

配电网监测数据的分布式 Map 压缩-查询技术

屈志坚,陈鼎龙,王群峰

(华东交通大学 电气与自动化工程学院,江西 南昌 330013)

摘要: 针对配电自动化中大量电力监测数据的处理问题,提出了回避聚合操作的配电网监测数据分布式 Map 压缩-查询新方法。通过将监测数据分布式 Map 压缩存储,利用 HQL 查询引擎及压缩接口将分布式 Map 压缩应用到连接查询的混洗阶段中,减小传递到查询聚合端的数据量,提高压缩数据的查询速度,并推导了时效性的相关公式。以北京某动车段 10 kV 电力运动监控系统的实测数据为例,搭建了四节点测试集群。压缩导入对比测试表明,分布式 Map 压缩速度快于分布式 Reduce 压缩,分布式 Map 的 Map_Deflate 压缩处理时间比分布式 Reduce_Deflate 减少了 45.3%;压缩-查询测试表明,当数据量为 2×10^7 记录级时,分布式 Map 的 Map_LZO 格式压缩-查询耗时大幅降低,比混洗阶段不压缩-查询时减少了 31.6%,验证了分布式 Map 压缩对加速查询的时效性。

关键词: 配电网; 大数据; 集群平台; 分布式压缩; 压缩-查询

中图分类号: TM 727; TP 311

文献标识码: A

DOI: 10.16081/j.issn.1006-6047.2017.12.027

0 引言

随着电网智能化和信息技术的快速发展,为使配电网更加安全可靠,配变电监测、用电信息采集、负荷控制、集中抄表等多类自动化系统投入运行^[1]。由于配电网处于电力系统末端,具有规模大、设备种类多、网络连接复杂、运行方式多变等特点,且监控数据采集点多、采集率较高,使系统产生大量的监测数据^[2-3],如江苏电网用电信息采集系统每日数据量可达数十 GB^[4],自动化系统信息长期运行的数据量级更是趋于 PB 级^[2,5-6]。目前配电自动化监控系统仍采用关系数据库,而关系数据库存取容量一般限制在 TB 级,越来越难以处理海量配电网监测数据,在某些地区的调度系统中已存在查询界面反应相当慢的问题,因此亟需研究大数据的快速查询和处理方法^[7]。数据压缩可以减少数据体量,有利于大幅缩小海量数据的查询规模,结合云计算集群环境下的分布式压缩处理和配电网大量监测数据需要连接查询的特点,研究分布式压缩-查询处理方法,可用于与配电网海量监测历史数据有关的数据分析,如历史趋势曲线和日月报表等,提供新的技术研究手段^[8]。

MapReduce 为 Apache Hadoop 生态链的一种分布式云计算模型,其可将计算任务分解为 Map 任务和 Reduce 任务两部分^[9],在实际应用中 Map 任务必

不可少,Reduce 任务可根据计算需要设置。对于大数据的分布式压缩,目前研究主要集中于数据的 Reduce 输出端压缩^[10-11],即分布式 Reduce 压缩,但该方法无法避开导入数据的 Reduce 任务,若能在 Map 环节直接进行分布式 Map 压缩,就可避免 Reduce 聚合过程的耗时。对于被压缩的大数据查询,由于不同应用需要选取不同的配电网监测数据关联属性,如历史曲线应用中有时选取监测对象编号,有时选取站地址等进行关联,同时由于 Hadoop 的主要查询工具 HQL (Hive SQL)不存在查询索引^[12],解压后的全表扫描量太大,使连接查询的性能问题尤为突出,制约了连接查询速度。而已有文献都只针对特定的查询语句,通过修改程序代码进行优化^[13-16],难以适应配电网时序数据应用中复杂多变的关联关系。当压缩数据连接查询时,若将 Map 压缩直接应用到连接查询的混洗阶段中,以大幅减小传递到查询输出端数据量的方式加快 MapReduce 的执行速度,就有可能提高压缩数据的连接查询效率。

针对配电网调度自动化应用中大量监测数据的处理问题,本文提出一种海量数据连接的分布式 Map 压缩-查询方法,即以分布式 Map 压缩方式将监测数据压缩存储并映射为 Hive 表,利用压缩接口将分布式 Map 压缩应用到连接查询的混洗阶段中,进行分布式 Map 压缩-查询,推导了压缩-查询时效性计算的相关公式,以 10 kV 配电网为算例进行测试,对分布式 Map 压缩-查询进行了验证。

1 配电网监测数据的分布式 Map 压缩

1.1 监测数据的分布式 Map 压缩映射

Sqoop 是 Apache 开放源码基金会组织的一个顶

收稿日期:2016-12-19;修回日期:2017-10-11

基金项目:国家自然科学基金资助项目(51567008);江西省杰出青年人才计划项目(20162BCB23045);江西省自然科学基金资助项目(20161BAB206156,20171BAB206044)

Project supported by the National Natural Science Foundation of China(51567008),the Outstanding Youth Talent Support Project of Jiangxi Province(20162BCB23045) and the Natural Science Foundation of Jiangxi Province(20161BAB206156,20171BAB-206044)

级项目,专用于 Hadoop 和关系库之间批量数据的传输,关系数据可导入 HDFS 并映射为 Hive 表。

配电网调度历史数据库中数量特别大的一类监测数据是时序数据,其典型格式定义为〈量测点标识,时间戳,量测值〉,以保存在调度监测数据库 hysd 中的历史遥测曲线表 TSTable 中的时序数据为例,用分布式 Map 压缩方式将其压缩映射为 Hive 表的流程如图 1 所示。

