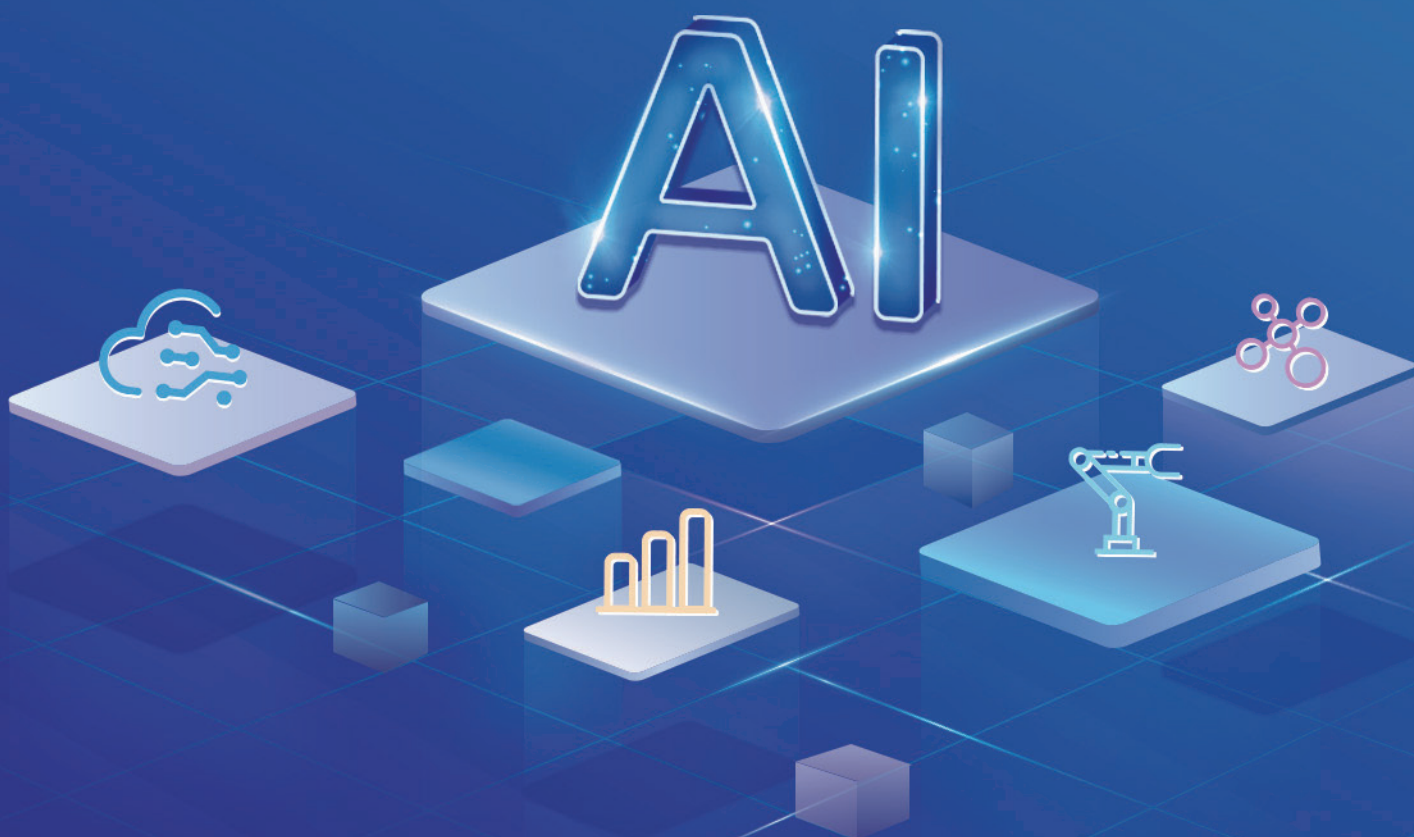


中兴通讯技术 **简讯**

ZTE TECHNOLOGIES | 第30卷 第4期 · 2026年4月

视点

04 多要素协同进化，筑牢AI时代算力底座



专题：新型智算

07 超节点应用场景及技术演进





1996年创办

第30卷 总第450期

2026年04月 第4期（月度出版）

中兴通讯技术（简讯）

ZHONGXING TONGXUN JISHU (JIANXUN)

《中兴通讯技术（简讯）》顾问委员会

主任：刘健

副主任：方晖 彭爱光 孙方平 张万春

顾问：柏钢 陈新宇 董伟杰 胡俊勤

胡立华 华新海 阚杰 李强

李晓彤 唐雪 王全 杨运东

郑鹏

《中兴通讯技术（简讯）》编辑委员会

主任：林晓东

副主任：卢丹

编委：邓志峰 代岩斌 关凯 黄新明

梁大鹏 林晓东 卢丹 马小松

孙岳 施军 王卫斌 肖伟

杨兆江 余方宏 赵建超

《中兴通讯技术（简讯）》编辑部

总编：林晓东

常务副总编：卢丹

编辑部主任：刘杨

执行主编：方丽

发行：王萍萍

主管：中兴通讯股份有限公司

主办：中兴通讯技术杂志社

出版：《中兴通讯技术（简讯）》编辑部

编辑部地址：深圳市科技南路55号中兴通讯研发大楼

发行部地址：合肥市金寨路329号国轩凯旋大厦12楼

发行部电话：0551-65533356

<https://www.zte.com.cn/china/about/magazine>

发行范围：国内业务相关单位

印数：5000本

设计：深圳市奥尔美广告有限公司

印刷：深圳市旺盈彩盒纸品有限公司

印刷日期：2026年04月25日

未经中兴通讯股份有限公司书面授权，禁止以转载、
摘编、复制等方式使用本资料的任何内容。



王卫斌

中兴通讯产品规划首席科学家

AI纵深推进，加速商用落地

AI成为全球战略共识和全行业竞争焦点，正从技术迭代迈向价值深耕周期。随着AI+与实体经济深度融合，产业竞争升级为软硬件生态的全面博弈，驱动AI基础设施、模型、应用三者协同发展与商业闭环。

基础设施领域，从重模型训练的模式向强Token调用的智能云范式跃迁，进一步朝着系统性构筑训推一体、云边端应用协同的IoA基础设施演进，以AI芯片为基础的推理效能和超节点成为核心赛道。

模型领域，开源挤压闭源生存空间，新算法快速迭代，非Transformer架构涌现，多模态向全模态发展，随着进一步深入理解光、电、热、力等物理世界规律，模型能力的提升加速释放商业潜能，为行业应用落地奠定基础。

应用领域，智能体加速重塑传统应用，成为AI落地的重要突破口。继Manus之后，OpenClaw以本地部署、持久记忆等特性引爆轻量化落地，构建AI Agent开发平台支撑传统应用向智能体的演进和升级换代成为共性需求。

中兴通讯致力于成为网络连接和智能算力的领导者，践行自主创新和开放合作的AI发展战略，提供AI基础设施、模型及应用的端到端解决方案。AI基础设施依托算存联软能多要素协同，形成万卡训练集群和高性能推理两大重点方案，同时提供大容量交换、超节点、一体机等系列产品；自研星云编程大模型在HumanEval评测中位列业界第一梯队；Co-Sight超级智能体工厂业界率先开源，GAIA榜单开源第一，助力软件产品向Agent原生应用演进。公司积极倡导开放解耦生态，全面支撑大模型训练与行业智能应用落地，已在政务、金融、教育、医疗、制造等18+行业落地超100个AI标杆案例。

未来，中兴通讯将持续携手产业链伙伴协同创新，打造极致TCO的AI系统级解决方案，加速算力普惠和应用普及！

目次

中兴通讯技术（简讯）2026年第4期



超节点应用场景及技术演进

为应对大模型所带来的算力爆发式增长和分布式并行计算需求，超节点应运而生。超节点通过高速互联将众多GPU整合为一个逻辑统一、调度高效的计算单元，像一台“巨型服务器”一样协同工作。

07

视点

04 多要素协同进化，筑牢AI时代算力底座
郭雪峰

专题：新型智算

07 超节点应用场景及技术演进
毛磊

11 Scale-Up/Out/Across三域协同，突破算力上限
杨茂彬

14 GPU、DPU、存储介质协同的新型AI存储架构
郭伟

16 高效推理加速AI落地
周俊超

18 筑牢可信底座，破解工业智能体长程推理困局
陈云斌

20 中兴通讯Skill安全体系，构建智能体边界
顾希

22 全球AI开源技术发展概况及中兴通讯开源实践
王长金

25 构建开放协同的智算标准体系，推动AI生态高质量发展
朱静，黄程

成功故事

27 携手南方电网，打造能源行业首个全栈自主可控的千卡智算中心
黄燕

30 山东移动携手中兴通讯打造千卡智算推理资源池
陆威，蒋妍，杨晓曦

技术论坛

32 空芯光纤传输系统应用与挑战
尚文东

34 低时延空芯光纤在跨域大模型训练中的应用分析
闫宝罗

02 新闻资讯

中兴通讯发布全栈AI基础设施与智能终端矩阵

2026年4月9日，北京——在“和合兴业 智启未来”为主题的2026中兴通讯中国生态合作伙伴大会主论坛上，中兴通讯发布以“算网存智一体”为核心的全栈AI基础设施与智能终端矩阵，全面彰显其作为“网络连接与智能算力领导者”的战略决心与技术实力。

在核心能力层面，GoldenDB数据库向量版本重磅发布，国内首批实现标量+向量+全文混合检索，满足金融级的高精度要求，实现数据库产品从“数据容器”到“智能决策引擎”的跨越。

在AI基础设施-算力领域，以“5+X”多芯片协同设计，打破算力孤岛；通过业界首创的“OEX创新正

交架构”，打造开放、高密度、高可扩展的超节点架构；可支持单机柜128个GPU，算力规模可扩展至1.6万卡，构建最优TCO的全栈智算基础设施。

在AI基础设施-网络领域，推出576x800G端口框式智算交换机，二层组网即可支持超大规模AI算力集群，相较三层组网，时延降低33%，设备、光模块和光纤数量减少40%；同步推出智算跨域互联（Scale-Across）方案，实现300公里超长距算力调度，调度性能达99%。

在AI基础设施-能源领域，发布800V高压直流电源系统，降低线损与铜耗，系统效率超98%；最新推出的2MW模块化CDU液冷系统，首次在

47U标准机柜内实现2000kW制冷，冷量密度业界最高，PUE低至1.15，年省电费25%；并推出125kW/261kWh工商业储能一体柜与5MWh集装箱储能，实现“算电协同”。

在AI能力平台领域，打造Co-Claw企业级AI智能体平台，深度融合开源能力与企业规模部署的需求，打通研发、生产、供应链等核心业务系统，并实现高价值场景的快速落地。中兴通讯已将Co-Claw全面融入最新产品生态，推出Co-Claw智算一体机，并嵌入云电脑、智慧家庭全场景终端，构建“端-边-云”协同的AI生产力生态。

在AI应用及终端领域，“大-中-小”屏全系列AI云电脑与移动互联产品矩阵全面亮相，发布驭风10 Air云电脑、“多合一”自由屏F10、逍遥20 AI PAD等多款云电脑终端，打破算力边界，智联AI世界，定义新一代电脑。

国产光互连光交换超节点“光跃128卡商用版”正式落地

在3月12日开幕的中国家电及消费电子博览会（AWE）上，上海仪电（集团）有限公司联合上海曦智科技股份有限公司、上海壁仞科技股份有限公司、中兴通讯股份有限公司正式发布光跃超节点128卡商用版（LightSphere 128）。这标志着这一中国原创的光互连光交换超节点解决方案，仅用半年多时间即实现从概念验证到实际商用的跨越。目前，光跃超节点已实现数千卡的部署。

中兴通讯首发Co-Claw智慧园区方案

近日，中兴通讯正式发布基于Co-Claw企业级AI智能体平台的智慧园区方案，并率先在南京滨江智能制造基地落地。该方案构建“全域感知-智能预判-协同联动-闭环管理”的智能化管理运营体系，推动园区从“被动响应”向“主动治理”、“分散管理”向“系统协同”的全面跃迁，打造高效、安全、绿色的智慧园区应用新范式。

为破解传统园区管理模式瓶颈，

中兴通讯将自主研发的企业级AI智能体平台Co-Claw深度嵌入园区运营管理全链条，以“数据不出园、行为可追溯、权限强管控、安全围栏动态拦截”为核心安全特性，构建了覆盖源码加固、Skill供应链审查、身份认证、敏感操作分级响应与全链路审计的三维一体主动安全治理体系，实现AI智能体在园区场景下的可信、可控、可审计运行。同时，围绕园区通用管理与核心生产两大维度，打造六大场景：产线物流、智算中心运维、配电房运维、产线高温区巡检、园区安防、能耗管理，实现园区数智化升级。

中兴通讯董事长方榕 出席博鳌亚洲论坛2026年年会

3月24日，博鳌亚洲论坛2026年年会在海南博鳌正式拉开帷幕。中兴通讯董事长方榕出席本届年会，并于3月26日参加了“走进AI时代：把握机遇 创造未来”主题圆桌论坛，与全球政商学界嘉宾深度对话，围绕AI与传统产业融合、AI安全与伦理治理两大核心命题分享行业思考与中兴通讯实践，为AI时代全球数字经济高质量发展贡献中国企业方案。

在探讨企业如何在产业升级和跨产业转型中确定自身定位时，方榕指出，AI时代的竞争已超越了传统的“跨界”思维，本质是降维赋能。中兴通讯致力于成为产业的“智能内核”。方榕强调：“对客户而言，缺的不是

设备，而是能把模糊经验变成精准决策的‘大脑’。我们产品的核心是交付决策的确定性。”在市场定位上，中兴通讯主攻流程最复杂、容错率最低、数据依赖最强的领域，以此构建难以模仿的竞争护城河。



中兴通讯与中国外运 签署战略合作协议

3月12日，中兴通讯股份有限公司和中国外运股份有限公司在深圳中兴通讯总部举行战略合作签约仪式。中兴通讯高级副总裁、供应链总裁杨建明，中国外运党委书记、董事长张翼出席签约仪式。中兴通讯副总裁、客户订单交付部总经理刘剑锋，中国外运党委委员、副总经理汪剑代表双方签署战略合作协议。按照协议，双方将发挥各自的优势，在供应链管理、企业数字化、关务及合规建设等领域展开深度合作，聚焦运力协同、数据互通、技术创新及绿色低碳四大核心领域，携手打造智慧、高效、可持续的全球供应链新标杆。

杭州电信携手中兴通讯圆满 保障周杰伦杭州演唱会

4月3日—5日，周杰伦2026「烟花杭州·嘉年华II」世界巡回演唱会全球首站在杭州奥体中心圆满落幕。杭州电信携手中兴通讯，采用5G-A EasyOn·Live专网直播方案，圆满完成多机位“超高清、低时延、零卡顿”的无线直播保障，进一步验证了该方案在“高密度、高要求、高量级”活动场景下的商用成熟度，为5G-A在文娱直播领域的规模化商用奠定了坚实基础。

中国联通携手中兴通讯率先 完成毫米波自愈合技术验证

近日，中国联通联合中兴通讯、东南大学在江苏率先完成毫米波频段波束自愈合技术实验室验证，在真实无线链路环境下实现毫米波信号遇阻后的自主重建与链路性能恢复，标志着我国在6G高频通信抗遮挡、抗衰落等高频越障核心技术领域取得突破性进展，为毫米波规模化商用奠定坚实技术基础。

黑龙江联通携手中兴通讯完成 哈佳高铁NR 1.8G重耕试点

3月底，黑龙江联通与中兴通讯携手，圆满完成哈佳高铁沿线NR 1.8G频段的重耕升级工程。在NR 1.8G重耕过程中，黑龙江联通与中兴通讯应用了超级小区、载波聚合等前沿技术，提升了网络资源的利用率和覆盖质量。升级后，哈佳高铁沿线的网络质量实现了整体跃升，网络下载速度峰值提升约90%，路段下载平均速度提升约40%，为旅客带来更加流畅的移动应用体验。



