

面向电力大数据的异构数据混合采集系统^①

孙超, 王永贵, 常夏勤, 陆鑫, 顾全

(南京南瑞继保电气有限公司, 南京 211102)

通讯作者: 孙超, E-mail: sunc@nrec.com

摘要: 随着智能电网的快速发展, 电力系统数据量的增长也非常迅速, 电力大数据急待开展深入研究. 电力数据产生的速率跨度大, 数据源众多且交互方式繁杂, 数据种类繁多等特点, 已有大数据采集方式难以适应多源异构数据的混合采集应用场景. 本文针对电力大数据提出了新的解决方案, 通过混合数据采集模型和采集集群实现了对异构数据源采集任务的混合调度和管理; 通过数据置信度标签技术, 在保留原始数据的同时, 标示数据的质量, 为后续大数据分析应用提供了便利; 通过 Sqoop、Kafka、文件传输等方式将采集与处理后的数据提交给大数据平台存储. 系统已经在用户现场部署并投入使用, 运行稳定, 效果良好.

关键词: 电力大数据; 异构数据源; 混合采集; 置信度; 数据监视

引用格式: 孙超, 王永贵, 常夏勤, 陆鑫, 顾全. 面向电力大数据的异构数据混合采集系统. 计算机系统应用, 2018, 27(12): 62-68. <http://www.c-s-a.org.cn/1003-3254/6667.html>

Mixed Heterogeneous Data Acquisition System for Power Big Data

SUN Chao, WANG Yong-Gui, CHANG Xia-Qin, LU Xin, GU Quan

(Nanjing NR Electric Co., Ltd., Nanjing 211102, China)

Abstract: With the rapid development of smart grid, the growth of power system data is also very fast. There is an urgent need of in-depth research for power big data. Because of the large span of the data acquisition speed, numerous data sources, complicated interaction interfaces, and various kinds of data type, the existing big data technologies are unable to adapt to power big data acquisition. In this study, a new solution for power big data acquisition is proposed. The system schedules and manages the heterogeneous data source acquisition tasks through mixed data acquisition model and collection cluster. With the technology of data confidence degree label, the system preserves the original data, indicates the quality of the data and provides convenience for big data analysis applications. The system submits collected data to the big data platform for storage by Sqoop, Kafka, file transfer, or other methods. The system has been deployed and puts into use in the user site. It runs stably and has a sound effect.

Key words: big data for power system; heterogeneous data sources; mixed data acquisition; confidence degree; data monitoring

随着特高压交流、柔性直流工程的建设, 电网的形态和特性发生重大变化, 西电东送规模不断扩大, 电网的联系愈加紧密, 电网运行方式更趋复杂, 未来电力供需平衡压力仍然巨大, 复杂大电网的潜在安全风险

将长期存在, 需要从多层级、大范围综合保障大电网的安全运行, 对电网运行人员驾驭大电网的能力、大范围、多目标资源优化配置的能力和电网运行的一体化运作水平提出了新的更高的要求. 现有电网监控类

① 基金项目: 国家高技术研究发展计划 (863 计划)(2015AA050201)

Foundation item: National High-tech R&D Program of China (863 Program) (2015AA050201)

收稿时间: 2018-04-29; 修改时间: 2018-05-24; 采用时间: 2018-06-05; csa 在线出版时间: 2018-12-03

系统获取的设备及电网运行的各种状态数据信息不能实现高度共享, 缺乏有效的管理, 孤立的数据难以形成有效的信息, 给电网的运行管理和科学决策带来了很大的盲目性, 已不能适应未来电网的发展要求. 近年来, 随着信息技术的发展, 全球数据量呈爆发式增长. 大数据的分析在国内外得到了迅速的发展和广泛的应用, 并取得了良好的社会效益. 随着我国电力行业信息化水平的快速发展, 电力系统数据量的增长也呈现出爆发的趋势, 电力大数据急待开展深入研究^[1-5]. 但是, 由于电力数据产生的速率跨度大^[6,7], 比如毫秒级广域向量测量实时数据, 秒级的稳态监视数据, 分钟级的微气象数据, 小时级的操作票流转数据和更长时间周期的设备实验数据等; 数据源众多且交互方式繁杂, 比如 WebService、电力专用规约、特殊文件格式等; 数据种类繁多, 比如实时数据、历史数据、文本数据、多媒体数据、时间序列数据等各类结构化、半结构化数据以及非结构化数据, 因此, 开展电力大数据分析的前提是开发多源异构数据混合采集系统.

