

# 测震波形数据存储和管理系统设计与实现

郭凯 彭克银 雷蕾  
(中国地震台网中心, 北京 100045)

**摘要:** 为了解决日益增长的海量测震波形数据在存储和管理方面存在的性能瓶颈问题, 基于测震波形数据管理的业务需求, 选取基于Hadoop大数据技术的分布式文件系统HDFS和分布式计算Spark架构进行数据的存储和计算, 并开发基于Web的测震波形数据存储和管理系统, 实现对海量测震波形数据的可视化管理和数据运行率检索。

**关键词:** 大数据; Hadoop; 测震数据; 分布式存储; 管理系统

中图分类号: P315.73

文献标识码: A

DOI: 10.3772/j.issn.1674-1544.2017.06.011

## Design and Implementation of the Seismic Waveform Data Management System

GUO Kai, PENG Keyin, LEI Lei

(China Earthquake Networks Center, Beijing 100045)

**Abstract:** In order to solve the performance bottleneck of seismic waveform data of the growing mass in the storage and management, based on the seismic business requirements, use the distributed file system of HDFS and distributed computing structure of Spark which based on the Hadoop big data technology to realize data storage and computation, and development of the seismic waveform data storage and management system based on Web, the realization of visualization management with massive seismic waveform data and query for data run rate.

**Keywords:** big data, Hadoop, seismic data, distributed storage, management system

### 1 引言

从1996年起, 中国地震局进行了大规模数字地震观测系统建设, 已建成由国家地震台网、区域地震台网和流动地震台网组成的数字地震观测系统<sup>[1]</sup>。中国地震台网中心承担着全国地震监测、地震中短期预测和地震速报以及各类地震监测数据的汇集、处理与服务的重要任务。仅测震波形数据, 现在每年汇集的数据量已经达到12TB, 目前汇集和管理的测震波形数据已经超过了100TB。

如此海量的数据, 受制于传统硬件IO、网络以及计算性能的限制, 如对一年的测震波形数据进行迁移和备份, 一般需要超过半个月的时间; 针对数据运行率计算, 主要作为评估数据质量的指标之一, 需要对每一条记录的波形数据进行读取、解压、计算完整性等一系列计算, 仅5年的计算结果就超过2亿条记录, 已经超过了传统的MySQL数据库的承受范围, 无法使用传统的数据库进行管理, 时间和效率上已经难以满足现在地震科学数据管理和科研的需求<sup>[2]</sup>。因此, 做好测震波形数据的汇集和管理就显得非常重要。作为

**作者简介:** 郭凯(1986—), 男, 硕士, 中国地震台网中心工程师, 主要研究方向: 测震数据处理与共享、大数据技术应用(通讯作者); 彭克银(1964—), 男, 博士, 中国地震台网中心数据服务部主任, 研究员, 主要研究方向: 科学数据管理; 雷蕾(1983—), 女, 中国地震台网中心助理工程师, 主要研究方向: 地震信息公共服务。

**基金项目:** 科技基础性工作专项重点项目“科技基础性工作数据资料集成与规范化整编”(2013FY110900); 地震行业科研专项“中国全球地震台网建设预研”(201508007)。

**收稿时间:** 2017年7月14日。

目前应用最广泛的开源大数据技术，Hadoop采用合理的框架分布存储和管理数据，将数据和计算资源进行有效的分配，基于集群的性能和数量，IO和计算能力可以线性增长。在其他领域处理海量数据上如医学基因<sup>[3]</sup>、电子商务<sup>[4-5]</sup>、气象<sup>[6]</sup>等领域采用基于Hadoop的技术架构均有较好的表现。因此，本文采用基于Hadoop的技术架构作为研发系统采用的主要技术，实现对测震波形数据的汇集、存储以及数据运行率计算。

## 2 测震波形数据的传输与汇集

经过多年的努力，中国地震局已经初步建成了地学的多学科、多门类的基础数据观测网，包括测震、电磁、形变、流体等四大观测台网和测线累计长度达15万公里的流动物理场观测系统<sup>[7]</sup>。按照地震科学数据分类，测震波形数据属于大类地震观测数据（D10000）下面中类的测震数据（D11000）。测震波形数据对于从事地震学研究的科研人员有着非常重要的意义，无论是进行区域速度结构、层析成像，还是大震后对地震序列特征等方面开展研究，都需要大量的地震波形数据的支持。