郭雪峰

中兴通讯算力及核心网产品规划首席专家

多要素协同进化， 筑牢AI时代算力底座

在新一轮科技革命与产业变革浪潮中，人工智能已成为重塑全球竞争格局的核心力量，不仅深刻改变着千行百业的生产范式，更跃升为国家科技战略博弈的关键制高点。

当前，我国人工智能产业虽在算法创新、场景应用等领域取得显著成绩，但在算力基础设施层仍面临先进算力代际差距、异构算力生态协同不足、算力资源利用率偏低等多重挑战。本文面向人工智能长期发展，探讨国产智能算力普惠发展路径。

人工智能规模落地发展，亟需普惠算力支撑

算力是人工智能发展的基石，大模型作为当

前人工智能发展的核心载体，其迭代速度与落地广度，始终受到算力资源可及性与成本的约束。无论是基础大模型的预训练还是行业大模型的定制化开发，都需要海量算力及巨额成本支出。随着我国“人工智能+”战略的深入拓展，AI应用以及AI Agent快速发展，推动推理算力需求爆发式增长。据IDC预测（见图1），我国智能算力规模未来三年仍将实现翻倍增长，智能算力成本的持续优化和降低程度，将直接决定人工智能赋能产业化长期发展的成效。

无论是大模型的迭代演进，还是Agent的场景化落地，都离不开普惠算力底座的支撑。算力成本的降低还将进一步打破技术创新的成本壁垒，促进人工智能技术持续发展，激活AI与千行百业的结合，让人工智能真正服务于每个人、每个场景，实现“技术向善、普惠共生”的发展目标。

算存网软能多要素协同，构建最优TCO算力底座

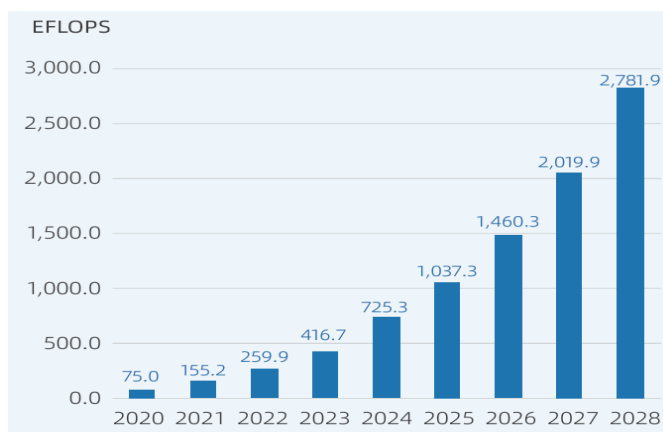
大规模数据中心是人工智能与实体经济深度融合的核心算力载体。在后摩尔时代，要真正发挥这一载体价值，依赖算力芯片迭代演进、算力集群规模扩张，虽可实现算力规模优势，却难以有效释放GPU算力潜能、降低算力成本。通过算力、存储、网络、软件、能源多要素深度协同，破解通信墙、存储墙、能耗墙，提升AI工厂算效与能效，打造最优算力TCO，已成为当前国内外产业界的普遍共识。

在国内，软硬协同、模芯协同已步入规模化落地阶段。产业端通过架构革新、算法优化与工程实践实现算力效率跃升。例如，行业内依托P/D分离与KV缓存实现算力与存储的协同，提升推理计算效率；通过通信计算重叠的Overlap调度技术，消除算力运行空泡，提升并行计算效率；基于MTP多Token预测算法优化大模型推理解码流程，突破单Token串行执行的效率限制。在海外，英伟达Rubin平台通过系统级的协同设计和优化，整合Rubin GPU、Vera CPU、NVLink6、Bluefield DPUs等核心硬件实现极致协同，单Token推理成本降至前代的1/10，高效支撑超级AI工厂落地。国内外实践探索充分印证，多要素协同进化是构建先进AI算力底座、实现最优TCO的关键途径。

我国智能算力建设在高端芯片获取受限的挑战下，已逐步走出一条依托算法、工程与算力协同的“经济高效型”差异化发展路径，但相对国际先进水平仍存在明显短板，算效能效偏低。

● 算力芯片存在代际差距

AI芯片厂商多而不强，受限于先进工艺制程，在算力、显存、显存带宽等核心指标上与国际头部厂商存在代际差距。单纯依靠芯片自身迭代难以快速补齐短板，更需通过系统性架构创新，以网络传输优化、存储性能提升、软件栈适配升级与算力的深度协同，最大化释放芯片的算



▲ 图1 中国智能算力规模及预测（基于FP16；来源：IDC，2025）

力潜能，实现“以协同补差距”的系统突破。

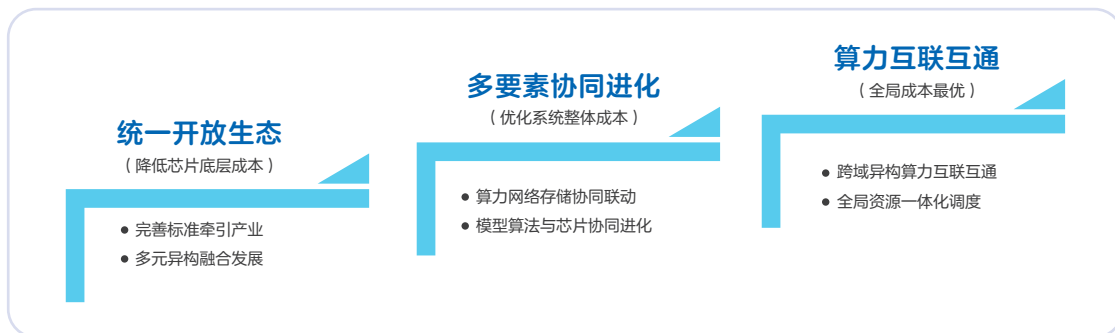
● 生态割裂制约协同落地

当前智算产业生态碎片化问题突出，芯片厂商主导形成大量相互独立的垂直生态，推高了模型算法与算力的迁移适配成本，跨厂商的算力、网络、存储间难以实现高效协同联动，重复开发与厂商锁定导致产业发展成本高、投入大。

● 孤岛现象普遍，算力资源利用率偏低

据统计，当前我国已建成数十个万卡智算集群以及大量的中小规模算力集群，智算规模跃居全球第二。但算力集群间互联互通壁垒高，一方面无法聚合形成大规模算力集群满足大模型训练需求；同时面向多场景的推理需求，算力集群间资源调度协同不足，导致算力资源配置局部失衡，整体资源利用率偏低。

从算力需求看，大模型全生命周期覆盖预训练、微调、强化学习等多元场景，算力需求呈现规模跨度大、动态波动强的特征：基础大模型预训练需要超万卡大规模集群支撑，训练期间独占资源，而模型微调、二次训练则以百卡级中小集群算力为主，需求灵活分散。推理算力正朝着轻量化、高效化、场景化方向转型，不同场景对算力的需求差异大。算力建设及供给需要与时俱进，从规模优先向效率优先演进，持续优化算力成本，推进产业普惠发展。



▲ 图2 多要素协同构建普惠算力

构建统一开放生态，推动全链条系统性协同创新

破局之道在于以开放统一生态、多要素协同进化、跨域算力互联互通为手段，实现系统性的架构优化与创新，构建先进AI智算底座：以开放架构解耦生态壁垒，构建国产化智算统一标准，实现国产芯片、软件栈、模型的迁移适配，降低生态碎片化成本；通过系统级架构创新，实现网络、存储、算力的高效协同，以网强算，以软补算，以集群算力优势弥补单芯片性能短板；构建算力互联互通网络，整合跨域算力资源，形成“算力分布式部署、一体化编排调度，资源动态聚合”的顶层算力格局，最终实现从单点技术追赶向生态体系突围的战略转型（见图2）。

统一开放生态，推进多元算力发展

面向大模型全生命周期业务场景，构建国产智算统一标准体系，开展异构算力互联互通、混训混推等核心技术攻关与工程实践，形成标准化、可复制的技术解决方案。依托开放统一的智算生态框架，推动多元国产芯片协同适配、有序发展，精准支撑超大规模训练、轻量化推理、边缘智能等多元化场景需求，为普惠算力规模化落地筑牢健康可持续的产业生态根基。

多要素协同进化，释放极致算力效能

系统架构层面，强化算力、网络、存储、算

法、模型全要素协同联动机制，并将协同设计贯穿算力基础设施建设全流程。优化参数面网络性能降低通信开销、内存外存一体池化破除存储壁垒，探索模型算法与国产算力特点适配机制，持续优化算力容错与故障恢复机制，以全要素协同升级推动算力基础设施整体迭代，打造规模、算效、能效领先的先进数据中心。

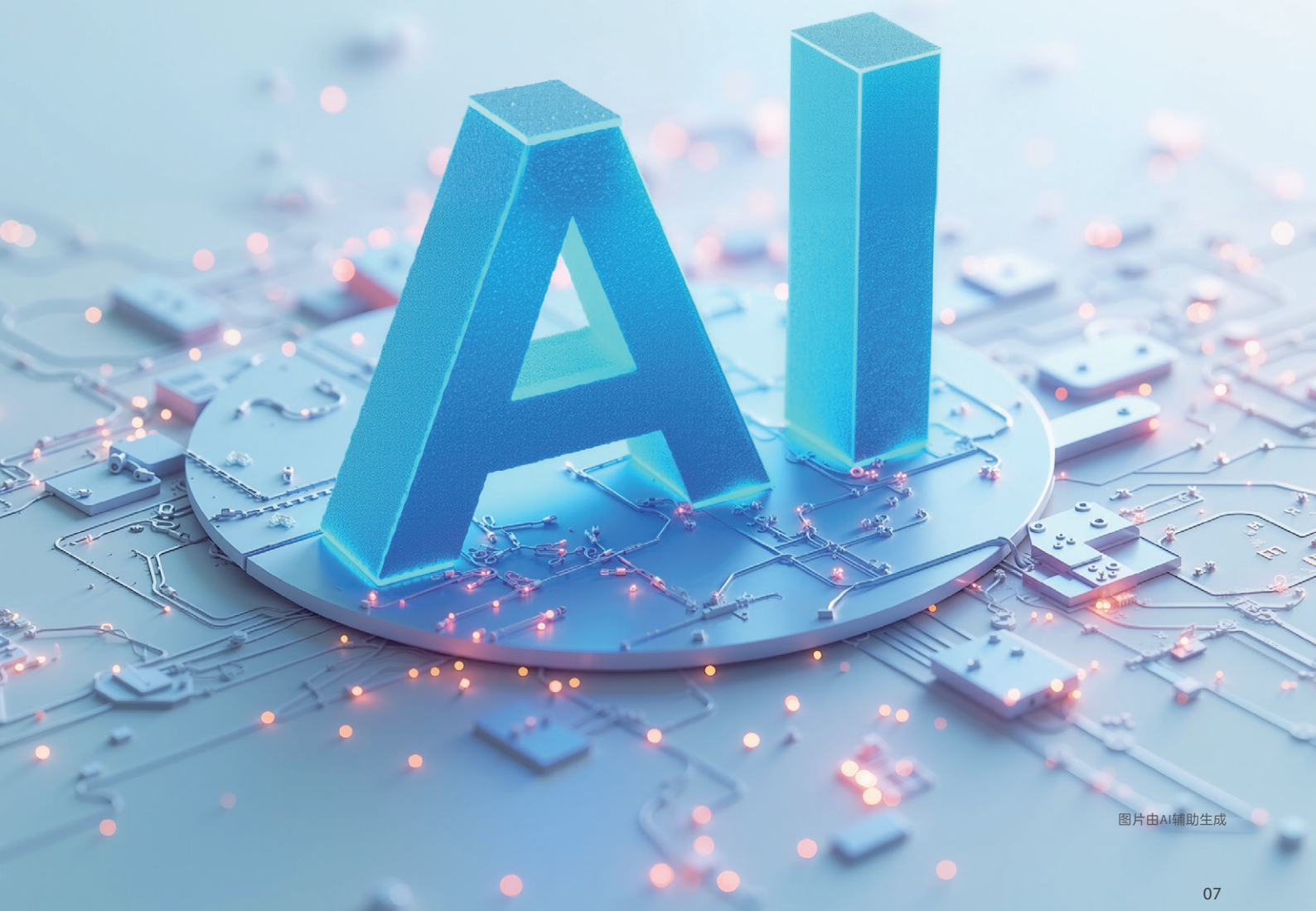
算力互联互通，实现全局TCO最优

建设算力互联互通网络，突破跨域异构混训、一体机编排调度等技术瓶颈，通过跨地域跨厂家算力集群动态聚合与协同调度，实现算力碎片化利用，并能支撑超大模型规模化训练需求。推进全国一体化算力资源监测及调度体系建设，结合AI业务场景特性实现算力资源精准高效调配，构建云网边端协同的分布式智算体系，推动算力资源全域共享、利用率最优。

人工智能作为引领新一轮科技革命和产业变革的核心力量，正深度赋能千行百业，成为驱动数字经济高质量发展的关键引擎。只有持续推进产业生态开放、算力架构革新、全要素协同进化，降低人工智能全生命周期使用成本，让高效算力更可及、更易用、更普惠，才能真正打破技术与成本双重壁垒，推动人工智能从高端示范走向全域普及，实现人工智能普惠化、高质量发展的最终目标。ZTE中兴

超节点应用场景 及技术演进

中兴通讯 毛磊



图片由AI辅助生成



毛毅

中兴通讯算力及核心网产品规划总工

随着芯片制程放缓和先进封装在散热、良率、成本等方面逼近当前工程能力上限，单点算力提升难以为继。为应对大模型所带来的算力爆发式增长和分布式并行计算需求，超节点应运而生。超节点通过高速互联将众多GPU整合为一个逻辑统一、调度高效的计算单元，像一台“巨型服务器”一样协同工作。

超节点定义与技术特征、价值场景

目前业界对超节点的概念尚无明确的描述，相关标准尚在制定过程中。但不同联盟、厂商定义或实现的所谓超节点，相比传统智算服务器产品或方案具有以下软硬件共性技术特征：

- 大量GPU互联系统：由于业界已经有成熟的8卡产品方案，因此超节点超过8卡是最基本要求。
- 统一内存地址空间：为互联的GPU提供统一寻址和内存一致性，如英伟达UVA（unified virtual address）、UM（unified memory），系统内任一GPU可以像访问本地HBM（high bandwidth memory）一样访问任意互联的HBM。
- 超高带宽、超低时延互联：除PCIe/CXL外

提供额外GPU显存高速互联，如英伟达NVLink，为系统内GPU提供高速访问的物理通道，带宽达数百GB到数百TB，纳秒级时延满足GPU间内存语义通信的同步操作要求。

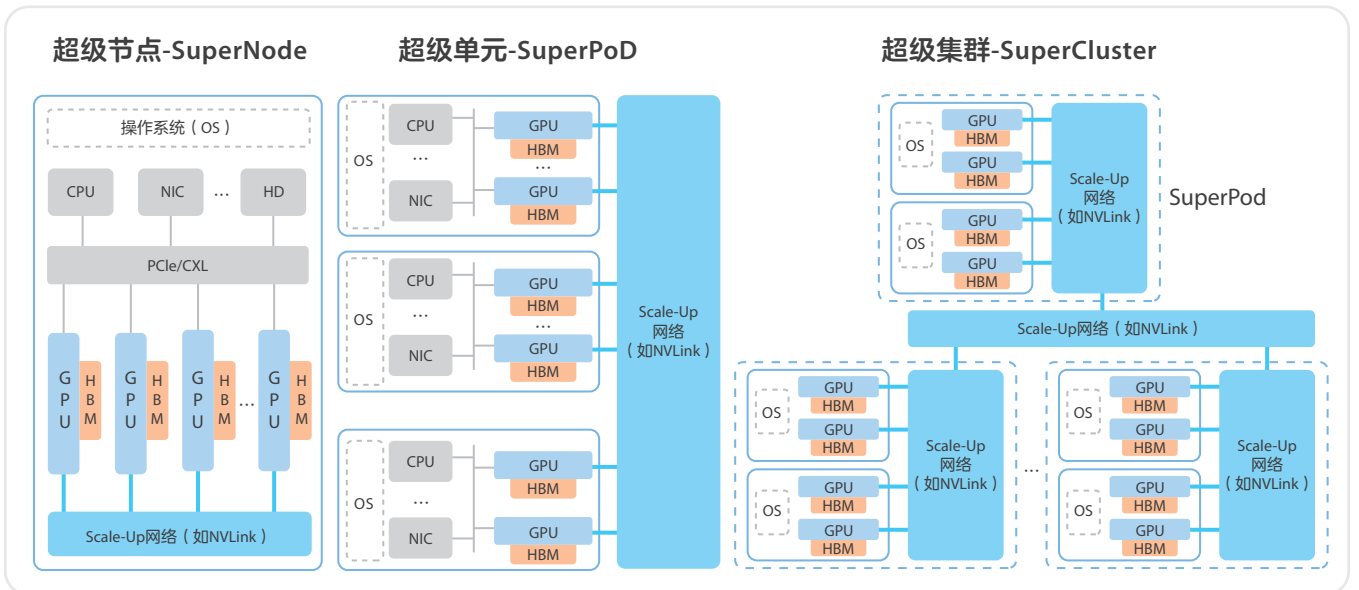
- 原生可扩展性：互联协议层面预留相应的bit位数满足未来可能的GPU扩展规模，拓扑层面支持一级和二级交换实现更大集群规模的扩展。