1 技术现状与研发目标

目前, 大数据领域有多种工具实现外部数据的采集采用和处理, 但他们都面向特定的应用场景, 部署和管理的机制也各不相同. 比如, Flume 是分布式日志采集技术, 可支持文本、数据库、console 输出等数据源, 将数据最终导入 HDFS 或 HBase 中; Kafka 是分布式发布订阅消息系统并不直接接触数据源^[8], 需要定制开

发数据源采集程序后经过 Kafka 分布式队列传输给数据消费端再存入 Hive、HBase、HDFS 等存储中; Sqoop 是面向关系数据库结构化数据的全量或增量数据采集并将数据存储到 HDFS 中的技术^[9], 以上三种技术都面向大数据海量采集场景, 实现了分布式横向扩展, 具有高吞吐量特征, 适用于采集和处理非实时或弱实时类数据. 当面对强实时数据的采集场景时需要更快速的数据处理技术^[10,11], 比如 Storm, 以便能够在内存中对毫秒级数据及时处理. 以上各类工具的数据采集面向的场景不同, 难以适应多源异构数据的混合采集应用场景; 各类工具的数据获取、转换和装载过程自成体系, 三个过程相互关联, 造成数据处理转换过程相对独立, 很难引入多源数据的交叉校验; 数据获取后直接入分布式数据库或分布式存储, 没有对对象进行统一的编码, 造成后续处理的困难.

针对这一现状, 研发面向电力大数据的异构数据混合采集系统, 实现异构数据源的混合接入和集群管理, 实现采集数据的高速缓存与刷新; 海量采集数据的数据质量校验与转化; 采集数据的统一编码、实时交换和数据接入情况的监视.

2 系统概述

2.1 系统逻辑架构

系统由数据接口层、数据采集与转换层和数据发布层三部分组成, 如图 1 所示.

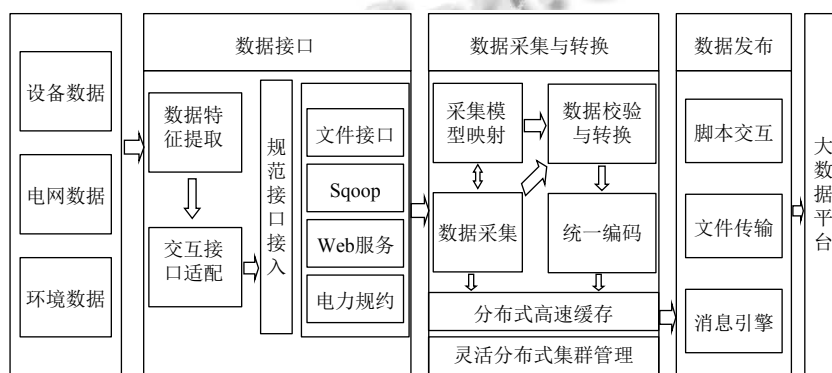


图 1 系统逻辑架构图

数据接口解决了不同类型采集数据接入方式的问题, 数据经过特征提取识别数据格式和交互方式, 适配对应的交互接口.

数据采集基于大数据分布式集成技术, 形成一个在线分布式的采集平台, 基于灵活分布式集群将异构系统的多源数据进行统一采集.

数据校核与转换基于分布式的内存数据库技术,实现了采集数据的高速刷新和处理.多样化的数据校核和转换,把对数据集的大规模操作分发给网络上的每个节点,实现海量数据处理的实时性和可靠性.引入了电力对象注册中心作为全局对象的统一管理设施.

数据发布基于高速实时总线技术,提供海量实时数据的消息总线,实现集数据的实时交换和发布.

2.2 系统存储架构

系统采集的数据最终提交给大数据平台,存储于HDFS分布式存储、HBase、Hive数据库中.图2是分类数据对应的存储模式.

