

文章编号: 1007-1423(2022)19-0091-04

DOI: 10.3969/j.issn.1007-1423.2022.19.016

基于“碳中和”的高性能计算集群组网建设方法探析

陈 阳, 陈坚泽

(广东液冷时代科技有限公司, 佛山 528000)

摘要: 围绕国家人工智能战略布局和发展需求, 基于“碳中和”政策目标, 介绍了高性能计算服务器集群建设方案, 包括高性能计算集群、节能建设、分布式存储集群、高性能计算网络等方案建设。分析高性能计算服务器集群管理需求, 并探析了高性能计算服务器集群管理方法。

关键词: 高性能计算; 服务器集群; 建设与管理

基金项目: 佛山市高新区高技术产业化创业团队专项(2020197000110); 佛山市“蓝海人才计划”创新创业团队项目(2030032000083)

0 引言

2021 年 3 月 13 日, 新华社公布了《中华人民共和国国民经济和社会发展第十四个五年规划和 2035 年远景目标纲要》(以下简称《纲要》), 《纲要》指出以人工智能为代表的新一代信息技术, 将成为我国“十四五”期间推动经济高质量发展、建设创新型国家, 实现新型工业化、信息化、城镇化和农业现代化的重要技术保障和核心驱动力之一^[1]。围绕国家人工智能战略布局和产业发展需求, 各地政府积极建设人工智能技术支撑平台, 如人工智能基础研究^[2]、智慧交通^[3]、区块链金融、物联网技术^[4]等, 同时启动智慧城市大数据平台^[5]等建设, 为人工智能研究和应用提供健全完善的基础平台服务。这些平台建设对基础计算算力提出了更高的要求, 亟需建设一个高性能计算服务器集群公共服务平台, 为人工智能技术支撑平台以及人工智能科学研究提供高性能、高通量的算力科研保障。同时, 高计算平台也意味着需要更高的能耗才能支撑高算力, 在国家“碳中和”的政策下, 如何利用新型技术, 解决高性能数据中心所带来的高能耗问题也是建立高性能数据中心的潜在研究问题。

本文以华南某科研机构建立高性能计算服

务器集群建设项目为例, 重点介绍高性能计算服务器集群建设中的高性能计算集群、节能建设、分布式存储集群、高性能计算网络等方案建设内容, 并讨论高性能计算服务器集群管理需求以及相应管理方法。

1 高性能计算服务器集群建设

华南某科研机构在前期投资建设项目中已基本完成实验室数据中心基础条件建设, 包括完成基础机房环境、电气系统、基础综合布线铺设等方面的建设, 现针对人工智能关键技术攻关的算力需求, 搭建十个小型计算集群, 每个计算小型集群配置算力 1PFLOPS 计算资源, 构建一个柔性的高性能计算集群, 同时采用液冷散热技术提升集群服务器散热问题, 减少空调机组装配, 实现“碳中和”节能目标。

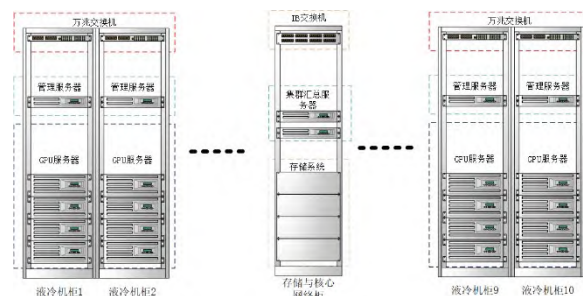


图 1 高性能计算服务器集群建设机柜置放示意图

下面从高性能计算集群、节能建设、分布式存储集群、高性能计算网络等方面阐述高性能计算服务器集群建设内容。

1.1 高性能计算集群方案

高性能计算集群建设基于高性能计算(High Performance Computing, HPC)技术基础构建计算集群平台,其中十台管理服务器,分别用于管理十个集群,同时配备两台管理服务器,用于做统一集群管理。高性能计算集群总体业务架构分为以下四层递进建设。