从2000年年底开始，中国地震台网中心接入全国48个测震台站，到2015年12月，中国地震台网中心已经实时接入了包括国家台网和区域台网共1029个测震台站。中国地震台网中心汇集的测震波形数据采用国际标准的Miniseed格式，以每个测震台站24小时的数据作为一个文件进行存储。测震台站采用了高精度的地震仪进行数据采集，一般采样率在100Hz左右，基于所在地的情况选择有线如SDH或者3G、4G网等方

式将数据汇集到台站所在的省级台网中心，然后省级台网中心采用流服务器的方式将测震波形数据汇集到中国地震台网中心，如图1所示。

测震波形数据最初采用磁带、光盘、硬盘等方式进行存储，可是存在数据量的急剧增长，使用磁带、光盘存储的数据随着时间的推移很难恢复，目前测震所有的数据又都汇集在NAS存储上，且受限于网络带宽以及NAS机头数量，传输速度难以超过100 M/s等问题。在面对TB级规模的测震波形数据，单纯采用文件方式存储，在数据汇集的速度、稳定性和安全性方面已经无法满足要求。

## 3 基于Hadoop的大数据体系结构及架构选择

Google公司于2003年提出了Google文件系统（GFS）的文件存储方法。当面对TB级或者GB级数据规模时，将所有的文件划分为若干块存储，每块大小64M，每个数据块在3个数据块服务器上冗余来保证数据的可靠性。2006年，Google公司又提出了面向结构化大数据的存储模型Bigtable。这是一个为管理大规模结构化数据而设计的分布式存储系统，可以扩展到PB级数据和上千台服务器。Hadoop开源实现了Google上述技术，具体实现了对应的分布式文件系统HDFS和分布式数据库Hbase<sup>[8]</sup>。

刘坚等<sup>[9]</sup>提出一种基于Hbase的地震大数据存储方法，通过搭建测试平台、Java语言开发测试程序，结果表明，Hbase存取地震数据耗时更低，在数据量较多时，其性能更加显著；李永红等<sup>[10]</sup>基于云计算环境下海量地震数据存储业务需

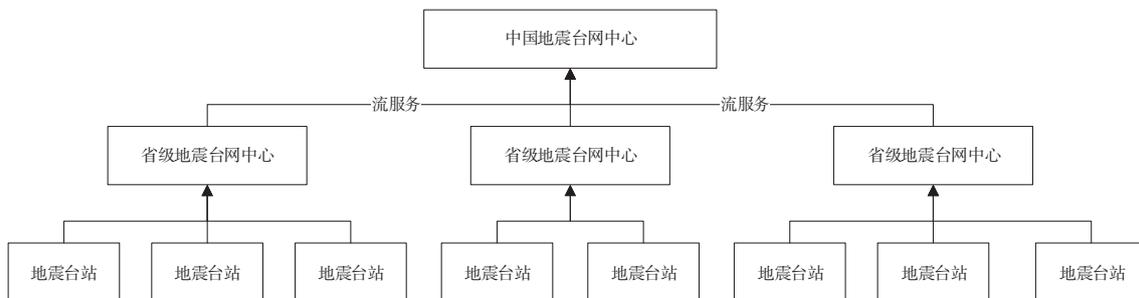


图1 测震波形数据的传输与汇集

求的分析, 采用NoSql数据库替代传统的Mysql进行数据的存储, 采用MapReduce进行计算; 王丹宁等<sup>[11]</sup>将测震波形数据解压后放入HBase进行测震波形数据的存储。从本文前述对地震业务的分析数据来看, 由于测震波形数据采用基于文件的方式进行存储, 且数据的完整性校验对于系统的吞吐量要求较高, 对每一个数据文件都要依次并发进行检索, 因此本文采用基于HDFS的方式对测震波形数据进行存储。目前开源的分布式计算架构主要为MapReduce和Spark, 由于MapReduce在Map和Reduce阶段产生的数据存储在硬盘上导致计算效率相对较低, 而Spark则将中间产生的数据放在内存中, 并采用RDD (Resilient Distributed Datasets, 弹性分布式数据集) 提高计算效率<sup>[2, 12]</sup>, 所以本文将采用Spark计算架构完成测震数据的运行率计算。