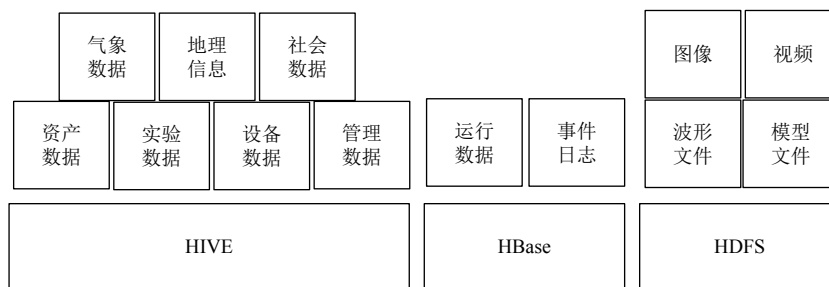


图2 系统数据存储架构图

结构化数据分为两种:一种是周期获取的非实时类数据,比如资产数据、实验数据、设备数据、管理数据、气象数据、地理信息和社会数据,此类数据具有固定的表结构,通常用SQL查询,存放在Hive数据库中;另一类是实时数据,比如电网运行数据,采集速率在毫秒级和秒级,事件日志类突发性强的数据,此类数据对数据吞吐性能要求较高,且访问方式较为单一,一般按时间序列和对象ID查询,采用键值对方式存放在HBase数据库中.半结构化数据,如波形文件、模型文件和非结构化数据如图像和视频以文件形式存放在HDFS分布式存储中.非实时类结构化数据通过Sqoop脚本定时从源系统增量抽取到Hive数据库中;实时类结构化数据由源系统发布到Kafka总线的实时数据主题中,采集端通过订阅相关主题数据存储到HBase库中;半结构化数据与非结构化数据通过文件传输协议存储到HDFS文件系统的按照文件类别和时间分类的目录中.

2.3 系统部署架构

系统采用PC服务器和虚拟化技术部署,主体功能部署在生产管理区,需要与互联网交互的功能,比如互联网上社会数据的获取、气象台预报数据的获取等,部署在DMZ区,系统部署图如图3.

前端采集集群负责与其他业务系统交互,采集各类数据,其中互联网数据需要通过DMZ区的互联网采集代理获取并缓存数据,再由前端采集集群发起二次

采集;数据转换集群负责采集数据的校验、转换和编码;最后,由数据发布集群按照数据类型特征将数据存储到大数据平台的Hive、HBase和HDFS中.系统与外部业务系统间通过生产管理大区的综合数据网交互;系统与互联网之间的数据采集通过DMZ区防火墙交互,数据交互只能由采集集群发起单向数据获取,从而保证内部系统与外部环境的安全隔离.

3 系统核心功能

3.1 混合数据采集模型

混合数据采集通过异构数据源模型智能映射技术实现采集模型的统一,它包含以下内容:(1)公共的模型信息;(2)公共对象信息和各异构数据源私有对象信息的映射关系;(3)公共数据服务和各异构数据库私有的数据服务的映射关系.

(1) 基于抽象容器的公共模型

异构数据源模型智能映射技术通过抽象容器将一般性的组织和事物中“包含”的关系理解为各级容器的从属结构,可以建立各种不同应用的模型,从而具有灵活性和扩展性.

如图4所示,在公共模型中,每种“容器”包含各自的“对象类型”及“容器描述属性”.当模型需要扩展时,只需在“对象类型定义”增加新的“对象类型”,在“容器描述定义”增加新的“容器描述属性”即可.