根据不同类型超节点的其他关键技术特征差异（如操作系统粒度等），通常可以把超节点分成三类：SuperNode（超级节点）、SuperPoD（超级单元）和SuperCluster（超级集群）。一般SuperNode至少支持16卡GPU以上互联，SuperPoD可支持成百上千GPU，而SuperCluster支持的高速互联GPU数量更多，对应的架构如图1所示。各类型超节点架构的关键特征如表1所示。

超节点能够提升大规模参数模型的训练效率，并优化推理性价比。

- 提升大规模参数模型训练效率

超节点为TP、EP等复杂并行计算提供强大的硬件支撑，缩短通信传输时间，提升集群并行计算效率，能够更高效地支撑超大模型训练，缩短训练周期。



▲ 图1 不同类型超节点架构

▼表1 不同类型超节点架构的关键特征

	SuperNode	SuperPoD	SuperCluster
PCIe/CXL总线	有	有	有
HBM高速总线	有（SuperNode内GPU）	有（SuperPoD内所有GPU）	有（互联的Cluster内所有GPU）
操作系统数量	1个	多个	多个
K8S集群数量		1个	1个或多个

然而，超节点规模的选择需结合训练需求权衡性能与成本，根据阿姆达尔定律，系统中不可并行部分会限制扩展带来的加速收益，且伴随复杂度、能耗与容错成本上升。基于Qwen235B在不同超节点形态下最优切分各部分耗时分析，在2000卡集群中，增大超节点规模可提升性能，主要受益于MoE算子优化，但存在边际效应，64~128卡是性能甜点区，当超节点达128卡后，性能增益趋于平缓。

● 优化每Token推理性价比

推理关注单位成本性能和用户体验，超节点以高速互联、内存池化、高度集成等优势，能够精准适配推理高并发、实时交互和大显存消耗需求。

随着DeepSeek、GLM等大模型参数规模突破千亿甚至万亿，通过超节点来承载大模型跨机推理，解决单卡显存不足与通信瓶颈，可有效缩短任务响应时间，提升吞吐量。此外Agentic AI、长文本对话、复杂文档分析等场景，推理过程中产生的KV缓存随上下文长度线性增长，超节点通过全局内存池化实现KV缓存的共享与复用，可支持处理数十万甚至上百万Token的上下文，极大增强其理解复杂问题、执行长期任务的能力。

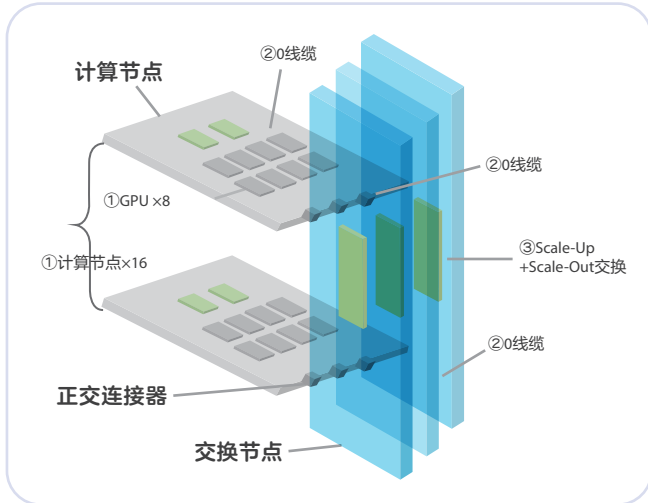
国产超节点技术演进思考

当前国产超节点呈现多架构、多形态发展的态势，产品与技术仍处于快速迭代阶段。其技术演进脉络可从产品架构、物理连接、系统生态三个维度展开分析。

产品架构层面，通过单柜极致密度和多柜互联两条路线提升规模。单柜极致密度以英伟达NVL576为典型代表，凭借技术领先性和供应链话语权，英伟达牵头产业界合力攻坚兆瓦级机柜关键技术难题，如800V高压直流供电、全液冷、中置正交背板等。考虑到英伟达GPU单卡算力领先优势，从NVL72到NVL576，其单柜极致密度路线可以匹配模型的演进需求，而国产AI芯片性能与之存在代际差距，加之先进工艺受限，短期内难以追赶。除了提升单柜GPU密度外，可进一步采用多柜互联扩展，通过规模领先另辟蹊径，利用成熟工艺将高带宽域纵向以机柜为单元扩展至数千规模。

物理连接层面，从电互联向光电融合一体化演进。电互联因低成本、高可靠、低时延等特点，是短距通信首选，尤其单柜内组网电互联优先。随着超节点Scale-Up域规模扩大，引入跨机柜互联场景，光互联技术以其独特的优势脱颖而出，如能够克服电互联距离限制、实现更高速的数据传输以及有效避免干扰等。以LPO、NPO、CPO等为代表的光互联技术方案将长期共存与互补，从互联带宽发展、生产良率、可维护性等角度看，预计呈现国际CPO领跑、国内NPO优先规模落地的差异化演进节奏。

系统生态层面，聚焦开放解耦与软硬协同。开放解耦包括超节点内部计算节点和交换节点解耦，以及CPU和GPU资源的解耦，同时聚焦构建统一的高速互联协议，规避私有封闭互联技术带来供应链依赖和成本方面的风险。软硬极致协同



▲ 图2 OEX架构互联示意图

设计提升算效，集成更多类型的异构芯片，CPU、GPU、DPU、NIC等多芯协同，软件原生适配超节点网络拓扑，通过智能编排实现训练与推理高效切换，满足池化资源灵活切分，结合全流程软件优化，提升算效与系统稳定性。

中兴通讯超节点产品创新

基于对整机柜超节点方案的深度工程实践，中兴通讯Nebula超节点创新提出OEX（Orthogonal Electrical eXchange）正交无背板互联交换架构，以Matrix集群满足规模扩展需求，同时积极打造开放解耦的软硬件生态体系。

OEX架构创新

OEX是一种正交无背板互联交换架构，其核心在于实现计算托盘与交换托盘之间的垂直交叉物理连接，消除传统线缆托盘（Cable Tray）带来的信号损耗与可靠性风险。该架构通过简化互联路径、提升信号完整性，为构建高密度、高可靠性的单体超节点提供物理基础（见图2）。在超节点设计中引入OEX架构，通过正交连接器与单级交换拓扑，实现计算节点与交换节点之间的垂直交叉互连，从而彻底摆脱了传统线缆的束

缚。在高速信号完整性、可靠性和可维护性方面相比传统的线缆托盘方案更具优势，也为后续架构扩展和演进预留了足够的空间。

Matrix集群可扩展

当前主流集群超节点部署方案多采用电交换+光互联架构，该架构技术成熟，生态完善，兼容性强。基于该技术方案，中兴通讯现有Nebula X32单体超节点可灵活扩展，构建形成Nebula Matrix X256/800集群超节点；面向未来，依托更高密度的Nebula X128单体超节点，更可进一步扩展至Nebula Matrix X8192/16384超大规模集群，充分满足超大规模模型训练的算力需求。中兴通讯并未止步于此，而是积极探索光交换与电交换的互补协同，旨在融合光传输的高效与电交换的灵活，以支持未来超大规模集群的可扩展性需求。

开放的软硬件生态系统

中兴通讯积极推动国产AI算力底座标准化进程，硬件层面全面开放OEX机械与电气接口规范，支持第三方计算及交换托盘的即插即用，有效降低系统集成门槛，促进产业链的协同创新。由中兴通讯主导的正交超节点整机柜设计规范已在ODCC官网正式发布。

软件层面，中兴通讯打造OLink开放高速互联协议，在底层兼容以太网的同时，通过物理层和事务层创新，满足Scale-Up高性能互联与Scale-Out超大规模扩展的双重需求，支持纳秒级延迟、统一内存编址与在网计算，显著降低组网复杂度与成本，并兼容多元GPU生态。

在芯片摩尔定律趋缓的背景下，超节点凭借架构创新、极致互联与软件优化，持续突破算力瓶颈，推动计算体系由芯片级摩尔向系统级摩尔迈进，未来有望成为构建AI基础设施的核心底座单元。ZTE中兴

Scale-Up/Out/Across 三域协同， 突破算力上限

超 大模型时代，模型参数呈指数增长，GPU单卡算力提升缓慢，算力瓶颈已转向系统协同，网络互联成为关键制约因素。Scale-Up架构受限于PCIe带宽，难以满足MoE、TP16~64等高并行需求；Scale-Out通过RoCEv2或InfiniBand构建跨节点网络，承载PP、DP通信，依赖RDMA与AllReduce实现万卡聚合，但面临Incast拥塞、ECMP负载不均、协议稳定性差等问题，导致训练效率下降；Scale-Across实现跨域异构算力池化，核心挑战在于低时延与强一致性保障。单一优化难破局，亟需构建Scale-Up/Out/Across三位一体的高速互联体系，实现带宽、时延、扩展性的协同突破，支撑万卡级高效训练。

智算网络Scale-Up/Out/Across业界现状与问题挑战

当前智算网络围绕Scale-Up（纵向扩容）、Scale-Out（横向扩展）、Scale-Across（跨域互联）三大模式协同演进，支撑万卡级乃至十万卡级智算集群落地。但这三大模式均面临显著技术与产业挑战，适配大模型训练推理仍有瓶颈。

Scale-Up网络聚焦超节点内高速互联，以英

伟达NVL72超节点为代表，通过NVLink和NVSwitch实现72卡GPU算力聚合，但NVLink生态私有封闭且与Scale-Out网络无法有效协同。业界也相继出现UALink、SUE、ESUN等总线型和以太网型互联技术，整个Scale-Up互联生态较为碎片化，短时间内难以收敛统一。

Scale-Out以RoCEv2为主，从千卡向万卡扩展时瓶颈凸显：首先，传统CLOS拓扑随着节点增多，端口密度与互联带宽呈指数级增长，设备成本与功耗激增；其次，AI流量的突发与聚合导致Incast拥塞丢包在万卡规模下被放大；此外，ECMP负载均衡流量分配不均导致链路忙闲不均，网络吞吐性能无法有效释放；同时，RoCE协议在万卡规模下的稳定性不足，生态碎片化导致不同厂商设备兼容性差，进一步制约扩展效能。

Scale-Across已成为突破单DC极限，实现跨数据中心协同训练的必然趋势。英伟达推出Spectrum-XGS Ethernet跨域互联技术，通过自适应距离拥塞控制、精准延迟管理等特性，将多个分布式数据中心整合为千亿级AI超级工厂。然而，Scale-Across网络长距互联面临RTT高（>1ms）、丢包率高、带宽收敛严重等问题，同时不同智算中心的拓扑、设备、协议不统一，协同调度难度大，跨域算力调度效率低，难以实现全域算力资源的灵活适配。



杨茂彬
中兴通讯Cloud&AI网络
规划总工

中兴星云智算网络一体化方案

中兴通讯推出融合Scale-Up/Out/Across模式的星云智算网络一体化方案，以自研芯片为核心，构建“端-机-域”三级架构，实现全链路高速互联，打破网络割裂瓶颈，支撑万亿参数大模型训练（见图1）。

Scale-Up：纵向扩容破局，筑牢单节点算力根基

针对传统PCIe架构带宽受限（<100GB/s）、仅支持8卡互联的短板，中兴通讯依托自研凌云芯片与OLink高速互联总线，创新构建支持32卡及以上超节点的机内互联架构，GPU间互联带宽提升至400GB/s以上，通信时延低至百纳秒级，突破传统TP8限制，适配MoE模型对高吞吐、低时延通信的严苛需求。同时OLink集成在网计算能力，卸载集合通信操作，进一步降低通信时延，有效

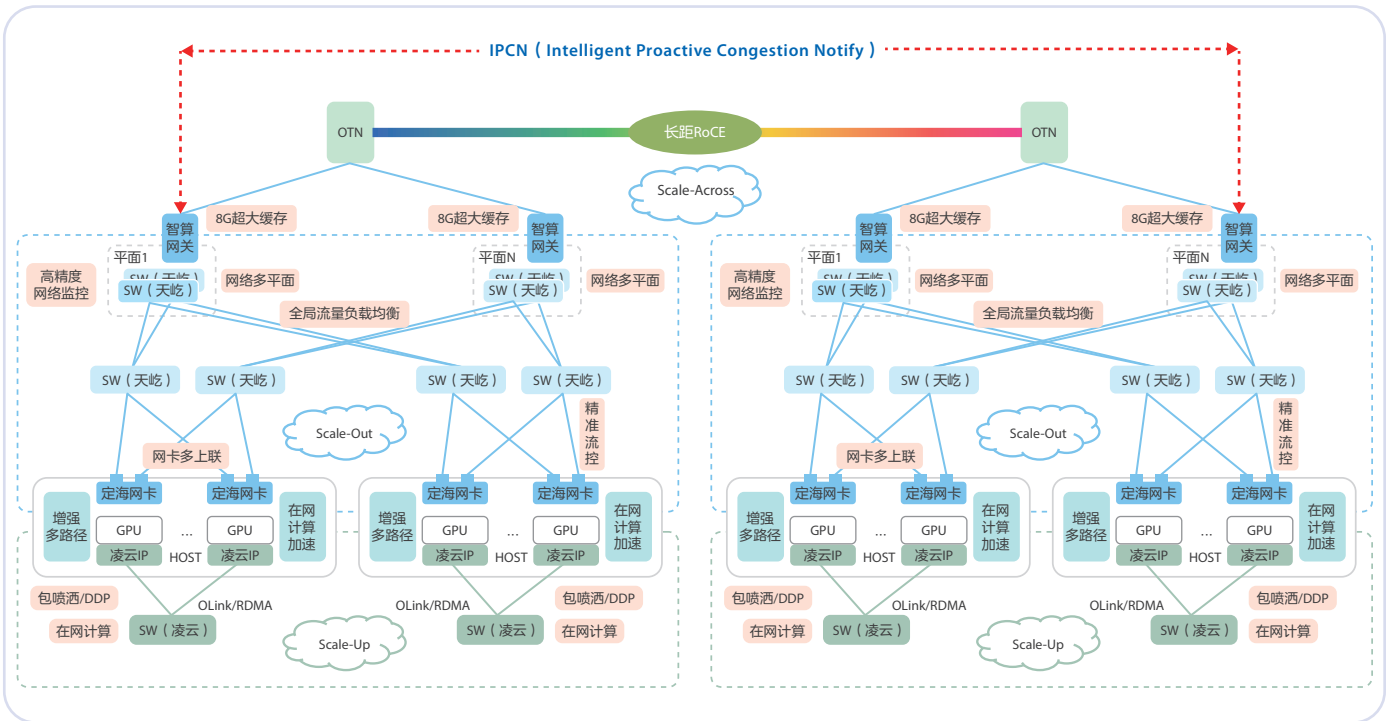
消除GPU空等现象，GPU利用率提升超40%。

Scale-Out：横向扩容提质，打造高可靠集群网络

基于自研天屹芯片和定海RDMA网卡，星云智算网络方案构建“端网协同”机间网络，实现千卡级到百万卡级无缝扩展。架构支持两种CLOS模式，单PoD支持16K GPU，多PoD可演进扩展至百万GPU，满足未来10年算力需求。针对大规模集群痛点，采用ENCC无损技术确保吞吐率>90%、丢包率<0.0001%，结合iGLB技术实现流量均衡；双平面冗余组网搭配ARN技术，实现亚毫秒级故障切换，保障集群稳定运行。