将遥测曲线表 TSTable 数据导入调度数据节点的语句为:sqoop import-connect url--direct--compress--

table TSTable--hive-import--split-by Num--hive-table TSTable。

具体步骤如下。

首先,将命令提交至 Sqoop 服务器,查询 TSTable 结构信息,生成导入数据的 MapReduce 作业代码文件 TSTable.jar,将代码文件提交至调度主机中的 Hadoop 作业客户端处理。

然后,作业客户端查询 TSTable 的序号字段的最大值和最小值,据(最大值-最小值)/Mapper 数平均分配表中的数据到 Mapper 映射器并行读取。以划

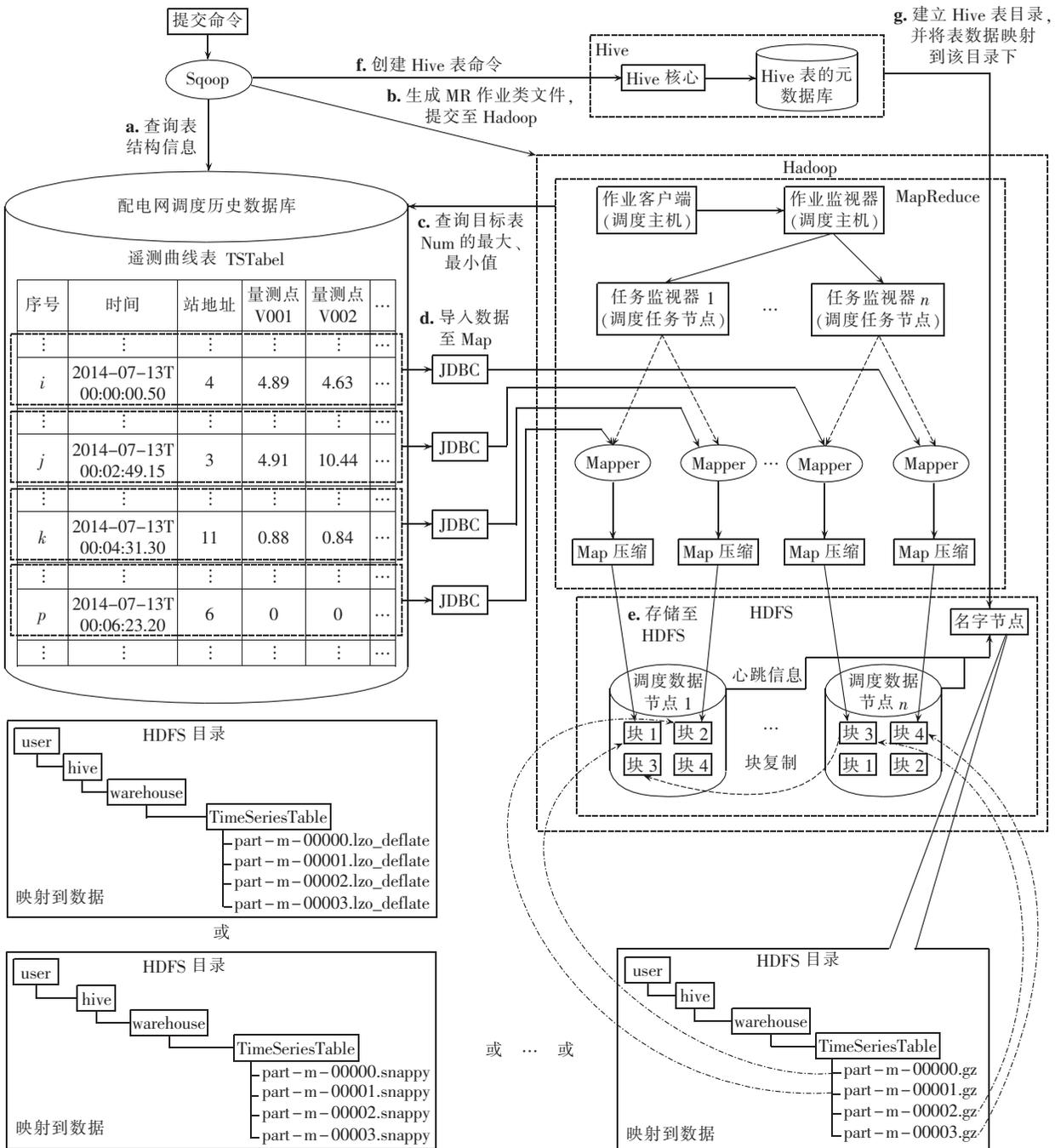


图 1 监测数据的分布式 Map 压缩流程图

Fig.1 Flowchart of distributed Map compression for monitoring data

分成 4 份为例,作业监视器将 4 份数据分配至 4 个 Mapper 映射器,通过 JDBC 接口从配电网调度历史数据库中并行读取数据,并将数据按指定压缩格式(Deflate、Gzip、Bzip2、LZO 或 Snappy),如 Gzip 格式,经分布式 Map 压缩后生成 part-m-00000.gz、part-m-00001.gz、part-m-00002.gz、part-m-00003.gz 共 4 个压缩文件,存储至 HDFS。

最后,利用 Sqoop 创建 Hive 表 TSTable,并将表结构信息存入元数据库,在 HDFS 的名字节点中创建一个与 TSTable 同名的目录,将 part-m-0000x 压缩数据文件映射至该同名目录。

