

# 基于共享式软件定义存储的 Hadoop

*结合利用 IBM Spectrum Scale 与 Hortonwork 的 Data Platform, 加速获取信息, 并提高 Hadoop 存储环境的价值*

**作者: John Webster**

**2017 年 6 月**



**Evaluator Group**

支持您作出最佳技术决策



## 数据是新资产

企业 CEO 和其他企业主管如今被不断高涨的数据浪潮所淹没，但是他们也从中发现了新的业务机会。随着消费者采用越来越多的方式来获取数据并与其他消费者保持互联，消费者本身也成为了一个数据源。每天都会有新的数据源带来“机器生成的”数据，通常指的是有线设备和无线设备生成的看不见的数据。目前，利用不断增加的数据源已经成为了 IT 部门的首要任务。随着数据分析技术的发展，企业高管意识到他们可以将企业拥有的数据转化为收益。数据是新资产。

企业的数据分析技术已经从交付周报、月报和年末报告的数据仓库发展到了能创造收入并提高运营效率和安全性的实时信息系统。大数据技术包含一系列低延迟、低成本且目前用于大规模数据分析的技术。2017 年，企业将从试用这些技术发展到大范围使用这些技术，以便在所有主要行业细分市场内支持特定的业务计划。示例包括：

### 工业 4.0

各种技术和数据源的融合使得制造业又经历了一场革命。计算能力将新的大数据分析功能与无线传感器、新的人机交互技术（触控界面和增强现实系统）、3D 打印和高级机器人整合一体。工业 4.0 计划大幅提高了制造效率、产品质量和员工安全性。

### Customer 360

现在，很多 CEO 都认为，通过深入了解客户，他们可以赢得或者至少维持竞争优势。他们希望尽可能地捕获和存储客户数据，并在客户与企业互动时，随时运用这些数据，因为他们相信这有助于他们提高客户忠诚度、满意度和收入。

### 保险科技

保险行业希望利用大数据分析改善从销售到承保的各个流程的运营。通过利用无线传感器、可穿戴设备、智能手机以及其他联网设备和可穿戴设备实时/近乎实时地获取和处理数据，保险公司能够更有效地管理风险，提高投保人的忠诚度，并优化销售机会。

### 近距离营销

零售商一直想要为实体店消费者提供类似在线购物的体验。如今，有远见的零售商正以战略方式在购物区部署无线电信标，这些信标已经成为了零售商采用的一项重要重要的辅助技术。客户通过简单、直观的智能手机应用与门店实时交互；这些应用允许客户浏览和购买商品，同时零售商也能自动交付有针对性的商业信息，并收集与店内消费者行为有关的重要数据。

## 数据货币化

CEO 意识到，他们可以出售或者“出租”其 IT 部门保管的数据。通过实施物联网 (IoT) 和其他上面提到的计划，CEO 还能够收集可为企业带来收益的数据。大型汽车制造商通过传感设备和其他车载数据源收集数据就是一个典型的例子。这些数据对于供应商合作伙伴、经销商和其他相关行业（比如轮胎制造行业）来说很有价值。

## Hadoop 与外部存储案例

如今，企业在实施这些新业务计划时，首要考虑的大数据分析平台是 Apache Hadoop。Apache Hadoop 是一个开源项目，起源于谷歌和雅虎的互联网数据中心，该项目以规模巨大、单位计算能力成本极低而著称。Hadoop 提供分布式处理能力，处理大量非结构化数据集。其本地存储环境 Hadoop Distributed File System (HDFS) 是一个基于 Java 的并行化分布式文件系统，旨在应用于目前可扩展至 200 PB 的 Hadoop 集群。HDFS 可以支持 4000 个节点的单一 Hadoop 集群。HDFS 支持多个应用和用户通过 Apache YARN (Yet Another Resource Negotiator) 同时访问数据。Hadoop 还具有容错性，这意味着，它能够承受磁盘和节点故障，同时不影响集群的可用性。用户可以根据需要替换出现故障的磁盘和节点。

在早期，用户利用 Hadoop 的 MapReduce 引擎，以批量处理模式运行分析应用。现在，数以千计的企业将大数据分析功能视为未来业务计划的关键驱动因素，比如上述业务计划，他们希望将 Hadoop 应用于各类分析应用。除了在同一集群内运行后续 MapReduce 作业外，他们还希望为不同类型的分析用户托管多个应用（参见下图 1）。其中包括 OLTP (Hbase) 和实时分析 (Storm 和 Spark)。



Hadoop 多应用处理环境 (来源: Hortonworks)

