

基于 AI 驱动的混合云 IDC 资源 自动化与成本优化策略

闫崇喆

(北京飞机维修工程有限公司, 北京 100000)

摘要 本文设计了基于 AI 驱动的大型混合云 IDC 规划方案, 聚焦五大核心目标: 自动化资源需求预测、优化成本结构、提升资源利用率与弹性、强化多租户安全合规及加速业务上线与容量扩展。通过分阶段、数据驱动等技术, 分析混合云 IDC 架构的负载资源需求动态化、峰均比高、资源利用率偏低等情况, 提出“数据可观测性建设—初步机器学习预测—深度 AIops 平台集成—多云联动全局优化”共计四个阶段演进路线, 推动混合云 IDC 资源规划从被动响应到主动预测, 手动式管理到自动化调度的转变。该方案基于核心目标提出具体实施策略, 同时提出团队能力建设与组织保障的建议, 以期对企业构建高效、敏捷、安全且成本可控的新一代混合云基础设施有所裨益。

关键词 混合云架构; AIops; 资源优化; 成本控制

中图分类号: TP18; TP2

文献标志码: A

DOI: 10.3969/j.issn.2097-3365.2026.01.024

0 引言

随着数字化转型的推进, 业务全球化下传统大型企业的 IT 基础设施建设方式日益复杂。混合云架构已成为既满足私有云的安全性、可控性, 又具有公有云的弹性、成本效益的主流 IDC 方案。但混合云架构仍需面临资源管理、成本控制和运维效率等方面的巨大挑战。为应对这些挑战, 并支撑企业业务的持续创新与增长, 本文聚焦五大核心能力: 一是自动化预测资源需求, 变被动扩容为主动规划; 二是优化成本结构, 通过精细化管理与智能调度降低 CAPEX 与 OPEX; 三是提升资源利用率与弹性, 打破孤岛并保障高峰服务质量; 四是强化多租户安全合规, 实现复杂共享环境下的数据隔离; 五是敏捷支撑业务上线与扩容, 缩短开发到上线周期。

1 大型混合云 IDC 规划面临的问题

1.1 工作负载特征复杂且动态

虚拟化与容器平台作为基础计算平台支撑着大量应用。在多租户环境下, 资源争抢可能导致性能变得异常波动。虚拟化层和虚拟机配置对性能有显著影响, 需要精细化设计。大数据分析的工作负载一般是周期性、计算密集型的, 如 MapReduce 任务, 其资源需求呈现明显的潮汐效应。AI 训练及推理是资源需求的“大

户”, 尤其消耗大量 GPU 等异构计算资源。且 GPU 需求通常是突发性的, 训练任务期间资源占用率极高, 而任务间隙则可能闲置。并且训练出的模型在转为推理使用后, AI 平台消耗的资源呈指数级降低。网络与安全设备等基础设施虽然自身计算需求相对稳定, 但其性能直接关系到所有业务的稳定性和安全性, 且会随时面临日益严峻的网络安全威胁。

1.2 普遍存在的资源管理挑战

云数据中心中的业务关键型工作负载, 其资源使用的峰值可能是平均值的 10 倍到 100 倍甚至更高。为了应对峰值, 传统规划方式往往导致资源在大部分时间里处于低利用率状态, 造成巨大浪费。工作负载具有显著的每日、每周的时间模式, 但不同资源类型 (如 CPU、磁盘 I/O) 的模式各不相同, 增加了手动规划的难度。

1.3 人工智能应用成熟度

市场调研显示, 多数传统大型企业在 AI 赋能资源规划领域仍处于调研阶段。传统行业 (如制造、能源) 的资源规划高度依赖历史经验、政策合规及多维数据整合, 受数据标准化、决策透明度等瓶颈制约, 难以应用 AI 技术。同时, 企业普遍缺乏现成模型、工具及专业团队。此外, AI 研发对计算、数据等资源需求巨大, 因此需制定清晰可行的渐进式路线图, 而非一步到位部署复杂系统。

作者简介: 闫崇喆 (1990-), 男, 本科, 工程师, 研究方向: IT 基础架构、混合云、智算中心。

2 AI 驱动混合云 IDC 规划方案与技术路线

本方案采用循序渐进的四阶段演进策略，从数据基础做起，逐步构建智能化、自动化的资源规划与管理能力。

2.1 阶段一：数据基础建设与全面可观测性建立

本阶段核心目标是构建混合云环境下统一数据采集与监控平台，实现对计算、存储、网络等全量资源及应用性能的全方位监控，为后续构建智能分析与预测模型提供数据基础。正如管理领域核心准则“无法量化，就无法管理”，精准的数据采集与监控是混合云高效运维的前提。本文将通过四大技术路线推进：第一，在物理服务器、虚拟机、容器及关键网络设备中部署统一监控代理^[1]（如 Prometheus Node Exporter、cAdvisor、Fluentd 等），实现数据全面采集；第二，构建含统一监控数据的时序数据库与日志聚合系统的中央数据湖，完成监控指标、日志及链路追踪数据的集中存储；第三，通过 Grafan 等工具构建多业务、多资源维度的多层次可视化仪表盘，支撑日常运维与工作负载行为分析；第四，系统性收集 CPU/ 内存利用率、I/O 性能、应用响应时间等关键指标，计算峰均比、波动性等统计特征，完成工作负载特征量化，为后续建模提供数据支撑。

