























该平台架构可以在使用 M2M 真实监测数据、不干扰在线系统正常运行的同时,对数据库中大数据进行任意操作,既可以作为主要的生产系统,亦可作为 X 集团大数据的备用平台,便捷安全。

在升级中,针对 X 集团单条数据量小、数据总条数多的特点,作者将传统关系型数据库中的长表存储改为宽表存储,节约存储空间,提高查询速度.M2M 平台利用大数据存储系统实现了海量级别的监测数据存储、查询和分析.经过改造的平台能够支持对 10 万台以上设备的数据进行管理,支持 300 亿以上监测数据的快速存储,支持每秒写入监测数据达到 10 万条以上.在多种工作场景下,表现出优于传统关系数据库的性能。

经测试,在 1 100 亿数据量级下,“大表查询”操作能在毫秒内返回结果.而在原有关系数据库中,需要耗时数秒.“监测数据写入”操作可以达到每秒十数万条监测数据的写入速度,满足当前以及未来长时间内的应用需求.“历史工况纠错”操作在存储百亿条记录的前提下,对错误数据的查询能在分钟级别内返回,对错误数据的修改能在毫秒内完成,而现有关系数据库则无法支持海量数据下的在线查询<sup>[14,16]</sup>.“起重机防锈保养”操作在存储百亿条记录的前提下,能在分钟级别内找出需要防锈保养的车辆信息,而现有关系数据库则无法支持海量数据下的在线查询.对于“故障预警”操作,M2M 平台利用 Map/Reduce 分析系统实现了在海量数据基础上的数据分析挖掘功能,对观察机械设备运转情况、及时发现故障具有重大意义.我们通过对工程机械行业特点的研究,推出了从时间、主机和工况的多维度分割合并方法,并给出了用于评价工况数据的关键指标体系,并完成了在大数据基础上的快速计算过程.在此基础上,设计了敏感工况分析工具、开机偏差热度图以及主机工况时序图等一系列可视化方法,为大数据技术在工业界应用提供了一套可操作的方法.平台通过对 590 亿条电子工况数据、380 亿条报警数据、120 亿条油耗数据和 10 亿多条故障记录进行决策分析,发现了 9 000 多个工程机械运行和操作异常现象,识别出 100 余种异常操作行为和安全隐患;通过基于智能传输协议的企业控制中心,减少 7 500 多名服务保障工程师的出勤次数达 60%以上。

用户表明:通过对工程机械易损件消耗特征进行统计,并利用远程监测技术结合定期维修计划预测关键易损件(例如泵管、切割环、眼镜板、输送缸等)的消耗和备件需求,减少呆滞库存 10%以上,每年可以有效节省呆滞库存约 9 300 万元.通过对 X 产品的历史工况数据进行分析,可帮助各个专业研究院寻找品质问题的产生原因,减少破坏性寿命实验所需投入的成本 20%。

## 5.2 面向天气预报的气象大数据应用示范

通用大数据平台技术是支持专业应用的基础平台,是支撑领域核心业务、优化服务质量的基础.为解决气象领域数据类型多、规模大、访问延迟低、业务逻辑复杂等数据管理难题,我们基于本文的一体化平台开发了专业化、定制化、一体化的气象大数据管理系统,针对气象大数据的特点提供了优化的数据解析、导入、存储和计算内核(如图 5 所示)。

