

大数据平台与系统集成架构优化研究

李烁

杭州西湖大数据运营有限公司, 浙江杭州, 310000;

摘要: 随着数据规模的爆炸式增长和企业信息化的深入发展, 传统的信息系统集成架构在处理海量、多源、异构数据方面逐渐暴露出瓶颈。构建高性能、灵活可扩展的大数据平台成为推动数据驱动业务转型的核心路径。本文从系统集成架构优化角度出发, 系统研究了大数据平台的技术体系、关键集成技术及应用优化路径。结合典型行业案例, 本文提出一种基于数据中台与服务网格的协同架构模型, 实现从数据采集、清洗、存储、计算、服务发布到业务闭环的全流程优化, 有效提升了平台资源利用率、数据服务响应效率和跨系统集成能力。研究成果可为政企单位在大数据平台建设、数据治理和业务集成中提供具有实践价值的参考方案。

关键词: 大数据平台; 系统集成; 架构优化; 数据中台; 服务网格; 异构数据处理

DOI: 10.64216/3080-1508.25.02.028

引言

本文聚焦“大数据平台与系统集成架构”的整体优化问题, 立足于平台基础架构演进趋势, 从技术体系、架构模式、集成机制与业务支撑等多角度展开分析, 提出一套可落地、可演进的系统架构优化策略, 涵盖从数据底座到服务治理的全过程, 力图解决当前平台建设中遇到的集成复杂性、计算压力与治理瓶颈等核心难题。

1 大数据平台架构体系的构建与优化基础

1.1 多层次大数据平台架构体系设计

当代的大数据平台往往会运用多层次的模块化设计方式, 这里面涵盖了数据采集层、数据处理层、数据存储层、数据分析层以及数据服务层等多个层面。在每一个层次当中, 都是依据其自身的功能以及性能方面的目标来开展相应的技术选型工作, 并且完成解耦部署操作。就拿采集层来讲, 其是借助 Kafka、F1ume 之类的组件来达成高吞吐数据流的接入任务的。而处理层, 则是运用 Spark、Flink 等分布式的流批一体框架去支撑起高性能的处理工作。存储层是以 HDFS、Hive、ClickHouse 等混合式的架构来构建起冷热分层存储的策略的。再看分析层, 它是把 BI 工具和 AI 算法模型相结合, 以此来实现对于业务的洞察以及预测等相关事宜。至于服务层, 其是通过 RESTful API 或者 GraphQL 来对外暴露服务接口, 进而让业务系统能够实现对其的调用操作。

1.2 异构数据的标准化与语义融合机制

大数据平台所面临的诸多重大挑战当中, 有一项极为关键, 那便是要对那些来源极为广泛且格式又相当复杂的数据展开统一化的处理工作, 并实现语义方

面的融合。为了能够妥善应对数据所存在的异构性这一状况, 平台特意引入了元数据管理体系以及数据标准化引擎, 针对接入的数据来开展诸如字段映射、数据清洗以及类型转换等一系列标准的处理操作。与此同时, 凭借着数据血缘以及语义建模机制, 去构建起一个统一化的数据资产视图, 以此来对数据的跨部门调用以及业务协同给予有力的支持。在此前提之下, 进一步引入数据标签体系以及知识图谱技术, 从而为数据挖掘以及 AI 分析赋予相应的语义方面的支撑。

1.3 平台资源调度与弹性扩展机制设计

在大数据平台不断持续运行期间, 计算资源的调度能力以及系统的弹性扩展性, 这二者无疑是决定该平台具备高可用性以及服务能保持稳定的关键所在。当面对业务呈现出高并发的状况、数据有着高吞吐的情况以及服务呈现高耦合的现实难题之时, 传统那种静态的资源分配方式已然没办法支撑起复杂多样的应用负载所出现的波动情况, 也难以应对动态任务进行切换的需求。本文所构建起来的大数据平台架构, 在针对资源调度这块内容方面, 引入了 YARN 和 Kubernetes 双引擎混合管理这样的模式。其中, YARN 比较擅长针对批量作业去做精细的资源分配工作, 并且在容错管理方面也有着不错的表现; 而 Kubernetes, 则是更契合于容器化微服务的灵活部署操作, 同时在面对相关情况时能够做到快速响应。在此基础之上, 该平台进一步构建了资源池化的策略, 把计算资源划分成诸如实时处理池、批处理池还有交互分析池等不同功能的区域。每一个资源池都会依据服务的优先级情况、并发方面的需求以及资源占用所设定的阈值等来开展独立的配置工作, 而且还能实现动态的伸缩调整。

为了确保系统在面对突发流量场景时,其业务能够保持连续不断地开展,平台精心规划并设计出了一种依托于负载感知的自动扩展策略。当系统监测到任务队列的长度、CPU的使用率或者内存的占用情况,一旦这些指标达到预先所设定的阈值之时,调度器便会即刻触发扩容方面的相关操作,会以动态的方式去创建计算节点,并且让这些新创建的节点能够自动地融入到任务调度网络当中。当各项任务得以完成,又或者是系统的负载出现下降的情况之后,系统能够依据相应的回收机制,逐步地去释放那些处于空闲状态的节点,以此来有效避免出现资源浪费的现象。在整个策略执行的过程之中,平台还将任务再平衡机制融入其中,其目的在于确保那些新扩展出来的节点可以承接一部分正在运行着的任务副本,进而达到减少排队所产生的延时效果,最终实现对计算资源进行高效且合理的再分配。