2.2 阶段二：基于统计与机器学习的初步预测

本阶段核心目标是基于前期积累的历史数据，初步开发资源需求预测模型，推动混合云运维从被动响应模式到主动预警模式转型。策略将通过四大技术路线落地实施：第一，针对大数据分析、常规业务应用等具有显著周期性的工作负载，采用 ARIMA、指数平滑法等时间序列预测模型，实现未来 24 小时、一周等时段的资源使用趋势预测；第二，为各应用及服务构建资源消耗“健康基线”，运用 3-sigma 原则^[2]等统计方法或孤立森林等基础机器学习算法，自动识别资源使用的异常波动并及时预警；第三，开发“What-if”分析工具，支持业务部门输入用户增长预期、大促活动等场景参数，同时模型应能够基于历史数据粗略估算额外资源需求，为手动扩容提供数据支撑；第四，强化团队能力建设，使工程师能够深度熟悉企业数据特征与业务场景，为模型落地与迭代提供人才保障。

2.3 阶段三：AI0ps 平台深度集成与智能调度

本阶段核心目标是构建并引入 AI0ps 平台，融合预测能力与自动化控制技术，实现混合云环境下资源智能调度、弹性伸缩及故障自愈，显著提升运维效率与资源利用率。策略通过四大技术路线推进：第一，开发自定义容器调度策略，调度决策既要考虑当前资源余量，又要结合预测模型结果，提前调度高负载 Pod 至

资源充足节点，或在资源紧张时驱逐低优先级 Pod；针对 AI 训练负载，调度器需感知 GPU 拓扑与显存状态^[3]，结合任务优先级及预计运行时长实现排队调度，最大化 GPU 集群利用率；第二，基于时间序列预测结果实现预测性自动伸缩，在业务高峰前提前扩容实例，高峰后平滑缩容，避免滞后响应导致的服务质量下降或资源浪费；第三，借助机器学习模型分析监控指标、日志与告警的关联性，在性能异常时快速定位根本原因，同步触发重启服务、隔离故障节点等自动化自愈预案；第四，融入成本感知调度策略，集成公有云 IaaS、PaaS 等定价模型，将可中断批处理任务优先调度至低成本实例，实现成本优化，该路线需要深入掌握公有云服务特性为基础。

2.4 阶段四：多云联动与全局资源优化

本阶段以最终目标为锚点，制定持续优化路线。该阶段核心是将私有云与多公有云平台整合为统一逻辑的“无限”资源池，使得应用负载在多云环境的自由智能流动，实现成本、性能与安全合规的动态最优平衡。策略通过四大技术路线推进：第一，部署或开发云管理平台，实现混合云资源的统一纳管，提供标准化服务目录与资源编排能力，夯实多云协同基础；第二，构建 FinOps 体系^[4]，培育 FinOps 文化并搭建工具链，实现云成本透明化、精细化分摊与持续优化，满足 AI 项目及公有云预算与成本管理的核心诉求；第三，优化智能工作负载放置策略，AI0ps 平台基于工作负载特征、数据主权要求、安全合规规则及各云实时成本，自动决策最优运行位置，如延迟敏感型业务部署于私有云，高弹性 Web 应用部署于公有云并实现多云平台智能比价；第四，攻克数据联动与一致性难题，解决跨云数据同步、一致性保障及网络延迟等问题，确保应用负载迁移后数据可访问性与性能达标。

3 针对核心目标的具体实施建议

3.1 自动化预测资源需求

从阶段一开始系统性收集数据，在阶段二应用时间序列模型进行趋势预测，并在阶段三引入更复杂的深度学习模型（如 LSTM）以提高预测精度，最终实现对计算、存储、网络需求的全面、精准预测。

3.2 优化成本结构（降低 CAPEX/OPEX）

成本优化是本研究核心价值之一，主要通过 CAPEX 与 OPEX 双维度实现精益管控。在 CAPEX 优化方面，依托精准的资源需求预测模型，规避数据中心过度采购问题。采购决策不再基于简单线性增长假设，而是以未来 6 ~ 12 个月的预测需求为核心依据，实现资源配置与实际需求的精准匹配。在 OPEX 优化层面，通过多

