

大数据技术在环境信息中的应用^①

李安增^{1,2}, 王 宁², 王常权³, 邱 燕⁴

¹(中国科学院大学, 北京 100049)

²(中国科学院 沈阳计算技术研究所, 沈阳 110168)

³(中国核示范电站有限责任公司, 荣成 264312)

⁴(德阳广播电视大学, 德阳 618000)

摘 要: 在“辽河流域水环境管理技术综合示范”项目中, 随着时间的累积, 环境监测数据处理系统采集到的数据量越来越大. 然而目前辽宁省环境监测数据处理系统无法有效处理日益增长的海量数据. 研究运用大数据技术, 改进环境监测数据处理系统中的数据中心. 利用 HDFS 强大的数据存储、管理功能, 以应对数据量的增长, 利用 MapReduce 及 Hadoop 其他相关工具, 快速处理海量数据, 降低数据规模, 最后将数据存储到数据库中.

关键词: 大数据; MapReduce; Hadoop; 并行计算; 分布式

Application of Big Data Technology to Environment Information

LI An-Zeng^{1,2}, WANG Ning², WANG Chang-Quan³, QIU Yan⁴

¹(University of Chinese Academy of Sciences, Beijing 100049, China)

²(Shenyang Institute of Computing Technology, Chinese Academy of Sciences, Shenyang 110168, China)

³(China Nuclear Demonstration Power Plant Limited Liability Company, Rongcheng 264312, China)

⁴(Deyang Radio and Television University, Deyang 618000, China)

Abstract: With the accumulation of time, data collected by environmental monitoring data processing system is more in the project of “Comprehensive Demonstration of Water Environment Management about Liao He River”. However, Liaoning province environmental monitoring data processing system cannot effectively deal with massive data growing. In this paper, we study and use big data technology. After that, we improve the data center of the environmental monitoring data processing system. We can use HDFS powerful function of data storage and management to cope with the growth of the massive data and take advantage of MapReduce and other Hadoop related tools, rapidly processing massive data and reducing the data size. Finally, the data will be stored in the database.

Key words: big data; MapReduce; Hadoop; parallel computing; distributed

近年来, 随着全球环境气候的日益严峻, 水污染事件的频发, 世界各国不断加大对环境保护的力度, 迫切需要通过环境信息化手段和环境监测能力的提高来为政府相关部门在保护环境、规划社会发展等方面提供决策所需要的信息支持. 随着科学研究、通信技术、IT 技术的迅速发展, 尤其是遥感、GIS、传感网和射频技术等现代技术的更是得到迅猛发展, 全面拓展了环境监测的时空尺度, 导致环境监测数据的种类和数量

呈现爆炸式增长, 而支撑政府相关决策科学性、准确性的基础依赖于对海量监测数据采集、传输和存储, 以及对海量数据的快速处理分析.

1 大数据技术的研究

1.1 大数据的定义

2011 年 5 月全球知名咨询公司麦肯锡发布报告《Big data: The next frontier for innovation, competition

① 基金项目: 国家水体污染控制与治理科技重大专项(2012ZX07505003)

收稿时间: 2014-04-21; 收到修改稿时间: 2014-05-19

and productivity》^[1]，报告中阐述了大数据的潜在巨大价值以及大数据技术的重要性。

目前，对大数据还没有明确的定义，但较为主流的认识是大数据应具有以下四个基本特征^[2-4]：数据规模大；数据种类多；数据要求处理速度快；数据价值密度低，但价值高。

大数据本身是一种现象而不是一种技术，但是伴随着数据的采集、传输、处理和分析的相关技术则是大数据处理技术，亦称大数据技术。

1.2 大数据处理分析流程

整个大数据的处理流程^[2]可以定义为：在合适工具的辅助下，对广泛异构的数据源进行抽取和集成，结果按照一定的标准进行统一存储，根据处理的数据类型和分析目标，采用合适的算法模型快速处理数据，并利用合适的数据分析技术对存储的数据进行分析，从中提取有益的知识并利用恰当的方式将结果展现给终端用户。图 1 所示为大数据的处理分析流程。

