

大数据技术栈

数据存储

HDFS
Hadoop Distributed File System
大量数据可以存在于数台计算机中，管理时可见的只有一个文件系统，而数据可以分布存储

计算能力

MapReduce
多个计算机协同处理任务的计算引擎（第一代）
计算过程简化到只有 Map 和 Reduce 两个中间用 Shuffle 进行串联
Map 将计算机中的数据取出，生成信息集合
Reduce 将信息集合进行计算处理再进行汇总

Spark
多个计算机协同处理任务的计算框架（第二代）
特点
降低 Map 和 Reduce 的区分界限，来达到更高的数据交换效率，使复杂算法能有更高的吞吐量
通过 Driver 把用户打包提交的 Spark 程序序列号后，发送计算任务到集群的工作节点进行并行计算
与 Hadoop Yarn 对标的 Cluster Manager 集群管理工具
集群部署方式
Standalone
Cluster 模式
Driver 运行在工作节点的进程中，完成任务后客户端进程立刻退出
Client 模式
Driver 一直运行在客户端进程中，并一直向 Console 输出运行信息
Yarn
Cluster 模式
Driver 直接运行在 Application Master 上 Yarn 直接管理，程序初始化后客户端进程退出
Client 模式
Driver 始终运行在客户端进程中与 Yarn 的 Application Master 通信获取节点资源
Mesos

高级抽象

Pig
模仿脚本方式描述 MapReduce 算法过程
Yahoo 开发，使用 Pig Latin 描述数据流
Hive
模仿 SQL 方式描述 MapReduce 算法过程
Facebook 主导，将结构化数据映射为一张数据库表可以将 sql 语句直接转为 MapReduce 任务
计算操作
代表每个计算阶段的数据集合
RDD 创建来源
Driver 推送
HDFS 推送
对标 Hadoop 的 Map
Transformations
不触发计算只产生 RDD 的中间结果
产生计算之前只在内存中生成 DAG（有向无环图）
Actions
对标 Hadoop 的 Reduce

交互引擎

计算引擎的计算速度难以接受，牺牲稳定性和通用性提升速度
交互 SQL 引擎直接读写 HDFS 越过计算引擎
Impala Presto Drill
Impal 大规模并行处理（MPP）式 SQL 引擎

流计算

在数据未存入 HDFS 仍在“流动”时就开始计算，达到实时计算
Spark Streaming
高吞吐、高容错的流式处理系统
Storm 是最流行的流计算平台
Twitter 开发，分布式容错实时计算平台

集中管理

Hadoop Yarn 提供任务调度和集群管理

独立组件

KV Store
根据“关键字”快速得到相关信息
MapReduce 需要扫描全部数据集合
用 KV Store 实现不计算、不处理高速读取
Cassandra HBase MongoDB
Cassandra 适合写入场景
MongoDB 适合读取场景
HBase 适应性强，效率中庸
分布式机器学习库
Mahout
数据交换编码库
Protobuf

ZoomKeeper
高一致性分布存取协同系统

Flume
高可用、高可靠的分布式日志采集聚合传输系统

Flume Agent Source 会处理这些数据事件并存放在 Flume Agent Channel
收集到的数据封装成事件发送给 Flume Agent 的 Source
Spooling Directory Source
Exec Source
logstash
可以监听追加内容
Flume Agent Sink 会从 Channel 中收集数据准备下一步的分发

Kafka
高吞吐、高容错的分布式发布订阅消息系统

业界实现
更类似于 Redis 的中间件系统，作为消息队列
数据->flume->kafka->HDFS
数据->flume->kafka->storm
数据->flume->kafka->spark streaming

sqoop

原理
用于在外部结构化数据处理于 Hadoop 导入导出数据的工具
在 RDBMS 和 HDFS 之间的数据的处理
数据导入
在内部将原数据集划分成不同的区块使用只有 map 的 MapReduce 完成数据传输每个 map 都对应一个区块写入 HDFS 中
数据导出
导出前提前创建表，基于目标表的元数据生成一个类从文本中读取数据并插入对应表
启动一个 MapReduce 任务从 HDFS 中读取文件用生成的类解析成记录执行导出

UUID

分布式环境下唯一元素识别码
由日期和时间、时钟序列、机器识别码（一般为网卡 MAC 地址）组成

ELK

ElasticSearch
可扩展的全文搜索分析引擎
基于 Apache Lucene
通过 RESTful API 隐藏 Lucene 的复杂
Logstash
Kibana
与传统数据库的类比
数据库 -> 表 -> 行 -> 列
索引 -> 类型 -> 文档 -> 字段