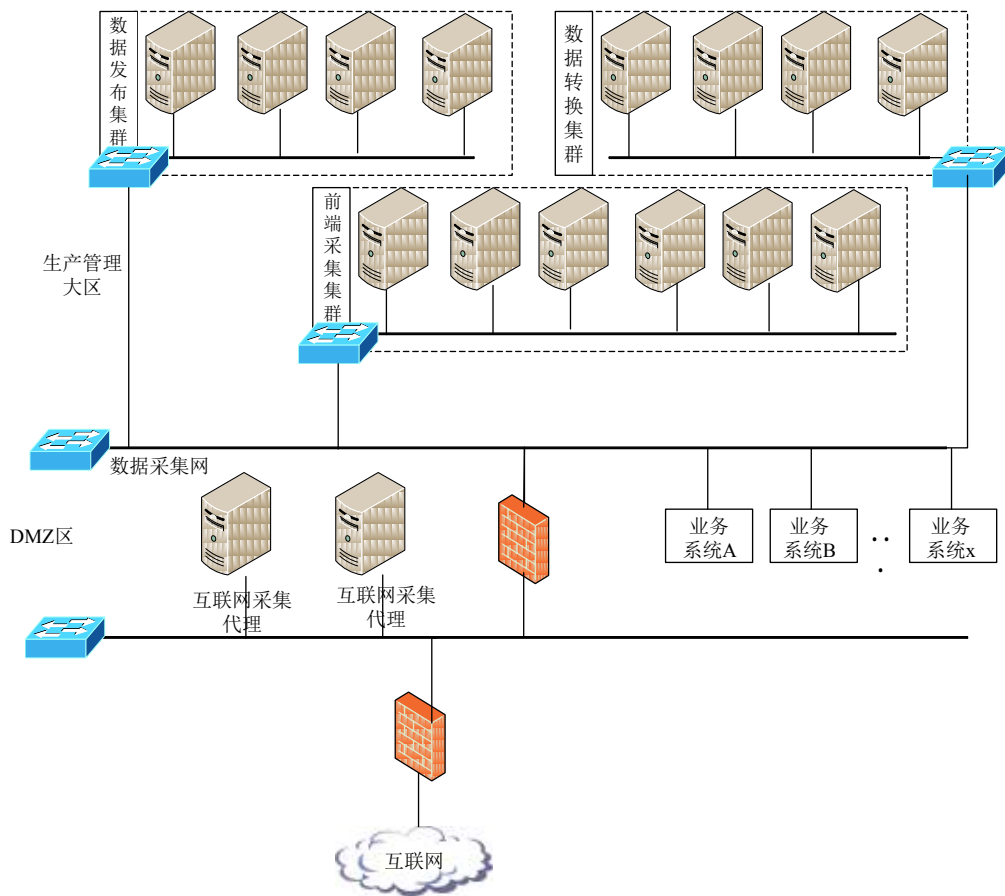


图3 系统部署架构图

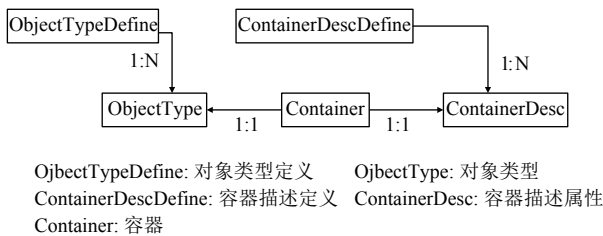


图4 容器及公共对象附加属性描述

(2) 公共对象信息和异构数据源私有对象信息的映射关系

公共对象信息是抽取了各异构数据源的对象信息部分, 通过公共命名方式对各系统的私有对象信息进行归纳, 提取出公共对象信息, 并将公共对象信息、私有对象信息以及它们的映射关系装载到已经建立好的具体的容器模型结构中。

(3) 公共数据服务和各异构数据源私有的数据服务的映射关系

异构数据源模型智能映射技术的管理数据服务分

为两部分: 公共数据服务管理、各异构数据库源的数据服务管理、以及两者之间的映射关系. 公共数据服务对混合数据采集提供统一的查询数据结构; 各异构数据库源的数据服务基于对各数据源 Agent 代理的子查询, 提供其内部具体的数据结构; 而公共查询到各数据源子查询是根据上述两者之间的映射关系。

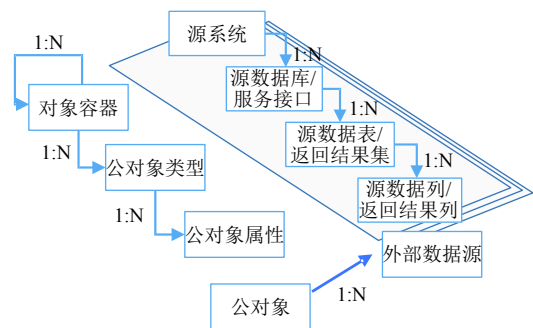


图5 公共数据服务和各数据源数据属性的映射关系

如图5所示, 公共模型的数据结构虽然与源系统

的数据结构不同, 但通过公共数据属性与源数据列之间的映射关系及源数据列与源数据表/服务接口返回结果集、源数据库/服务接口、源系统之间的层级关系, 即可获取源系统的各项信息.