(1) 基础设施层:利用已建成的模块,方便快捷组装,本项目在现有机房基础设施上部署服务器集群、搭建高性能存储集群、搭建内部集群网络。

(2) 平台核心集群模块:搭建核心GPU计算集群、分布式存储集群、高性能运算网络。计算方面,搭建高性能GPU计算服务器集群,支撑高密度算力运算需求;存储方面,搭建分布式海量存储集群,支撑海量数据高吞吐访问及高容量存储;网络方面,搭建高速运算网络、高速存储网络、管理网路。

(3) 硬件上构建核心软件中台:主要实现集群管理与作业调度功能,其中,集群管理通过安装部署集群管理软件,构建服务器集群体系,实现算力资源虚拟化管理;作业调度则在集群基础上部署作业调度管理软件,构建算力资源调度体系,实现算力资源弹性调度管理。

(4) 算力业务应用层:通过构建高性能计算开放平台,对各用户提供适用于各种运算业务的运算资源和数据,如大数据计算、基因测序、多模态数据模型、图像识别等。

1.2 节能建设方案

从最新的国家政策导向可知,数据中心低碳节能已经是数据中心建设和运营很重要的一个指标,数据中心节能降耗成为国家“碳中和”“碳达峰”战略的重要一环。因此,基础计算平台建设应该满足IDC能耗政策要求。

本次高性能计算服务器集群节能方案采用液冷管理节能技术,该技术利用超高导热系数的液冷导热模组将服务器高热流密度的核心芯片热量通过液冷模块带到服务器外,进一步通过水循

环冷却系统排到室外,液冷循环采用自然冷却无需压缩机,实现节能效果。图2为本次建设采用的间接液冷系统热管理架构图,间接液冷系统由导热液冷模组、快速接头、液冷分配单元、冷量温控单元、一次冷却环路、二次冷却环路、封闭气冷通道、自然冷却单元等构成,具有低能耗、高功率密度、高可靠性等优点。

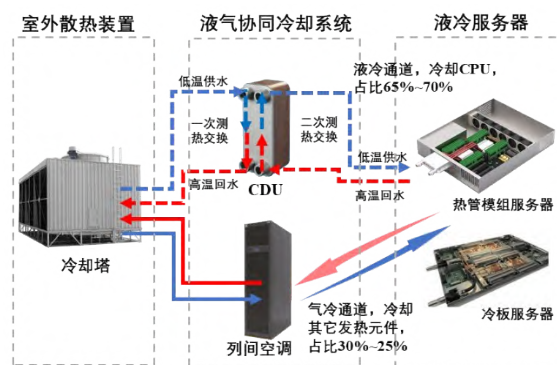


图2 间接液冷技术原理图

1.3 分布式存储集群方案

为构建分布式高性能存储集群,集群建设采用基于BeeGFS并行文件系统分布式存储方案,其主要优势表现在:①分布式文件内容和元数据,有效避免架构瓶颈,一方面可跨多个服务器的条带化文件内容,另一方面则可使文件系统的元数据存放于多个元数据服务器内。②兼容性好,BeeGFS存储服务基于横向扩展(Scale-Out)设计。每个BeeGFS文件系统实例可以具有一个或多个存储服务组件,方便提高性能与空间。一个存储服务实例具有一个或多个存储服务组件。③缓存优化能力强,由于BeeGFS自动使用存储服务器上的所有可用RAM自动进行缓存,因此它还可以在将数据写入磁盘之前将较小的IO请求聚合到较大的块中。④优化高并发访问,BeeGFS用于在高I/O负载的情况下提供最佳的稳健性和性能,优化解决简单的文件系统(比如NFS)在高并发访问的情况下存在严重的性能问题,以及在多个客户端写入同一个共享文件时会损坏数据等典型问题。

