算力网络应用平台研究与设计

许俊东¹ 李兆滨¹ 宋德华¹ 田 韧¹ XU Jundong LI Zhaobin SONG Dehua TIAN Ren

摘要

针对算力网络 (CPN) 因算力多样泛在、多要素融合、需求差异化的特点导致的用户使用困难、应用管理难等问题,从纳管与接入、决策与调度、应用与开放三个层面研究算力网络应用的管理以及异构算力融合和跨域编排关键技术,并设计一个算力网络应用平台。通过提供一种即开即用、按需付费的零感知算网应用服务,使用户关注于应用服务需求,无需关心算力资源需求和算力网络复杂环境,从而实现降低用户的算力网络使用门槛,提升算网应用的使用和管理效率的目标。

关键词

算力网络应用平台;零感知算网应用服务;应用管理;异构算力融合;跨域算力编排

doi: 10.3969/j.issn.1672-9528.2024.02.006

0 引言

算力网络是我国信息通信领域历经多年积累提出的原创性概念,其基本内涵是一种根据业务需求,在云、网、边之间按需分配和灵活调度计算资源、存储资源以及网络资源的新型信息基础设施^[1]。算力网络是云网融合发展的新阶段,其核心内容是通过高速网络连接整合云、边、端、智算、超算等多层次的算力资源,同时结合 SDN、云原生、大数据、AI、区块链、零信任、隐私计算等新兴技术,提供一个比传统云计算更泛在、更灵活的算力基础设施,主要有以下特点。

- (1) 算力更加多样化。在算力种类上包括通用算力、智能算力、超算算力以及量子算力等,其中通用算力是基于CPU 芯片的服务器提供的计算能力;智能算力是基于GPU、GPFA、ASIC 等芯片的加速计算平台提供人工智能训练和推理的计算能力;超算算力是基于超级计算机等高性能计算集群所提供的计算能力^[2];量子算力是基于量子计算机所提供的利用量子态的性质来执行计算的能力。在算力分布上包括跨地域的各级算力节点,如中心节点、边缘节点以及智能终端设备,当前主流的云、边、端即构成了多层级的分布式算力架构。
- (2) 多要素深度融合。在算力网络中除算力外还包括 其他要素,如网络、大数据、AI、区块链、隐私计算等。其 中网络是算力网络的根基要素,通过无处不在的网络来实现 对算力、数据、用户的连接,特别是随着5G网络、专线/专网、 全光网络、确定性网络的发展,对于无线传输和跨地域的长 距离传输提供了保障;大数据和AI则是算力网络向自动化、 智能化发展的关键要素,通过数据分析、深度/强化学习、

数字孪生等技术提供一个更加智能高效的算力网络环境;区 块链、隐私计算、零信任等技术则是保障算力网络安全可信 的关键要素。

(3)需求差异化突出。在算力网络中的资源环境具有跨地域、多层级、多种类的特点,传统以资源为中心的"资源式"服务模式已不适用算力网络,用户也很难在复杂多样的算力间进行资源选择。服务模式需要从"资源式"向"任务式"转变,为用户提供智能、极简、无感的算网服务^[3]。任务式服务模式通过屏蔽底层算力资源的复杂细节,使其专注于业务服务本身,然而各行各业的需求差异较大,如对同一应用可能面临着云化部署、私有部署、混合部署、跨域部署等不同的部署要求。

在政策层面,为推动国家数字经济的发展,国家陆续出台具备战略性、引领性和创新性的系列政策和举措,加快构建以算力和网络为核心的新型基础设施体系。2021年5月,国家发展和改革委员会、中央网信办等四部委联合印发了《全国一体化大数据中心协同创新体系算力枢纽实施方案》^[4],明确提出了建设八大国家枢纽布局方案,推动算力、网络、能源、数据、应用的一体化协同发展与创新。2022年2月,"东数西算"工程正式启动^[5],通过构建数据中心、云计算、大数据一体化的新型算力网络体系,将东部算力需求有序引导到西部,优化数据中心建设布局,促进东西部协同联动。自此,各企业和组织积极响应国家战略加入算力网络的建设中,促进了算力网络的快速发展。

