

基于分级存储环境的数据即时应用产品^①

何翔¹, 张明明², 黄高攀²

¹(国网电力科学研究院, 南京 210003)

²(国网江苏省电力公司, 南京 210024)

摘要: 目前很多大型企业的核心业务采用二级部署模式, 此模式下总部对全公司范围内业务数据进行全面即时统计分析时会存在数据量大、存储分散、需求变化响应不及时以及手工统计填报等问题。基于以上问题, 设计了一种基于分级存储环境的多维数据分析软件。该产品能够让用户从不同视角分析业务数据, 快速形成满足业务工作需要的数据视图, 最后通过报表的形式进行展示。该产品已成功应用于国家电网公司, 具有一定的经济效益和理论价值。

关键词: 分布式; 数据即时应用; 模型构建; 内存计算

Real-Time Data Application Product Based on Hierarchical Storage Environment

HE Xiang¹, ZHANG Ming-Ming², HUANG Gao-Pan²

¹(China State Grid Electric Power Research Institute, Nanjing 210003, China)

²(JiangSu Electric Power Company, Nanjing 210024, China)

Abstract: Currently, many large corporations adopt secondary deployment pattern in their core business. However, when head office analyzes the company-wide data timely and fully in this model, it may cause problems of large amount of data, scattered data storage, not timely responding to the changes in requirements and handiwork statistics. In this paper, based on the above problems, a multidimensional data analysis technology and software tool which can support the end users on-demand self-service data analysis are proposed. The software can analyze data according to different requirements of users, form a data view quickly, and show data through report forms finally. Furthermore, this flexible, efficient, practical and real-time data analysis tool has been successfully applied to State Grid Corporation of China (SGCC), which shows that it has some value in theory and practical applications.

Key words: distributed; real-time application; model building; memory computing

随着“SG186”工程的建设 and “SG-ERP”的深化应用, 信息系统在国家电网公司各业务领域快速拓展, 建成了纵向贯通、横向集成的一体化信息平台, 覆盖总部及省(市)公司的两级数据中心, 奠定了公司信息化管理的基础^[1]。

在信息化建设初期, 考虑到集团化公司在各业务领域信息系统众多, 不同系统差异化程度较大以及管理模式不一致等因素, 公司在部署重要业务系统时, 一般采用二级部署和本地存储的建设模式, 以保证公司的经营管理和生产运行更加稳定。但随着大数据时代的到来, 数据分析已不仅限于高层管理者的决策之

用, 也日益成为企业员工日常操作必需。现实业务中对全公司数据进行综合性分析统计的即时数据应用类需求越来越多。

目前 IT 界没有相对成熟的数据分析应用工具对同类业务、跨多个地域分级存储的数据进行即时应用的能力^[2]。导致集团化公司在处理分级环境的数据时还需要通过系统线下的方式进行下发、填报、收集和汇总离线数据表格。这样不仅给基层业务人员增加了大量的工作, 同时手工填报数据的准确性也得不到保证, 无法满足公司总部对全局数据进行多维度、实时分析的需求^[3]。

① 收稿时间:2015-07-30;收到修改稿时间:2015-10-19

基于业务需要和现实考虑，迫切需要研发出基于分级存储环境的数据即时应用平台，通过分布式计算技术，提升数据分析效率，实现跨地区部署系统业务数据的即时应用。从而改变公司针对此类多级部署、数据分级存储的业务系统的数据分析应用现状，为公司经营决策提供更加准确的分析数据支撑。