2 系统集成模式的演进与优化策略

2.1 从 ETL 流程到数据中台的演进路径

传统的 ETL(即 Extract-Transform-Load)模式,以往在数据仓库建设方面确实起到了颇为重要的作用。不过呢,随着企业业务的复杂程度不断攀升,再加上数据来源也越发多样,这 ETL 模式的种种弊端就渐渐凸显出来了。在实际的操作环节当中,ETL 所依靠的是那种强耦合的作业流程,还有静态的数据模型。当碰到要对多系统并发的情况,以及要处理大规模异构数据的整合任务的时候,常常就因为作业链路太过复杂、执行所花费的时间很长、对变更的响应又特别慢等这些状况,进而变成了业务推进过程中的瓶颈所在。而数据中台理念被引入之后,就打破了 ETL 那种以“管道”作为核心的串行加工逻辑,转而去构建起了一个以“数据资产服务化”当作目标的统一数据平台。处在数据中台的架构之下,平台会在采集到原始数据之后,马上就对其展开标准化的处理操作,并且使其沉淀成为主题数据域。

然后,通过服务目录、API 调用、数据标签以及权限机制等等这样的一些方式,来支持不同的业务系统按照自身的需求去进行调用、去做组合以及相互之间的交互,如此一来,就明显地提高了数据流转的效率,同时数据的可复用能力也得到了提升。还有一点很关键的是,数据中台在和 AI 建模平台相互融合之后,不但能够对数据可视化分析以及实时监控起到支撑的作用,而且还可以凭借模型训练、预测以及反馈闭环等这些操作,去推动智能决策的达成。在这样的一种模式之中,数据可就不再仅仅是那种单纯用来“存取与加工”的对象了,而是成了业务流、决策流当中能够起到动态驱动作用的因子,这样就从根本上

实现了从“数据加工流水线”到“数据智能枢纽”的转变。

2.2 跨平台系统集成中的接口适配机制

当下,企业所面临的系统集成相关任务,常常得要连通几十甚至多达上百个异构信息系统。在这些系统彼此之间,存在着诸如技术栈有差异、通信协议无法达成一致、数据结构难以实现统一等诸多问题,这对集成的效率以及数据的一致性都产生了极为严重的影响。在这样的一种背景情形之下,平台务必要去构建起一种接口集成机制,该机制得具备高度的兼容性以及动态适配方面的能力。通过把 API 网关引入进来,让其充当统一的入口,如此一来,平台既能保证其安全性,又能够达成对服务流量加以控制、进行路由转发以及实现负载均衡等目的。再结合数据适配器与协议转换模块,系统首先可在内部将数据标准予以统一,之后再依据需求把外部接口数据转换成为平台能够识别的格式,进而实现像 JSON、XML、Thrift、SOAP 等多种协议的无缝适配。

对于那些老旧的系统或者是已经处于封闭状态的系统而言,借助于“接口镜像”技术来构建起一个虚拟交互层,这样一来,这些系统便无需经历大规模的重构操作,就能达成现代化集成的目标。与此同时,服务目录体系能够清楚明晰地将每个接口的功能、数据结构、调用方式以及版本管理等内容都标示出来,再配合上身份验证以及权限控制机制,以此来确保接口调用既安全又具备可审计性以及可追溯性。这样的一种适配机制,在促使系统互联效率得以提升的同时,也为接口治理以及变更管理构筑起了一个稳定的基础,从而构建起了关于系统集成平台化、服务化的技术闭环。

2.3 服务网格在微服务架构中的集成优势

随着微服务架构在企业信息系统里面得到广泛运用,服务的数量一下子激增起来,调用链也变得复杂了,而且依赖关系还很难去追踪,这些问题就渐渐暴露出来,如此一来,服务治理也就成了系统集成当中极为关键的一个挑战。服务网格技术,其实就是为了应对这样的复杂性而专门设计出来的,它通过引入特别设置的“数据平面”以及“控制平面”,对服务之间的通信展开抽象化的管理操作。这么做呀,一方面解耦了服务开发和运维的逻辑关系,另一方面还让服务的可观测性以及自治能力都得到了提升。在构建大数据平台的时候,像 Istio、Linkerd 这样的服务网格,在微服务的治理层面被广泛地应用起来,它们能够支持流量调度、熔断限流、检查以及可追踪调用等这些非常重要的核心功能。