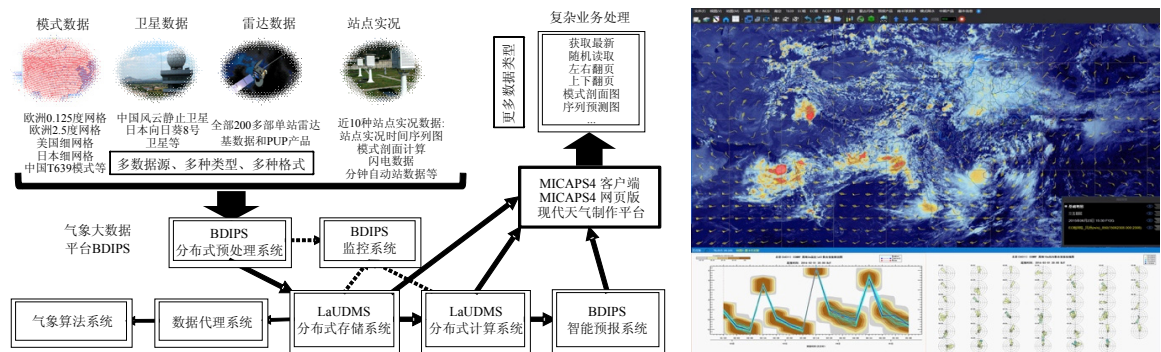


Fig.5 Meteorological big data synthesis processing framework

图 5 气象大数据综合分析处理系统架构

为解决气象预报中的大数据问题,采用跨层优化的先进软件工程理念,研制了气象大数据的近实时数据分

析与处理平台.平台对气象数据的存储组织进行了优化,提升了数据检索效率,同时能够兼容原有存储系统;通过一体化数据处理流程,提升了数据使用效率;通过对标准格式的支持,简化了数据处理流程;通过分布式的数据处理,提升了数据分析速度;通过按地域对数据切片,提升了数据的传输效率;通过分布式存储,提升了数据加载及传输的速度;针对气象数据多维度的特性,进一步演化出多维数据空间模型,进一步提高了系统性能.

针对气象实时数据多类型、高维度、弱模式的特性,我们利用多源异构数据的自由表模型,实现全部海量气象数据从文件系统至非结构化分布式数据库的迁移,采用弹性可扩展服务端架构,应对海量高并发行业及公众用户对于气象实时数据的访问.采用气象大数据分布式实时流式解析技术,实现数据产生即可见的高速加工流水线,同时给出了气象大数据分布式存储解决方案,确保全部数据毫秒级写入与查询.系统服务器端强大的容灾备份能力,高度可靠性保障,完备的系统监控,自动化、智能化的数据处理流程.该系统目前已应用到国家、省、市三级,支撑每日天气预报业务.新系统具有高稳定性、高容错性和良好的可扩展性,支持海量、多类型气象数据的快速解析、导入、存储和访问,数据规模峰值达 16TB/天,数据检索速度达到毫秒级.

## 6 总 结

在大数据时代,软件系统和工程面临的机遇挑战体现在互为依赖的两方面:一方面,软件系统与工程应针对大数据处理的需求,研究如何开发支持大数据处理各个环节的软件技术与系统,形成面向大数据的软件工程——面向大数据生命周期的一体化集成设计开发环境;另一方面,软件系统与工程实施过程中会涉及大量具有大数据特征的系统运行过程数据,因此有必要对这些多维数据进行充分的关联挖掘和机器学习,发现数据驱动的开发和运行规律,形成大数据的软件工程方法学,指导大数据软件系统的开发——面向软件生命周期的大数据应用系统运行分析工具.

大数据应用系统是个万花筒,覆盖数据的采集提取、存储、计算、分析、可视化等大数据全生命周期的多个技术环节,而各个环节都涉及多种解决方案,涉及到的各类系统有几百种之多,这给面向领域的大数据应用系统构建带来了极大的挑战.图灵奖得主 Stonebraker 提出了“one size does not fit all”的理念<sup>[15]</sup>,认为大数据就应该面向特定领域和问题进行定制和优化.然而,这必然导致大数据生态圈的碎片化,特别是在大数据技术从消费互联网向产业互联网渗透过程中,开发人员从复合型极客转换为产业领域型人才,凸显出大数据应用系统选型困难、系统配置难以预设、维护管理代价大等难题.