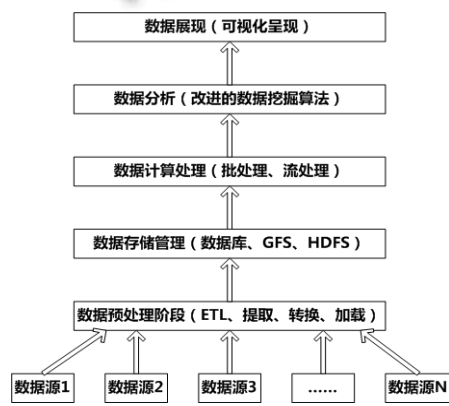


图 1 大数据处理分析流程

1.2.1 数据预处理.

在进行存储和处理之前，一般会采用 ETL 技术将数据源中的数据进行抽取(Extract)后经过数据转换(Transform)加载(Load)到数据仓库或者数据中心。

1.2.2 存储管理.

数据经过 ETL 处理后存放在数据库或者利用谷歌文件系统和 Hadoop 的分布式文件系统，这类分布式存储系统采用了分布式架构，既能达到较高的并发访问能力又具有较高的可扩展性、灵活性。

1.2.3 计算处理.

对大数据进行处理需要消耗大量的计算机资源，这对机器的运算速度和成本都提出了更高的要求，显

然采用分布式并行处理技术是最佳选择，MapReduce 框架可以用一系列廉价的机器构成，在成本、可扩展性和处理速度上有着巨大的优势。

1.2.4 数据分析.

数据分析是指从大量数据中发现规律提取新知识，它是大数据处理流程的核心，是大数据的价值体现。在计算处理后，对传统数据挖掘算法改进优化，建立数学模型，进行数据分析。

1.2.5 数据展现.

以更直观和互动的方式展示分析结果，便于人们理解。大数据分析系统应该提供数据分析、查询机制等一系列功能，并以可视化的方式呈现给最终用户。可视化技术可采用与 Web 技术相结合，以图表或图像的形式呈现。

2 环境监测数据处理

2.1 环境监测系统

目前，辽宁省环境监测系统示意图如图 2 所示。

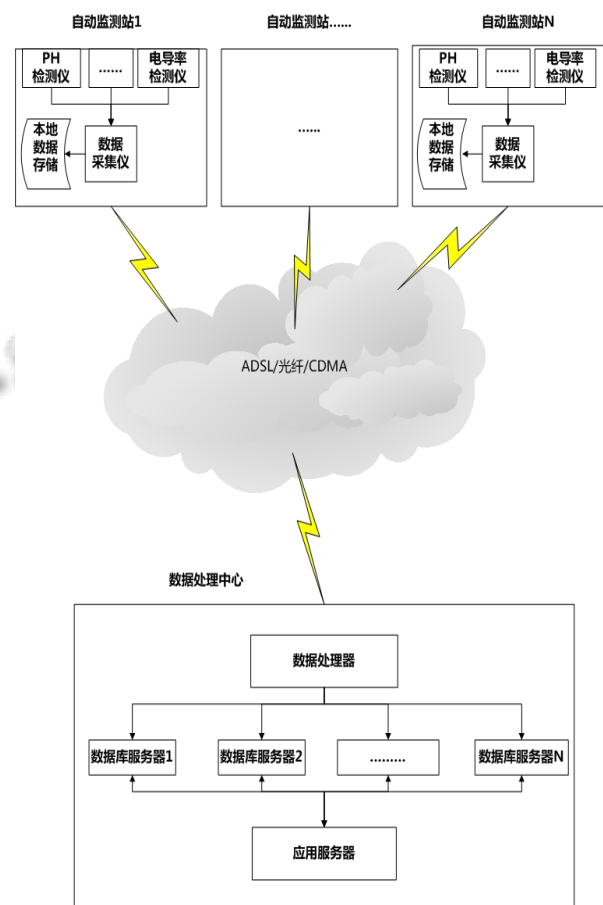


图 2 环境监测系统

自动监测站,负责采集环境监测数据,例如COD、氨氮、PH值、电导率等,采集到的数据通过网络发送到远程数据中心,另外还需要进行本地存储,防止出现网络中断,造成数据丢失的情况。如果出现网络中断,则下次网络恢复时,由数据补送机制将期间未发送给数据中心的数据从本地取出发送给数据中心。