服务网格还能够支持那种细粒度的安全策略,就

好比 mTLS 加密通信、服务身份认证以及访问控制等方面，这就使得跨服务通信的安全等级得到了提升。在集成的场景之中，服务网格还会提供面向服务的“入口网关”和“出口控制”，以此来实现内外部数据流的统一接入以及分发，进而为不同系统之间的集成给出一条可控可观测的通信路径。通过服务网格所进行的这种统一管理，平台不但降低了运维的复杂程度以及部署的风险，而且还明显地增强了系统的容灾能力和业务的弹性，从而为构建那种高可靠、高可用的大数据系统集成架构筑牢了十分坚实的基础。

3 架构优化的应用实践与综合效益分析

3.1 电信行业数据平台优化案例分析

在数字化转型压力不断加剧这般背景之下，某省级电信运营商于其原有的数据平台架构当中暴露出了为数不少的各类问题。比如说，数据更新的周期显得颇为漫长，接口响应起来总是迟滞不前，而且系统之间接口冲突更是频繁发生，另外还很难满足在多元化业务场景里面对实时数据处理以及高频调用所提出的相关需求。面对这样一系列的种种挑战，该运营商着手开展了全面的数据平台重构方面的相关工作，其核心要点主要涵盖三个方面的内容。其一，着手建立起统一的数据中台，以此来达成对运营、客服、计费、网络等诸多系统数据加以整合与施行治理的目的，进而打破信息孤岛的这种局面。其二，引入服务网格以及微服务架构，把原有的那种粗粒度服务重新构建成为灵活且能够进行调节的轻量化服务单元，从而促使系统弹性以及调用效率得以提升。其三，全面推进容器化部署工作，借助 Kubernetes 编排来实现服务的弹性扩展以及具备自动修复的相应能力。

在架构完成升级之后，平台成功达成了将近 200 个系统的统一接入这一成果。就接口调用方面而言，其平均响应时间大幅下降，下降幅度超过了 70%。与此同时，数据共享率也有了明显的提升，已经提升至 82% 的程度。再者，运维效率方面也获得了很不错的提升，而系统故障处理时长更是缩短到了原来时长的四分之一。这样的一系列优化举措，一方面强化了企业数据中台释放价值的的能力，另一方面也极为显著地提升了业务系统的运行效率，并且还让用户满意度得到了提升，从而使其成为了电信行业信息架构优化方面极具代表性的典型案例。

3.2 智慧城市场景下的集成应用模式

智慧城市的建设涉及城市运行的诸多重要领域，像交通方面的管控工作、能源的调度安排、环保的监测事宜以及城市的安防事务等等，在其建设过程中，大量信息系统之间非得达成深度的互联状态以及实

现实时的协同效果才行。就在这样的一种背景之下，某地所打造的智慧城市大数据平台运用了融合式的集成架构，把大数据中台、城市信息模型（CIM）还有物联网边缘网关都整合到了一起，由此搭建起一个能对数据进行统一感知且能针对业务做出响应的框架。借助中台来对源自交通信号、公共摄像头、气象设备、车辆 GPS 等多种不同来源的数据展开统一的标准化处理操作，并且在 CIM 平台里完成空间坐标的映射处理以及语义的融合工作，最终能够对路网运行预测、突发事件智能调度、区域能源负荷调节等这些核心场景给予有力的支撑。

3.3 架构优化的经济效益

从投资以及运维方面来讲，系统架构予以优化的话，可不单单是让性能获得提升、促使技术实现升级这么简单，它还能够更为直接地转变成企业实实在在的经济效益，并且达成资源节约的效果。在诸多行业的实际运用当中，经过优化之后的集成架构，借助容器化部署的方式、依靠资源池化调度的手段以及通过对微服务进行拆解等举措，在很大程度上削减了物理服务器的数量，同时也让软件授权成本得以大幅降低，如此一来便使得初期硬件方面的投资有所减少。在多业务同时并行运行的场景之下，平台的模块复用率出现了颇为显著的提升情况，其平均开发周期更是缩短了大约 30%，这就极大地提高了新业务上线的效率，也让系统具备了更强的灵活性。与此同时，因为服务都被统一封装成为了标准的 API 接口，所以接口复用率提升的幅度超过了 40%，进而减少了重复开发以及维护所需要花费的成本。

4 结语

大数据平台的构建与系统集成架构的优化，已成为数字化转型的重要基石。本文在分析技术演进与架构实践的基础上，提出了一套以中台化理念为核心，融合服务网格与智能调度技术的优化路径。实践表明，系统架构的合理升级不仅提升了数据处理与服务能力，也为企业带来了显著的运营收益与业务弹性。未来，随着 AI 技术、边缘计算与多云部署的发展，大数据平台的架构优化将更趋智能化、自治化，为数字生态体系的持续演化提供持续动力。

参考文献

- [1] 马晓勇, 张娟. 人工智能赋能玻纤行业质检创新研究和实践[J]. 江苏通信, 2025, 41 (01): 97-101.
- [2] 李小龙. 铜阳极板机器视觉质检系统应用实践[J]. 湖南有色金属, 2025, 41 (01): 42-45.
- [3] 陈桃. AI 视觉赋能工业质检的研究和实践[J]. 江苏通信, 2024, 40 (04): 91-95.