1.4 高性能计算网络方案

在高性能计算场景下,由于集群之间需要相互通信,所以对网络的带宽和时延要求比较

高(应用之间带宽>40 Gbps, 时延<10 us 微秒), 现有的 TCP/IP 软硬件结构无法满足该需求, 因此需要使用 RDMA(Remote Direct Memory Access) 技术远程直接内存访问, 构建 IB(Infiniband) 网络实现高性能场景下高速度、高吞吐网络传输需求。RDMA 模式对数据包的加工都在网卡内完成。因此就跳过了操作系统, 直接把数据发送到网卡内, 少了应用内存与内核数据之间的交互, 所以速度上更快, 时延更短。IB 网络: 基于无限带宽技术, 这种网络有很高的带宽(100 Gb/s 以上)和非常低的时延(毫秒级)。

2 高性能计算服务器集群管理方法

2.1 高性能计算服务器集群管理需求分析

本次高性能计算服务器集群管理通过集群管理软件, 构建服务器集群体系, 实现算力资源虚拟化管理。其中集群管理需求主要表现在如下方面:

- (1) 满足对多种深度学习、机器学习及大数据任务的资源调度和管理需求, 要求提供大规模 GPU 集群调度、集群监控、任务监控、分布式存储等功能。
- (2) 实现集群资源调度与服务管理统筹, 提供针对 GPU 优化的调度算法, 实现集群资源调度高效管理。
- (3) 提供面向用户的可视化接口或应用接口, 网页端可视化界面、客户端 SDK、集成开发环境(IDE)拓展接口等。

(4) 提供丰富的用户管理, 集群、任务监控, 任务调度, 任务错误分析, 任务监控等服务功能, 提高运维人员的工作效率。

(5) 实现容器化和微服务化, 使得运行环境可以在开发和运维达到统一。软件需支持任何形式的计算任务以及大部分计算框架, 包括各种深度学习框架和机器学习框架(如 PyTorch、Tensorflow)等。

2.2 高性能计算服务器集群管理方法探析

针对以上分析的高性能计算服务器集群管理需求, 设计图 3 的高性能计算服务器集群管理架构图, 共分为用户管理、集群管理、业务管理三大方面。

2.2.1 用户管理

及对用户组群集资源参数进行管理如用户组的 GPU 数量、存储配额、组名、最大运行作业数、等待作业数等; 支持同时在多个已分配资源的租户空间执行任务, 各用户资源互不影响。设置不同的资源分配和服务访问权限; 不同用户组间实现数据隔离。

2.2.2 集群管理

包括分布式管理、并行训练作业管理、集群总计显示等内容。分布式管理是集群管理的重点, 包括如下方面:

(1) 分布式计算集群监控: 包括集群资源总体监控人 GPU 资源监控。集群资源总体监控, 监控和显示群集 GPU、内存和存储总体使用情况。

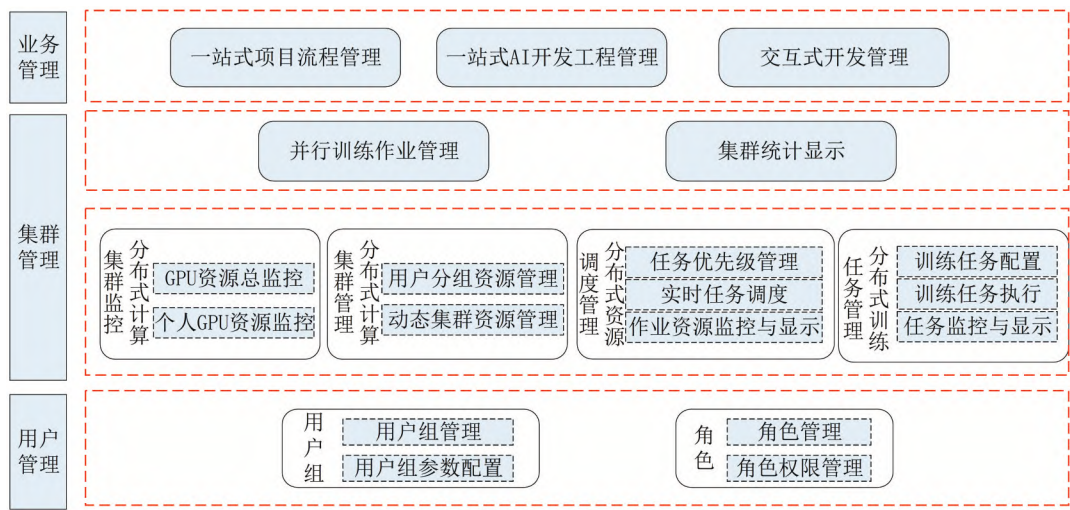


图 3 高性能计算服务器集群管理架构图

(2) 分布式计算集群管理: 对用户组进行集群资源配额管理, 对组内用户进行集群资源配额管理; 集群资源管理, 动态添加基础资源, 动态分配和管理集群资源。

(3) 分布式资源调度管理: 执行任务优先级管理; 实时任务资源分配和调度管理; 以任务方式根据优先级分配计算资源, 任务完成进行计算资源回收; 计算作业资源监控和执行情况显示。

(4) 分布式训练任务管理: 训练任务配置管理, 基础任务参数管理; 训练任务执行管理; 训练任务监控, 查看模型训练和资源使用情况。

2.2.3 业务管理

基于高性能计算服务器集群应用的一站式业务流程管理, 包括一站式项目流程管理、一站式AI开发工程管理、交互式开发管理等。

(1) 一站式项目流程管理: 实现项目流程构建、项目流程表单设计、任务分配、团队管理、项目流程可视化、项目检索等项目流程管理支持。

(2) 一站式AI开发工程管理: 包括数据集管理、数据集推荐、模型训练、模型部署、API调用示例、关联用户训练任务与部署任务等内容。

(3) 交互式开发管理: 支持用户通过平台内置AI镜像进行创建交互式开发环境, 环境实例可以使用CPU资源也可以使用GPU资源; 平台支持开发实例的持久化。

3 结语

在国家人工智能战略布局和产业发展大背景下, 高性能计算服务器集群公共服务平台建设需求日渐突出。本文提供一种切实可行、低碳高效的高性能计算服务器集群建设方案, 并探析高性能计算服务器集群管理方法, 借助自主研发集群管理软件, 构建服务器集群体系, 实现算力资源虚拟化管理。相关管理方法有待在实践中进一步优化和深化。

参考文献:

- [1] 《国家及各地区国民经济和社会发展第十四个五年规划和2035年远景目标纲要》[J]. 中国信息界, 2022(5):110.
- [2] 朱磊明. 云计算平台基础服务[J]. 信息与电脑(理论版), 2018(4):30-31.
- [3] 杨燕艳, 朱春燕. 基于大数据技术的智慧医疗平台设计[J]. 信息记录材料, 2022, 23(6):173-176.
- [4] 上海建立类脑芯片与片上智能系统研发与转化功能型平台[J]. 华东科技, 2019(3):11.
- [5] 崔昌云. 智慧城市大数据平台的研究与应用探究[J]. 中国信息化, 2021(4):112-113.

作者简介:

陈阳(1984—), 男, 广东潮州人, 硕士研究生, 高级工程师, 研究方向为计算机应用技术研究

陈坚泽(1983—), 男, 广东潮州人, 硕士研究生, 研究方向为电气自动化、计算机应用技术研究

收稿日期: 2022-05-04 修稿日期: 2022-06-27

Analysis of High Performance Computing Server Cluster Construction and Management Methods

Chen Yang, Chen Jianze

(Guangdong Liquid Cooling Technology Co., Ltd., Foshan 528000)

Abstract: Based on the national artificial intelligence strategy layout and development needs, and the policy goal of “carbon neutrality”, the HPC server cluster construction scheme was introduced in this paper, including HPC cluster, energy saving construction, distributed storage cluster, HPC network construction scheme. The requirements of high performance computing server cluster management were analyzed and the methods of high performance computing server cluster management were discussed.

Keyword: high performance computing; server cluster; construction and management