1 总体架构设计

从整体来说，本文设计的是一个分布式架构的系统，包括总部侧和分公司侧。总部侧包括 Web 展现交互服务和逻辑服务，分公司侧包含计算代理。具体架构图如图 1 所示。

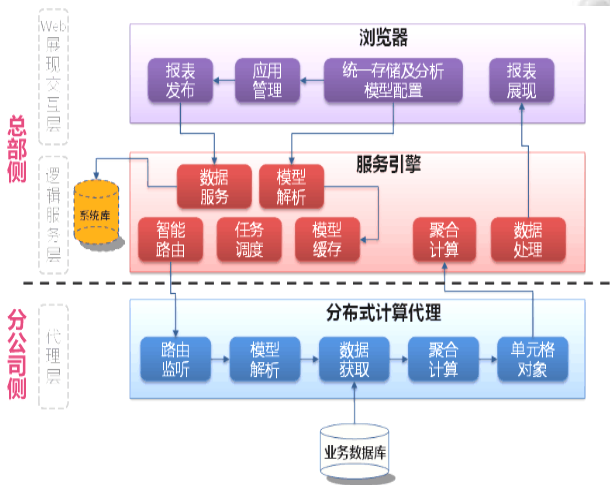


图 1 系统架构图

总部到分公司通过智能路由连接在一起。业务人员在浏览器中配置统一存储及分析模型，配置数据即时应用模型。服务引擎在任务调度的调度下，解析模型，生成路由指令，执行智能路由。计算代理利用分公司的计算资源，是分布式计算的一个节点。当路由到一个计算代理节点时，计算代理解析路由指令，动态拼接 SQL 语句，从分公司数据库提取明细数据，根据路由指令这包含的算法执行计算。计算完毕，把计算结果封装到路由指令中的单元格对象中，然后返回到总部。总部接收到计算代理返回的单元格对象后执行动态聚合计算，生成最终结果，并推送到浏览器展示。

从现实情况分析，该系统应用到生产环境中，必须解决分级存储环境下业务数据数据结构差异性，存储位置不确定性以及快速完成数据即时分析报表构建等技术难题。

2 系统模型设计

本文采用了面向对象的方式来描述不同位置、相同类型的业务数据的系统模型，从而屏蔽跨地区分级部署系统的业务数据个体的差异性，达到对业务数据存储结构和应用方式的抽象，为用户自助式构建数据即时应用场景打下基础。

该系统模型是基于分级存储环境的业务数据统一存储与分析模型。其中统一数据存储模型主要描述各级业务数据的存放位置、基础业务数据类别及数据、不同业务范围数据特征分类、不同业务数据存储时的物理结构、业务数据间的关联性描述、业务数据的全局一致性和各地差异性处理等信息。统一数据分析模型主要描述不同业务数据的分析视角、业务数据计算指标识别、指标关联性识别与描述、数据分布处理与聚合规则模型、数据展现布局模型、数据分析过滤条件模型、图表展现模型、数据展现分析结果持久化模型等。

在该产品的设计和实现过程中，最终展示给用户的是报表，因此将系统模型又分为基础数据模型、报表数据模型、报表展现模型三个部分。其中统一数据存储模型是指基础数据模型，统一数据分析模型分为报表数据模型、报表展现模型两部分。具体如图 2 所示。

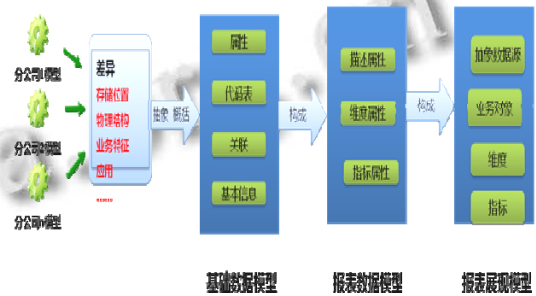


图 2 系统模型图

2.1 基础数据模型

基础数据模型是以面向对象的方式描述各级业务数据的数据源信息、存储结构、基本属性、关联属性、代码表等信息。其中基础属性是描述数据的物理存储结构，主要属性包括模型名称，数据来源，约束条件等。代码表用于存放业务对象用到的编码信息，主要属性包括代码表名称，约束条件，代码表物理字段名称，中文映射名称，数据类型，显示名称等。关联属性

