

文章编号:1673-8411(2019)01-0065-04

Hadoop 在气象数据存储中的应用

任晓炜, 黄志, 詹利群

(广西区气象信息中心, 广西 南宁 530022)

摘要: 利用虚拟资源池搭建基于 Hadoop 的大数据存储架构, 将海量自动站文本数据、数字化历史图片以及二进制雷达基数据按照自定义 ETL 存储规则进行数据清洗之后存入大数据框架, 在并发读取效率测试中取得了良好的效果, 为应对海量气象资料增长在扩展性和系统性能方面提出的挑战提供解决思路和基本模型。

关键词: 大数据; Hadoop; ETL; 存储

中图分类号: P49

文献标识码: A

The application of Hadoop in the storage of meteorological data

Ren Xiaowei, Huang Zhi, Zhan Lique

(Guangxi Regional Meteorological Information Center, Nanning Guangxi 530022)

Abstract: A large data storage architecture based on Hadoop is constructed by using virtual resource pool. The massive text data, digitized historical pictures and binary radar-based data are cleaned according to the customized ETL storage rules and stored in the large data frame. Good results have been achieved in the concurrent reading efficiency test, which provides solutions and basic models to meet the challenges posed by the growth of mass meteorological data in terms of scalability and system performance.

Keywords: big data; Hadoop; ETL; storage

引言

随着气象现代化建设不断发展, 各类气象观测数据的数据量每年以 TB 级增加。本部门从 2009 至今开展的气象历史资料数字化工作一直有条不紊地进行, 陆续完成了地面、高空等各类历史观测报表数字化扫描归档工作, 积累了大量的数字图片, 目前文件个数已达 4-500 万张, 数据容量已达 3T; 全区 2500 多个自动站从建站至今的历史数据文件个数约为 3.5 亿个, 数据量约 1.5T; 全区 9 部多普勒天气雷达基数据从建站至今数据量约为 3T。以上数据目前存储在不同的服务器磁盘或其他外置存储介质中, 且每年不同程度地剧增。气象数据归档具有数据量大、小文件多、数据处理复杂、实时吞吐量大的特点^[1], 其零散

低效率的存储方式已不能满足部门内部对其数据进行处理和获取的时效需求。

以上提及的 3 类数据分别为 jpg、txt 以及压缩二进制文件格式, 内容较为复杂, 数量巨大, 且数据类型多样化, 呈现出比较明显的大数据特征。传统的对大规模数据处理大多使用分布式的高性能计算、网格计算等技术, 需要耗费昂贵的计算资源, 而且对于如何把大规模数据有效分割和计算任务的合理分配都需要繁琐的编程才能实现, 而 Hadoop 分布式技术的发展正好可以解决以上的问题^[2]。目前, 以 Hadoop 为代表的基于云计算环境的水平扩展分布式架构成为大数据存储服务应用未来发展趋势和热点之一, 本文通过对 Hadoop 大数据核心技术框架的研究和分析, 结合日益增长的气象业务数据的特点, 设计并搭建基

收稿日期: 2019-01-10

基金项目: 国家档案局科技项目《Hadoop 技术在气象数字档案馆存储中的研究与应用》(编号: 2016-X-06) 资助。

作者简介: 任晓炜 (1965-), 女, 黑龙江哈尔滨人, 大学本科, 高级工程师, 从事气象信息系统设计与管理工作的。

制文件, 频次为每6min一次体扫, 每部雷达一天约240个文件, 日数据量约100多M。考虑到数据实际情况, 可以将每部雷达每日的数据进行压缩合并为一个日数据文件约30-150M, 全区共9部雷达的数据需要处理。

由于FastDFS 文件管理系统相比较于HDFS 在于其存储块更灵活, 特别适合雷达基数据存储需求。使用FastDFS 的nginx 模块可以将文件的下载请求直接转交指定nginx 处理, 能有效的降低Portal 网站的并发压力, 提高客户端的响应速度, 并辅以MySQL 数据库作为雷达基数据的索引信息的存储介质。

雷达基数据FastDFS 储存设计可在每台DataNode 上(hadoop01-hadoop06) 的Volume 里只建立一个Storage Server。一共建立6 个volume。处理步骤如下:

步骤 1, 在hadoop01-hadoop06 服务器完成安装部署FastDFS, 将hadoop01 设置为tracker 存储跟踪器提供目录检索和管理功能, 将hadoop02-06 设置为存储服务器storage server 用于数据存储。

步骤 2, 将storage Server 的存储路径统一配置为/***/fastdfs/storage, 保证其路径与hadoop 系统的hdfs 路径同为根目录, 严禁配置成父子目录。

步骤 3, 搭建nginx fastdfs-nginx-module 提供fastdfs 的WEB 前端访问, 主要用nginx 直接下载fastdfs 文件, 本文采用FastDfs 上API 接口进行下载处理。

步骤 4, 在MySQL 数据库中按月建立索引表T_File_yyyyMM, yyyyMM 表示年月。

步骤 5, 系统将本地雷达文件按日合并为压缩文件逐个上传到FastDFS 文件存储系统, 假设本地待上传的物理文件为F1。

步骤 6, 接收FastDFS 返回的服务器存储路径信息 P1。

步骤 7, 接将F1 按规则拆分识别出相应的业务属性之后连同P1 插入到T_Fi-le_yyyyMM 表中, 并根据站点名字段建立查询索引。如: <RADA_ST_DOR_L2_CAP_Z9771_20131201.tar.gz> 对应存储的索引表应该是T_FILE_201312, 建立存储索引数据表之后的如下图 4 所示, 雷达基数据在服务器上fastdfs 物理存储目录效果如图 5 所示。