1.2 分布式压缩的编码格式

Hadoop 集群支持 Deflate、Gzip、Bzip2、LZO 和 Snappy 压缩格式,其中 LZO 改进了 LZ 算法压缩格式和字符匹配搜索方式,与基于 LZ77 算法的 Deflate 和 Gzip 等相比,具有更快的压缩处理速度,其压缩编码格式如图 2 所示。图中, D 为指回距离,即当前待编码监测数据与已编码数据的最大匹配部分所间隔的字符数; L 为重复长度,即当前待编码监测数据与已编码数据的最大匹配部分的长度; N 为字节数。

2 分布式 Map 压缩加速查询的原理

2.1 分布式 Map 压缩的连接查询

将监测数据进行分布式 Map 压缩存储,并向 Hive 映射时,若将分布式 Map 压缩应用到连接查询的混洗阶段中,就可大幅减小传递到输出端的数据量,从而加快 MapReduce(简记为 MR)任务的执行。任意

的多表连接 HQL 语句为:

SELECT 目标列表表达式 1,目标列表表达式 2,目标列表表达式 3,...

FROM 表 1 JOIN 表 2 ON 表 1.连接字段=表 2.连接字段 JOIN 表 3 ON 表 1.连接字段=表 3.连接字段...

[WHERE <条件表达式>]。

以 Hive 中压缩的遥测曲线表 TSTable 与遥测定义表 Ycdef 通过站地址做等值连接查询为例说明加速查询的原理,对应的 HQL 查询语句为:

```
SELECT TSTable.DateTime,TSTable.
StationAddr,Ycdef.StationName,
Ycdef.OrderName,
TSTable.V001,Ycdef.ID,Ycdef.Unit
FROM TSTable JOIN Ycdef
ON TSTable.StationAddr=Ycdef.StationAddr
AND Ycdef.OrderName='V001'。
```

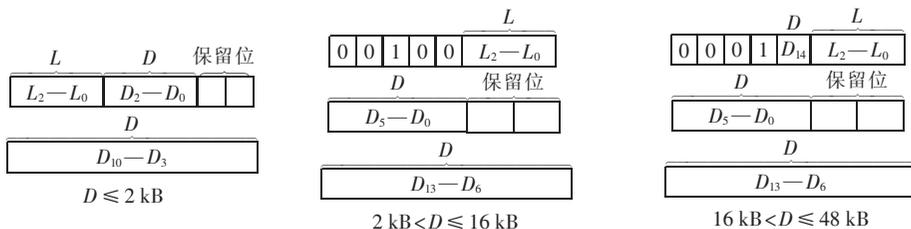
具体加速处理的流程如图 3 所示。

a. 分布式读取数据,添加表标记。

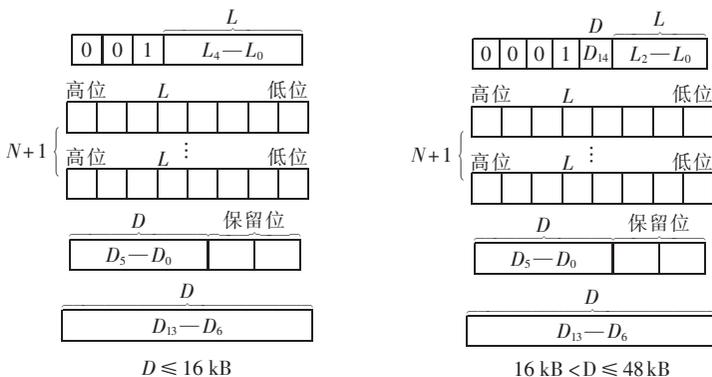
Mapper 负责从 HDFS 中读取数据块,并将其解压,以站地址作为键,以整条记录作为对应的值输出,并在值的首部添加表标记,表示此记录的表来源,如表标记“1”表示来自表 TSTable,表标记“2”表示来自表 Ycdef。

b. 分布式压缩数据。

对 Mapper 输出的监测数据进行压缩,减小传输到 Reduce 端的数据量,加快 MR 任务的执行。以



(a) 重复长度 $L \leq 8$



(b) 重复长度 $L > 8$

图 2 压缩编码格式

Fig.2 Compression code format

图 3 中 Mapper₁ 输出的监测数据压缩为例说明压缩过程。若 Mapper₁ 输出的监测数据按键、值进行二次排序后的结果为“3,1,j,2014-07-13T00:02:49.150,3,4,91;4,1,i,2014-07-13T00:00:00.500,4,4,89”,则记为 T_{S1} ,同理 Mapper₂ 二次排序结果可记为 T_{S2} ,依此类推,Mapper_k 二次排序结果可记为 T_{Sk} ,对此监测时序数据进行 LZ0 格式的分布式 Map 压缩的流程如图 4 所示(以 Mapper₁ 为例)。

T_{S1} 压缩后的监测数据为“(00)_H(0F)_H,3,1,j,2014-07-13T00:02:49.150,3,4,91,4,1,i,(2C)_H(74)_H

(00)_H, (0A)_H,0:00.500,4,4,89, (00)_H(00)_H”,即以 3 个字节表达“2014-07-13T00:0”。由于监测数据中的时间值、监测值重复较多,压缩前监测数据经过排序处理,相邻数据冗余量大,因而可以获得较高的压缩率。

c. 混洗和排序。

在用混洗排序处理压缩数据时,按键的哈希值将键值对发送到各个聚合输出端,以使站地址相同的记录汇集到同一个查询 Reduce,为在 Reduce 端进行等值连接查询做好准备。