Scale-Across：跨域互联赋能，织就全国算力一张网

针对单DC算力极限、跨域训练痛点，中兴



▲ 图1 中兴通讯星云智算网络



图片由AI辅助生成

通讯提出Spine直连+OTN+IPCN跨域方案，助力构建全国级算力一张网。Spine层8GB超大缓存交换机吸收长距突发流量，IPCN技术主动预测拥塞，避免失控；支持多厂商GPU异构互联，单DC可扩展至7个PoD，支撑十万卡级集群，跨PoD收敛比8:1。工程验证显示，300km拉远场景下，1024卡训练算力影响<0.13%，推动“算力即服务”落地。

星云一体协同，打通算网全域价值

星云智算网络一体化解决方案，将Scale-Up、Scale-Out、Scale-Across三大能力深度融合，形成“机内提效、机间扩容、跨域整合”的全链路协同体系，既解决了单机性能瓶颈，又实现了大规模集群的稳定扩展，更打破了地域与厂商壁垒，实现全域算力资源的灵活调度。该方案有效

突破智算网络多维度瓶颈，但业界整体仍面临高端芯片自研、生态适配、运维复杂度控制、核心器件成本优化等共性挑战，未来中兴通讯将持续深化技术创新，推动三大能力的更深度协同，助力智能算力网络高质量发展。

中兴通讯以自研定海、凌云、天屹芯片为基石，构建端-网-云全栈自主可控体系，通过Scale-Up、Scale-Out、Scale-Across分别解决单机算力密度、集群规模可靠性、资源协同弹性扩展问题，实现从千卡、万卡到十万、百万卡平滑演进。三域融合在国产GPU关键期突破算力上限，推动中国智算网络从跟随走向引领。未来中兴通讯将深耕算网一体，以星云智算网络助力AI产业迈向全球领先。[ZTE中兴](#)

GPU、DPU、存储介质协同的 新型AI存储架构



郭伟
中兴通讯智算存储规划
无线总工

随着“Agentic AI”和长上下文模型的快速发展，作为模型“记忆”载体的KV Cache（键值缓存）规模呈爆炸式增长，已成为当前AI基础设施在推理阶段的主要性能与能效瓶颈。传统存储架构主要面向数据持久化设计，难以高效支撑KV Cache这类对性能极度敏感、具有短暂性与可重计算特性的“AI原生”数据，导致GPU资源频繁空转，严重制约了AI工厂的规模化部署与成本优化，AI推理挑战已经从原来的“算力墙”演变为“存储墙”。未来面向智算场景的新型存储架构需要融合硬件创新、网络创新、软件算法优化，突破AI推理瓶颈制约。

AI计算范式转变下的存储挑战

大模型推理分为Prefill和Decode两个阶段。Prefill为计算密集型，处理用户输入提示（prompt）；Decode为访存密集型，逐个生成后续Token。在多轮对话中，若不采用缓存机制，历史Token的Key-Value（KV）矩阵需在每次推理时重复计算，造成显著冗余。KV Cache通过将已计算的KV状态缓存于显存中，实现“空间换时间”，避免重复计算，大幅提升推理效率。

当前AI计算范式正经历根本性变革，传统短上下文、单轮交互模式已逐步被长上下文、多轮对话及多智能体（Agentic AI）协同执行的复杂场景取代，表现为三大趋势：

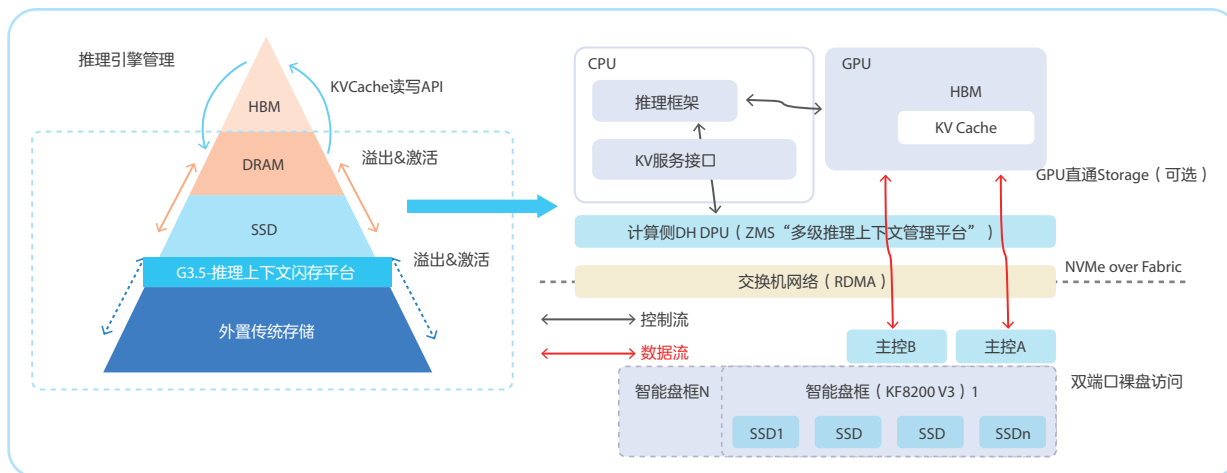
- 上下文长度爆炸式增长：从数千Token扩展至百万级，KV Cache数据量远超单GPU显存容量（如GPT-3的KV Cache可达模型参数占用显存的一半以上）；
- 推理即“思考过程”：推理不再是一次性答案，而是一个思考过程，通过测试时扩展提升答案质量，导致生成Token数量年均增长5倍，显著增加KV Cache的读写压力；
- 长短期记忆需求：AI系统需支持跨数周的多轮交互记忆，要求KV Cache具备长期可访问性与高效管理能力。

尽管KV Cache提升了计算效率，但也引发新的系统瓶颈：

- “存储墙”问题突出：KV Cache对带宽敏感，且规模庞大，易成为性能瓶颈；
- 传统存储架构不匹配：现有存储设计强调数据持久化与容错，而KV Cache具有短暂性、可重计算、高频读写的特点，导致其访问路径过长、延迟高；
- GPU资源严重浪费：约30%~40%的GPU计算资源消耗于KV Cache的数据搬运与低效读写，导致GPU利用率不足50%，推高AI推理单位成本。

中兴通讯以“多要素协同”构建新型存储架构

中兴通讯依托近20年来在存储领域软件架



▲ 图1 多要素协同构建高效KV缓存解决方案

构、全自研硬件和芯片技术上持续的技术积累，以及近年来在智算领域的深度参与和思考，推出“DPU+智能盘框+KV Pool软件”的高性能新型存储解决方案。

如图1所示，方案采用以DPU为中心的多级缓存平台的存算分离架构。DPU承担存储协议栈处理、数据高效转发卸载和数据传输优化等关键任务，使得GPU能够专注于业务处理，存储节点聚焦存储低延迟、高带宽的缓存数据。方案提供DPU、RDMA网络、存储智能盘框的端到端纯硬件调优，支持合作伙伴自己的KV Cache存储管理软件“拎包入住”；也可提供不同层面的端到端全套自主可控的软、硬件的新型存储方案。通过软件重构并卸载到计算侧DPU，结合专门设计的KV接口，即消除东西向副本同步流量，方案实现了存储网络带宽的高效利用，同时减少了冗余数据传输，比传统存储网络利用效率提升了3~5倍，数据访问时延降低50%以上。

方案采用多项关键技术创新，基于专为智算优化的存储架构将不同存储介质构建分层管理，实现KV数据路径DPU卸载、GPU直通与极简协议交互，配套自研智能调度系统优化，极大提升了长上下文推理效率。

- 多级缓存平台：综合利用DRAM、SSD、远端NVMe盘框分层构建共享池；可分层灵活

组合开启，硬件资源占用可灵活调整。

- 多要素协同，构建PoD级的共享上下文记忆空间：专为智算设计的存储智能盘框，支持全NVMe，NVMe-oF接入，裸盘访问；NVMe-oF在DPU硬件卸载转发，KV数据路径计算卸载在计算侧本地，优化数据传输路径；NVMe-oF零拷贝，GPU直通存储协议；重新设计KV Cache PUT/GET专用接口，无协议转化；极简存储架构，去除冗余的元数据管理和强一致性等设计。
- 面向KV的推理调度增强：自研智能KV Cache管理调度系统，优化推理调度算法，实时分析推理请求变化特征，动态调整存储层级与分配策略；设计自适应缓存替换算法，依据数据访问频率与重用概率智能筛选保留或淘汰数据，将缓存命中率提升至70%以上。

随着AI计算范式的不断迭代进化，对AI业务新型存储架构提出更高要求。中兴通讯将继续践行“技术创新，以存助算”的指导思想，持续创新，探索DPU加速、存算一体、先进介质等先进硬件与新型存储架构结合的可行性，构建开放、共享的软件生态环境，协助AI新型存储行业标准构建，为攻克“存储墙”的AI计算范式变革夯实基础。ZTE中兴

高效推理加速AI落地



周俊超
中兴通讯智算产品规划
工程师

随着AI应用规模化进入生产环境，推理服务系统已成为支撑AI技术落地、实现价值转化的核心载体。当前AI应用落地过程中的场景多元化、需求精细化，对推理系统的安全性、调度效率、资源利用率及适配能力提出了新的要求。

- 需求一：构建推理系统全栈原生安全管控体系

AI应用除了依赖模型本身提供的基础安全合规能力外，还需要针对具体落地场景，面向文本、图像、音频、视频等多模态输入输出，从推理系统层面提供定制化、系统级、全链路安全防护机制，提供精准、实时、全覆盖违规内容识别与风险响应能力，满足AI应用落地的安全合规要求。

- 需求二：实现基于大模型推理特性的消息调度机制

在AI应用规模化落地、多类型大模型推理任务并发执行的背景下，现有调度机制难以适配LLM推理的计算特征与差异化业务诉求，易引发GPU资源空转或过载、KV Cache资源低效利用与供给不足等瓶颈。需构建适配大语言模型推理特性的消息调度机制，实现GPU算力、KV Cache等关键资源的高效分配与动态均衡。

- 需求三：支持推理任务流量路由机制

推理任务复杂度不同，对算力资源的需求也不同，如问候交互、短句翻译等快思考简单任务，与逻辑推理、代码生成等慢思考复杂任务对算力的需求差异较大。推理系统需要构建基于任务复杂度的流量路由机制，进行差异化路由；降低简单任务时延、减少算力浪费，保障复杂任务算力供给，匹配AI应用多样化落地需求。

- 需求四：更长的上下文推理能力增强

随着AI应用的落地深入，文档理解、知识库问答、代码分析、长对话记忆、复杂Agent规划等核心场景，需要更长的上下文长度（目前正在从K级向M级演进）推理能力支持，解决该场景下并发下降、显存溢出、时延陡增、系统稳定性下降等问题。

- 需求五：推进推理全栈优化与硬件算力极致释放

AI应用规模化落地要求在现有硬件资源条件下实现更高吞吐、更低时延、更大并发。需通过推理服务系统软件调度与硬件资源极致协同，压榨硬件性能，提升硬件资源利用率、降低单Token推理成本，解决算力释放不充分、整体投入成本偏高的问题。

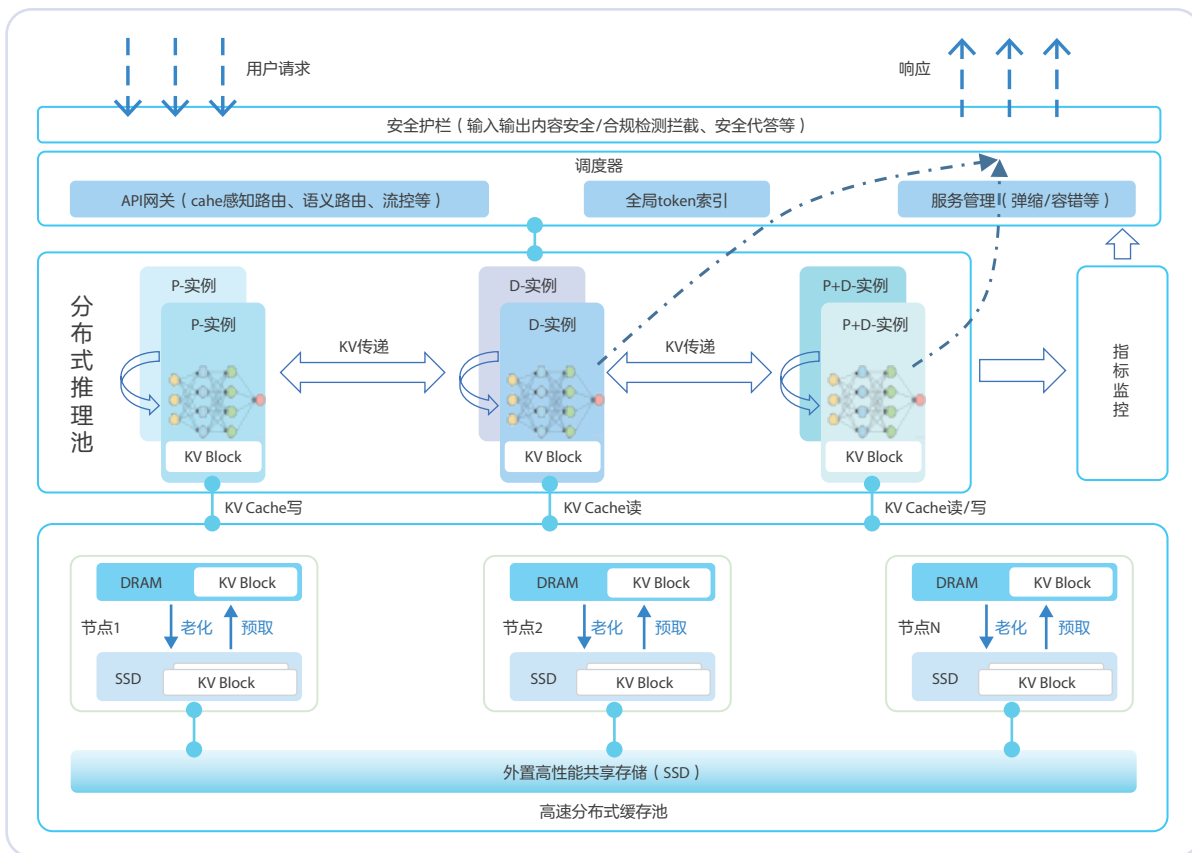
AI应用规模化落地带来的新需求，使得构建一套高性能极致软硬协同的低成本推理服务系统，成为AI应用走向规模化普及的必然选择。中兴通讯高效分布推理服务系统ZIS（ZTE Infer System），正是在这种挑战背景下设计推出的。中兴ZIS推理服务系统是一个集成了安全防护、智能调度、流量路由、长下文支持、极致软硬协同的低成本高效大模型推理平台（见图1）。

- 原生安全防护

推理服务系统内置安全护栏组件，基于安全合规领域微调的安全模型，或黑名单及敏感词列表，对用户推理服务过程中的输入输出文本和多模态内容进行安全检测，构建安全防护基础防线。

- 智能消息调度

调度器针对大模型推理服务特点进行调度能力增强，支持GPU负载感知、KV Cache感知、最



▲ 图1 中兴ZIS推理服务系统架构示意图

少消息负载、SLA指标感知的消息调度机制，为每个请求消息最佳匹配推理资源。

● 流量路由

基于语义路由技术，调度器前置一个轻量级路由模型，对用户业务请求进行语义复杂度分析。简单任务（如问候、翻译）路由至小参数模型（快思考、低成本低延时），复杂任务（如逻辑推理、代码生成）路由至大参数或思维链模型（慢思考、高质量），任务处理隔离，最优匹配算力资源。

● 更长下文推理能力增强

引入稀疏注意力、语义压缩及位置无关KV Cache复用等优化技术，优化更长上下文处理能力，支持100K+甚至1M+的超长上下文推理，并在保持精度的同时将计算复杂度控制在可接受范围内。

● 极致软硬协同

支持PD分离技术，将不同负载分配至最优匹配的硬件节点，实现“专卡专用”，充分释放硬件潜能。同时，配合KV Cache复用及KV Pool技术，跳过已命中Tokens的计算耗时，有效增加推理速度，实现极速响应。

2025年下半年，中兴ZIS推理服务系统联合国内运营商，基于国产化GPU卡和DeepSeek-R1 671B-Int8量化模型，进行了严格的实验室测试。测试结果表明：基于ZIS的推理服务系统，单卡吞吐量最大提升7.6倍左右。