表示基础模型之间的关联关系，主要属性包括关联名称，父类型，子类型，关联方式，父类型属性，子类型属性。基础属性对代码表是一种引用关系，能够自动实现代码转换。通过关联模型，可以实现对多个有关联关系的业务对象的数据即时应用。

基础数据建模主要分成以下三个步骤：

- ① 数据源定义，定义不同数据源连接，可实现总部、网省不同数据源连接；
- ② 对象反向映射，通过选择表模式下某数据表、部分数据表或全部数据表，可拖拽方式反向映射成业务对象模型；
- ③ 对象模型定义，反向映射后的对象基本信息、属性、关联关系进行定义。



图 3 基础数据模型

2.2 报表数据模型

报表数据模型负责描述不同业务数据的分析视角、处理与聚合规则以及过滤条件等，主要包括描述属性、维度属性和指标属性。报表数据对象的主体来自待形成报表的核心业务数据表，同时将与之关联的所有可能用于统计分析的对象属性拉平到同一报表数据对象中。

在多维分析中维度是人们观察数据的角度^[3]。本文的维度模型是公共维度，是指单独的、共享的维度表。它主要有公共维度名称，对应基础数据模型，主键域，显示域，约束条件等属性。在电网业务中，例如电压等级等属性可以在多个报表中应用，就可以定义成公共维度。

指标模型是指一个可以被多个应用共享的指标集^[4]。指标集包含多个指标，以及多个适用维度。指标是由指标名称、指标类型、计算公式(求和、求平均、最大值、最小值、记录数等)等构成。适用维度定义该指标

集包含的维度，可以引用公共维度，也可以把指标集对应的基础数据模型的属性定义成私有维度^[5]。私有维度仅在此指标集中有效，不能复用。



图 4 维度定义

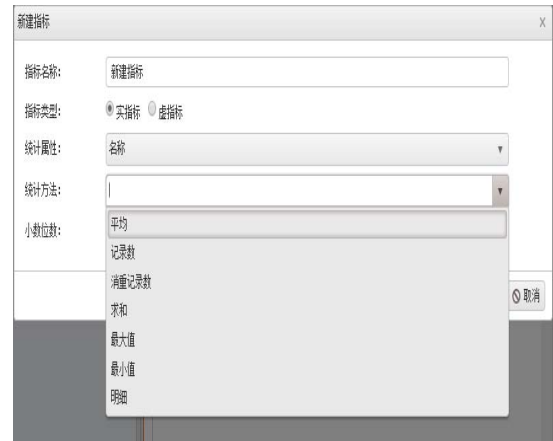


图 5 指标定义

2.3 报表展现模型

报表展现模型是在报表数据模型的基础上进行构建，以拖拽报表数据模型上相关属性为手段，快速定义报表展现方式，最后的结果展示给用户。其构建的步骤如下：

- ① 以拖拽方式设计报表模型，组合维度、指标形成表格，设置行高、列宽等显示样式
- ② 为增强设计阶段交互体验，设计时表格数据可用模拟数据(类似 VS 的表格)，而非即时获取。
- ③ 以报表为单位设置总部本地获取还是从各省市公司获取。
- ④ 维度的值放置于表格中可设置动态或静态取值之分。定义报表时，首先获取静态维度数据，可以拖拽静态维度数据或维度属性

⑤ 每个维度均可设置是否出现“合计”行(列)

⑥ 多个指标属性作为一个整体, 仅能出现在行或列, 且作为整体移动与维度间的属性(最前或最后, 不夹在中间)

⑦ 对于统计表格, 可设置是否可显示明细表及明细表定义. (区分主表是否为统计表的依据主要是是否选择了描述属性)

最后生成如图 6 所示的效果图.