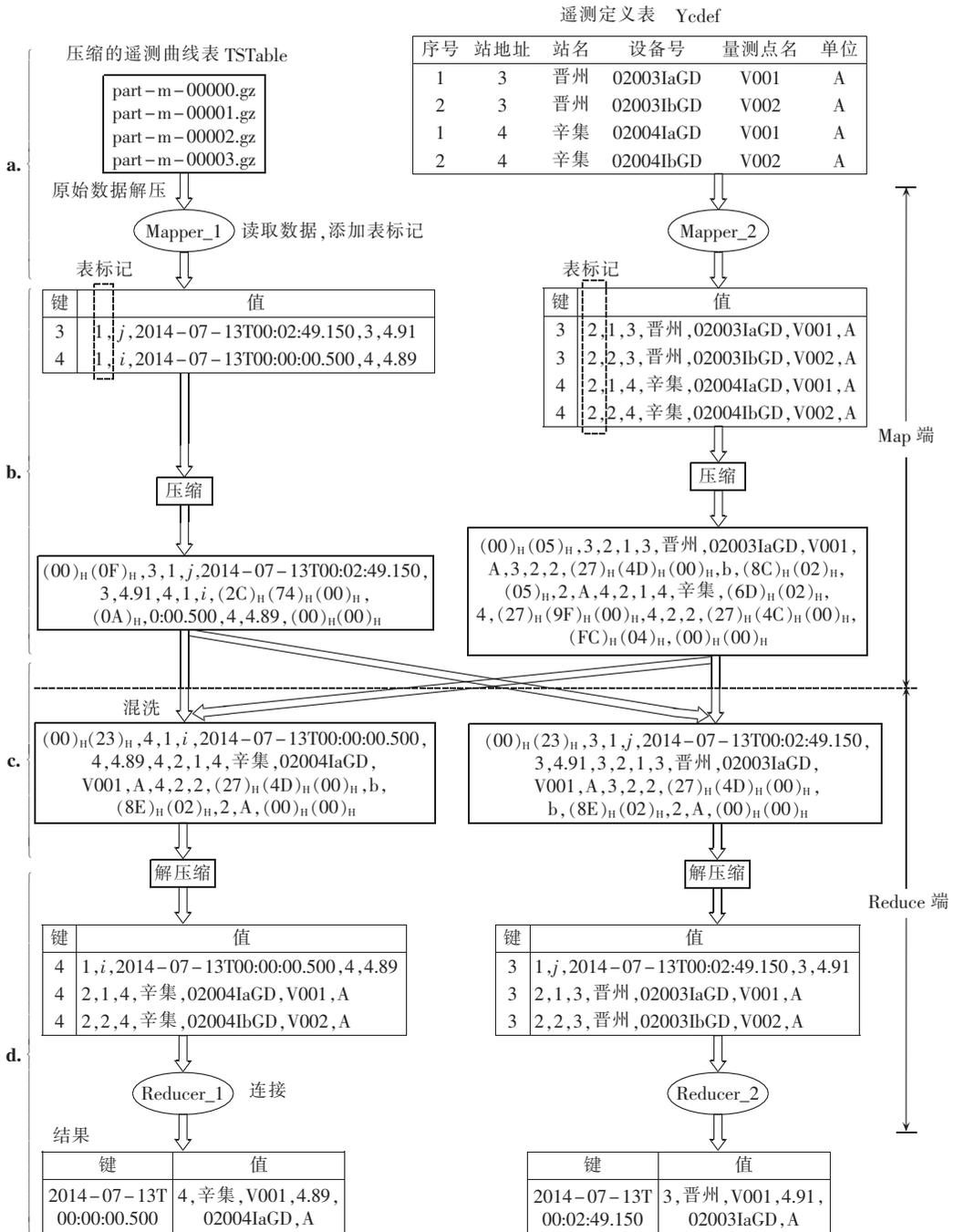


图 3 分布式 Map 压缩加速连接查询的原理

Fig.3 Principle of distributed Map compression to speed up join query

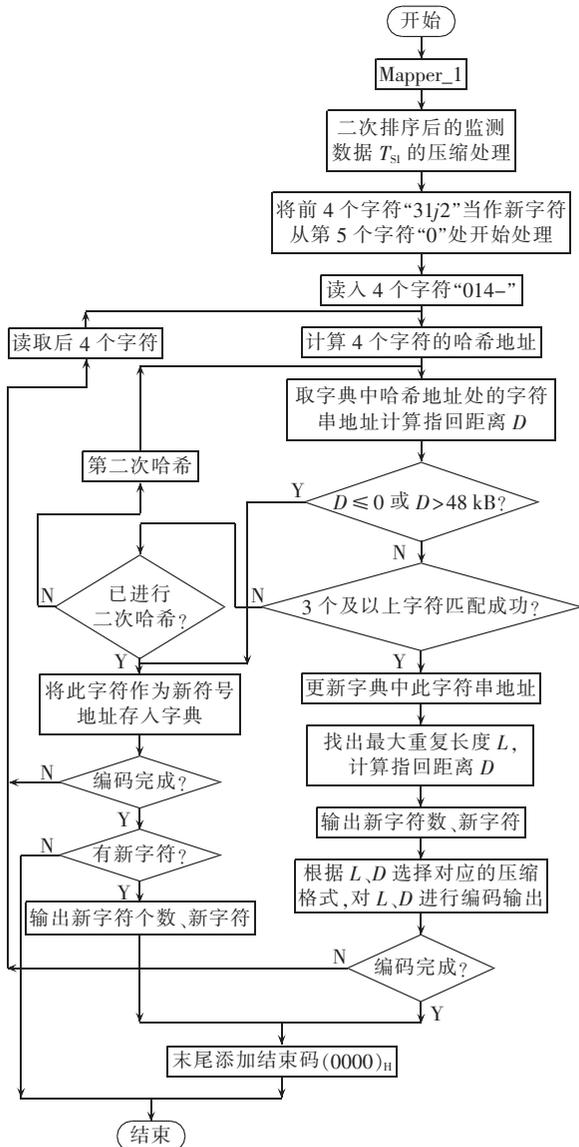


图 4 LZ0 格式的分布式 Map 压缩流程图

Fig.4 Flowchart of distributed Map compression based on LZ0 format

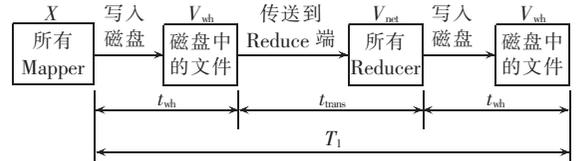
d. 解压缩及连接查询输出。

查询 Reduce 端收到监测压缩数据后,按相应的解码格式进行解压,根据步骤 a 中添加的表标记,将键值相同、来自不同表且测量点为 V_{xxx} 的记录连接,如连接查询语句中的 $Ycdef.OrderName = 'V001'$ 条件,指定测量点为 V001 的记录进行连接,如图 3 中得到的为辛集站测量点 V001 对应设备 02004IaGD 的监测值及其单位,最后将结果作为值输出至 HDFS。

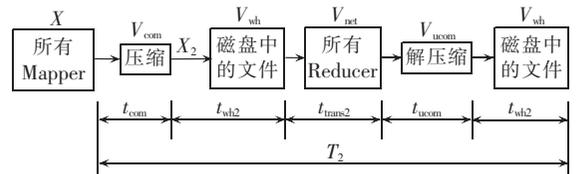
2.2 混洗阶段压缩-查询的时效性

对于一个经混洗处理的查询,令未采用分布式 Map 压缩时整个查询作业中所有 Map 输出的数据大小为 X , t_{wh} 为写磁盘时间, V_{wh} 为写磁盘的速度, t_{trans} 为数据由 Map 端传送到聚合输出端所用时间, V_{net} 为网络传输速度, T_1 为混洗阶段所需时间。

与此对应,令分布式 Map 混洗压缩时整个查询作业中所有 Map 输出的数据压缩后大小为 X_2 , V_{com} 为压缩速度, t_{com} 为压缩时间, t_{wh2} 为写磁盘时间, t_{trans2} 为数据由 Map 端传送到聚合输出端所用时间, t_{ucom} 为解压缩时间, V_{ucom} 为解压缩速度, T_2 为混洗阶段所需时间,未列出的变量及取值可参照未采用分布式 Map 压缩时的定义。压缩及解压缩速度单位相同,且与数据大小的单位相对应,如 X 的单位为 MB,则所有速度的单位为 MB/s。各环节中变量的定义如图 5 所示。



(a) 未采用分布式 Map 压缩时混洗阶段变量的定义



(b) 采用分布式 Map 压缩时混洗阶段变量的定义

图 5 混洗阶段变量定义

Fig.5 Variable definitions of shuffle stage