因为 HDFS 专为 Hadoop 量身定制，所以它无法为 Hadoop 用户提供现代化存储平台所提供的功能集，这些是专为在多个 IT 生产数据中心用例内存储和管理数据而提供的功能集。通过以下方式，数据中心级存储系统的属性对于 Hadoop 用户来说极具价值：

### 增强数据保护和灾难恢复功能

HDFS 依靠在摄取数据时创建的克隆数据副本（通常为三个副本），从磁盘故障、数据丢失场景和相关停机中恢复过来。尽管该流程允许集群在不停机的前提下出现磁盘故障和替换，但是它会放缓数据摄取操作，对信息获取速度产生负面影响，并且无法覆盖包含数据损坏的数据丢失场景。同时，这还会导致用户无法高效使用存储媒介。用户希望在集群内将数据保存 7 年以满足可能出现的合规性要求时，他们会非常担心这个问题。

现代化存储平台利用纠删码和外部数据保护功能（同步和异步复制）提供大规模的自动化外部数据保护，并且不需要在摄取数据时创建三个数据副本。内部和外部的自动化分层能让用户在多年的操作中高效使用存储资源。

### 支持 Hadoop 集群更高效地使用计算和存储资源

HDFS 整合计算与存储，以便将数据处理位置和数据存储位置之间的“距离”最小化，从而大规模提升性能。但是，将 HDFS 作为长期永久存储环境使用会带来一些意想不到的后果。若要以数据节点的形式添加存储容量，管理员还必须增加处理和网络资源，不论他们是否需要。计算和存储的这种耦合会限制管理员运用自动化存储分层以大规模利用固态硬盘的能力。

通过集成现代化存储平台与 Hadoop，用户能够根据需要独立扩展存储，无需过度配置计算和网络资源。通过在同一系统中使用 SSD 和 HDD 获得分层的性能，用户能够平衡存储资源，大幅减少对过度配置存储资源的需求，进而提高性能。

### 用更少的时间完成更多分析

将 Hadoop 应用于分析应用的另一个优势在于，用户能够针对大量非结构化数据运行查询。正因为此，Hadoop 也经常的定位成“大数据湖”。这里的理念是复制活动数据存储中的数据，将副本传输至数据湖。但是，这个流程相当耗时（从几小时到几天不等）且占用大量网络资源，具体取决于数据量。

解决这一问题的途径之一是利用 Hadoop 存储面向生成数据的多个应用整合存储，不再需要耗费大量时间生成和跟踪数据以及在网络中传输数据副本。支持行业标准访问协议（如 NFS、SMB、iSCSI、S3 和 Swift）的现代化存储系统能做到这一点。通过集成现代化存储系统与 Hadoop，其他业务应用用户能够即时使用 Hadoop 应用生成的数据。通过使用同时提供多个用例的多用途存储环境，用户不需要修改企业可能依靠的事务性数据架构。

### 降低复杂性

开源社区创建了一些附加项目来增加新功能并解决不足之处。在 Hadoop 内，DistCp 可用于在 WAN 距离内实现集群的周期性同步，但是随着时间的推移，集群之间会出现不一致的现象，这时需要用户手动协调差异。Falcon 支持数据生命周期和管理。但是，站在 Hadoop 管理员的角度，他们通常是将 Hadoop 作为独立实体来学习和管理 Hadoop。每个 Hadoop 都有自己的生命周期，需要单独进行跟踪、更新和管理。在这方面，企业 Hadoop 管理员会越来越重视简化，这是自然而然的事情。通过使用拥有内置功能的存储环境，比如提供低成本存储的分层功能（包括 IBM Cloud）、Snapshot 集成式数据保护和利用通用命名空间实现全球数据共享等功能，用户能简化管理，减少出错几率。

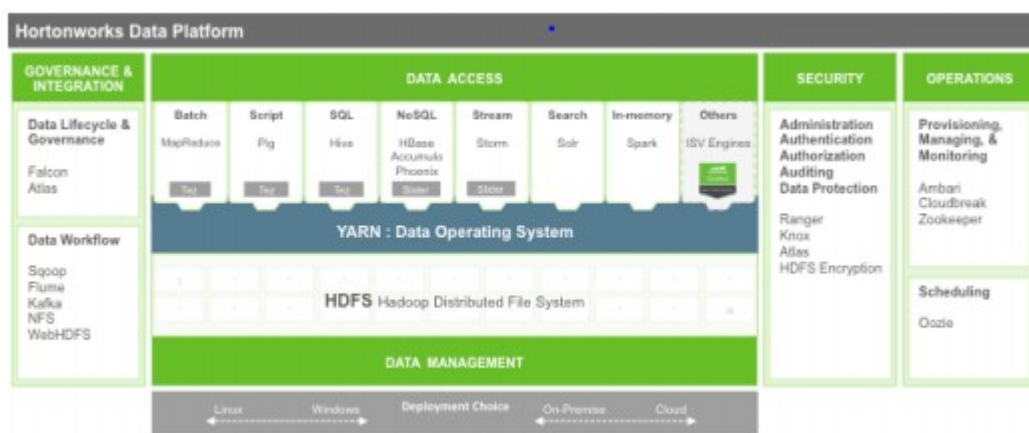


通过集成向外扩展的现代化共享存储平台与 Hadoop 集群，用户能够额外获得这些价值。在该场景下，用户无需增加、集成和管理更多项目，即可将存储平台的信息共享和数据持久性功能应用于 Hadoop，取得同样的成果。