## 4 系统设计与实现

### 4.1 数据存储

测震波形数据文件按照每个台站的每个测向 (一般为BHZ、BHE、BHN) 以天为单位进行存储, 采用Miniseed格式。它由512个字节组成, 分为数据头段和数据体两部分。其中, 数据头段长度为64字节, 包含了该条波形数据的属性, 如所属的台网、台站、经纬度信息以及波形数据产生时间等; 数据体为448字节的压缩过的地震波形数据, 数据体中的数据根据Steim2算法进行

解压缩。

采用HDFS存储测震波形数据。首先将测震波形数据从NAS、硬盘、光盘等汇集到Hadoop集群的存储上, HDFS提供了许多Shell命令来实现访问文件系统的功能, 而这些命令构建在HDFS File System API之上, 通过Shell命令将测震波形数据推送到HDFS上, HDFS的Block设置为256MB。Hadoop提供了基于Web的数据显示界面, 汇集到HDFS上的测震波形数据如图2所示。

完成数据的汇集后, 需要对数据进行运行率计算, 对数据进行完整性校验, 这是评估汇集数据质量的重要参考指标之一。测震波形数据基于Hadoop技术的存储流程如图3所示。由于入库的数据量较大, 为了避免入库中断而导致重复数据集

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
drwxr-xr-x	hadoop	supergroup	9.02 MB	2017/07/25 14:27:45	3	256 MB	AH-LANG-00-BHE-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	9.16 MB	2017/07/25 14:27:48	3	256 MB	AH-LANG-00-BHN-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	8.11 MB	2017/07/25 14:27:50	3	256 MB	AH-LANG-00-BHZ-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	11.91 MB	2017/07/25 14:27:53	3	256 MB	AH-LAS-00-BHE-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	11.91 MB	2017/07/25 14:28:00	3	256 MB	AH-LAS-00-BHN-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	11.94 MB	2017/07/25 14:28:06	3	256 MB	AH-LAS-00-BHZ-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	13.98 MB	2017/07/25 14:28:13	3	256 MB	AH-LER-00-BHE-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	12.98 MB	2017/07/25 14:28:27	3	256 MB	AH-LER-00-BHN-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	12.93 MB	2017/07/25 14:28:30	3	256 MB	AH-LER-00-BHZ-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	9.62 MB	2017/07/25 14:28:44	3	256 MB	AH-LZY-00-BHE-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	9.88 MB	2017/07/25 14:28:48	3	256 MB	AH-LZY-00-BHN-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	9.72 MB	2017/07/25 14:28:51	3	256 MB	AH-LZY-00-BHZ-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	5.92 MB	2017/07/25 14:28:53	3	256 MB	AH-LHZ-00-BHE-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	6.61 MB	2017/07/25 14:28:54	3	256 MB	AH-LHZ-00-BHN-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	6.09 MB	2017/07/25 14:28:56	3	256 MB	AH-LHZ-00-BHZ-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	7.83 MB	2017/07/25 14:28:58	3	256 MB	AH-LYN-00-BHE-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	8.33 MB	2017/07/25 14:29:02	3	256 MB	AH-LYN-00-BHN-20150910000001.miniseed
drwxr-xr-x	hadoop	supergroup	7.86 MB	2017/07/25 14:29:04	3	256 MB	AH-LYN-00-BHZ-20150910000001.miniseed