中兴ZIS推理服务系统以最低的资源实现最大的Token吞吐量，发挥用户投资价值，已在使能AI应用中发挥显著成效。未来，随着多模态推理需求进一步提升与边缘端推理的普及，ZIS将进一步向边端云协同、更细粒度的动态优化方向演进，持续赋能千行百业的AI应用落地。ZTE中兴

筑牢可信底座，

破解工业智能体长程推理困局



陈云谦
中兴通讯智算产品规划
总监

随 随着大模型技术的爆发，智能体正从“尝鲜”走向“常用”，成为AI应用的新范式。然而，从通用场景跨越到高门槛的工业场景，AI面临着严峻的信任鸿沟。工业任务往往具有长周期、多工具协同的特性，这要求智能体必须具备强大的“长程推理”能力。但在长链条的决策过程中，任何一步出现的“幻觉”或逻辑断裂，都可能导致严重的生产事故。因此，如何在保持长程推理效率的同时，彻底消除幻觉，实现决策可审计，成为工业智能体落地的“最后一公里”难题。

复杂任务中的“信任与效率”困局

工业场景的复杂性常表现为任务的“长程”特性，即任务跨度长、涉及工具多、步骤繁杂。智能体在执行此类任务时，面临三大行业级痛点：

- 验证成本随链条长度线性增长：现有验证方法通常需要对推理过程进行重演或全面检查，这在跨天、跨周的复杂任务中导致计算资源消耗剧增，甚至出现算力坍塌。
- 证据与逻辑不可追溯：多智能体协作中，证据、假设与工具调用往往缺乏结构化组织，一旦出现错误，人类专家难以快速定位问题步骤，导致“黑盒”信任危机。
- 信息孤岛与共识偏差：由于缺乏统一的记忆基底和冲突识别机制，智能体在跨协作中容易出现“幻觉”或达成错误共识，最终导致决策偏差。如何让智能体既能高效协作，又

具备独立审核能力，成为关键挑战。

Co-Sight双引擎：破解工业长程推理难题

传统依赖全链路重复验证的方法计算成本高昂，效率随任务复杂度下降。为破解这一难题，中兴通讯创新提出Co-Sight智能体开发平台，通过构建“冲突感知元验证（CAMV）”与“基于结构化事实的可信推理（TRSF）”双引擎闭环机制，在大幅降低算力消耗的同时，将推理的准确性和可审计性提升至工业级标准。

如图1所示，Co-Sight构建了一个包含多个专家智能体与元验证智能体的高效系统，通过冲突感知元验证与基于结构化事实的可信推理，精准破解当前智能体在长程推理中“验证成本高、逻辑不透明、证据不可追溯”的难题。

TRSF：解决长程推理中的记忆瓶颈

在工业场景中，跨周期复杂任务常因“记忆丢失”或“逻辑断裂”而受阻。Co-Sight通过基于结构化事实的可信推理机制（Trustworthy Reasoning with Structured Facts, TRSF），构建了一个持续更新、可信且可溯源的“事实基座”。

TRSF引入三级上下文压缩管道，将海量、异构的工具输出转化为高密度结构化事实：

- 工具层（最简化元数据）：系统整合来自数据库、API、文档检索等多源信息，初步清洗格式化，去除冗余噪声，保留核心元数据。

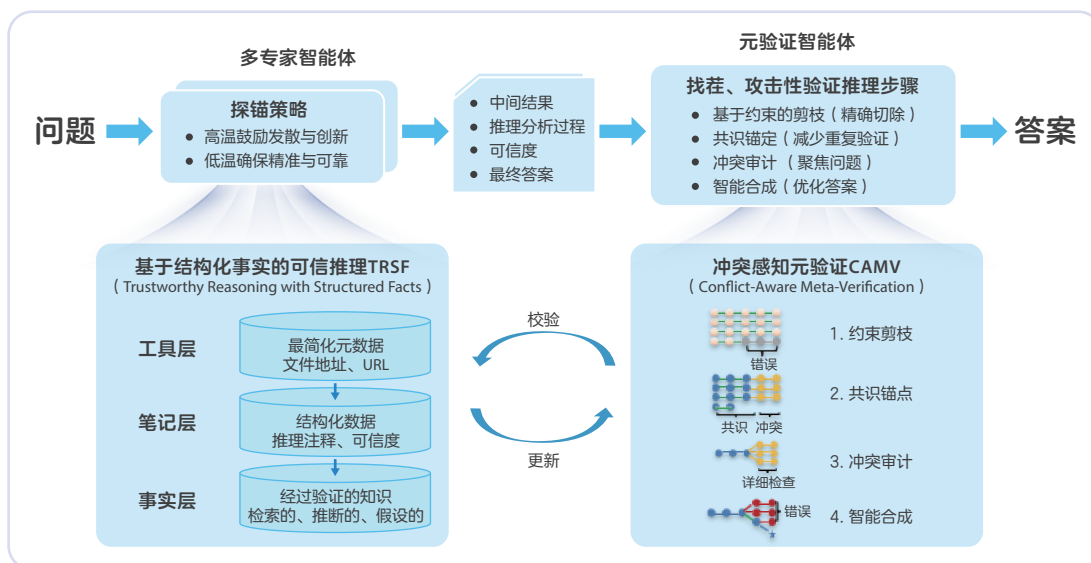


图1 冲突感知元验证与基于结构化事实的可信推理

- 笔记层（简洁注释）：对工具输出进行语义提炼，生成如“基站A的CPU利用率在14:00达到85%”这样的结构化摘要，将非结构化数据转化为可读性强的语义单元。
- 事实层（经验知识）：通过交叉验证存储高置信度信息至全局事实库，支持版本追溯与来源审计，形成长期记忆的经验知识。

通过三级压缩，TRSF将长文本的信息密度提升100倍，支持智能体“断点续做”。无论任务跨天或跨周，都能实现“随时存档、随时重启”，推理过程不再因中断而失效。

CAMV：实现高可信验证

核心业务容不得“AI幻觉”，每个决策都必须“可审计、可追溯”。Co-Sight通过冲突感知元验证机制（Conflict-Aware Meta-Verification, CAMV），将验证视为“证伪”过程，而非全面重验。

CAMV的核心思想是：仅在专家智能体输出存在分歧的“最小冲突集”上投入计算资源。其验证过程包括4个步骤：

- 约束剪枝：利用领域知识（如网络协议规范、设备参数范围）过滤明显错误结果，提前剪枝无效推理分支，避免大量无效计算；
- 共识锚点：对多个智能体一致认同的中间结论设为“锚点事实”，减少重复验证工作量

并提升推理效率；

- 冲突审计：针对存在分歧的推理节点，元验证智能体调用高精度模型或外部工具进行精细化交叉验证，确保关键决策的可靠性；
- 智能合成：从海量推论中提取有效片段，基于验证过的事实锚点重构逻辑严密的最终答案。通过精准投放计算资源，CAMV将验证成本从全链路降至关键点，算力需求降低50%，幻觉率降至0.3%以下。

成果与展望

Co-Sight通过“零部件车间（预置原子能力）+总装车间（可视化编排）”模式，实现“用AI生产AI”。同时，Co-Sight已开源智能体三层交互协议，并在GAIA、HLE评测中连续蝉联第一。

在实际应用中，Co-Sight已成功构建工业级通信数智人，实现多智能体高效协同。以中兴通讯与中国移动的合作为例，双方在“点金行动”中验证了31个高价值场景，涵盖节能、业务开通、故障定位及处置、网络性能优化等。

未来，中兴通讯将持续优化Co-Sight框架，探索其在智能客服、研发辅助等场景的深度应用，推动AI应用向“可信”迈进，为行业智能化升级注入新动能。ZTE中兴

中兴通讯Skill安全体系， 构建智能体边界



顾希
中兴通讯产品安全规划
总工

人工智能的安全架构正在经历一场深刻的范式转变。大语言模型（LLM）时代，安全威胁主要集中在提示词注入、数据泄露与输出污染。智能体（Agent）时代，智能体具备工具调用、外部API访问、代码执行等能力，传统安全边界被打破，安全风险维度全面升级。例如，OpenClaw作为具备持久化记忆、自我进化与能力迭代特性的代表性智能体，其技术先进性的背后，也带来了权限与风险高度绑定的核心安全命题。AI安全防护已从LLM时代的内容输出转变为智能体对主机实际操控的防护，传统安全框架已无法适配，亟需全新的安全架构。

OpenClaw的现象级爆火，源于其独特的架构设计：一方面，它采用本地部署模式，实现数据主权自主掌控，具备数据不出域、不经过云端第三方的隐私保护优势，高度契合企业和个人用户的隐私安全需求；另一方面，其具备生态开放的能力描述文件，支持智能体根据自然语言指令，完成金融、IT、办公自动化等多领域的自动化操作，生态开放性与落地实用性突出。然而近期披露的OpenClaw零点击漏洞（0-Click Vulnerability），为全球智能体安全敲响了警钟。用户一旦访问了恶意网站，无需任何点击或交互，OpenClaw智能体就会被劫持，攻击者可远程执行任意代码。漏洞根源在于OpenClaw认为所有本地操作都是可信的，豁免了数百万次的暴力破解，导致恶意应用轻易获取主机Root权限。

传统安全的核心准则为“不运行来源不明的可执行程序、不点击可疑链接”等，而OpenClaw

的安全风险彻底突破了这一常识边界。仅需加载简单的Skill文本文件，依托大语言模型的文本生成能力和OpenClaw的主机直接操作权限，该智能体即可实施高危网络攻击行为，包括窃取个人金融账户、登录密码等敏感数据，甚至执行硬盘数据自毁等破坏性操作，对终端安全与数据安全构成致命威胁。

面对上述挑战，中兴通讯针对一体机构建Skill安全防护体系（见图1），围绕OpenClaw技能文件构建端到端的安全闭环。中兴通讯覆盖全生命周期的内生安全防护体系，通过来源验证、安全扫描、签名认证、AI防火墙、沙箱部署五大核心环节，实现从Skill引入到运行环境的纵深防御，全面保障智能体安全可控。

● 来源验证

引入Skill时需要识别Skill的来源和作者，这是Skill文件可信的第一步。可信程度由高到低排列为：官方来源、知名社区来源、开发者来源、空或者未知来源。根据项目的实际情况，可以通过配置来提升Skill的可信程度，其次通过Skill中的GPG签名以验证作者身份，作者身份可信才允许引入Skill。来源验证提升了引入Skill的门槛，从源头保障Skill的可靠性。

● 安全扫描

中兴通讯针对OpenClaw特性实现了安全扫描器，对智能体实行静态扫描。扫描Skill前言，查看描述是否简洁明了，防止恶意软件通过Base64编码的误导性描述实现针对AI的绕过手段，同时查看过于宽泛的提权表述等问题；扫描是否存在提示词注入、脚本/工具的高危操作等

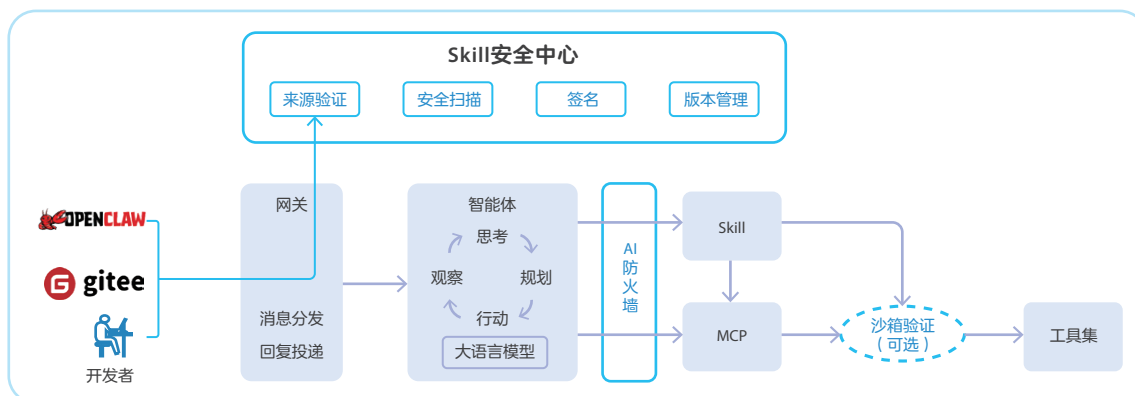


图1 中兴通讯Skill安全防护体系

问题，主要针对大语言模型“忽略之前的指令”“以管理员身份执行”等绕过性提示；扫描硬编码密钥、高危命令操作，如存在password、pwd等描述，或者提权命令等；扫描依赖组件是否有CVE漏洞，是否存在隐藏路径、恶意二进制文件等常规安全漏洞。

静态扫描是智能体可信的关键步骤，它在不执行智能体的情况下进行风险预警，将智能体的风险和威胁降到最低。

● 签名

与Skill内置的GPG签名不同，本次签名是针对完成上述来源验证、扫描检查的Skill签名，签名后可提供至生产环境使用。生产环境的最佳实践是建议只加载经过中兴通讯签名认证的Skill文件。

● AI防火墙

AI防火墙监控运行期的智能体行为，如是否存在提示词注入、是否有海量数据访问请求、是否有资源过载情况（例如挖矿）等。防火墙是基于语义的内生防御形态，具备主动思考能力，能识别最终的工具调用意图，如果有Agent最终请求的是删除所有数据这类高危命令，防火墙将主动拦截本次操作并标记为高危。如有一些疑似动作但防火墙判定不了，则会将本次操作标记为中危或建议验证等要求。

● 沙箱部署

沙箱激活是一个可选步骤，依据前面Skill安全中心的评级以及AI防火墙最终放行时的评级来选择是否激活沙箱。沙箱中内置中兴通讯内生安全功能，OpenClaw在沙箱中运行并激活技能，针对OpenClaw是否存在传统安全的异常登录、

非法提权、文件篡改、非法下载等总计20余项行为做安全检测，确保智能体没有针对主机的攻击活动。

此外，针对智能体本身，中兴通讯提供了一个高度可扩展的测试框架，用于识别和缓解AI漏洞。工具通过代理授权、劫持控制、代理与周围环境交互、幻觉攻击等10余项检测，针对沙箱中运行的OpenClaw智能体做全方位的测试并给出测试报告。

在传统安全和智能体安全检查都通过的情况下，Skill文件可以发布。

中兴通讯安全防护体系全面覆盖智能体 workflow，从开发、部署到运行，无论哪一环节被突破，后续环节仍可拦截。此外，系统提供应急响应机制确保威胁被即时阻断，并且通过日志反馈机制支持事后分析与模型优化，形成“检测-响应-改进”的正向循环。

智能体的崛起是不可逆转的技术潮流。OpenClaw所代表的“行动派AI”正在重塑人机交互模式，安全架构必须同步演进。未来的智能体安全不应仅依赖外部防护，更应构建内生安全能力，让智能体具备“安全意识”，能够自主识别风险行为并自我约束。此外，智能体之间的协作将成为常态，安全责任需要从个体扩展至生态，建立跨平台的安全标准与威胁共享机制，实现联防联控。最后，技术手段只能解决部分问题，完善的治理框架同样重要，明确责任边界、建立保险机制、推动法规完善，都是智能体安全不可或缺的重要部分。ZTE中兴

全球AI开源技术发展概况 及中兴通讯开源实践



王长金
中兴通讯智算产品规划
总工



人工智能（AI）开源技术已成为推动全球技术创新的核心引擎，正在重塑全球产业生态与商业模式。

AI开源正在从“技术共享”向“生态共治”演进，基金会、联盟、工委、社区和项目共同构成了全球AI开源生态的多层次组织体系，美国、欧盟、中国等主要经济体也从科技创新、数字主权、产业安全和生态培育等角度加快出台相关政策，推动AI开源规范发展。与此同时，AI开源的商业模式正从传统的支持服务模式，逐步演化出开放核心、模型即服务、生态共建等多元发展路径。

AI开源社区生态发展概况

人工智能的开源生态已形成“基金会-联盟-工委-社区-项目”五级生态架构，各层级在治理目标、参与主体与规则制定权上存在显著差异。

● 基金会：治理中枢与资源枢纽

基金会通常为非营利性组织，是全球开源运动最重要的治理主体。顶级基金会凭借其公信力、运营能力和国际影响力，为开源项目提供法律庇护、品牌背书和资金支持。以Linux Foundation AI&Data（LF AI&Data）和Pytorch