在应用层面,当前以5G、人工智能、大数据、云计算等为代表的新一代信息技术飞速发展,与各个行业深度融合。各企业已经认识到数字化转型的重要性,逐渐将数字化能力融入到生产经营中,提升产能和效率。然而数字化转型的成

^{1.} 浪潮通信信息系统有限公司 山东济南 250101

熟度在企业中差异较大,大型企业不论从顶层规划还是从数字化投入都远远领先于中小型企业。很多企业的数字化转型升级之路还很漫长,一方面是算力基础设施缺失,小型企业一般无法承担算力基础设施的建设和管理成本,另一方面是数字化方案缺失,市场上虽然有各种各样的应用软件和解决方案,但很多小型企业无法获得满足自身需求的最佳方案。

基于此,本文通过算力网络应用平台提供一种即开即用、按需付费的零感知算网应用服务,使用户关注应用服务需求,从而降低用户的算力网络使用门槛,助力企业数字化建设。

1 关键技术分析

算力网络应用平台包括纳管、调度和应用等环节,在关键技术方案上本文从纳管与接入、决策与调度、应用与开放等三个方面进行研究分析,如图 1 所示。



图 1 总体技术架构

1.1 纳管与接入层

纳管与接入层采用分类标准化对各类异构算力资源进行接入适配、统一纳管,然后采用融合编排的方式实现基础算力环境的融合,对云、边、端、智算和超算算力资源的统一编排,为跨域编排和协同调度提供基础。

在分类标准化方面,根据不同的算力节点类型抽象统一的标准模型,形成对一类算力资源的标准规范,包括统一资源对象、统一数据模型、统一服务接口等,如表 1 所示。同时因各类算力在异构特点和服务开放接口上存在差异,采用插件化架构,在满足标准接口规范的前提下,支持对不同的纳管算力类型进行快速组装和扩展,避免由于大量的定制导致系统架构的腐化。

| 序号 | 类型 | 资源对象 |
|----|----|-----------------------------------|
| 1 | 云 | 云主机、云硬盘、VPC、子网、安全组、弹性网卡、 快照、备份 |
| 2 | 超算 | 集群、节点、存储、队列、作业、软件 |
| 3 | 容器 | 应用、服务、容器、部署、Pod、任务、存储、镜像 |
| 4 | 智算 | 集群、节点、文件、环境、任务、镜像、模型 |

表 1 统一资源对象

在容器融合编排方面,基于 Kubernetes、K3s、Singularity 在基础算力环境上构建统一的算力融合层,将容器能力延伸到边缘侧、端侧、智算和超算侧,实现对云、边、端、智算和超算算力资源的容器化编排管理能力。Kubernetes 是一个开源的容器编排引擎,用来对容器化应用进行自动化部署、扩缩和管理,并可通过插件方式支持对 GPU 的调度 ^[6],在算力网络应用平台中主要负责对云和智算算力的编排管理。 K3s 是经 CNCF 一致性认证的轻量级 Kubernetes 发行版,针对边缘计算、物联网等场景进行了高度优化,并支持 ARM 架构 ^[7],在算力网络应用平台中负责对边缘节点和端计算节点进行管理。Singularity 专门用于高性能计算场景的容器技术,基于可移植性进行虚拟化,更加轻量级,部署更快,目前主要被应用在高性能计算中心 ^[8],在算力网络应用平台中主要负责对超算算力的编排管理。