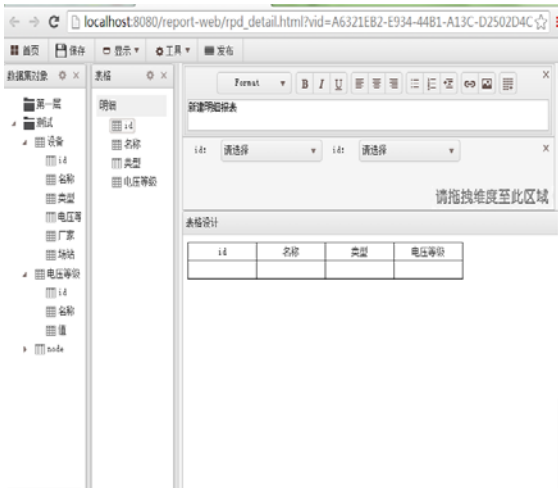


图 6 报表展示模型效果图

3 计算引擎设计

计算引擎是数据即时应用的核心功能模块, 通过计算引擎, 可以把各种模型(如路由模型、维度模型、指标模型等)转换成内存对象, 并把各种对象按照既定逻辑组织起来, 最后把对象化的计算结果转换成前端所需的 JSON 对象输出到前端^[6]. 计算引擎算法工作流程如图 7 所示.

① 加载模型: 根据要展示的数据即时应用编号, 查找是否有数据即时应用模型缓存, 如果没有, 加载该数据即时应用模型.

② 加载维度成员: 维度成员包括静态维度成员和动态维度成员. 静态维度成员定义在数据即时应用模型中, 是不再变化的. 动态维度需要到维度表中查询出维度成员.

③ 生成数据集对象: 根据数据即时应用关联的统一存储及分析模型, 以及约束条件, 生成 SQL 语句、数据源标识以及统一存储及分析模型编号并生成数据集对象.

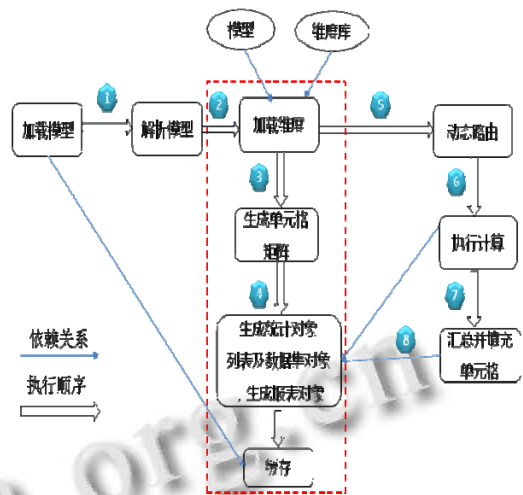


图 7 报表展示模型效果图

④ 生成单元格矩阵: 根据数据即时应用模型中维度的层级关系、维度所在坐标轴以及指标数目, 遍历 XML 节点, 生成单元格对象矩阵. 单元格对象分为三种类型: 表头、维度和度量.

⑤ 生成统计对象: 首先判断是否有“地区”维度, 如果有, 选取第一个“地区”维度对应的度量单元格并生成对象列表; 如果没有, 把所有度量单元格生成对象列表. 对于实指标对象, 每个对象要关联到哪个业务对象哪个属性; 对于虚指标对象, 要关联到哪些实指标对象. 对象列表中对象编号和单元格矩阵中度量单元格对象的编号一致, 便于后期的数据回填(把对象数据回填到度量单元格). 对象列表分为实指标对象列表和虚指标对象列表.

⑥ 智能路由: 根据维度成员信息查找路由表中是否存在记录, 如果有记录, 提取路由路径. 如果没有记录, 广播所有路径.

⑦ 计算代理执行计算: 计算代理首先解析数据集对象, 把明细数据取到内存中. 然后计算每个对象. 对于实指标对象, 通过过滤条件和属性以及算法, 从内存数据集中取得数据并计算. 实指标对象计算完毕后再计算虚指标对象.