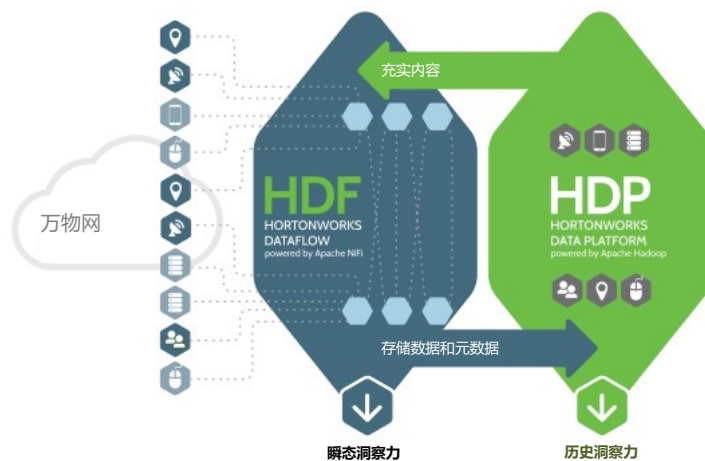
我们以 Hortonworks 发行版 Apache Hadoop 为例，并着重研究延迟敏感型用例，比如前面提到的用例。然后，我们还发现 Hortonwork 刚刚宣布支持 IBM Spectrum Scale 存储平台，以应对不断增长的面向 Hadoop 的应用场景。

## 面向当今业务计划的 Hortonworks HDP 和 HDF

Hortonworks 的 Apache Hadoop 发行版在处理大数据时，会以两种不同的方式来分析静态数据和动态数据。基于 Apache Hadoop 构建的 Hortonworks Data Platform (HDP) 提供了一组重要的 Hadoop 功能，这些功能通常能满足企业生产 IT 环境的需求，其中包括访问和管理数据、安全性、简化运营以及遵守治理实践和法规（参见下图 2）。HDP 针对的是静态数据。



Hortonworks Dataflow (HDF) 基于 Apache NiFi 技术构建。NiFi 能满足用户加速和简化系统之间数据流的需求。HDF 是一个单一整合型平台，用于数据获取、简单事件的处理、传输和交付；它非常适合用于处理和整合不断激增的实时数据源（如智能手机和传感器）生成的多样化且复杂的数据流。通过与 HDP 集成一体，HDF 能够捕获这些数据源生成的且通常为瞬态的数据，并将这些数据传输给 HDP。HDP 和 HDF 的集成（参见下图 3）将生成一个单一的集成式平台，用于数据获取、数据处理和持久性数据存储，这是实施 Customer 360、物联网和工业 4.0 等当今业务计划的必备要求。HDF 支持用户实时收集和处理的瞬态数据，而 HDP 则能处理和存储瞬态数据，并用历史数据充实瞬态数据。



HDP 与 HDFS 的结合旨在在这些新业务计划的情境中和常见的 Hadoop 批量处理型应用中加速信息和可执行洞察力的交付。因此，这是一个能利用可扩展、多用途的数据中心级存储平台来增强的多用途环境。

为实现这一目标，Hortonworks 和 IBM 联合宣布了一个支持协议，企业可利用 IBM Spectrum Scale 软件定义存储平台来增强 Hortonworks HDP/HDFS，进而提高用户生产力，通过提高存储效率来降低总体拥有成本，简化数据管理，并提高业务连续性。

## IBM Spectrum Scale for Hadoop

IBM Spectrum Scale 存储软件是一个企业级平台，用于文件和对象存储以及数据管理。该软件基于并行文件系统构建，以便在向外扩展的架构中提高可扩展的性能和容量。它为 Hadoop 的 HDFS 存储层提供透明的支持。通过与 Hortonworks HDP 相集成，IBM Spectrum Scale 存储软件作为持久性存储环境（用于生产级 Hadoop 数据存储库）的主要属性包括：