图2 测震波形数据的在HDFS上的存储显示界面

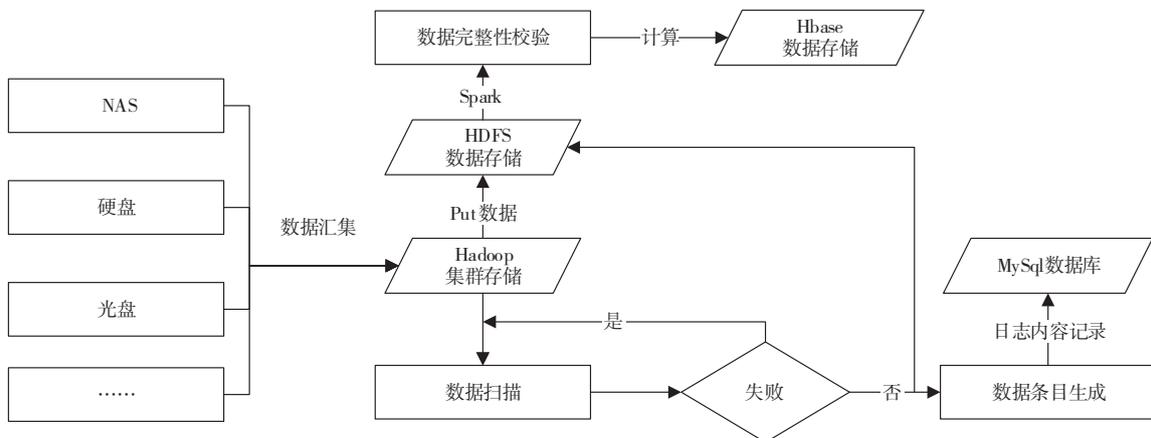


图3 测震波形数据的汇集和管理流程图

的入库情况，还需要对入库的数据提供日志功能，这里采用MySQL数据库记录已经入库的数据。

#### 4.2 数据运行率计算

基于Hadoop搭建的集群配置环境，采用Spark进行测震波形数据运行率的分布式计算。计算流程如图4所示，首先数据处理模块的主线程获取HDFS上的数据目录集合，并将目录集合以任务集的方式提交至计算节点进行计算。计算数据运行率的单位为小时，计算节点每次取出一天的数据进行计算，计算完毕后将解析获得的台网代码、台站代码、当前小时值以及24个小时的连续率一起推送到HBASE进行存储。由于该任务是基于Hadoop的分布式计算，其效率得到了较大提升。一般地，采用6台机器构成的Hadoop集群计算效率相对单机可以提升6倍<sup>[2]</sup>，并且该效率将随着集群规模的增大呈线性增长。

#### 4.3 系统功能设计和展示

在完成基于Hadoop技术的数据存储和运行率计算后，为了及时掌握测震台网的数据汇集状态，系统提供了以下两个功能：一是基于Web的展示功能，能够动态监控台站、台网的数据传输链路；二是为有效地了解全国各测震台站历史波形数据汇集和实时运行情况，能够以日、月、周、年等方式自动生成测震波形数据统计报表，包括运行率、异常台站等信息。

Web应用主要采用了以下技术：前端采用基于HTML5技术的BootStrap框架，后端采用Spring MVC框架，图表控件采用Echarts，地图控件采用百度地图以及Openlayers。如图5所示，该系统运行界面分为上下两部分进行展示：在上面部分，左侧以动态地图汇聚的形式展现全国各台站实时的运行状态，以省级台网为单位进行数据的汇集显示；右侧展现近一小时内全国各台站连续率运行率排名，排序方式为连续率由高到低进行排序。在下面部分，左侧显示系统的任务调度情况；右侧为整体台网每个小时的连续运行率情况。

系统提供了按时间范围查询统计全国所有台网、省级台网以及单个台站的运行率，基于查询目标在查询中选择统计类型为“全部台网数据查询”“台网内各台站运行率”“单一台站运行轨迹”进行统计。统计结果包括：全国台网整体运行率，省级台网的整体运行率，单个台站的运行率。其中，“单一台站运行轨迹”统计结果以阴影图的形式进行显示，显示内容为所选台站的整体运行率的变化状况，查询结果如图6、图7、图8所示。