根据以上定义可知,未采用分布式 Map 压缩时混洗阶段所需时间 T_1 为:

$$T_1 = 2t_{wh} + t_{trans} = 2\frac{X}{V_{wh}} + \frac{X}{V_{net}} \quad (1)$$

采用分布式 Map 压缩时混洗阶段所需时间 T_2 为:

$$T_2 = t_{com} + 2t_{wh2} + t_{trans2} + t_{ucom} = \frac{X}{V_{com}} + 2\frac{X_2}{V_{wh}} + \frac{X_2}{V_{net}} + \frac{X_2}{V_{ucom}} = \frac{X}{V_{com}} + 2\frac{XR}{V_{wh}} + \frac{XR}{V_{net}} + \frac{XR}{V_{ucom}} \quad (2)$$

其中, R 为压缩比率,表示压缩后文件大小与压缩前文件大小之比。

T_1 与 T_2 的差值为:

$$\Delta T = T_1 - T_2 = X \left[\left(\frac{2}{V_{wh}}(1-R) + \frac{1}{V_{net}}(1-R) \right) - \left(\frac{1}{V_{com}} + \frac{R}{V_{ucom}} \right) \right] = X(t_{save} - t_{add}) \quad (3)$$

其中, t_{save} 为采用分布式 Map 压缩后节省的磁盘 I/O 和网络 I/O 时间; t_{add} 为采用分布式 Map 压缩后增加的时间。

根据式(3)可知,对于一个确定的 Hadoop/Hive 集群, V_{wh} 、 V_{net} 近似为常数,当选定压缩编码格式后,将 R 、 V_{com} 、 V_{ucom} 代入式(3)使 $\Delta T > 0$ 时,才能加快分布式 Map 压缩的查询速度,且当 $t_{save} > t_{add}$ 时, X 越大,

ΔT 越大,性能就越好。因此,需要通过实验确定不同压缩格式的分布式压缩-查询性能。

3 算例测试

3.1 实验平台及实验数据

以北京某动车段 10 kV 电力远动监控系统为例进行实验,其监控对象主要包括配电所、检查库变电所、检修库变电所、生产调度变电所、信号楼 10 kV 变电所和转向架变电所。

为了验证利用分布式 Map 压缩时的连接查询性能,搭建了五节点分布式实验平台,其中 4 个节点构建 Hadoop/Hive 集群,1 个节点作为数据服务器,Hive 配置在主节点上作测试机。