阶段技术落地层层递进：借助阶段三 AIOps 自动化能力减少日常运维人力投入；通过智能调度与预测性伸缩，将资源平均利用率从行业低位提升至 60%~70% 的目标水平，降低闲置成本；阶段四则通过多云成本套利策略，智能选用公有云实例等低成本资源，进一步压缩运营开支。需强调的是，AI 项目的云资源预算管理是成本优化成功的关键支撑。

3.3 提升资源利用率与弹性

提升资源利用率与弹性是本文核心价值体现。在利用率优化方面，通过智能调度器实施工作负载混合部署策略，将日间高峰在线业务与夜间批处理任务等错峰负载^[5]部署于同一物理集群，填补资源使用“波谷”，实现全天候资源高效利用。在弹性增强层面，以公有云无限资源池作为私有云弹性扩展补充，当 AIOps 平台预测流量洪峰超出私有云容量时，自动触发“云爆发” workflow，将部分流量或应用实例无缝迁移至公有云，充分发挥其快速扩展能力，保障服务稳定性。

3.4 加强多租户安全与合规性

AIOps 的价值不仅限于效率提升与成本优化^[6]，其数据分析能力可延伸至安全领域，核心支撑异常行为检测与合规性监控自动化。在异常行为检测方面，通过分析资源调用、网络流量与 API 访问日志模式，AIOps 平台能识别偏离正常基线的行为，为安全漏洞、恶意攻击及内部违规操作提供早期预警。在合规性监控自动化层面，将数据存储地理位置等合规策略编码为调度规则，确保工作负载部署自动满足合规要求。构建完善安全框架，是数字化解决方案落地的必要前提。

3.5 加速新业务上线与容量扩展

借助 AIOps 与基础设施即代码 (IaaS) 融合，新业务上线的资源申请、审批、配置及部署流程实现全自动化，交付周期从数周压缩至数小时乃至数分钟。并且利用 AIOps 平台的精准容量预测支撑“容量即服务”模式，助力基础设施团队提前规划扩容，保障业务部门资源按需供给，有效避免因资源短缺导致的新业务上线延迟问题。

4 团队建设与组织保障

4.1 团队技能图谱构建

组建跨职能的“云智能中心”团队，成员应包括：数据科学家/AI 工程师（负责开发和优化预测模型与调度算法）；应用运营工程师（负责 AIOps 平台的搭建、运维和与现有系统的集成）；云架构师（负责设计混合云和多云架构）；FinOps 分析师^[7]（负责云成本的分析、优化和预算管理）。

4.2 工具选型

在自研与采购之间取得平衡。核心的预测和调度算法建议自研，以贴合传统大型企业实际情况的独特业务模式。而监控、日志、告警等基础平台可优先考虑成熟的开源或商业解决方案。

4.3 文化变革

推动从传统的 IT 运维文化向数据驱动、主动预测的 AIOps 文化转型。鼓励实验、容忍失败，并建立清晰的量化指标（如资源利用率、平均故障恢复时间、单位业务成本）来衡量项目成效。

5 结束语

在数字化转型持续推进的背景下，承载复杂工作的传统大型企业正面临传统 IDC 资源规划模式的瓶颈，静态化、人工化管理方式已难以适应新业务快速迭代与弹性伸缩的需求。为突破企业运维困境、重构资源管理体系，人工智能驱动的 AIOps 模式提供了创新性解决方案，这种模式已成为降本增效与提升业务敏捷性的核心方向。本文构建的四阶段演进路线，为 AI 资源规划能力的从零构建提供了可落地的实施框架。通过稳步筑牢数据基础、融合应用机器学习与自动化技术，最终建成全局优化的智能混合云管理体系，既能高效应对当下业务挑战，又能夯实适配未来技术变革与市场波动的运维根基，推动 IT 基础设施实现从“成本中心”到“价值创造中心”的关键跨越。

参考文献：

- [1] 佚名.AI 与可持续性能力解决云转型难题[J]. 软件和集成电路, 2023(01):64-66.
- [2] 孙磊. 人工智能推动交通银行数据中心向智算中心转型[J]. 中国金融电脑, 2025(06):13-16.
- [3] 丁建仁, 吴俊, 张学亮, 等. 多架构计算与存储需求下 IDC 云平台资源优化配置与高效调度研究, 2025(03): 69-72.
- [4] 熊俊. 人工智能算力中心网络建设探究[J]. 中国安防, 2025(08):44-47.
- [5] 刘伟, 李朝阳, 史海超. 人工智能技术推动能源变革的政策体系和应用挑战研究[J]. 信息通信技术与政策, 2025, 51(06):27-32.
- [6] 柳雨晨, 李涛, 王志佳, 等. 基于大语言模型的电信运营商数字化系统智能运维(AIOps)的设计与挑战[J]. 广东通信技术, 2025, 45(04):37-41.
- [7] 周春云. 高性能多微云成本优化调度系统设计与实现[J]. 通讯世界, 2025, 32(01):38-40.