数据处理器,进行一些简单的验证机制,例如:验证数据来源的安全性,防止非法数据来源干扰甚至破坏整个系统的稳定性和可靠性,影响数据分析结果的准确性,造成政府相关政策的误判;验证环境监测数据取值是否在合理的范围等,并将不合理的取值丢弃或者将数据进行标准化处理(根据国家相关标准规则或业务规则计算处理),然后将处理后的数据存储到分布式数据库服务器中。

应用服务器,对客户端进行请求响应,然后从数据库中获取相关数据,在计算处理后将结果返回给客户端。

据统计,辽宁省空气质量监测数据自2013年起到目前大约有1TB,现在每年以700G左右的速度增长,而且增长速度会不断增加。而水环境监测数据与空气质量监测数据相仿。随着数据规模的不断扩大,上述数据中心难以适应海量数据的处理。而大数据技术无疑是最好的选择之一。

2.2 环境监测数据处理系统设计

应用Hadoop及其他工具改进数据中心的数据获取、处理和存储过程^[5,7]。改进后的数据中心数据处理系统如图3所示。



图3 数据处理系统

改进的数据中心主要由四个模块组成:

1)数据过滤模块

自动监测站采集数据,经过网络传输存储到数据中心的数据过滤器,该数据过滤器进行数据的简单预处理,主要负责将采集到的数据进行逻辑判断,去除格式错误、不符合常理的数据,并将处理后的数据保存到HDFS中,充分利用HDFS强大的数据存储、管理功能;

2)数据转换模块

主要有数据转换器Hive和数据转换处理器Hadoop组成。Hive负责转换规则的制定与输入,根据指定的业务规则,编写SQL脚本,然后将SQL脚本翻译成复杂的MapReduce任务。利用Hadoop MapReduce^[6,8]强大的数据处理能力,对HDFS中的数据进行转换处理,生成业务所需的Hive数据表;

3)数据加载模块

如果直接从Hadoop HDFS中获取数据,过程非常缓慢,因此利用Sqoop将上述生成的Hive表加载到关系型数据库中,而从关系型数据库中获取数据的效率较高。另外,经过上述的数据转换过程,数据规模大大下降,完全可以存放到关系型数据库中。Sqoop屏蔽了底层数据库的细节,可以选择加载到不同的数据库中,例如MySQL、SQLServer、Oracle等,大大适应了不同的系统需求;

4)数据应用模块

应用服务器接受客户端请求响应,直接从数据库中获取相应请求所需要的数据,然后进行计算处理将最终的结果返回给客户端。

改进后的数据处理系统,可以根据不同的业务需要编写相应的SQL语句,对存储到HDFS的数据进行转换处理,生成相应的数据表存储到数据库中。同时本系统还可以适应不同的环境监测需求,例如水质环境监测系统,采集到的数据仍旧存放到HDFS中,只需重新定义业务规则,编写相应的SQL语句,即可重新生成该系统的数据库表,为了提高效率,最后数据加载模块可以根据不同的系统要求,将数据库表加载到不同的数据库中,以获得较好的性能。

3 模块设计

3.1 数据过滤模块

自动监测站采集数据,经网络发送到数据中心,并写入到数据中心部署的HDFS上,此处采用简单的写入TXT文件,每一行代表一个自动监测站的采集数

据, 格式为: 自动监测站编号、采样时间(yyyy-MM-dd HH:mm:ss)、监测数据值等, 各个数据之间采用”Tab”键间隔, 数据内容可以根据业务需要动态添加, 只需在数据转换模块编写 SQL 语句重新生成所需数据表即可. 假设采样频率为每 30 秒采样一次, 则一天采集的数据条数为 $24 \times 60 \times 2 = 2880$ 条, 完全可以存放在一个 TXT 文件中. 为了便于数据的管理, 各个文件采用分层管理, 具体如图 4 所示.