3.2 数据采集集群

分布式采集集群通过异构系统模型智能映射获取各个数据源的元数据信息, 将每个数据源按照采集量横向分片形成不同的采集任务, 每个采集任务对应一个或多个冗余采集通道, 采集集群通过按节点分配或 NAT 映射等技术将采集通道在集群内各节点上的均衡分散运行, 提高集群的整体并发性.

(1) 节点分布式负载均衡策略

集群具备按节点分布式运行的能力, 不同数据源的采集通道具备按节点动态负载均衡技术分组集群并行处理能力. 采集应用在 M 个采集节点上运行, 采集通道有 N 个, 按照最理想化的负载均衡效果, 当前时间, 每台采集节点上将有 $\lceil N/M \rceil$ 个通道在运行, 这样保证了通道采集在每台采集节点上运行的负载均衡性. 算法如图 6 所示.

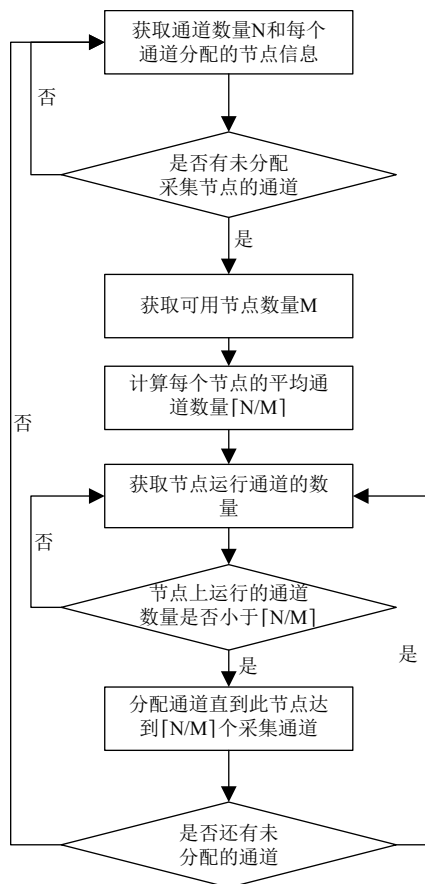


图 6 采集通道负载均衡算法

假如某个采集集群在四台采集节点上运行, 采集通道共有 96 个, 在正常情况下 (所有采集通道通信正常, 所有采集节点运行正常), 每个采集节点将有 24 个通道正常运行, 在某一时刻, 如果某个采集节点发生故障, 这在该节点上运行的 24 个采集通道, 将会按负载均衡算法转移到另外正常的三台采集节点, 即每个采集节点运行 32 个采集通道. 同样, 一旦故障采集节点恢复正常, 采集通道将恢复到四节点运行状况.

在前端采集集群中, 各通道在采集节点上运行都具有优先级指数, 该优先级指数在配置各采集通道时派生, 派生方式可以人工设置, 也可以通过程序按照当前通道总数和采集节点总数关系自动派生.

采集管理采用竞价机制竞选同一通道在不同节点上获取资源权限, 优先级高的节点优先获取资源, 进入启用状态, 在一定时间内通道连通后正式获取资源, 其他节点处于候选状态. 如该节点一定时间内不能连通则把该通道权限移交该优先级低的节点启用.

3.3 置信度评估

对于大数据分析而言对数据质量的要求与常规数据挖掘中对数据的质量要求不同, 大数据分析侧重弱关联关系, 需要保留更多的原始数据, 而常规的策略往往将数据按照业务的要求清洗异常数据. 系统通过在数据校核和转换的处理过程中加入置信度评估方法, 将经过校核转换处理的数据和原始数据分别给与不同的置信度标签, 而不是将数据直接清洗删除掉, 在保留数据的同时也起到了数据辨识的作用, 为后续的数据分析类应用提供了更多的选择权. 例如, 当数据校验发现数据 A 跳变, 并将其按照平滑算法计算出合理值 B, 系统会为数据 A 打上原始值的标签, 数据 B 打上处理值的标签, 并置数据 B 可信度高于数据 A. 后续的分析应用如果分析奇异点或跳变现象则会选择使用原始值 A, 如果分析数据变化的合理趋势则会选择置信度较高的处理值 B.

如图 7, 置信度评估的过程是在评估引擎中对采集的数据根据评估规则通过评估引擎的计算, 对数据的标签置位, 以表示数据的质量. 评估规则包括数据校验的各项规则和数据处理过程中对数据的修改规则, 数据标签包括数据溯源标签和数据可信度标签.