在容器融合架构上通过异构容器的集成适配形成多容器 集群的管理模式,如图 2 所示。API 服务主要对外开放接口, 如容器创建、资源绑定、配置变更等,供决策与调度层和内 部组件使用。调度器主要根据策略约束和可用资源来确定满 足要求的可用集群,并对每个可用集群进行打分排序(如健 康度、规格、利用率等),将资源绑定到最合适的集群,进 行跨集群的调度。控制器主要监控集群中各个容器对象,并 通过 API 服务与 Agent 交互,在对应的集群中创建和管理容 器对象。Agent 负责管理分布的集群,将本地集群注册到系 统中进行集群的全生命周期管理,同时将本地集群及其资源 的运行状态上报到系统,并接收 API 服务的请求进行本地集 群资源对象的管理。Agent 主要面向不同的容器编排技术实 现差异化的适配,并形成标准接口向 API 服务注册能力以供 控制器和调度器进行调用。

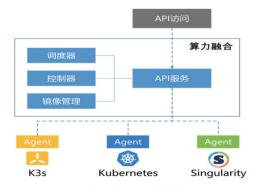


图 2 容器融合架构

1.2 决策与调度层

决策与调度层分为需求调度和动态调度两方面。其中, 需求调度主要对用户需求进行识别和调度,动态调度主要通过感知运行状态进行实例运行优化。

(1) 需求调度由用户触发,对用户需求进行识别解析生

成应用需求、算力需求、网络需求、安全需求等约束,采用实 训模型进行多要素的分析决策,在已纳管的跨区域、跨层级的 算力节点之间选择合适的资源进行调度和任务分配,充分利用 算力多样化优势提高应用计算效率和保证业务安全。

首先接收用户提出的应用需求,一般包括使用模式(本地应用、云应用、融合应用)、性能要求(时延、效率、SLA)等。其中,本地应用是指将所有应用包部署到本地可管控的算力节点上,具有高安全、高可控的特点;云应用是通过多租户的方式直接开通云上部署的应用实例,具有成本低、方便灵活的特点;融合应用是通过分析应用的不同模块算力和性能要求,并结合安全限制将应用的各个模块分布式部署到算力网络中,具有高性能、高可靠的特点。

然后通过需求解析将用户需求转换为具体的需求描述, 如应用解构方法、部署模式、算力要求、网络时延要求、可 用性要求、模块安全要求、数据安全要求等条件。

最后在需求解析完成后,结合各算力节点的资源状态指标和相关决策模型(成本优先、性能优先、安全优先、资源优先等),通过多要素决策算法为用户匹配最佳的算力方案,通过接入层 API 服务完成算力请求分发,实现在全域范围内的算力需求指标和算力资源指标的优化调度。

(2) 动态调度由感知事件触发,通过实时感知算力资源和应用实例的运行数据,根据规则和模型的分析识别生成感知事件,以修复隐患、平衡功耗、最大化效益为优化目标进行各个实例的动态调整,保证算力网络应用平台的健康度。

首先感知功能获取各算力节点、资源、平台、应用的运行数据作为原始数据,如资源利用率、资源容量、实例副本数、业务请求量、CPU利用率、磁盘利用率、网络带宽、上下行速率、网络时延、网络抖动等监控指标。监控指标可以分为服务层面和资源层面,服务层面关注服务质量和用户体验,资源层面关注资源的运行情况。如Google的"4个黄金指标"^[9]中的延迟、通讯量、错误、饱和度四个指标,可以在服务级别帮助衡量终端用户体验、服务中断、业务影响等层面的问题。RED方法^[10]中的速率、错误、耗时三个指标,适合于云原生应用以及微服务架构应用的监控和度量,可以有效帮助用户衡量云原生以及微服务应用下的用户体验问题。USE方法^[11]中的使用率、饱和度、错误,可用于分析系统资源问题,指导用户快速识别资源瓶颈以及错误。

然后根据原始感知数据和 RED/USE 等方法进行事件的分析和定义,这主要有两种方式,一种是通过门限的方式设置静态或动态的阈值,如 CPU 利用率在大于某个阈值时定义为 CPU 利用率过高,在小于某个阈值时定义为 CPU 利用率不足;另一种方式是通过构建智能体,采用机器学习、深度学习等算法,通过分析感知数据,发现数据间的时空关联关系,发现隐含事件和事件预测。