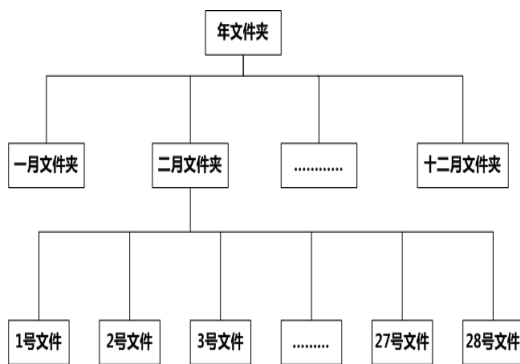


图 4 文件层次结构

在将环境监测数据写入 TXT 文本前, 需要进行简单的预处理. 例如判断监测数据值是否在合理的取值范围内, 像电导率一般情况下都为正数, 如果数据不在合理范围内, 数据直接丢弃; 将部分数据进行标准化处理(根据国家相关标准规则或业务规则计算处理)等.

3.2 数据转换模块

数据转换模块主要就是根据业务需要, 利用 MapReduce^[8]强大的并行处理能力, 对数据进行转换处理生成 Hive 表. MapReduce 的详细工作流程如图 5 所示.

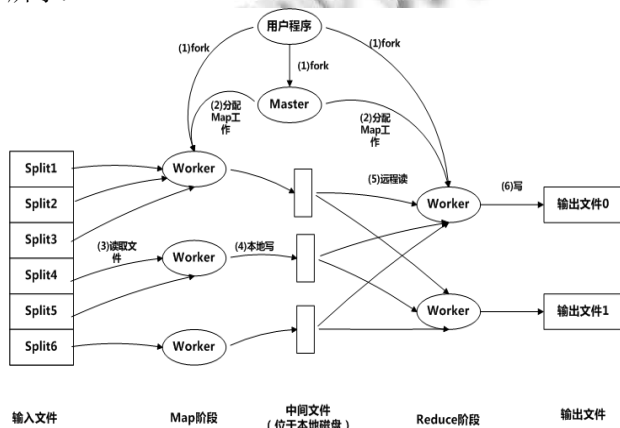


图 5 MapReduce 工作流程

①首先, 用户程序中的 MapReduce Library 会将要处理的文件切分成大小为 16M 到 64M 之间的 M 份, 一般来说, 这 M 份数据是放在了集群中的多台机器上. 然后使用 fork 在整个集群的机器上启动用户进程的多份拷贝, 如图所示, (1).

②其中一份程序拷贝是 master, 其余拷贝是 workers. master 会将任务分配给 workers. 任务包括了 M 个 map 的任务和 R 个 reduce 的任务, master 会找到空闲的 worker 然后将 map 或者 reduce 的任务分配给它. 如图所示, (2).

③当一个 worker 被分配到 map 任务, 那么它就会处理 M 份中的一份 split. 输入数据中的数据也是以很多的 Key/Value 对形式存在, Map 任务的 worker 会将这些 Key/Value 对一个一个传递给 Map 函数, Map 函数产生的很多新的 Key/Value 对会被缓存在内存中. 如图所示, (3).

④缓存中的 Key/Value 对会周期性的写入到本地磁盘中. 如图所示, (4). 由于有 R 个 reduce 任务, 所以磁盘中的 Key/Value 对会按照一定规则分为 R 个区, 每个区的数据都特定的给某一个 ReduceWorker. 在图中的每个 Map Worker 所对应的中间文件其实被分成了 R 个区. 这些区的位置会被传递到 master 上, 它负责将这些位置交给 reduce workers. 可以看到, master 在这里起到了一个“承上启下”的作用, Map 函数和 Reduce 函数之间的连接是由 master 完成的.