通过置信度评估引擎, 根据评估规则库里定义的评估规则对生产数据、资产数据、环境数据、实验等数据进行评估, 由评估结果给原始数据置上数据标签. 这样不仅能够完整的保留原始数据, 还可以通过数据标签给后续数据使用者决定数据的取舍.

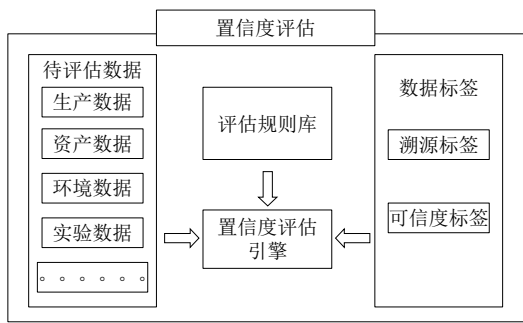


图7 海量数据置信度评估图

3.4 数据交换

数据采集与提交需要解决数据交互过程中交互机制多样化与规范化的矛盾。系统根据不同数据的特点在转发给大数据平台时使用不同的数据交换方案,数据主要通过以下三种形式提供给大数据平台。

(1) Sqoop 脚本

对于结构化的数据,通过类 SQL 语句的 HQL 快速实现映射为一张数据库表,通过编写 Sqoop 脚本实现数据的全量导入或者增量导入。Sqoop 导入的主要是商用关系库中的数据,如设备资产信息,设备缺陷信息,设备试验数据,设备跳闸数据,电压监测数据等。

(2) Kafka

Kafka 是高并发性的分布式消息系统,基于 Kafka 消息订阅的实时数据发布服务按照实时数据区域属性分成多个主题并发传输,在消费端同样采用并发策略,将订阅的实时数据并发写入 HBase。系统可以根据实时数据的变化弹性扩展发布主题的数量和服务数量,从而充分利用 Kafka 和 HBase 并发吞吐量大的特性提高对实时数据存储的响应能力。采用 Kafka 交互的数据主要有实时量测数据、在线监测直采数据、广域向量测量数据、告警事件等。

系统定义了消息服务报文格式,其中,报文头采用 JSON 格式统一定义, body 定义也采用 JSON,可根据类型不同而不同,表 1 为报文头各属性定义及 body 定义举例。

```
{
  "header": {
    "systype": "<<system_type>>",
    "sysname": "<<system_name>>",
    "sysdomain": "<<system_domain>>",
    "msgtype": "<<message_type>>"
  },
  "body": [
    ...
  ]
}
```

属性	定义
header 报文头	systype 来源系统类型 sysname 来源系统名 sysdomain 所属的子控区或主控区 msgtype 消息类型
body 消息体	自定义内容 由各接口类型自行定义

表 1 实时消息服务报文头各属性定义及 body 部分定义

告警消息的 body 消息体定义如表 2,其采用 JSON 格式。

表 2 告警消息体定义

Key	Value 格式	说明
id	long	安管告警 id
time	longlong	告警/事件时间
type	Char[32]	告警的类型
level	Int	等级
object_gid	Char[34]	告警对象通用 ID
object_name	Char[128]	告警对象名称
sub_type	Char[32]	告警子类型
device_gid	Char[34]	故障发生节点的 GID
status	int	告警处理状态
handle_time	Long long	告警处理时间

(3) 文件传输

系统通过 ftp 方式将获取的非结构化或半结构化文件数据传输给大数据平台,存储到 HDFS 分布式文件系统中。这部分数据包括模型文件、录波文件、图像文件、视频文件等。

以上三种数据发布的任务均由系统提供的集群管理软件负责调度,可混合调度三类任务,对于 sqoop 和文件传输任务采用批处理定时启停模式,对于实时数据发布服务采用在线分布式弹性扩展方式调度。

4 系统运行与测试

面向电力大数据的异构数据混合采集系统在用户现场实际部署在 16 台虚拟服务器节点上,具体软硬件配置如表 3 所示,部署图如图 3 所示。

表 3 软硬件配置

配置项	参数
CPU	XeonE7 2.1 GHz 12 核*2 颗
内存	32 GB
操作系统	RedHat EL6.5
第三方类库	ACE5.6、Java JVM1.7、Python2.7