最后在生成事件的同时,根据事件分析结果选择合适的调度方案对当前实例的运行参数进行调整优化,以修复隐患、

平衡功耗、最大化效益为目标,保证算力网络的运行最优。

1.3 应用与开放层

应用与开放层面向最终用户提供应用平台服务和算网开放服务,为用户打造一个零感知的算网服务体验。

应用平台服务对应用全生命周期进行管理,可分为应用发布、应用审核、按需订购、应用开通、应用变更、应用下线等环节。其中,应用发布提供对应用基本信息创建和程序包管理功能。应用审核对发布提交的内容进行检查,并在隔离的沙箱环境中进行模拟部署测试,确保应用的程序包、镜像和配置等方面的安全性、完整性、可用性。按需订购是用户可以根据自身需求选择订购的应用并确认应用使用模式、性能要求等内容。应用开通对用户需求进行应用的自定义构建部署和账号的开通,并在使用完毕后进行退订操作,完成即开即用的目标。应用变更包括信息变更和应用升级,信息变更是对应用的基本信息进行维护,应用升级是指对应用程序包和镜像升级的管理。应用下线是指应用从应用商城下线,不再提供对外服务的功能,已订购的还可以继续使用,并在一段时间后完全销毁(默认一个月),完全销毁时已订购的服务同时关闭。

算网开放服务将系统能力通过标准的 OpenAPI 开放服务接口进行开放,支持与其他平台或系统进行集成,释放算力网络的开放共享、泛在调度的价值,助力数字化转型升级。

2 应用平台功能架构

综合上述分析,在应用平台功能架构中分别在接入、调度、应用等三层进行算力网络应用平台的分析设计,如图 3 所示。



图 3 算力网络应用平台功能架构

2.1 接入层

接入层负责云边端等异构算力的接入和容器化融合,打造统一的融合基础设施,包括算力接入和算力融合两部分。

算力接入通过算力注册功能,用户可以将本地自有数据中心、边缘服务器、公有云、智算中心等作为算力节点注册 到系统中,其中公有云、智算、超算等算力可以以租户方式 进行注册,将租户的算力资源作为一个算力节点接入。算力 适配主要对不同的算力节点按分类标准化规则进行适配,形成可调度的标准算力资源和接口。资源管理负责接收上层调度请求对算力资源进行开通、变更、删除等操作,这里的算力资源主要是指虚拟机、云主机、GPU服务器或超算节点等,作为容器的宿主机。运行监控主要是对已开通算力资源的运行状态进行监控,包括性能指标和告警事件。算力开放是指对外提供的OpenAPI接口,可以被上层模块调用。

算力融合基于算力类型和级别将不同的算力资源加入不同的容器集群中进行统一编排管理。API服务负责对外暴露开放接口,供决策与调度层和内部组件使用。调度器负责根据策略约束和可用资源来确定满足要求的可用集群,将资源绑定到最合适的集群,进行跨集群的调度。控制器负责监控集群中各个容器对象,并控制在对应的集群中创建和管理容器对象。元数据管理主要对各个容器集群上报的元数据信息进行管理和维护。镜像管理负责对全局镜像和应用镜像的管理,并分发到各个集群镜像环境。

2.2 调度层

调度层负责解析应用层需求同时感知运行状态进行资源 和应用的优化调度,保障业务监控运行,包括感知中心、决 策中心和调度中心三部分。

感知中心是对算力节点、容器集群、应用服务等对象的 运行状态数据进行实时监控分析并生成相关感知事件通知调 度中心处理,包括感知接入、数据管理、感知分析、感知事 件和数据共享。