⑤ReduceWorker 从 master 获取了这些区的位置, 然后将这些 Key/Value 对从 mapWorker 处通过 RPC 读取. 注意, Reduce Worker 只读那些对应于自己的区的数据. 如图所示, (5). 读取完数据后需要对于这些 Key/Value 对按照 Key 进行排序, 如果需要, 可能还采取外排的方式. 之所以需要排序是因为对于一个 Key, 会有很多组不同的 Key/Value 对, 通过排序可以将它们聚合在一起.

⑥ ReduceWorker 迭代整个被排序后的 Key/Value 数据, 将每个独一无二的 key 和它所对应的 value 序列送入 Reduce 函数中. Reduce 函数对于这样的输入产生的结果往往是 0 或者 1 个值. Reduce 函数对于 Reduce 函数的结果会被添加到最终的输出文件里. 如图所示, (6).

⑦ 当所有的 map 和 reduce 任务结束后, master 会唤醒用户的程序. 这时候, 用户的 MapReduce 调用返回.

⑧最后根据业务要求,利用 Hive 工具中的有关命令创建 Hive 表并加载数据。

3.3 数据加载模块

如果直接从 Hadoop HDFS 中获取数据,过程非常缓慢,因为 Hive 结果表在被查询时会临时发起 MapReduce 任务,这对快速返回查询结果不利;如果结果表本身比较大、并且查询方式本身还有点复杂的话,那响应就更慢了,因为 Hive 对查询的优化能力跟关系型数据库管理系统没法比。所以最好把 hive 结果表再通过 Sqoop 导入到关系型数据库管理系统的一张表中。根据数据转换模块中已经生产的 Hive 表,然后利用 Sqoop 工具的相关命令完成导入关系型数据库管理系统中。下面以导入 mysql 数据库为例进行演示:

①连接 mysql 并列出数据库中的表

```
./sqoop list-tables --connect jdbc:mysql://IP 地址:端口号/数据库名称--username 用户名--password 密码
```

②将 hive 中的表数据导入到 mysql 中

```
./sqoop export --connect jdbc:mysql://IP 地址:端口号/数据库名称--username 用户名--password 密码--table 表名--export-dir/user/hive/2013/8/3/--inputfields-terminated-by '\t'
```

4 结语

本文设计改进的数据处理系统实现了面对海量的环境监测数据快速处理的需求,有效的为业务系统提

供了快速数据处理能力,同时在面对类似系统时无需大规模改动即可适应新系统要求,扩展性、灵活性高。改进后的系统主要利用 Hadoop 及其他工具对海量数据处理的强大能力,配合业务需要,快速处理海量数据,灵活生成业务所需要的数据。

参考文献

- 1 Manyika J, Chui M, Brown B, Bughin J, Dobbs R, Roxburgh C, Byers AH. Big data: The next frontier for innovation, competition, and productivity[Technical Report]. McKinsey Global Institute, May 2011.
- 2 孟小峰,慈祥.大数据管理:概念、技术与挑战.计算机研究与发展,2013,50(1):146-169.
- 3 王珊,王会举,覃雄派,周烜.架构大数据:挑战、现状、与展望.计算机科学,2011,10:1741-1752.
- 4 马建光,姜巍.大数据的概念、特征及其应用.国防科技,2013,4:1671-4547.
- 5 夏秀峰,张亮,石祥滨,徐蕾.一种改进的分布式 ETL 体系结构.计算机应用与软件,2010,4:174-176.
- 6 李建江,崔健,等.MapReduce 并行编程模型研究综述.电子学报,2011,11:2635-2642.
- 7 莫荣强,艾萍,吴礼福,岳兆新,冯鹏.一种支持大数据的水利数据中心基础框架.水利信息化,2013,6(3).
- 8 郝树魁.Hadoop HDFS 和 MapReduce 架构浅析.邮电设计技术,2012,7.