集群配置为:1 个主节点,CPU Intel Core i5-4590、3.30 GHz、4 核,内存 8 GB、DDR3;3 个从节点,从节点 1 的配置与主节点相同,从节点 2,CPU Intel Core T7600、2.33 GHz、双核,内存 3 GB、DDR2,从节点 3,CPU AMD Athlon 5200+、2.7 GHz、双核,内存 4 GB、DDR2;操作系统均为 64 位 Ubuntu 14.04。数据服务器的配置与集群的主节点相同,操作系统为 64 位 Windows7。

将系统现场记录的 2014-11-20 至 2014-12-19 时段遥测表 Yc、遥测设备表 AnalogDev 导入本实验平台的数据服务器,作为测试数据。其中,表 Yc 主要包括动车段变配电所的三相电压、三相电流、有功功率、无功功率、频率、变压器绕组和铁芯温度等模拟量的时序数据,30 d 共有 3.2×10^7 条记录,平均每日记录数约为 10^6 条,平均采样周期为 80.84 ms;表 AnalogDev 记录数为 2 171 条。其格式分别如表 1、表 2 所示。

表 1 遥测表 Yc 格式

Table 1 Telemetry table based on Yc format

时间	站地址	设备号	遥测值
2014-11-20T00:00:00.196	3	02003k3-Ia	14.4
2014-11-20T00:00:00.196	3	02003k3-Ic	13.9
2014-11-20T00:00:00.196	3	02003k3-P	230.2

表 2 遥测设备表 AnalogDev 格式

Table2 Telemetry equipment table based on AnalogDev format

序号	设备号	站地址	设备名	单位	参比因子
37	02003k3-Ia	3	馈出三-测量 Ia	A	0.002 75
38	02003k3-Ib	3	馈出三-测量 Ib	A	0.002 75
39	02003k3-Ic	3	馈出三-测量 Ic	A	0.002 75
43	02003k3-P	3	馈出三-P	kW	0.098 90

3.2 分布式 Map 压缩导入性能测试

为了测试分布式 Map 压缩对配电监测数据的压缩性能,将关系库中的遥测表 Yc 分别用分布式 Map 压缩、分布式 Reduce 压缩方法导入 Hive。导入遥测表 Yc 的 3.4×10^6 至 3.1×10^7 (3.4×10^6 、 6.8×10^6 、

1×10^7 、 1.3×10^7 、 1.6×10^7 、 2×10^7 、 2.3×10^7 、 2.6×10^7 、 2.9×10^7 、 3.1×10^7) 条监测记录到 Hive,测试用不同压缩格式时 2 种压缩方法的导入时间,如图 6 所示。

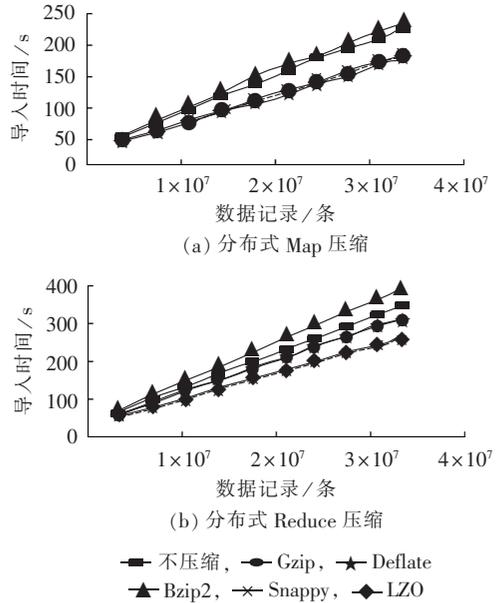


图 6 分布式 Map/Reduce 压缩导入时间

Fig.6 Import time of distributed Map/Reduce compression

对比图 6(a)、(b)可知,5 种压缩格式的分布式 Map 压缩均明显比同种压缩格式的分布式 Reduce 压缩速度快,如当记录数为 3.1×10^7 条时,采用分布式 Map 的 Map_Deflate 压缩导入时间比分布式 Reduce_Deflate 压缩减少了 45.3%。

从图 6(a)还可知,除 Map_Bzip2 格式外,采用其他 4 种压缩格式的分布式 Map 压缩后,与未采用压缩时相比,导入时间明显减少,当记录数大于 2×10^7 条时趋于平缓;采用 Map_Deflate 格式效果最佳,当记录数为 3.1×10^7 条时,比未采用压缩减少了 24.5%。

3.3 分布式 Map 压缩-查询测试

对存储在 Hive 中的压缩的遥测表 Yc 与遥测设备表 AnalogDev,进行混洗阶段未采用分布式 Map 压缩和采用不同压缩格式的分布式 Map 压缩的连接查询对比实验。两表通过设备号进行等值连接查询,对于同一组数据,压缩格式分别取 Map_Deflate、Map_Gzip、Map_Bzip2、Map_LZO、Map_Snappy,记录每组实验的查询耗时、Map 输出字节数。Map 输出字节数对比如图 7 所示。

由图 7 可知,采用分布式 Map 压缩后,Map 输出字节数明显大幅减小,大量减小了混洗阶段的磁盘 I/O 和网络 I/O 时间。如当记录数为 3.1×10^7 条时,采用这 5 种压缩格式的分布式 Map 压缩后,均能使 Map 输出字节数由 1 870 MB 减小到 400 MB 以下,减小了约 79%。