⑧ 动态聚合计算并填充单元格: 计算引擎根据路由信息, 判断所有节点数据都返回完毕(对象数据返回到总部采用同步机制), 开始执行动态聚合操作. 聚合时, 如果没有“地区”维度, 累加相同编号对象的值并赋值于相同编号的度量单元格对象. 如果有“地区”维度, 把“地区”添加到对象编号的前面作为三段式编

号的前缀,然后查找相同编号的单元格对象并赋值.聚合计算后计算合计,最终形成报表展示对象.

4 智能路由召唤及响应设计

本文设计的产品主要是针对分级存储环境下业务数据的统计和分析,解决目前数据来源不定、手工填报问题,提升数据分析的即时性和准确性.本文提出了一种针对分级存储环境的基于业务特征的数据智能路由召唤与响应技术.



图8 智能路由形成图

由上图所示,总部和各个分公司是通过智能路由技术连接在一起的.总部侧会部署产品应用服务,主要完成了业务场景构建和报表数据的展示.其中业务场景的构建就是本文中基础模型和报表数据模型的构建过程,而最终展现结果则是通过报表展现模型进行展示.分公司侧会部署计算代理服务,计算代理负责完成各节点的数据计算任务,根据主站发送计算请求实体,完成数据查询并异步发送主站消息队列.在总部和分公司进行数据交互就需要通过智能路由召唤和响应技术.因此解决在分布式环境下数据实时、安全的交互是该产品实现过程中技术难点.

分级环境下的数据即时应用与传统应用模式下的数据分析相比,单个数据分析结果视图(如一张统计表格),甚至单个目标统计结果数据其细粒度的原始基础数据可能来自多个不同地域的数据存储节点,需要即时从分散的多个远端获取数据并聚合.智能路由机制包括路由模型的构建和路由规则的构建.路由模型由地区维度、路由表、路由节点组成.三者之间的关系如图9所示.

路由节点也叫分布式计算节点,是分级存储环境中的一个计算代理节点,通过在路由节点部署计算代理,实现下级数据的本地化高效处理,从全局来看实

现了大量数据的分布式计算.

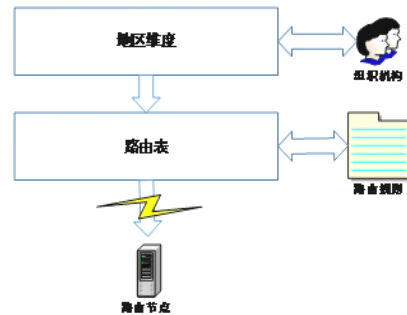


图9 智能路由构成图

路由表由路由规则组成,关于路由规则在下面描述.路由表记录业务数据应用场景和路由节点之间的映射关系.由于路由规则根据业务特征可以动态变化,所以路由表的内容也是动态变化的.路由表主要属性包括数据应用场景名称、路由节点名称、路由节点服务地址等.

地区维度用于创建路由节点和组织机构之间的映射关系,如把组织机构意义上的某省公司和逻辑意义上的该省公司这一计算节点关联起来,业务人员只需要关注组织机构上的该省公司,不需要关注技术领域的计算节点的细节.

通常一个数据即时应用场景会有多个路由规则,每个路由规则根据数据业务特征生成一条维度数据和路由节点之间的关联关系,如1000Kv电压等级的缺陷数据到哪些省公司提取数据,而不是到所有的节点提取数据,因为并不是所有的省公司都有1000Kv的数据.路由规则根据自学习功能自动维护.所谓自学习是指,第一次路由时如果路由表中没有记录,就广播到所有的路由节点,当所有的路由响应回来的时候,就能够知道哪些省公司有数据,哪些省公司没有数据,由此就生成了相应的路由规则,下次就根据路由规则不再广播,以此可以提高整体应用的效率^[7].路由规则有一个定期刷新机制,利用系统空闲时间,会再次广播路由,如果发现新的路由路径或原来的路由路径已经失效,就更新路由表.