本文研究了面向大数据生命周期的软件工程方法学和面向软件生命周期的大数据软件系统,特别是两者之间的本质依赖关系,通过对大数据系统的开发与运行的共性模式进行抽象,提出了一体化开发运行平台框架,具体包括两个方面:第一,研制一体化集成设计开发环境,支持高抽象层次的、面向领域的需求描述框架,具有用户友好的可编程性、并能描述复合业务需求的过程模型;第二,研制了大数据应用系统运行分析工具,通过有效处理、分析软件运行生命周期中生成的运行数据,如运行日志、系统配置、部署脚本、性能指标等,从中提取有用信息、做出优化决策,帮助软件开发者以数据驱动方式进行大数据应用软件的运行优化.

## References:

- [1] Jagadish HV, Gehrke J, Labrinidis A, Papakonstantinou Y, Patel JM, Ramakrishnan R, Shahabi C. Big data and its technical challenges. *Communications of the ACM*, 2014,57(7):86–94.
- [2] Beatty J, Wieggers K. Forward thinking for tomorrow's projects requirements for business analytics. Seilevel Whitepaper, 2015.
- [3] Chen H, Chiang RHL, Storey VC. Business intelligence and analytics: from big data to big impact. *MIS Quarterly*, 2012,36(4): 1165–1188.
- [4] Chung L, Nixon BA, Yu E, Mylopoulos J. On non-functional requirements in software engineering. LNCS 5600, Springer-Verlag, 2012. 363–379.
- [5] Computing Community Consortium, Computing Research Association. Challenges and Opportunities with Big Data: A Community White Paper Developed by Leading Researchers Across the United States. White Paper, 2012.

- [6] Hu H, Wen Y, Chua TS, Li X. Toward scalable systems for big data analytics: A technology tutorial. Access, IEEE, 2014,2:652–687.
- [7] Huang XD, Wang JM, Zhong Y, Song SX, Yu PS. Optimizing data partition for scaling out NoSQL cluster. Concurrency and Computation: Practice and Experience, 2015,27(18):5793–5809.
- [8] Kambatla K, Kollias G, Kumar V, Grama A. Trends in big data analytics. Journal of Parallel and Distributed Computing, 2014,74(7): 2561–2573.
- [9] Khouri S, Bellatreche L, Jean S, Ait-Ameur Y. Requirements driven data warehouse design: We can go further. In: Proc. of the Int'l Symp. on Leveraging Applications of Formal Methods, Verification and Validation. Berlin, Heidelberg: Springer-Verlag, 2014. 588–603.
- [10] Long M, Wang J, Sun J, Yu PS. Domain invariant transfer kernel learning. IEEE Trans. on Knowledge and Data Engineering, 2015, 27(6):1519–1532.
- [11] Long M, Wang J, Cao Y, Sun J, Yu PS. Deep learning of transferable representation for scalable domain adaptation. IEEE Trans. on Knowledge and Data Engineering, 2016,28(8):2027–2040.
- [12] Manyika J, Chui M, Brown B, Byers AH. Big data: The next frontier for innovation, competition, and productivity. Report, McKinsey Global Institute. 2011.
- [13] Russom P. Data analytics. TDWI Best Practices Reports, Fourth Quarter. 2011. 1–35.
- [14] Song S, Zhang A, Wang J, Yu PS. Screen: Stream data cleaning under speed constraints. In: Proc. of the 2015 ACM SIGMOD Int'l Conf. on Management of Data. ACM Press, 2015. 827–841.
- [15] Stonebraker M, Cetintemel U, Zdonik S. The 8 requirements of real-time stream processing. ACM SIGMOD Record, 2005,34(4): 42–47.
- [16] Wang J, Song S, Lin X, Zhu X, Pei J. Cleaning structured event logs: A graph repair approach. In: Proc. of the 31st Int'l Conf. on Data Engineering (ICDE). IEEE, 2015. 30–41.
- [17] Zhang QL, Li SL, Li ZH, Xing YJ, Yan Z, Dai YF. CHARM: A cost-efficient multi-cloud data hosting scheme with high availability. IEEE Trans. on Cloud Computing, 2015,3(3):372–386.



王建民(1968—),男,吉林磐石人,博士,教授,博士生导师,CCF 高级会员,主要研究领域为大数据与知识工程,软件工程.