**单一全球命名空间**，通过向集群添加新节点，在单一 Spectrum Scale 环境内支持各种扩展 Hadoop 部署项目，包括小型项目和大型项目。



**统一的存储环境** — 支持文件和基于对象的数据存储。数据访问方法包括 POSIX、NFS、SMB、S3 和 Swift。

文件系统或文件集级别的**快照**，同时还支持备份到外部存储目标（备份设备和/或磁带）。

LAN、MAN 和 WAN 距离的**异步和同步数据复制**，同时保证事务一致性

**自动化云存储分层**，利用存储层之间自动且由策略驱动的数据移动，对基于云的对象存储或公有云存储进行透明的云存储分层。用户也可利用磁带作为额外的归档存储层。

支持**无中断运营**，无需系统停机，即可在全球命名空间层级优化活动文件管理和文件放置。同时，还支持用户在系统不停机的前提下进行滚动升级。

在每个文件的基础上实施**由策略驱动的数据压缩**（即，由系统管理员控制何时以何种方式压缩哪个文件），这种方式能够将存储效率提高两倍，并减少 Hadoop 集群节点的处理负载。

**基于存储的安全功能**，包括可选静态数据加密、安全擦除和用于身份验证的 LDAP/AD。还支持通过 Active Directory 或其他 LDAP 源进行身份验证和授权。

以简单的基于 GUI 的方式管理存储环境，其中包括自动化资源配置和存储系统性能监控

**IBM zSystem 集成**面向的是希望将 IBM zSystem 大型机数据与 Hadoop 集成一体的用户

Spectrum Scale 为 Hadoop 集群带来了企业数据管理收益，但是您并不需要用 Spectrum Scale 完全取代 HDFS。它能够与部署在 HDFS 上的 Hadoop 集群共存，因为它提供了一个单一命名空间，该命名空间覆盖了 HDFS 和 Spectrum Scale 管理的所有存储。

---

### *Evaluator Group 评价:*

*如今在很多细分行业中，越来越多的数字企业计划采用了 Hadoop，这推动了 Hadoop 在企业内知名度的提升。企业内有一些群体希望开展新型数据分析，他们听说竞争对手已经在 Hadoop 上开展此类分析。这种情况下，这些群体的压力越来越大。*

*Hadoop 存储环境是企业实施这些计划时的一个关键考量因素，尽管 Hadoop 用户和管理员通常并不这样认为。我们认为，IBM 的 Spectrum Scale 在以下方面对于整个 Hadoop 存储环境来说价值巨大：能够整合来自不同数据源的数据；为 Hadoop 和其他基于 Hadoop 的分析应用相关的系统提供一个连续的数据层。利用事务性一致的异步和同步复制功能支持的数据保护和业务连续性功能可帮助 Hadoop 满足企业生产 IT 数据中心部署项目的要求。因此，我们看到 IBM 与 Hortonworks 建立了合作，旨在完成 Hortonworks 在 IBM Spectrum Scale 上的认证，而这一目标也将会在未来实现。*

*企业 IT 人员必须能够从容地将 Hadoop 当作面向生产分析应用的平台加以管理。与 HDFS 存储环境相集成的企业数据中心级存储系统能够满足 IT 人员的很多要求。我们已经介绍过，IBM 的 Spectrum Scale 软件能够以企业 IT 人员熟悉的方式简化、保护和管理 Hadoop 数据资源，同时遵守企业现有的数据管理策略和实践。*

---

## 关于 Evaluator Group

*Evaluator Group Inc. 致力于帮助 IT 专业人士和供应商制定并实施战略，使他们能够从存储系统和数字化信息中获取最大价值。Evaluator Group 的服务能为 IT 专业人士深入、公正地分析存储体系结构、基础架构和管理做法。自 1997 年以来，Evaluator Group 已向成千上万的最终用户和供应商专业人员提供了产品和市场评估、竞争力分析以及教育方面的服务。[www.evaluatorgroup.com](http://www.evaluatorgroup.com) 请关注我们的 Twitter 帐号 @evaluator\_group*

### **Evaluator Group, Inc. 版权所有。保留所有权利。**

*未经 Evaluator Group Inc. 明确的书面许可，无论出于何种目的，均不得以任何形式或者通过任何电子或机械方式（包括影印和录制）复制或传播本出版物的任何部分，或者将其存储于数据库或检索系统中。本文档中包含的信息如有更改，恕不另行通知。Evaluator Group 对错误或遗漏不承担任何责任。在本文档中，Evaluator Group 对文中所述产品的使用或操作不作任何明示或默示的保证。Evaluator Group 对因本出版物的任何方面引起或与之相关的任何间接、特殊、非继发性或附带损害概不负责，即使已被告知可能会发生此类损害也是如此。Evaluator Series 是 Evaluator Group Inc. 的商标。所有其他商标均为它们各自的公司所有。*