面向电力大数据的异构数据混合采集系统已接入

22 个业务系统和数据源,其中包括 3 个互联网数据源,涵盖电网运行、设备和环境信息等 239 类数据,数据年处理量超过 1 PB。系统具备 7×24 小时连续运行能力,现场实际运行超过一年,未出现因系统故障造成的数据采集中断。图 8 是数据采集实时监视界面。

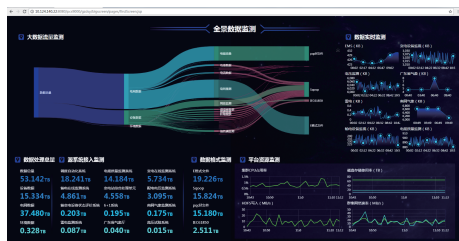


图 8 大数据接入监视界面

针对电力行业实时数据采集场景,系统经过性能压力测试。测试结果表明采集集群单机可接入实时数据点数量为 4 305 576 点,集群整体可实现多机横向扩展千万点以上的实时数据采集;实时数据交换与发布的响应时间平均为 11 毫秒;数据交换和发布服务的准确率为 100%;30 分钟 20 万点实时数据每秒连续变化压力试验丢包率为 0;系统完全能够满足电力行业实时数据采集对系统容量、可靠性和实时性的要求。

5 结论与展望

本文介绍了一套面向电力大数据的异构数据混合采集系统,通过混合数据采集模型和采集集群实现了对异构数据源采集任务的混合调度和管理;通过数据置信度标签技术,在保留原始数据的同时,用合理的方式标示数据的质量,为后续大数据弱关联分析提供了便利;通过 Sqoop、Kafka、文件传输方式将采集与处理后的数据提交给大数据中心。系统已经在用户现场部署并投入使用,运行稳定,效果良好。在接下来的工作中将进一步研究容器技术和微服务框架,增强前端

采集集群和后端数据发布服务弹性扩展灵活性,以便实现更大规模,更多类型数据的采集和处理。

参考文献

- 闫龙川,李雅西,李斌臣,等. 电力大数据面临的机遇与挑战. 电力信息化, 2013, 11(4): 1-4. [doi: 10.3969/j.issn.1672-4844.2013.04.001]
- 岳阳,张晓佳,高一丹. 基于 Hadoop 的电力大数据技术体系研究. 电力与能源, 2015, 36(1): 16-20.
- 陈超. 电力大数据质量评价模型及动态探查技术研究. 现代电子技术, 2014, 37(4): 153-155. [doi: 10.3969/j.issn.1004-373X.2014.04.043]
- 彭小圣,邓迪元,程时杰,等. 面向智能电网应用的电力大数据关键技术. 中国电机工程学报, 2015, 35(3): 503-511.
- 张沛,杨华飞,许元斌. 电力大数据及其在电网公司的应用(英文). 中国电机工程学报, 2014, 34(S1): 85-92.
- 蒋湘涛,贺建飏,李楠. 电力信息采集的通用型通信规约解析系统研究与设计. 电力系统保护与控制, 2012, 40(9): 118-122. [doi: 10.3969/j.issn.1674-3415.2012.09.021]
- 谢大为,杨晓忠. 调度自动化系统中运动技术网络化的实现. 电网技术, 2004, 28(8): 34-37. [doi: 10.3321/j.issn:1000-3673.2004.08.008]
- 王岩,王纯. 一种基于 Kafka 的可靠的 Consumer 的设计方案. 软件, 2016, 37(1): 61-66. [doi: 10.3969/j.issn.1003-6970.2016.01.015]
- Chen CA, Jiang S. Research of the big data platform and the traditional data acquisition and transmission based on sqoop technology. The Open Automation and Control Systems Journal, 2015, 7(1): 1174-1180. [doi: 10.2174/1874444301507011174]
- Wang J, Wang WH, Chen RF. Distributed data streams processing based on flume/kafka/spark. Proceedings of the 3rd International Conference on Mechatronics and Industrial Informatics. 2015. 167-171.
- 王润华,毋建军,侯佳路. 分布式实时计算引擎——Storm 研究. 中国科技信息, 2015, (6): 68-69. [doi: 10.3969/j.issn.1001-8972.2015.06.027]