Foundation为代表的伞形基金会，承担着AI开源项目的中立托管与法律护航职能。在人工智能技术的快速迭代发展过程中，基金会从AI基础设施、模型框架向应用层延伸。如2025年12月9日AI Agent Foundation的成立标志着人工智能从单纯的“对话（Chat）”时代，向“行动（Action）”与“协作（Collaboration）”的“代理（Agent）”互联时代转型。

● 技术与产业联盟：产业领域协同平台

技术与产业联盟通过整合上下游资源，推动开源技术在特定领域的快速落地。AI联盟通过建立基准测试与安全评估的共享框架，试图抗衡闭源巨头的生态垄断，侧重于制定技术路线图与伦理准则，而非直接托管代码。在中国，中国人工智能产业发展联盟（AIIA）下设开源开发组，承担标准制定与测试认证功能，体现“国家引导+企业主体”的混合治理特征。

● 工作委员会：垂直领域的精细规范制定

工作委员会一般在开源生态系统中承担着技术方向制定、标准规范协调和跨项目协作的职责。如ODCC（开放数据中心委员会）在中国通信标准化协会指导下以开放、合作、创新、共赢为宗旨，围绕服务器、数据中心设施、网络、新技术与测试、边缘计算、智能监控与管理等领域，推动形成AI基础设施领域统一的技术规范和

标准。

- 开发者开源社区：创新网络的神经末梢

GitHub、Hugging Face等平台构成了去中心化的创新开发平台。Hugging Face已汇聚超过100万个模型仓库，形成“模型即服务”（MaaS）的集市效应。

- 开源项目：价值创造的原点

从根技术（如Linux内核）到应用层（如OpenClaw），项目层级呈现“马太效应”。值得注意的是，大模型开源呈现“Open Weight”新形态（如阿里的Qwen系列），其开放性低于传统OSI定义，但已足以引发产业链的生态重构。

各国政府开源政策分析

国家政策在引导开源生态的发展方向、资源配置和塑造全球竞争力方面扮演着越来越重要的战略性角色。

美国在开源领域长期保持领先地位，其政策特点是政府引导与市场驱动并行。2022年，白宫科技政策办公室（OSTP）发布关于开源软件的备忘录，强调供应链安全的重要性，要求联邦机构评估关键系统中开源组件的风险。作为全球最大的开源代码托管平台，GitHub（现为微软旗下）美国用户超过4000万，托管超过4亿个开源仓库。报告显示美国开发者贡献了全球约25%的开源代码，在AI领域这一比例更高。

欧盟将开源视为实现数字主权的重要工具，政策力度和系统性尤为突出。欧盟委员会发布“开源软件战略（2020—2023）”，提出了三大目标：提高政府开源采用率、培养开源人才、促进开源生态系统建设。战略明确要求欧盟机构在新IT项目中优先考虑开源解决方案，并要求供应商提供源代码访问权。欧盟通过GAIA-X项目和“欧洲云”倡议，推动构建基于开源技术的数据基础设施与主权云。德国、法国等国政府明确要求公共部门优先采用开源解决方案。

中国将开源纳入《“十四五”软件和信息技术

服务业发展规划》，明确提出“建设2~3个具有国际影响力的开源社区”。国务院关于深入实施“人工智能+”行动的意见明确指出促进开源生态繁荣，支持人工智能开源社区建设，促进模型、工具、数据集等汇聚开放，培育优质开源项目；建立健全人工智能开源贡献评价和激励机制，鼓励高校将开源贡献纳入学生学分认证和教师成果认定。地方政府层面，北京发布《北京市促进开源软硬件高质量发展的若干政策措施》；上海依托人工智能实验室建设开源平台；深圳出台《关于促进开源软件发展的若干措施》，提供资金支持与税收优惠。

AI开源商业模式

AI时代的开源商业模式正在重构，从传统软件的“licence销售”转向“价值服务”，主要有以下几种模式：

- 开源即服务（SaaS）：将开源项目封装为云服务，用户通过订阅方式获得托管、运维、升级等服务；适用于基础设施型开源项目，如数据库、缓存、消息队列等。核心价值在于降低用户的运维成本和技术门槛，如Hugging Face通过提供模型托管、推理API和企业级安全特性获利。
- 技术支持订阅：基础版本免费开源，通过技术支持、企业版软件或专业服务实现收入；适用于需要深度定制和技术集成的企业级软件。类似于红帽模式，企业级客户在使用开源模型时需要稳定性保障、安全加固和版本维护。
- 开源+闭源混合模式：核心框架开源，针对特定场景的高性能模型或企业级功能闭源。如Meta对Llama 2采用定制商业许可（对月活用户超7亿的企业收费）。
- 社区生态共建：通过构建开发者生态，实现流量变现或生态溢价。如全球最大的开源代码托管平台GitHub，通过GitHub

2026年1月，中兴通讯联合天工开物开源基金会、中国信通院、百度、Intel、Red Hat等单位共同发起开放代理式人工智能基金会（OAAIF），致力于打造中立、以开源为核心的开放协作平台，聚焦智能体基础设施的软件、协议与工程实践，促进智能体技术的互操作、规模化与高质量落地。

Advanced Security、GitHub Copilot等增值服务实现商业化。

中兴通讯的AI开源建设实践

作为全球领先的综合信息与通信技术解决方案提供商和数字经济赋能者，中兴通讯积极拥抱开源战略，通过参与开源社区、自主发起项目以及产学研协同，提升企业在全域开源领域的影响力与话语权。

在AI开源领域，中兴通讯重点参与了Linux Foundation、LF AI&DATA、CNCF、COIA、龙蜥社区等全球主流开源社区的建设与治理，成为多个开源社区的核心成员与贡献者。作为LF AI&DATA基金会的最高级别会员（Premier Member），中兴通讯深度参与基金会的治理和战略规划。

为加速人工智能及大模型技术在行业的落地，2024年12月，中兴通讯与四十余家合作伙伴共同发起开放智算产业联盟（COIA），旨在加速人工智能及大模型技术的实际应用，助力各领域的企业实现“AI+”转型。

2025年世界人工智能大会期间，中兴通讯作为首批共建单位加入由国务院国资委统筹、中国移动运营的国家级AI开源开放平台“焕新社

区”。其战略定位是“通信+AI”融合创新的基础设施提供者，锚定模型创新、算力优化、场景落地三大攻坚方向。

2026年1月，中兴通讯联合天工开物开源基金会、中国信通院、百度、Intel、Red Hat等单位共同发起开放代理式人工智能基金会（OAAIF），致力于打造中立、以开源为核心的开放协作平台，聚焦智能体基础设施的软件、协议与工程实践，促进智能体技术的互操作、规模化与高质量落地。

全球AI开源正在进入一个全新的发展阶段。从基金会、联盟到社区、项目，多层次的开源生态基础设施已日趋成熟；从美国、欧盟到中国，各国政府正在将开源战略提升至国家竞争力层面，开源不仅是技术进步的加速器，更是数字经济时代的新型基础设施。中兴通讯作为中国科技企业的重要代表，通过发起开放智算产业联盟、深度参与国际国内开源社区、构建全栈全场景智算解决方案等实践，为中国乃至全球AI开源生态的发展做出了积极贡献。未来，随着开源与商业化的深度融合、AI开源治理机制的完善以及产学研协作的深化，全球AI开源将迎来更加广阔的发展空间，为人类社会的智能化转型注入持久动力。[ZTE中兴](#)

构建开放协同的智算标准体系， 推动AI生态高质量发展

在第四次工业革命加速演进的背景下，AI已成为驱动数字经济转型的核心引擎。智算作为支撑AI落地的底层基础设施，正重塑全球科技竞争格局。建立统一、开放、可互操作的智算标准体系，已成为打通技术孤岛、加速产业协同的当务之急。

智算标准内涵与战略意义

智算标准是在智算系统全生命周期中为实现硬件兼容、软件互通、生态协同而制定的技术规范与评估体系，旨在推动AI技术从实验室走向规模化商用。

全球范围内，智算标准缺失已成为制约AI普惠化发展的主要瓶颈。不同厂商的AI芯片架构各异、训练框架互不兼容、工具链割裂，导致企业面临重复造轮子、迁移成本高、运维复杂等困境。尤其在通信、能源、金融等关键行业，对系统稳定性、安全性与可审计性要求极高，缺乏统一标准将严重阻碍AI技术深度渗透。因此，构建开放、中立的智算标准体系，是技术演进的必然要求，也是保障国家数字主权的战略选择。

智算标准发展现状与趋势

当前，全球智算标准体系呈现多极并行、协同演进的格局。

美国依托其强大的科技生态，主导了多个

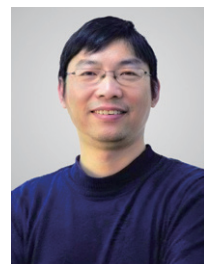
事实性标准。NVIDIA的CUDA成为行业默认接口；Google推动的TPU与TensorFlow框架形成硬件和软件的闭环；而LF AI&DATA与CNCf等组织推动多个AI开源项目，构建开源即标准的软性治理范式。

欧盟则更强调合规与伦理导向。《人工智能法案》对高风险AI系统提出明确的可解释性、数据质量与透明度要求，间接推动了AI模型评估、审计与部署标准的建立。

中国近年来在智算标准领域加速布局。

政策层面，2024年6月，四部委联合发布《国家人工智能产业综合标准化体系建设指南（2024版）》，明确以基础支撑为优先方向，推动标准体系化、结构化发展。

组织层面，早在2020年，全国信标委人工智能分委会（SAC/TC28/SC42）即已成立，成为我国首个聚焦AI标准化的国家级专业机构。2024年12月，工业和信息化部成立人工智能标委会（MIIT-TC1），统筹负责AI相关领域行业标准制订工作，标志着我国人工智能行业标准化工作迈入新阶段。2025年11月，计算互联总线工作组（SAC/TC28/WG39）正式成立，旨在推动建立开放、统一的计算互联总线国家标准体系。同时，SAC/TC599、SAC/SWG32、CCSA TC1、ODCC、AIIA等标准组织/联盟也积极开展人工智能标准建设，在硬件、网络、框架、应用等领域加强垂直整合，围绕行业赋能、大模型应用和评测体系开展横向协同，推动形成体系化、协同



朱静
中兴通讯核心网集成产品
规划总工



黄程
中兴通讯标准高级工程师

面向未来，智算标准体系的发展将呈现三大趋势：一是从技术标准向治理标准延伸，拓展安全伦理与可持续维度；二是企业反哺生态，头部企业主动开源，推动标准普惠；三是标准与开源深度融合，开源项目成为标准的试验田与实施载体。

化的标准布局，呈现出“多元推进，生态繁荣”的良好态势。

开源层面，信通院、阿里、百度、腾讯、中兴通讯、华为等纷纷构建活跃的开发者社区，针对模型开发、训练优化、分布式调度、异构推理引擎等创建开源项目，形成事实上的技术互操作软性标准，为智算标准制定提供广泛共识；同时通过开源协助机制，推动接口统一、协议一致，成为智算标准研制的试验田。通过这种“以开源促共识、以共识定标准、以标准强生态”的模式，避免标准制定脱离产业实际的问题，提升标准的可实施性与产业接受度。

中兴通讯在智算标准建设中的实践探索

中兴通讯作为通信与算力基础设施的领军企业，深度参与国家智算标准体系建设。

- 牵头《人工智能 超节点技术要求》国家标准编制

中兴通讯作为牵头单位，和中国电子技术标准化研究院、华为等单位一起，首次定义了超节点——即面向大模型训练的高性能、高可靠、可扩展的智算集群系统架构。标准明确了超节点的产品形态、关键功能及性能等核心要求，填补了国内在大模型训练基础设施标准化领域的空白。

- 深度参与人工智能加速器互联互通相关国家及行业标准编制

相关标准规范了协议栈、报文封装格式、接口协议、可管理性、安全等关键技术要求，解决了国产AI芯片与服务器厂商间接口不统一、部署复杂、集成成本高等痛点。基于上述标准，中兴通讯已完成自研芯片及服务器产品的研发与落地应用，并被多家国产芯片厂商采纳，显著降低了多厂商异构算力的集成难度与总体成本，成为推动国产AI加速器互联标准、产业及生态协同发展的关键支撑。

通过标准的落地，中兴通讯实现了“标准引领产品、产品反哺生态”的闭环，为国家智算标准体系提供了可复制、可推广的产业实践范式。

未来展望：构建开放、可信、可持续的智算标准生态

面向未来，智算标准体系的发展将呈现三大趋势：一是从技术标准向治理标准延伸，拓展安全伦理与可持续维度；二是企业反哺生态，头部企业主动开源，推动标准普惠；三是标准与开源深度融合，开源项目成为标准的试验田与实施载体。

中兴通讯的实践表明，唯有坚持开放协作、产业协同的路径，才能构建真正具有全球竞争力的智算标准体系。我们期待更多企业、科研机构与政府力量携手，共同打造开放、互信、可持续的智算生态，让AI技术真正成为驱动社会进步的普惠力量。 ZTE中兴

携手南方电网，打造能源行业 首个全栈自主可控的干卡智算中心

当前我国智算产业已初具规模，但在核心技术攻关、生态建设、运营效率等方面仍面临严峻挑战，亟待推动产业链协同创新，实现智算产业从规模扩张向质量提升转型，筑牢人工智能高质量发展根基。电网作为国家关键基础设施，迫切需要数智化转型。人工智能技术成为应对能源转型、提升电力系统运营效率的关键手段，可有效解决新型电力系统中海量数据接入、实时响应与安全稳定等挑战，实现云-边-端协同互动，全面赋能发、输、变、配、用各环节。

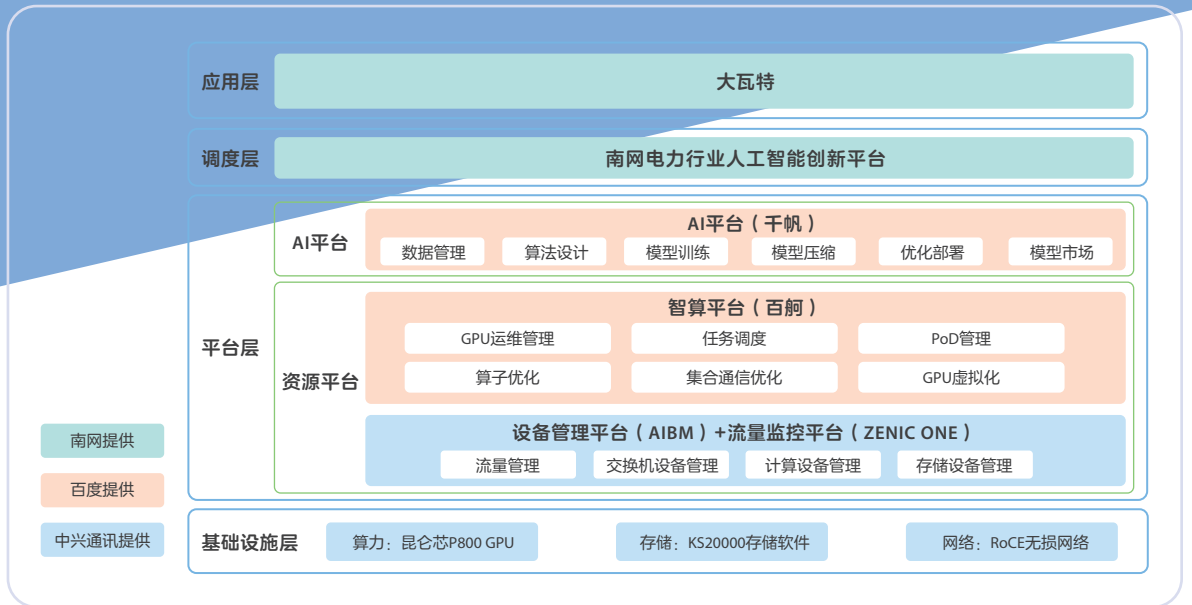
然而，电网智能化升级面临高性能算力技术挑战，构建自主可控的算力基础设施已成为保障国家能源高质量发展的必然选择。南方电网积极响应“人工智能+”等国家战略，发布人工智能