决策中心根据感知共享的数据通过机器学习、强化学习 等手段进行建模并训练,形成智能化的模型和推荐方案,基 于优化目标、用户需求为调度中心和用户提供最优的推荐参 数和配置,包括模型管理、知识管理、训练中心、业务规则等。

调度中心则主要监听用户需求和感知事件,根据调度策略和决策中心的推荐参数进行任务的调度执行,实现算力节点、容器集群、应用服务等对象的优化运行,包括调度方案、需求识别、调度任务、调度策略等。

2.3 应用层

应用层负责为最终用户提供一站式的应用即开即用和生 命周期管理服务能力,同时提供与其他系统集成的能力,包 括应用商城、应用管理、能力开放三部分。

应用商城提供算网应用的订购、使用和退订等功能。用户可以在应用商城中选择需要的应用进行订购,订购时可以指定应用的部署方式和性能等要求,调度层则根据用户需求推荐最优的使用方案供用户确认。

应用管理主要对应用的全生命周期进行管理,包括发布、审核、开通、变更、下线等。

能力开放通过将系统能力进行封装和开放,提供能力目录、能力订阅等功能,方便外部系统的集成访问。

3 总结

算力网络应用平台通过以应用服务为中心的方式,使最终用户优先关注需求和应用,屏蔽了底层的算力网络细节,解决用户对算力资源进行选择、维护和评估困难的问题,降低了算力网络的使用门槛。同时基于容器融合编排技术实现在各类算力节点间的跨域融合调度能力,通过分布式的应用部署实现应用模块的优化分布,提高了业务性能和安全性。但是该平台仍然存在一定的问题和局限性,在智能化、服务体验、安全等方面需要进一步提高,同时针对大规模场景还需要进一步研究和验证。

参考文献:

- [1] 雷波, 刘增义, 王旭亮, 等. 基于云、网、边融合的边缘计算新方案: 算力网络[J]. 电信科学, 2019,35(9):50-57.
- [2] 中国信息通信研究院. 中国算力发展指数白皮书 [R/OL]. (2023-09-01)[2023-09-15].http://www.caict.ac.cn/kxyj/qwfb/bps/202309/P020230914584614752938.pdf.
- [3] 中国移动通信集团有限公司. 中国移动算力网络白皮书 [R/OL]. (2021-11-03)[2023-06-21].https://cmri.chinamobile.com/wp-content/uploads/2021/11/ 算力网络白皮书.pdf.
- [4] 中华人民共和国国家发展和改革委员会. 关于印发《全国一体化大数据中心协同创新体系算力枢纽实施方案》的通知[EB/OL].(2022-01-12)[2023-05-24].https://www.ndrc.gov.cn/xwdt/ztzl/dsxs/zcwj2/202201/t20220112_1311853.htm-1?code=&state=123.
- [5] 中华人民共和国国家发展和改革委员会. 东数西算 [EB/OL]. [2023-01-09].https://www.ndrc.gov.cn/xwdt/ztzl/dsxs/.
- [6] Kubernetes[EB/OL]. [2023-06-09].https://kubernetes.io.
- [7] K3s[EB/OL]. [2023-06-19].https://www.rancher.cn/k3s/.
- [8] Singularity[EB/OL]. [2023-07-01].https://sylabs.io/singularity/.
- [9] NIALL R M, BETSY B, JENNIFER P, et al. Site Reliability Engineering: How Google Runs Production Systems[EB/OL]. [2023-06-30].https://research.google/pubs/pub45305/.
- [10] GREGG B. The USE Method[EB/OL].[2023-07-05]. http://www.brendangregg.com/usemethod.html.
- [11] WILKIE T. The RED Method: key metrics for microservices architecture[EB/OL].[2023-07-06]. https://www.weave. works/blog/the-red-method-key-metrics-for-microservices-architecture/.

【作者简介】

许俊东(1984—), 男, 硕士, 研究员, 高级工程师, 研究方向:分布式计算与系统、算力网络、软件工程。

(收稿日期: 2023-12-01)