### 5 结论

本文利用大数据技术进行了测震波形数据存储和管理系统设计与实现，取得以下主要成果。

(1) 基于Hadoop的HDFS进行了测震波形数

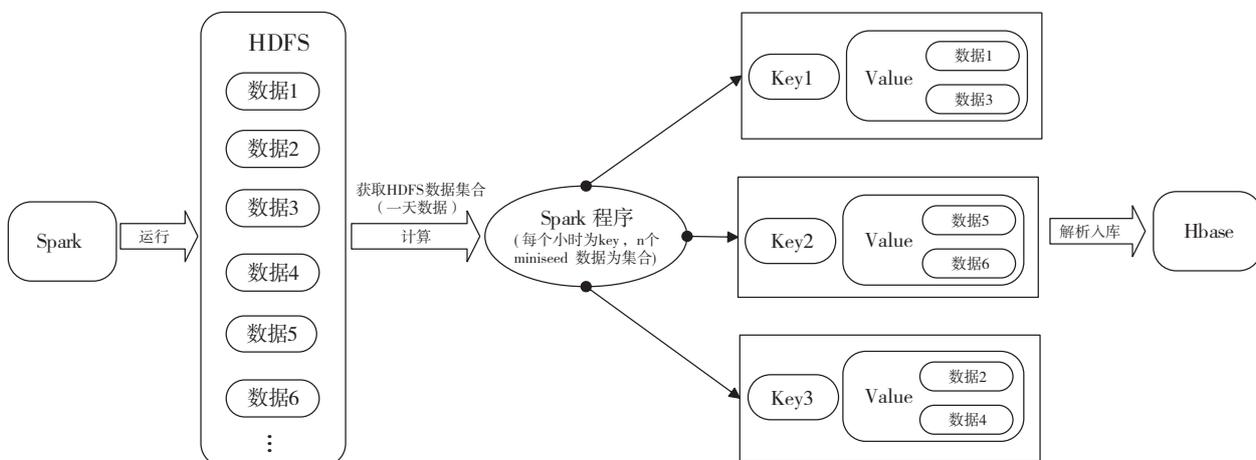


图4 测震波形数据运行率计算流程图



图5 系统运行界面



图6 系统查询省级台网



图7 省级台网所属台站

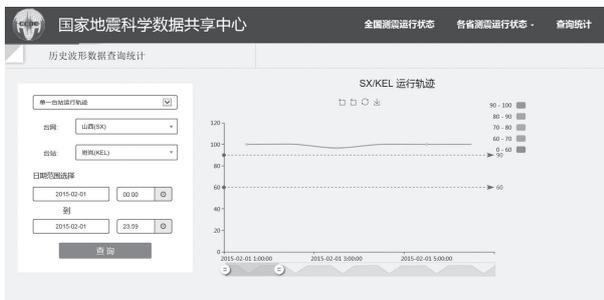


图8 单个台站的运行率结果

据的分布式存储，采用Spark技术对测震波形数据运行率进行计算，计算效率基于集群的数量可以呈线性增加。

(2) 实现了基于Web的测震波形数据汇集和管理系统，以小时为单位的颗粒度实现对全国测震台网波形数据汇集情况的可视化，实现了对省级台网、单个台站波形数据运行率的查询统计。

(3) 基于大数据技术的多副本技术，相对于将数据放入NAS可以提高数据的安全性和完整性，同时在数据处理效率上也有了较大的提升，便于对海量测震波形数据的管理。

### 参考文献

- [1] 刘瑞丰. 中国地震台网的建设与发展[J]. 地震地磁观测与研究, 2016, 37(4): 201.
- [2] 郭凯, 黄金刚, 彭克银, 等. 数据技术在海量测震数据中的研究应用[J]. 地震研究, 2017(2): 317-323.
- [3] 齐向明, 郑帅, 魏萍. 基于Hadoop的微阵列数据两阶段并行K近邻基因提取[J]. 计算机工程, 2016(5): 54-59.
- [4] 李克然. 基于云计算的电子商务数据管理模式研究[D]. 西安: 西安电子科技大学, 2011.
- [5] 陈吉荣, 乐嘉锦. 基于Hadoop生态系统的大数据解决方案综述[J]. 计算机工程与科学, 2013(10): 25-35.
- [6] 黄妍. 基于Hadoop的气象信息云存储系统设计与实现[D]. 成都: 电子科技大学, 2016.
- [7] 刘瑞丰, 蔡晋安, 彭克银, 等. 地震科学数据共享工程[J]. 地震, 2007(2): 9-16.
- [8] TOM White. Hadoop权威指南[M]. 北京: 清华大学出版社, 2014.
- [9] 刘坚, 李盛乐, 戴苗, 等. 基于Hbase的地震大数据存储研究[J]. 大地测量与地球动力学, 2015, 35(5): 890-893.
- [10] 李永红, 周娜, 赵国峰, 等. 云计算环境下地震数据管理与服务应用研究[J]. 震灾防御技术, 2015, 10(Z): 811-817.
- [11] 王丹宁, 柴旭超, 王文青. Hadoop平台下的地震波形数据存储与应用规划[J]. 软件工程, 2016, 19(1): 48-49.
- [12] 陆宏治, 邹时容. 一种基于SSD的高性能Hadoop系统的设计与应用[J]. 科技资讯, 2015, 13(29): 1-2.