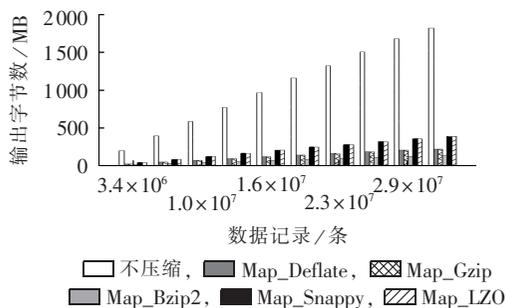


图 7 Map 输出字节数对比

Fig.7 Comparison of Map output bytes

对于压缩率,采用 Map_Deflate、Map_Gzip 时均为 12.2%,采用 Map_Bzip2 时为 7.4%,采用 Map_Snappy 和 Map_LZO 时均为 21.2%,其均不随数据量变化。结果表明,当这 5 种压缩格式的分布式 Map 压缩应用于连接查询的混洗阶段时,Map_Bzip2 的压缩率最大,Map_Deflate 和 Map_Gzip 的较大,Map_Snappy 和 Map_LZO 适中。查询耗时对比如图 8 所示。

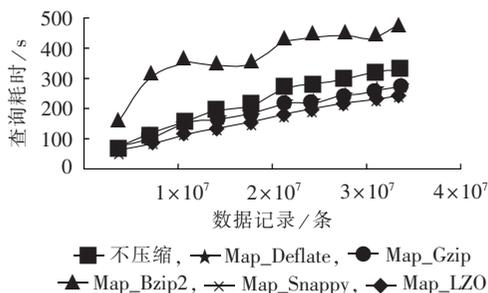


图 8 查询耗时对比

Fig.8 Comparison of query time

由图 8 可知,采用分布式 Map 压缩后,当选用合适的压缩格式时,两表连接的查询耗时会明显减少,加速效果明显;分布式 Map 的 Map_LZO 格式压缩效果最好,略优于 Map_Snappy,当记录数为 2×10^7 条时,比未采用分布式 Map 压缩时耗时减少了 31.6%,原因在于 Map_LZO、Map_Snappy 格式的压缩率适中,压缩、解压缩速度快,使减少的磁盘 I/O 和网络 I/O 时间大于由于压缩和解压缩而增加的时间,即使得式(3)大于 0。

对于分布式 Map 的 Map_Default、Map_Gzip 格式压缩,在数据量大于 1×10^7 条记录数后,开始有优化效果,并在达到 2×10^7 条记录数时效果明显,而 Map_Bzip2 由于压缩率太大,虽可使数据量大幅减小,在一定程度上节省了磁盘 I/O 和网络 I/O 时间,但其压缩及解压缩的处理时间却大幅增加,超过了因数据体量降低而节省的时间,即使式(3)小于 0。

4 结论

a. 针对海量数据压缩及其快速查询问题,在云计算集群基础上提出了一种分布式 Map 压缩及混洗查

询的快速处理方法,通过压缩-查询有关变量的定义,确立了影响查询时效性的压缩比、压缩速度和解压缩速度之间的函数关系,推导了时效性计算的公式。

b. 以北京某动车段电力运动监控的实测数据为例,研究了配电网监测数据的分布式 Map 压缩编码及接口。集群环境下分布式压缩导入测试结果表明,分布式 Map 压缩比分分布式 Reduce 压缩速度快。

c. 压缩-查询测试的结果表明,当数据量达到 2×10^7 条记录数时,分布式 Map 的 Map_LZO 格式压缩所需的查询时间,相比混洗阶段未压缩而直接查询的时间大幅降低,减少了 31.6%,验证了所提方法的有效性,可提高海量监测时序数据的查询性能,不仅减少了数据体量,而且加快了查询处理的速度。

参考文献:

- [1] 王金丽,盛万兴,王金宇,等. 中低压配电网统一数据采集与监控系统设计和实现[J]. 电力系统自动化,2012,36(18):72-76,81.
WANG Jinli, SHENG Wanxing, WANG Jinyu, et al. Design and implementation of a centralized data acquisition and supervisory system for medium-low voltage distribution network [J]. Automation of Electric Power Systems, 2012, 36(18): 72-76, 81.
- [2] 葛磊蛟,王守相,瞿海妮. 智能配用电大数据存储架构设计[J]. 电力自动化设备,2016,36(6):194-202.
GE Leijiao, WANG Shouxiang, QU Haini. Design of storage framework for big data of SPDU [J]. Electric Power Automation Equipment, 2016, 36(6): 194-202.
- [3] 薛振宇,胡航海,宋毅,等. 基于大数据分析的县公司综合评价策略[J]. 电力自动化设备,2017,37(9):199-204.
XUE Zhenyu, HU Hanghai, SONG Yi, et al. Comprehensive evaluation based on big data analysis for county electric power company [J]. Electric Power Automation Equipment, 2017, 37(9): 199-204.
- [4] 郑海雁,金农,季聪,等. 电力用户用电数据分析技术及典型场景应用[J]. 电网技术,2015,39(11):3147-3152.
ZHENG Haiyan, JIN Nong, JI Cong, et al. Analysis technology and typical scenario application of electricity big data of power consumers [J]. Power System Technology, 2015, 39(11): 3147-3152.
- [5] 李济沅,钟一俊,王东举,等. 基于 CIM 和 ESB 的企业级电网准实时资源中心研究与应用[J]. 电力自动化设备,2015,35(8):148-155.
LI Jiyuan, ZHONG Yijun, WANG Dongju, et al. Research of enterprise-class quasi-real-time grid resource center based on CIM & ESB and its application [J]. Electric Power Automation Equipment, 2015, 35(8): 148-155.
- [6] Marie-Luce PICARD,潘旭阳. 从法国公共电力企业的视角看大数据带来的挑战和机遇[J]. 电网技术,2015,39(11):3109-3113.
Marie-Luce PICARD, PAN Xuyang. Big data challenges and opportunities for a utility as EDF [J]. Power System Technology, 2015, 39(11): 3109-3113.
- [7] 王德文,肖磊,肖凯. 智能变电站海量在线监测数据处理方法[J]. 电力自动化设备,2013,33(8):142-146.
WANG Dewen, XIAO Lei, XIAO Kai. Processing of massive online monitoring data in smart substation [J]. Electric Power Automation Equipment, 2013, 33(8): 142-146.
- [8] 邵振国,吴瑾樱,苏文博. 面向海量历史监测数据的谐波污染用

