的资源供多个用户使用。使用者可以请求（租借）资源，而非拥有资源，并可随时调整使用配额，多余资源可被终止收回。这降低了使用者使用的经费，可以更加高效地使用并可按需配置服务提供者的资源。

## 1. 从部署模型的角度分类

从部署模型的角度分类，云计算可以被分为如下四种类型。这个分类标准最初来源于美国国家标准技术研究所 (NIST)，现在已被广泛采纳。

### (1) 公用云 (Public Cloud)。

公用云服务通过网络及第三方服务供应商开放给客户使用。“公用”意味着服务广泛的使用范围，而不一定代表免费和无限制使用。在亚马逊、VMware 和 Ubuntu 社区提供的公用云服务中，云供应商通常会对用户实施访问控制机制。公用云作为资源管理的解决方案，不仅可以保障不同用户（不同付费水平）的服务质量 (QoS 和 SLA)，而且使得资源管理具有成本效益。弹性的云服务是公用云的一个显著特征，按需使用服务降低了用户的系统运维成本。国外著名的公用云提供商有亚马逊、IBM 和微软等，国内著名的公用云提供商有新浪、百度和阿里巴巴等。

### (2) 私有云 (Private Cloud)。

私有云与公用云的差别主要在于提供的服务的管理模式不同。私有云中服务的数据与程序皆在组织内管理，而且往往对资源限制较为宽松（如网络带宽、安全规范、法规影响等）。此外，私有云服务通过机构/公司内部的虚拟化资源，用户更细粒度、高权限地掌控基础架构、安全策略和使用模式的设计与调整。著名的私有云产品主要有开源产品 Eucalyptus 和 OpenStack，以及收费的商业产品 VMWare vCenter 和微软 Windows Server 产品线。

### (3) 社区云 (Community Cloud)。

社区云由众多利益相仿的组织掌控及使用，如特定安全要求、共同宗旨等。社区成员共同使用云数据及应用程序。社区云也可以看成公用云的一种特例，更多的是跨组织的行



业视角的云服务应用。社区云虽然不是新概念，但它的发展确实是当前最为新兴的一类。例如，国际赌博科技公司（International Game Technology）最近推出了 IGT Cloud，面向赌博公司赌场的业务管理。又如，联合健康集团（United Health Group）推出了 Optum Health Cloud（Optum 健康云），面向医疗保健行业的公司客户，旨在充分利用全美保险服务的云资源。此外，弗吉尼亚社区学院（VCCS）与企业合作推出的教育社区云，为 40 个校园、23 所高校提供了服务，达到了提高效率和改善服务的双赢效果。

#### （4）混合云（Hybrid Cloud）。

混合云是介于公用云及私有云之间的云服务使用模式，在私有云的私密性和公用云的灵活低廉性之间进行权衡。用户通常将非关键信息外包在公用云上处理，同时将企业关键服务及数据放在机构内部处理。相对而言，混合云的案例较少，但也有相关的商业公司提供了完善的解决方案。例如，Amazon 提供的 VPC（Virtual Private Cloud，虚拟私有云）将 Amazon EC2 的部分计算能力接入企业的防火墙内；VMware 提供的 vCloud，可将自动化业务连续性与支持虚拟化的安全性和合规性相结合，针对 IT 服务的访问、安置和成本提供个性化的控制力。

### 2. 从使用模式的角度分类，云计算通常存在如下三种服务模式。

从使用模式的角度分类，云计算通常存在如下三种服务模式。

#### （1）软件即服务（SaaS）。

用户使用应用程序，但并不掌控操作系统、硬件或运作的网络基础架构。这是一种基础服务模式，软件以付费或免费租赁的方式被提供。典型的操作模式是，服务提供商提供一组账号和密码，用户根据用户信息认证和按配额使用。例如，客户资源管理的服务提供商 Microsoft CRM 与 Salesforce.com 提供了相关的 SaaS 服务。

#### （2）平台即服务（PaaS）。

用户使用虚拟平台操作应用程序，掌控运作应用程序的环境，包括部分主机的操作权，但并不掌控操作系统、硬件或运作的网络基础架构。平台通常是指应用程序的管理环境，包括应用服务器和 Web 服务器等基础环境。例如，Google 公司提供的 App Engine 和新浪公司的 Sina App Engine。

#### （3）基础架构即服务（IaaS）。

用户使用一台或数台机器的基础计算资源，如处理能力、存储空间、网络和中间件，但并不掌控整个虚拟化环境。用户部署操作系统，使用存储空间，使用已部署的应用程序及防火墙和负载平衡器等，可以按自己的需求开发配置适合的应用，属于最底层的云服务。例如，Amazon AWS 和 Rackspace 就是典型的这类服务。

用户通常希望商业化的产品能够满足服务质量（QoS）的要求，并且一般情况下要提供服务水平协议。同时，开放标准对于云计算的发展是至关重要的，开源软件可对众多不同需求的客户提供快速、高效的部署和实施过程。自从 Google 公司提出了著名的三大件，即无共享架构便捷编程的 MapReduce 框架、针对数据极少修改的大文件提供持续存储的分布式文件系统 GFS，以及压缩的高可扩展性的非结构化存储 Big Table，相关开源相关