2.4 气象数字化历史图片数据存储设计

数字历史图片文件个数为百万数量级, 每个扫描件约 700K, 数字历史图片文件个数不多且文件大小适中, 同样适合存储到FastDfs 上, 处理步骤与上述雷达基数据大致相同, 同样辅

以MySQL 数据库作为其索引信息的存储介质。唯一区别在于对FastDfs 的文件索引保存到不同的表上, 数字历史图片按年建立索引表, 在MySQL 上数字图片数据的存储索引表名为T_PIC_yyyy (yyyy 表示年份), 数据库表结构与雷达基数据索引表基本相同。

3 系统测试

测试工具是采用美国Mercury 公司产品LoadRunner8.1, 它可以模拟上几百上千万用户实施并发访问测试。本次测试的主要是使用虚拟用户来模拟实际用户对读取 2M 以上文件(以雷达基数据作为参考)进行读取效率测试, 测试由于受制部署环境和软件版权, 采用单机模拟测试, 测试结果如下图 6 所示。图中“度量”说明为: Action_Transaction: 用户从登陆到退出系统的时间跨度; vuser_end_Transaction: 用户退出界面用时; vuser_init_Transaction: 输入账号/密码登陆系统用时; 读取: 雷达页面选择查询条件用时; 雷达查询和读取: 点击查询到结果输出用时。

图 6 中显示负载测试的每秒期间执行Vuser (虚拟用户)脚本的Vusers 的数目及其状态和每秒内执行事务所需的平均时间, 在半分钟时集合所有用户同时并发, 服务器的负荷随着用户的增加而增加, 响应的时间也会增加。并发 100 用户, 事物“雷达查询和读取”平均响应时间最大值 5.44s, 最小值 1.045s, 平均值 3.239s, 本次测试文件大小 7M, 文件越大读取所花的时间越多,

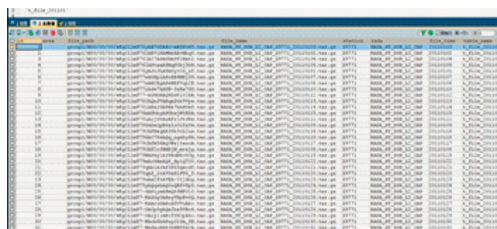


图 4 雷达基数据存储索引数据表

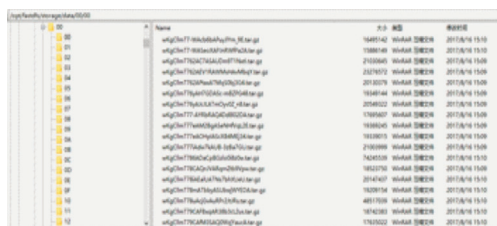


图 5 雷达基数据在 FastDFS 中的物理存储目录效果图

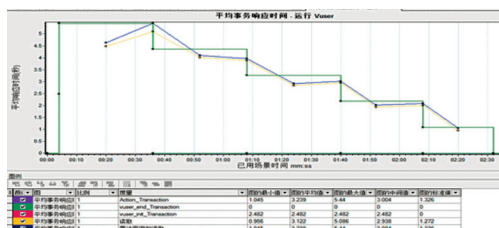


图6 LoadRunner 模拟测试结果

按正式环境7台分布式服务器同时运作, 响应时间会更快, 综合分析雷达基数据读取2M以上文件用时约1-2s, 读取效率较以往服务器磁盘等方式有明显提升。

4 结论

本文搭建以Hadoop为代表的基于虚拟池的水平扩展分布式架构开展3类代表性气象数据的应用, 结合系统的总体测试的效果, 认为基于Hadoop大数据框架下访问大文件的分布式存储系统, 能够有效解决了海量数据文件的存储问题, 并提供了高并发、高扩展性、高效率的解决方案; 为后续气象数据的共享和服务奠定基础, 为应对海量气象资料增长在扩展性和系统性能方面提出的挑战提供解决思路和基本模型。

参考文献:

- [1] 张金标, 张恩红. 基于多线程流水线的光盘自动刻录技术研究 [J]. 气象研究与应用, 2018, 39(2): 94-97.
- [2] 崔杰, 李陶深, 兰红星. 基于Hadoop的海量数据存储平台设计与开发 [J]. 计算机研究与发展, 2012, (S1): 12-18.
- [3] 王军, 刘文化, 于伟东, 等. 一种基于Hadoop的纺织海量生产数据存储设计 [J]. 微型电脑应用, 2013, (6): 53-57.
- [4] 王海荣, 刘珂. 基于Hadoop的海量数据存储系统设计 [J]. 科技通报, 2014, 30(9): 127-130.
- [5] 薛胜军, 剑寅. 基于Hadoop的气象信息数据仓库建立与测试 [J]. 计算机测量与控制, 2012, 20(4): 926-932.
- [6] 尹颖, 林庆, 林涵阳. HDFS中高效存储小文件的方法 [J]. 计算机工程与设计, 2015, 36(2): 406-409.