

第8届 Al+ Development Digital Summit

Al+研发数字峰会

拥抱AI重塑研发

11月14-15日 | 深圳





EDEAI+ PRODUCT INNOVATION SUMMIT 01.16-17 · ShangHai AI+产品创新峰会



Track 1: AI 产品战略与创新设计

从0到1的AI原生产品构建

论坛1: AI时代的用户洞家与需求发现 论坛2: AI原生产品战路与商业模式重构

论坛3: AgenticAl产品创新与交互设计

2-hour Speech: 回归本质



用户洞察的第一性

--2小时思维与方法论工作坊

在数字爆炸、AI迅速发展的时代, 仍然考验"看见"的"同理心"

Track 2: AI 产品开发与工程实践

从1到10的工程化落地实践

论坛1: 面向Agent智能体的产品开发 论坛2: 具身智能与AI硬件产品

论坛3: AI产品出海与本地化开发

Panel 1: 出海前瞻



"出海避坑地图"圆桌对话

--不止于翻译: AI时代的出海新范式



Track 3: AI 产品运 AI 产品运营与智能演化

从10到100的AI产品运营

论坛1: AI赋能产品运营与增长黑客 论坛2: AI产品的数据飞轮与智能演化

论坛3: 行业爆款AI产品案例拆解

Panel 2: 失败复盘



为什么很多AI产品"叫好不叫座"?

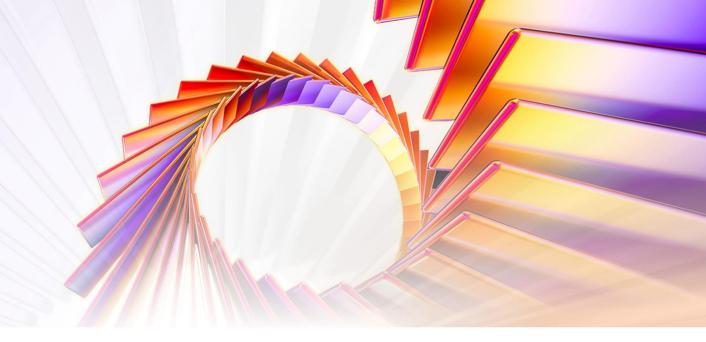
--从伪需求到真价值: AI产品商业化落地的关键挑战

智能重构产品数据驱动增长



Reinventing Products with Intelligence, Driven by Data





智算集群故障诊断算法研究与实践

陈文潇 | 华为技术有限公司





陈文潇

华为天才少年 智能体软件专家

华为天才少年,智能体软件专家,清华大学NetMan实验室博士,一直从事AIOps,人工智能,网络自动运维等工作。在WWW、INFOCOM、FSE等国内外学术会议上发表演讲,担任SIGMETRICS, AAAI, WWW等多个国际会议期刊审稿人。担任openFuyao核心技术负责人,与科大讯飞、蚂蚁集团、清华大学共同打造昇腾生态竞争力。