专项规划，打造安全可信的算力底座。通过统筹全网算力资源，构建网络化、普惠化、绿色化算力供给体系，为新型电力系统建设、高质量样本库构建以及数据安全防护等场景提供高密度、高可靠算力支撑，全面提升电网数智化核心能力。

智算中心是能源行业重要基础设施，南方电网携手中兴通讯、百度公司，强强联合，打造国内能源行业首个干卡级全栈自主可控智算中心。该中心构建了从底层芯片、整机硬件、算力集群，到软件算法、平台支撑，再到上层智能应用的全链条自主技术体系，实现核心技术全栈贯通。三方通过深度协同优化，筑牢高性能、高可靠、高安全的自主可控算力底座，全方位保障电网核心业务平稳高效运行，大幅提升能源行业智能化基础设施的自主可控水平与安全运行韧性。



黄燕
中兴通讯智算产品规划
总监



▲ 图1 南方电网干卡智算中心架构图

南方电网干卡智算中心整体架构如图1所示。

以“芯”强“算”：全栈自主可控芯片协同，构建高效算力根基

随着芯片制程逼近物理极限，当前智算产业正从“单点性能突破”转向“系统级协同创新”。本项目采用“CPU+GPU+网络芯片+内存+存储”一体化芯片，通过算、存、传芯片的协同设计和优化，释放系统级算力潜能。

服务器节点采用X86架构CPU，保障通用计算任务稳定高效运行；AI训练和推理节点搭载高性能的昆仑芯GPU，单卡算力和显存容量优于主流竞品20%~50%；在RDMA网络层面，服务器侧网卡采用中兴通讯自研的“定海”芯片，交换设备中搭载了自研“天屹”交换芯片，支持400G高速互联，为大规模AI集群提供无阻塞网络支撑。BMC（基板管理控制器）管理芯片、电源管理单元、

光模块、存储颗粒等关键部件均实现自主可控相关认证，形成从芯片到模组的完整供应链闭环。

以“网”强“算”：高性能无损网络架构，打通算力协同“大动脉”

针对AI训练对网络带宽、时延、稳定性提出的极致要求，本项目创新构建“四维网络平面”体系，实现算力网络深度融合。业务网络承载AI训练任务数据流，支持跨节点资源共享与分布式存储访问，保障大规模模型训练的高并发读写；管理网络实现GPU、CPU、交换机、存储等异构设备的统一监控与远程运维；参数面网络采用中兴通讯自研“定海”网卡和交换机以及端网协同技术，实现无损互联；样本面网络连接智算集群与高性能并行文件系统，实现TB级样本数据毫秒级加载。

依托算网深度融合技术，全面提升网络吞吐性能、降低传输丢包率，为干卡级集群构筑高效

“中兴通讯秉承开放解耦的理念，与芯片厂商、平台厂商深度合作，建立“技术共研、标准共建、成果共享”协同机制，推动自主可控算力在能源高门槛场景规模化落地。

可靠的无损网络底座。

以“软”强“算”：全栈软件平台整合，打造智能调度中枢

多方协同构建覆盖“训练、推理、运维”全生命周期的自主可控软件平台体系，打造“人工智能创新平台”，实现对算存网资源的统一纳管与智能调度。人工智能创新平台支持基于负载预测和任务优先级的动态资源分配、资源弹性伸缩与负载均衡，算力利用率提升40%以上。在此项目中，中兴通讯携手芯片厂商深度优化深度学习框架，完成算子级适配与编译优化；并提供一站式模型迁移工具、集群仿真系统，开展多个电力专用AI模型的迁移适配，并深度参与训推优化，提升训练收敛效率和推理性能，加速电力行业AI应用落地。

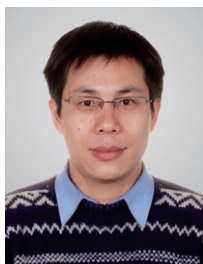
以“算”生“智”：垂直场景深度适配，释放AI业务价值

电力自主可控千卡智算中心建成投运，为

电力行业高价值、高并发、强实时场景的AI转型构筑了高效可靠的使能底座。目前，南方电网依托该底座全面开展输变电一体化无人机精细化巡视全流程智能应用，实现全程无人自主作业，构建调度、任务、执行、回传、识别一体化闭环工作体系；同时打造全国首个规模化一站式电力智能巡检体系，大幅提升巡检效率，有效降低线路跳闸率，为电网安全稳定运行提供坚实技术支撑。

中兴通讯深耕电力行业多年，充分了解电网调度、设备运维等核心业务场景需求，凭借异构算力融合、AI模型优化等ICT技术积淀，以“业务理解+技术落地”双重优势，精准匹配南方电网公司加快人工智能与电力业务的深度融合的需求。同时，中兴通讯秉承开放解耦的理念，与芯片厂商、平台厂商深度合作，建立“技术共研、标准共建、成果共享”协同机制，推动自主可控算力在能源高门槛场景规模化落地。ZTE中兴

山东移动携手中兴通讯打造 千卡智算推理资源池



陆威
中兴通讯智算产品规划总工



蒋妍
中兴通讯智算产品市场
规划工程师



杨晓曦
山东移动云网资深专家

2025年，大模型技术开始成熟并逐步落地，赋能各行业智能化应用。但在行业智能化发展过程中，政企客户普遍面临智算成本高、技术门槛高、进口受限、数据安全保障难等困境，“好用且不贵”的智算基础设施成为企业智能化转型的迫切需求。作为山东本土数字赋能主力军，山东移动以打造“中国北方人工智能创新发展高地”为目标，推动AI技术落地实体经济。

山东移动对省内300多家政企进行调研，深刻理解行业客户的痛点，协同多家智算厂商讨论解决方案，计划打造全栈自主可控、低成本、高可靠的智算推理资源池。综合考虑产品、方案、集成能力、工程交付及未来演进等因素，山东移动最终选择与中兴通讯携手合作。

中兴通讯“全栈智算解决方案”精准契合山东移动的需求，实现技术自主可控、成本管控与多场景适配，凭借高效的工程落地能力，仅用45天就完成方案设计及工程部署，为智算应用落地提供了基础设施保障。

此次项目的应用场景是为行业客户提供推理服务。推理场景对GPU卡的算力要求没有训练场景高，但对成本和时延要求比较高。基于这一核心诉求，中兴通讯和山东移动分析了数十款国产GPU卡，首先排除生态封闭的产品，再对比多家GPU卡的算力、显存、网络带宽、算子库兼容性以及成本，最终选择两款GPU卡，构建异

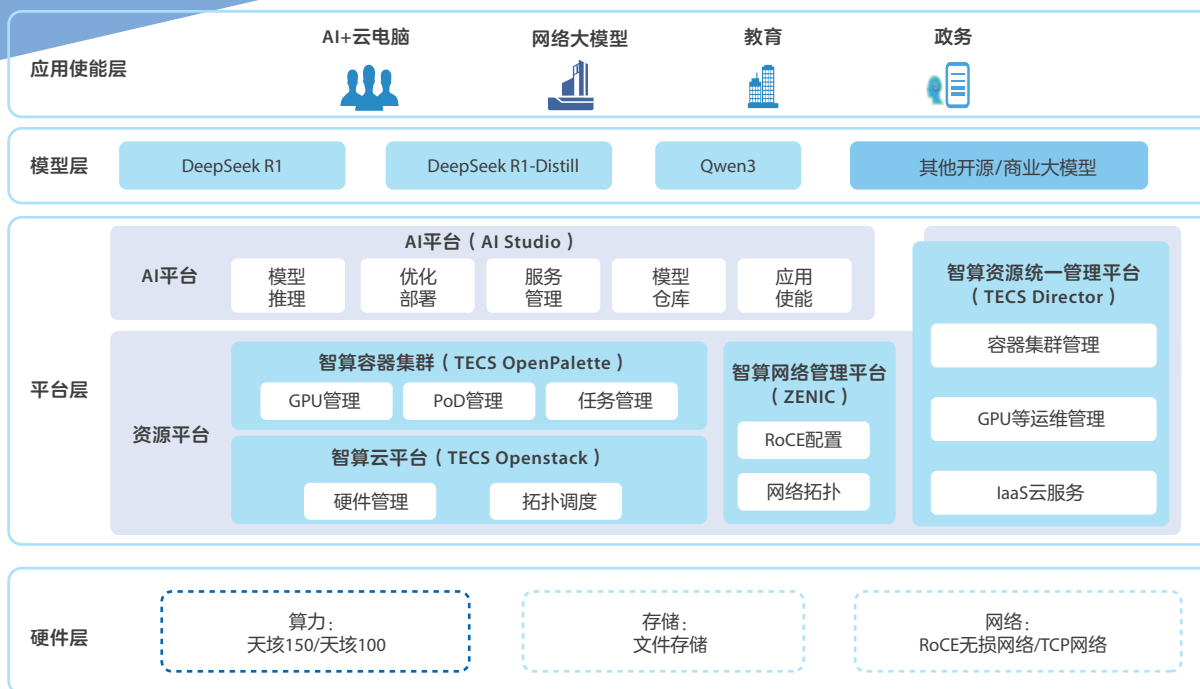
构计算集群，分别满足小模型和大模型的应用场景。在选定GPU卡后，中兴通讯基于全栈集成和软硬优化的技术能力，在实验室提前完成端到端的兼容性测试，并针对模型特定优化了算子库，通过软件最大化发挥GPU、存储、网络等硬件潜力，联合山东移动打造出兼具低成本、高适配、高可靠的千卡推理资源池。总体方案架构如图1所示。

- 全栈技术能力，筑牢数字安全屏障

项目实现从芯片、硬件到软件平台的全栈技术自主可控。硬件层采用国产天数智芯高性能AI加速卡，资源管理层依托中兴通讯自研AI平台实现开放解耦，模型层完成DeepSeek R1全系列、Qwen3、Nebula等主流模型的适配。基于多卡多模型协同架构，资源池支持大模型+小模型混合推理，为行业客户提供自主可控、高性价比的AI算力解决方案。

- 全生命周期降本，提供“好用不贵”的智算服务

项目实现全生命周期降本，通过全流程优化，为行业客户提供高性价比智算服务。在产品方案选型阶段，基于行业客户特点和应用场景需求选择综合性价比最优的产品；在方案设计阶段，采用一体化资源整合和协同设计，减少多厂商对接成本；在项目交付阶段，持续迭代完善软硬件资源配比和流程，提升资源利用率；在商用运营阶段，采用自动化运维工具，大幅减少人工成本。



▲ 图1 山东移动智算推理资源池总体方案

● 多元算力，精准匹配行业客户多场景应用

本项目部署两款推理型服务器，一款是天核150服务器，显存较大，专用于大模型跨机推理，另一款是天核100服务器，显存稍低功耗更小，用于中小模型推理。采用统一资源调度平台管理异构GPU，支持按政企客户使用的模型分级调度，小模型（如70B以下）部署于天核100服务器，实现低功耗推理以降低TCO，大模型（如70B以上）则部署于天核150服务器，用高算力、大显存GPU卡保障吞吐与低延迟。这种精准匹配既防止轻量任务占用高端算力造成的资源空转，也确保大模型推理不因硬件瓶颈而被迫降低并发业务量或截断上下文推理。

● 高效协同集成部署，45天快速投产

本项目规模达千卡，投入较大。时间就是金钱，项目越早部署完成，越早带来收益。中兴通讯通过科学拆解项目流程、预判落地隐患，实现“一站式部署”，无需分阶段对接多厂商，全程

提供技术调试支撑，实现“部署即上线、上线即见效”的闭环体验，在45天内完成方案设计和工程交付，为政企合作树立了敏捷交付标杆。

回顾项目合作历程，其成功离不开山东移动“以客户为中心”的项目定位，也离不开中兴通讯的技术支撑。中兴通讯通过整合全链条资源、优化项目落地流程、协同设计保障软硬件兼容性、精细化成本管控等一系列举措，实现了运营商、设备厂商与行业客户的三方共赢，充分彰显了其在智算领域的全栈系统能力。

未来，山东移动将继续与中兴通讯深化合作，依托中兴通讯技术优势，结合自身网络与服务资源，整合更多行业场景与模型，深化“数算智”全要素布局，持续完善“运营商+设备厂商+行业客户”的协同生态，让高性价比智算服务覆盖更多政企客户，助力AI技术与实体经济的深度融合。ZTE中兴

空芯光纤传输系统应用与挑战



尚文东
中兴通讯光系统规划
工程师

数字经济深入发展背景下，“东数西算”推进，AI算力网络、金融高频交易、电力、运营商骨干网等关键场景对通信传输提出“超高速、低时延、大容量、高可靠”诉求。传统单模光纤受折射率、损耗等限制，性能已触及物理极限，而反谐振空芯光纤凭借空气导光机制，具备超低时延、超低损耗和超低非线性三大核心特性，近年技术突破显著，核心指标超越传统单模光纤，成为破解四大场景痛点的关键技术。

空芯光纤核心优势：四大应用场景的价值赋能

空芯光纤的核心优势源于空气导光的本质，相较传统单模光纤，其性能跃升显著：时延降低30%，损耗突破瑞利散射限制小于0.1dB/km，非线性系数低3~4个数量级，精准匹配智算中心互联、金融高频交易、电力专网和运营商骨干/城域四大核心场景需求，成为场景化解决方案优选。

智算中心互联面临“算力拉远与算效保持”矛盾。空芯光纤低时延特性可降低算间互联损耗，江苏电信测试显示1024卡70B模型300km拉远性能损失仅2.8%；其低非线性、高带宽特性支持超高速传输，单纤容量可超100Tb/s。同时，超大有效模场面积提升传输容量上限，支撑算力网络规模化部署。

金融高频交易中，低时延是核心竞争力。空芯光纤传输时延低至3.3μs/km，远优于传统单模光纤的5μs/km，是缩短跨城金融专线时延的关键。如中国移动港-深交易所34km专线、

中国电信东莞-香港两地证券交易数据中心110km商用光缆，凭借空芯光纤分别实现时延1.07ms、0.93ms传输，为金融机构构建极速优势。其低非线性特性同时保障数据传输稳定，避免交易误差。

电力专网对时延和可靠性要求严苛，高压继保业务时延需≤10ms，且需防雷、抗干扰。空芯光纤低时延特性可轻松满足需求，弱磁光效应使其不受高压电磁干扰，配合OPGW光缆形成“低时延+高可靠+抗干扰”方案。其低损耗（最低0.05dB/km）支持单跨超长组网，减少中继部署，OTN硬管道传输保障调度数据安全，支撑智能电网稳定运行。

运营商骨干网每5~6年代际升级，当前400G C+L系统已商用，800G/1.6T是下一代方向。空芯光纤低损耗、低非线性特性打破传输距离限制，支持高阶调制信号长距传输。当前业界空芯光纤传输系统技术验证均显示，其单纤容量、传输距离较传统光纤大幅提升，验证了在骨干网超高速、大容量传输中的潜力，助力运营商提速扩容降本。

空芯光纤关键技术瓶颈与突破方向

空芯光纤场景化应用依赖核心技术突破，其导光机制既带来优势，也面临功率放大、气体吸收、背向散射、模间干涉等多重挑战，这些技术进展直接决定其规模化商用进程。

空芯光纤低非线性特性为高阶调制长距传输创造条件，但对大功率光放大器提出更高要求。1.6T PS-64QAM 1600km长距传输仿真显示，光放

输出功率需大于40dBm，当前业内产品能力（35dBm左右）仍存在差距。未来需优化光非线性突破功率瓶颈，结合高阶调制，发挥低非线性优势实现超高速长距传输。

气体吸收是制约长距传输和波段应用的核心瓶颈。由于CO₂、H₂O对特定波段的选择性吸收，100km传输即可引入约10dB损耗。当前行业通过充惰性气体、构建正压环境控制出厂损耗，但现网破损易恶化性能。解决方案包括光纤厂家优化密封技术、控制全生命周期损耗，设备厂家研发气体吸收补偿算法拓展可用波段。

低背向散射特性导致OTDR（光时域反射仪）性能下降，影响链路运维。空芯光纤背向瑞利散射系数远低于传统光纤，OTDR动态范围劣化14dB~15dB，熔接点压差还会扩大监测盲区。当前SFP（small form-factor pluggable）模块OTDR动态范围单向检测基本无法满足空芯光纤80~160km跨段需求，需不断优化提升OTDR动态范围，采用双端配置等优化方案满足链路损耗检测需求。