由于广域网通信有许多不确定因素,常常会有路由召唤及响应失败现象的发生^[8].分析故障发生的位置,如果发生在总部侧,总部侧会对失败的路由指令保持缓存,并不断尝试重新路由,直到路由成功.如果发生在计算代理侧,计算代理对响应结果保持缓存,

并不断尝试响应总部,直到响应成功。

5 工程应用情况

基于分级存储环境的数据即时应用产品已在国家电网调度领域得到应用。

国家电网总部 OMS 系统中采用就是分级部署和本地存储的模式。目前国家电网总调和华北分局部署了本文设计的数据即时应用产品,江苏省公司和山西省公司作为试点单位部署了计算代理服务。总部用户通过该产品,在线自助设计数据分析模型,配置多种应用场景,随时获取各个分公司数据,满足了多角度、频繁变化的数据应用需求。

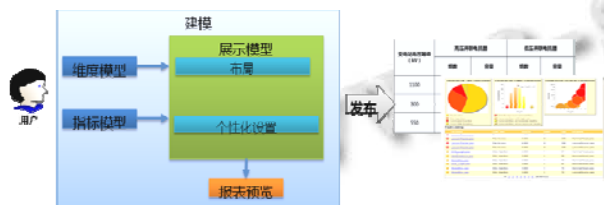


图 10 自助式多维分析报表图

国家电网公司在分级存储环境下,进行全公司数据统计时,信息系统还无法做到完全支撑,通常还需要通过一些系统线下的方式下发、填报、收集和汇总很多离线数据表格,如 OMS 系统通过 Excel 由各地向总部报送数据等方式。本文设计的产品很好的解决了这些问题,并已取代了人工填报的模式。

目前该产品正在向国家电网 PMS 系统、电能质量系统进行推广,提升公司现有生产管理系统、调度数据管理、电力营销、电能质量分析等业务系统的数据即时应用能力,进一步提升此类系统对资产、运行等业务数据分析处理的能力、降低人工处理的周期和成本,在提高数据分析质量的基础上取得良好的经济效益。

6 结语

本文提出并实现了基于分级存储环境的数据即时应用技术,以分布式计算等信息化手段实现跨层级业务数据的即时获取、智能分析与多维自助展现,取代

了传统的以离线、手工为主的多级数据报送、汇总与统计工作模式,大幅提升了业务数据分析工作的执行效率,减轻了基层业务人员的手工统计、填报数据的工作量,同时能够快速响应业务需求的变化。本文实现的产品在大型企业信息化建设中具有积极的作用。由于该产品侧重于数据的抽取、聚合和计算,在报表展现、可视化方面实现的相对粗糙,也是后期需要改进的地方。

参考文献

- 曹占峰,赵强,刘海涛.国家电网公司数据交换监控模块研发及应用.电力信息化,2011,9(7):6-11.
- Mateescu R, Serwe W. Model checking and performance evaluation with CAPP illustrated on shared-memory mutual exclusion protocols. Science of Computer Programming, 2013, 78(7): 843-861.
- 王宇飞,朱伟,刘丹.基于 OSGi 的分布式 Web 应用结构.计算机系统应用,2015,24(8):73-78.
- Kayaaslan E, Cambazoglu BB, Aykanat C. Document replication strategies for geographically distributed web search engines. Information Processing & Management, 2013, 49(1): 51-66.
- 唐云善,缪巍巍.一种对象化并行计算框架.计算机系统应用, 2015,24(7):35-40.
- 李培军,吕立,李喜旺,马存,于喜清.电网业务中的海量数据存储技术.计算机系统应用,2014,23(4):70-76.
- Gui Z, Yang C, Xia J, et al. A performance, semantic and service quality-enhanced distributed search engine for improving geospatial resource discovery. International Journal Of Geographical Information Science, 2013, 27(6): 1109-1132.
- 杨阳,王品,杜少华.NFS 在分布式数控系统中的应用与改进.计算机系统应用,2015,24(6):202-206.