- 户统计建模方法[J]. 电力自动化设备,2016,36(8):110-114.
- SHAO Zhenguo,WU Jinying,SU Wenbo. Statistical modeling based on massive historical monitoring data for harmonic pollution customer[J]. Electric Power Automation Equipment, 2016,36(8):110-114.
- [9] GUNARATHNE T,WU T L,QIU J,et al. MapReduce in the clouds for science[C]//Cloud Computing Technology and Science, 2010 IEEE Second International Conference on Digital Object Identifier. Indianapolis, USA:IEEE,2010:565-572.
- [10] 屈志坚,郭亮,陈秋琳,等. Hadoop 云构架的智能调度无聚集群压缩技术[J]. 电力系统自动化,2013,37(18):93-98.
- QU Zhijian,GUO Liang,CHEN Qiulin,et al. Intelligent dispatching lossless cluster compression technology based on Hadoop cloud framework[J]. Automation of Electric Power Systems,2013,37(18):93-98.
- [11] 屈志坚,郭亮,刘明光,等. 智能配电网量测信息变断面柔性压缩新算法[J]. 中国电机工程学报,2013,33(19):191-199.
- QU Zhijian,GUO Liang,LIU Mingguang,et al. New variable section flexible compression algorithm for measurement information in intelligent distribution network[J]. Proceedings of the CSEE,2013,33(19):191-199.
- [12] 屈志坚,陈阁. 容错存储的电力系统监测数据查询优化技术[J]. 电网技术,2015,39(11):3221-3227.
- QU Zhijian,CHEN Ge. Query optimization for power system monitoring data with fault-tolerant storage[J]. Power System technology,2015,39(11):3221-3227.
- [13] 黄山,王波涛,王国仁,等. MapReduce 优化技术综述[J]. 计算机科学与探索,2013,7(10):885-905.
- HUANG Shan,WANG Botao,WANG Guoren,et al. A survey on MapReduce optimization technologies[J]. Journal of Frontiers of Computer Science and Technology,2013,7(10):885-905.
- [14] 张首正,周凯东. 基于 MapReduce 的 SQL 查询优化分析[J]. 计算机应用,2014,34(增刊 2):63-65.
- ZHANG Shouzheng,ZHOU Kaidong. MapReduce based query optimization for SQL analysis[J]. Journal of Computer Applications,2014,4(Supplement 2):63-65.
- [15] 赵辉,杨树强,陈志坤,等. 基于 MapReduce 模型的范围查询分析优化技术研究[J]. 计算机研究与发展,2014,51(3):606-617.
- ZHAO Hui,YANG Shuqiang,CHEN Zhikun,et al. Optimization of range queries and analysis for MapReduce systems[J]. Journal of Computer Research and Development,2014,51(3):606-617.
- [16] 王梅,邢露露,孙莉. 混合存储下的 MapReduce 启发式多表连接优化[J]. 计算机科学与探索,2014,8(11):1334-1344.
- WANG Mei,XING Lulu,SUN Li. MapReduce based heuristic multi-join optimization under hybrid storage[J]. Journal of Frontiers of Computer Science and Technology,2014,8(11):1334-1344.

作者简介:



屈志坚

屈志坚(1978—),男,江西南昌人,副教授,博士,主要研究方向为智能监控理论与信息处理技术(E-mail:08117324@bjtu.edu.cn);

陈鼎龙(1988—),男,广西贵港人,硕士研究生,主要研究方向为调度自动化监控信息大数据处理技术;

王群峰(1992—),男,江西南昌人,硕士研究生,主要研究方向为智能配电网大数据技术。

Distributed Map compression-query technology for distribution network monitoring data

QU Zhijian,CHEN Dinglong,WANG Qunfeng

(School of Electrical and Automation Engineering,East China Jiaotong University,Nanchang 330013,China)

Abstract: In view of the considerable quantity of monitoring data in distribution automation,an innovative method of distributed Map compression-query for distribution network monitoring data is proposed,in which the aggregation operation can be avoided. The monitoring data is compressed and stored through distributed Map compression. Then,the distributed Map compression is applied to the shuffle stage of join query by HQL query engine and compression interface,which reduces the data quantity transferred to the reducers,and improves the query speed of the compressed data. Besides,the formulas of time effectiveness are deduced. The measured data of a Beijing EMU(Electric Multiple Unit) depot 10 kV power remote monitoring system is taken as an example,where a four-node cluster is constructed to carry out experiments. The compression importing comparison test results demonstrate that the distributed Map compression speed is faster than the distributed Reduce compression,and the Map_Deflate compression processing time decreases by 45.3% compared to distributed Reduce_Deflate. Meanwhile,the compression-query test results demonstrate that the compression-query time of Map_LZO greatly reduces when the amount of data is 2×10^7 records,which decreases by 31.6% compared to uncompressed-query in shuffle stage. Hence,these results verify the efficiency of distributed Map compression for accelerating query.

Key words: distribution network; massive data; cluster platform; distributed compression; compression-query