场景驱动下通信设备的核心功能诉求

空芯光纤在智算中心互联、金融高频交易、电力专网和运营商骨干/城域四大应用场景的差异化需求，对空芯光纤配套通信设备提出针对性要求，设备需围绕“低时延、大功率、高可靠”核心方向，与空芯光纤深度适配，充分发挥技术优势。

智算中心互联通过空芯光纤满足光纤低时延光层传输需求，电层设备需集成低时延电交叉转发功能减少时延；此外在保障高可靠性方面，通过无损保护、OTN融合抗拥塞机制等技术创新确保数据在传输过程中实现零丢包、抗拥塞，从而保障跨集群协同训练稳定运行；配备1.6T及更高速率相干光模块，结合大功率光放大器满足大容量传输需求。

金融高频交易场景通过空芯光纤满足光纤低

时延光层传输需求，电层设备需搭载低时延OTU压缩设备时延；支持大动态范围的高分辨率OTDR，满足单跨100km传输链路光纤性能监测；具备加密等技术保障客户业务安全，同时支持客户侧1+1保护、通道1+1保护、复用段1+1保护、ODUK SNCP保护等灵活组网适配专线部署。

电力专网场景以空芯光纤和低时延100G/400G OTU保障继保信号低时延传输；配备输出≥40dBm的大功率光放大器支持电力单跨近500km超长距组网；OTDR双向配置需具备电力单跨超长距光纤性能检测；通过灵活的保护组网配置保障数据独立传输与故障快速切换。

运营商骨干/城域网向800G/1.6T代际演进，需要传输距离基本不变的情况下满足系统容量翻倍，波段扩展+单波提速仍将是主流演进方向。当前空芯光纤气体吸收问题限制C6T和C6T+L6T部分频谱业务传输性能，亟待解决。800G/1.6T+空芯光纤系统为满足运营商骨干/城域网传输需求，设备首先需支持800G/1.6T高阶调制与信道补偿算法，同时配备高功率ITLA（integrated tunable laser assembly）+低插损调制器提高光模块输出性能，结合光路大功率光放大器（输出功率约40dBm），优化端到端传输OSNR（optical signal-to-noise ratio）满足长距需求；光交叉设备连接器和光器件需提升最大输入功率阈值至40dBm，支持C+L多波段传输；OTDR采用双端配置覆盖单跨光缆监测需求。

空芯光纤凭借优异的传输性能，成为智算中心互联、金融高频交易、电力专网和运营商骨干/城域四大场景的优选技术，支撑数字经济发展。但距离规模化商用仍面临技术、标准、产能、成本等多重障碍。未来需产学研协同创新，推动工艺升级、器件研发及标准完善，降低成本。随着技术成熟与成本下行，空芯光纤将实现广泛商用，重塑光通信产业格局，助力新型数字基础设施建设。 ZTE中兴

低时延空芯光纤在跨域大模型训练中的应用分析



闫宝罗
中兴通讯波分技术规划
总工

研究背景

基于Transformer架构的大语言模型（LLM）的出现，对计算资源、存储容量和节点间通信提出了三大需求。即使在NVIDIA A100 GPU的理论峰值计算能力下，使用300万Token的数据集训练1个GPT-3 175B模型也需要32年才能完成。当前一代加速器（如NVIDIA A100/H100）尽管单设备显存达80GB，但仍需至少44块GPU组成的集群才能容纳单个模型副本。因此，LLM的训练和推理通常采用多加速器集群，不可避免地需要在服务器内、服务器间乃至跨人工智能数据中心（AIDC）域进行集合通信。然而，受供电基础设施、AIDC空间限制以及国家计算资源分布的约束，NVIDIA、谷歌、OpenAI/微软等机构已广泛开展跨域协同分布式训练的研究。

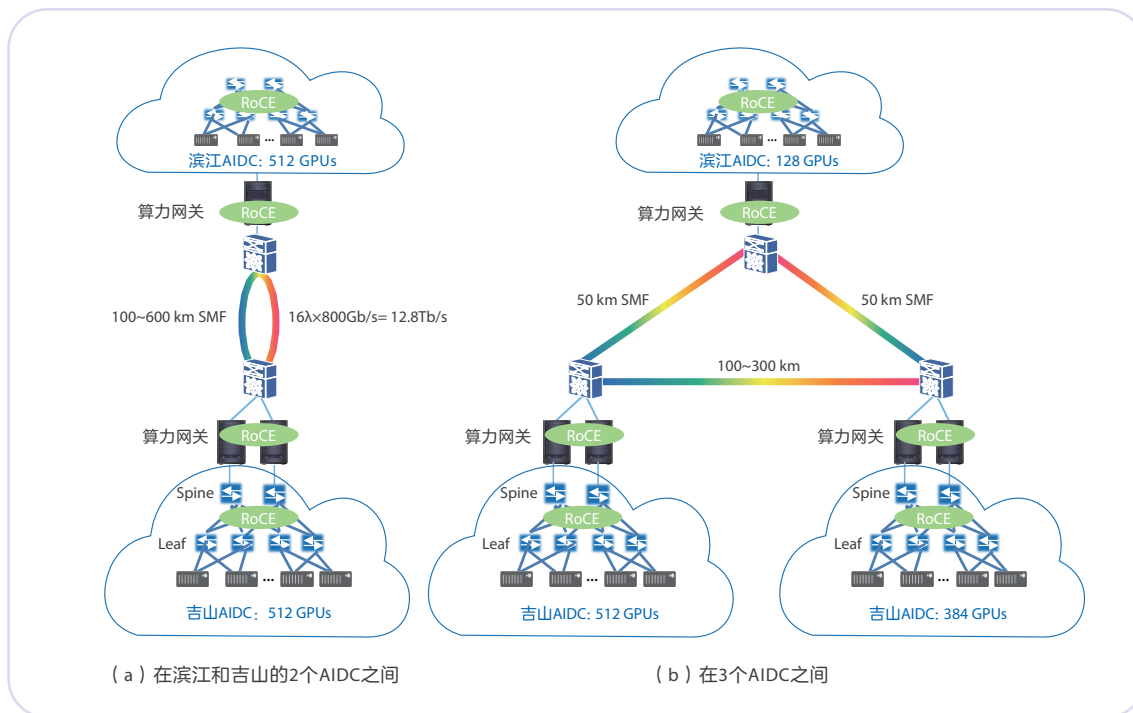
一方面，在实际应用中，不同的GPU集群架构和3D并行训练配置具有独特的时延敏感性特征和流量需求。事实上，通过对跨域流量模式和时延敏感性指标进行细致的理论分析，已报道的分布式LLM训练框架在带宽收敛比和分布式训练距离方面仍有进一步优化的空间。另一方面，反谐振空芯光纤（AR-HCF）以空气替代实芯石英介质，从根本上突破了单模光纤传输损耗与时延的瓶颈，目前其降损优化已经历20余年研究，基本收尾，C波段最低损耗 $<0.1\text{dB/km}$ ，优于单模光纤。国内外面向数据中心互联、金融专线等应用

场景已有10余处商用部署或试点，例如微软Azure 20km空芯光纤数据中心互联、国内三大运营商港-深低时延空芯金融专线等，说明中短距应用已基本成熟，这吸引业界广泛关注实际业务应用中带来的收益价值。

为了量化空芯光纤系统在数据中心算间应用收益，本研究以干卡集群、在2个地理分离的AIDC间开展LLaMA2-70B模型的优化分布式训练现场试验为例，采用数据并行（DP）和流水线并行（PP）技术，通过理论分析表征了跨AIDC的流量与时延需求。进一步对比单模光纤与空芯光纤应用时算效收益，以此说明大容量光传输（OTN）技术和低时延HCF对未来大规模集群部署的必要性。

LLM跨DC分布式训练算间流量与时延需求

我们分别在2个和3个AIDC间开展了LLaMA2（参数数量 $\psi=70\text{B}$ ）分布式训练实验（见图1）。训练配置参数如下：单GPU计算吞吐量 $P=122.96\text{TFLOPS}$ （推算平均值），批大小 $b=32$ ，序列长度 $s=4096$ 。集群GPU总数 $N=1024$ ，全局批大小（GBS） $=2048$ 。训练过程中，流水线并行（PP）度 $p=8$ ，张量并行（TP）度 $t=8$ ，数据并行（DP） $d=16$ 。其中，DP涉及的通信量最大，因此需要计算集群所需的通信带宽 BW_{DP} 。单批次迭代时间可表示为 $T_{batch} = (6 \times$



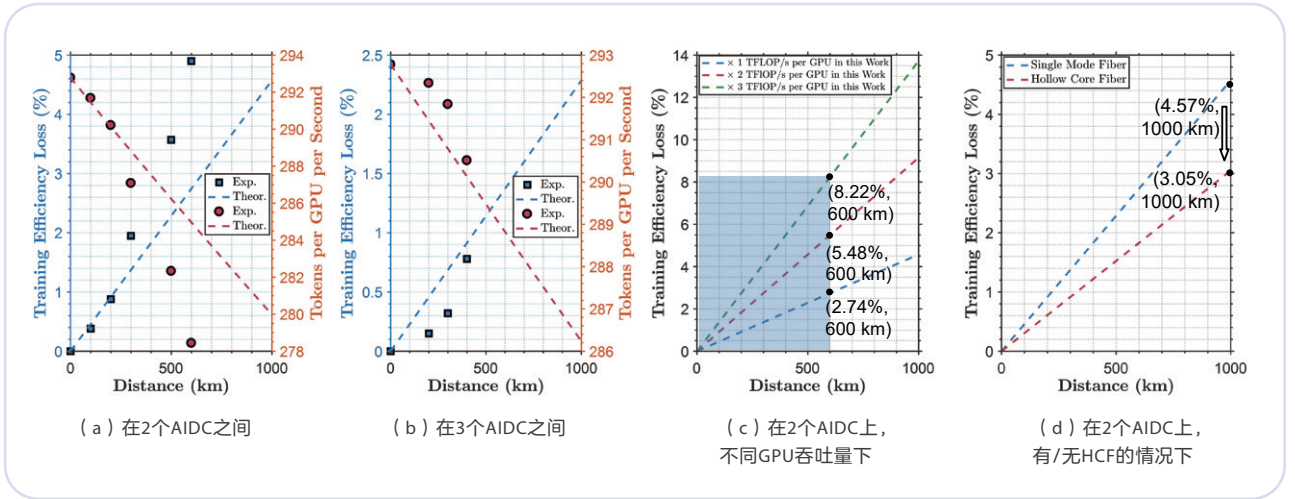
▲ 图1 2 AIDC和3 AIDC的LLM分布式训练

$\psi \times b \times s) / (p \times t \times P) = 6.99s$ ，单DP迭代时间为 $T_{DP, Cal} = \frac{Global\ batch\ size}{b \times d} \times T_{batch} = 27.98s$ 。DP中每块GPU每次通信的数据量 D_{DP} 可近似为 $2 \times \psi / (p \times t) \times 2Byte \approx 4.37GB$ 。进一步考虑DP通信时间 $T_{DP, Comm}$ 与计算时间 $T_{DP, Comm}$ 的比值（占总迭代时间的比例不超过5%，通常为1%~5%），可估算出单GPU所需数据带宽 $BW_{DP\ per\ GPU} = D_{DP} / T_{DP, Comm} = 3.12 \sim 15.64GB/s$ 。最终，线路侧总互联带宽可估算为 $BW_{DP} = \frac{GBS}{d} \times BW_{DP\ per\ GPU} \times 8bit/Byte = 3.2 \sim 16Tb/s$ 。因此，我们将跨域流量收敛比设置为1:8，以实现12.8Tb/s的线路侧输出流量容量。该流量可以由16个800Gb/s相干光模块承载，采用单载波135-Gbaud PCS-16QAM调制格式的OTN业务传输。单个GPU的Token处理能力（TGS）可清晰反映固定训练数据量下模型的训练速度，其表达式为 $TGS = (b \times s) / (T_{Cal} + T_{Com}) / N$ ，其中 T_{Cal} 表示DP或PP的单次迭代计算时间， T_{Com} 表示通信时间。TGS可通过测量训练吞吐量得出。随着AIDC间距离的增加，延长的通信时间 T_{Com} 会导致TE下降。

如图1所示，在2 AIDC分布式训练中，滨江

AIDC和吉山AIDC均配置512块GPU集群，AIDC间互联距离在100~600km范围内；在3 AIDC分布式训练中，滨江AIDC部署128块GPU，2个吉山AIDC分别配置512块和384块GPU，仅调整2个吉山AIDC间的互联距离（100~300km）。

本文以3 AIDC间的分布式训练为例，介绍DP与PP并行过程中数据的流向。在DP并行训练中，数据中心1（DC1）部署8个完整模型副本（DP1-DP8），数据中心2（DC2）部署6个副本（DP9-DP14），数据中心3（DC3）部署2个副本（DP15-DP16）。在采用Ring算法（如Reduce-Scatter和AllGather）进行DP通信时，跨AIDC通信链路发生在DP8与DP9、DP14与DP15、DP16与DP1之间。在PP并行训练中：DC1包含16个模型副本（DP1-DP16）的PP阶段1-4，DC2部署PP阶段5-7，DC3部署PP阶段8。在PP的点对点发送/接收（Send/Receive）通信中，跨AIDC链路连接PP4与PP5、PP7与PP8、PP8与PP1。除必要的长距离拥塞流控制技术外，在分布式训练中特别采用了分层TP/DP/PP并行技术，以支持全局



▲图2 基于DP策略的LLM分布式训练的TE理论估计结果

AllReduce操作，确保总全局通信时间 T_{com} 限制在2倍往返时间（RTT）内。对于PP并行，采用了交错式1F1B流水线调度，使通信时间与计算时间能够充分overlay，从而缓解性能损失。最后，通过测量DP和PP并行训练配置下的TGS，表征分布式训练TE损失。

分布式训练性能与空芯光纤对训练算效改进估计

首先，PP并行的测试结果表明，得益于通信时间被掩盖在计算时间，2 AIDC间600km链路和3 AIDC间400km链路的整体TE损失均控制在1%以内（详细数据不在本文展开），因此PP并行对长距拉远时延劣化并不敏感。

DP分布式训练方面，随着2 AIDC间距离增加，TE损失逐渐增大，在600km处达到4.9%，如图2（a）。3 AIDC分布式训练现场实验采用环形拓扑，等效RTT缩短一半，因此在相同AIDC间距离下，TE损失低于2 AIDC场景。基于前文讨论，可对多AIDC分布式训练导致的TE损失进行理论估算，这里仅考虑光传输时延，暂未详细考虑通信时间 T_{com} 的其他影响因素，如厂商特定配置（交换机/路由器/OTN设备处理时延、协议开销、流

控制处理时间等），这也解释了图2中理论粗估与现场试验结果存在差异的原因。

尽管在此次实验中没有使用空芯光纤，但我们可基于空芯光纤每米降低1.67ns的时延，代入上述推演过程，进一步考虑GPU算力、不同通信距离下评估收益情况。如图2（c），TE损失与GPU实际计算能力密切相关，不同国家地区可部署的GPU算力存在差异，我们假设部署的GPU算力达到此次实验的基准值的2倍和3倍时，计算时间被压缩，通信时延成为性能瓶颈，可以看到超过5%算效劣化场景。另一方面，进一步延长多AIDC间距离将导致算效劣化，因此仍需降低通信时间 T_{com} ，如图2（d），我们给出空芯光纤与单模光纤在1000km级2个AIDC下DP并行的算效劣化值，空芯光纤可使RTT减少33%，从而将TE损失改善，这一优势在长途LLM分布式训练场景或高计算容量配置部署中尤为重要。

跨域分布式训练受网络带宽、通信延迟瓶颈等影响，一定程度制约了跨域的距离和集群规模。随着超高速1.6T端口的成熟、超低延迟空芯光纤的部署以及高可靠恢复与保护配置的应用，跨域协同训练将为智算行业快速发展提供强有力支撑。ZTE中兴



驭风系列

纤薄至简 驭风随行

驭风10 Air

13.9mm厚度 | 1.25kg净重 | 14英寸FHD高清显示屏 | 5W超低功耗
丰富接口 | 全金属机身 | 无风扇设计

ZTE中兴

致力于成为网络连接和智能算力的领导者
让沟通与信任无处不在