目录 CONTENTS

- I. 背景
- II. 问题/痛点
- III. 解决思路/整体方案
- IV. 具体实现/技术实践
- V. 总结与展望



PART 01

背景: 讯飞大模型的发展



▶ 从0到1,讯飞星火大模型实现从快速追赶到自主创新





大模型评测体系发布



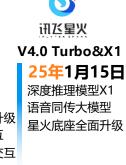




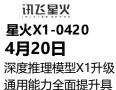












6

通用能力全面提升具 备深度推理模式



▶ 讯飞持续探索全栈国产化无人区



讯飞星火实现了训练和推理的全国产化,星火大模型持续引领国产平台发展 "飞星一号"平台2024年全年平均使用率达到95%

500+次

解决基础软硬件问题

30+项

新增框架和平台特性

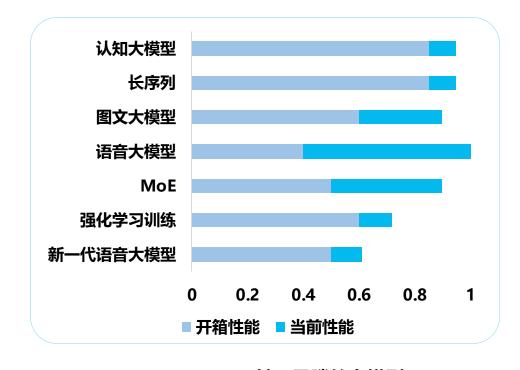
150+个

基础、通信和融合算子优化

75%

模型训练适配优化平台耗时优化

2024年 "飞星一号" 集群 优化过程中创新积累



基于昇腾的大模型 训练效率持续优化 讯飞&华为联创探索"开车换车轮" 规模无关的断点续训新方案 可支持**集群规模再次倍增**

> 智算集群规模的 再次跃迁

"飞星二号"达到万P

未来基于国产算力的大模型自主技术创新,探索新模型新算法的持续适配及智算集群规模的再次跃迁



PART 02

痛点: 智算集群维护

▶大模型的能力及效果受多种因素影响,集群算力是驱动大模型へiDD ỗ號 创新的基础

数据集质量

更好的数据质量更长的训练时间

新训练范式

强化学习决定模型对齐程度

Agent

决定使用工具、上下文记忆能力

多模态

决定模型功能丰富度

MOE

决定特定复杂任务处理能力

"算力规模≈参数量*数据量/训练时长"

集群大算力是支撑更高质量大模型创新落地的最关键基础

算力需求与数据量成正比

多模态数据训练需求是文本数据的320倍

算力需求与序列长度成正比

32K+序列的模型已商用, 1M+序列已出现

算力需求与模型参数成正比

百亿百卡、千亿千卡

精度需求要匹配合适的浮点运算

综合使用FP32/BF16/FP16以快速模型收敛

Grok4算力分配:通过RL压强投入获得模型能力大幅提升

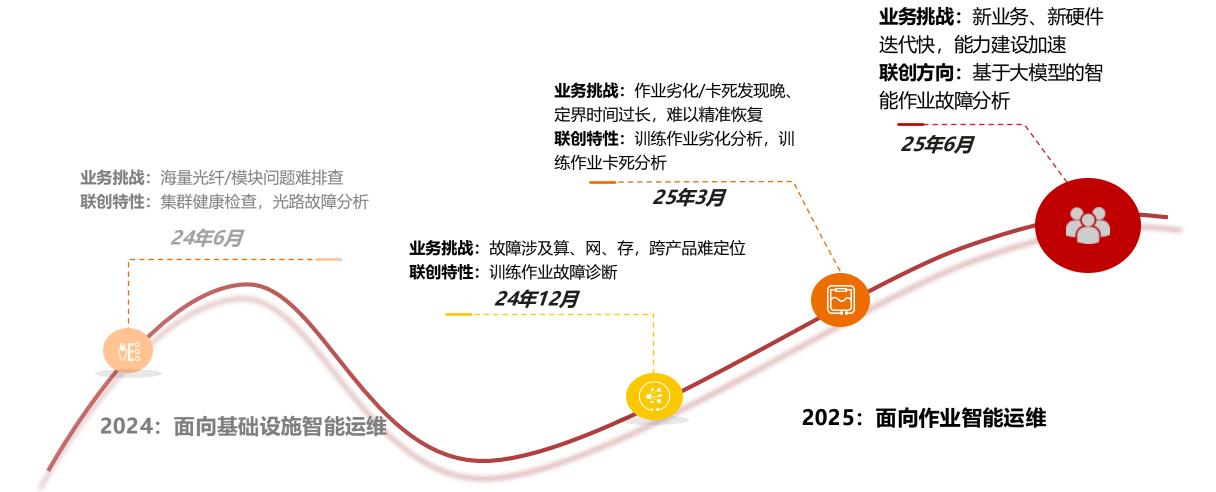


国内TOP企业计算规模较国外存在近10倍差距



▶ 从0到1, 讯飞星火大模型实现从快速追赶到自主创新







PART 03

整体方案: 运维智能化

▶ 将华为CCAE融入集群日常运维流程,提升运维智能化水平





飞星一号: 高频次使用CCAE特性

训练作业保障:

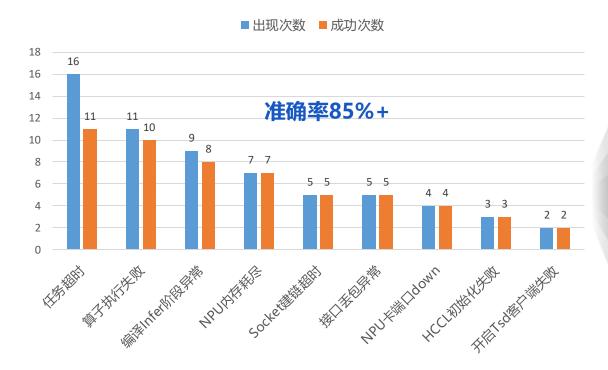
作业故障诊断,月均自动调度次数150+, 准确率85%+

基础设施故障诊断:

覆盖算网存基础设施故障,累计识别关键 故障180+

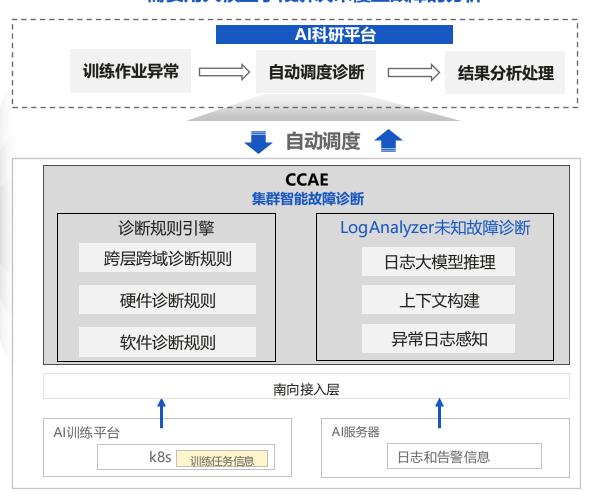
▶故障模式库已覆盖场景诊断准确率85%+,需要加强未知故障へiDD ᠍ 定位准确率

现网故障诊断TOP问题及诊断情况



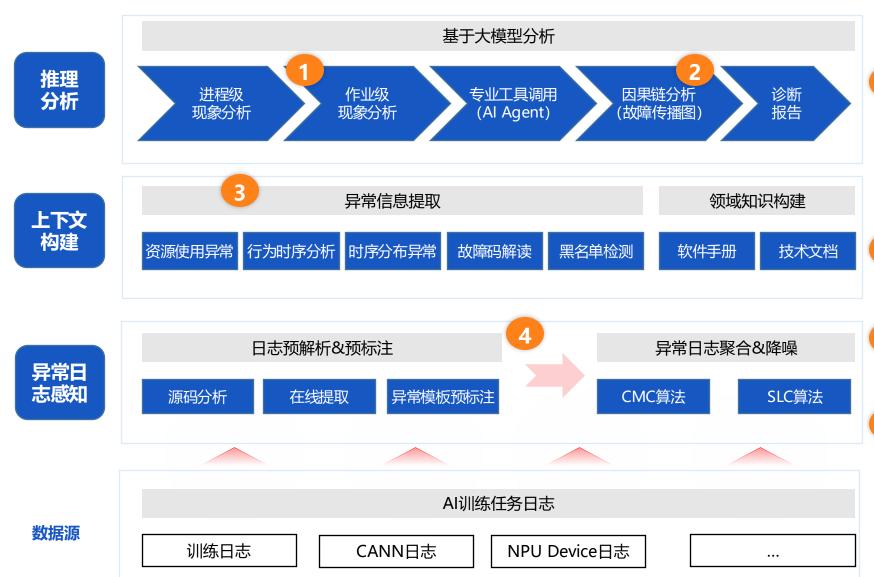
- ✓ **已知故障无法覆盖所有问题**:随着新硬件上市、新业务应用,故障模式 库逐步增长,基线-配套固定3个月周期,故障覆盖率呈上下波动特征。
- ✓ 未覆盖故障分析依赖各域专家手工分析,故障分析需要天级-周级不等

需要用大模型手段解决未覆盖故障的分析



▶ 首次引入基于日志大模型底座的AI辅助分析引擎





作业级故障分析技术

- 结合提取的异常信息和CANN、NPU等领域 知识, 生成训练进程故障时间线和故障现象
- 汇总进程级现象,结合HCCL通信知识分析导 致集群作业中断的故障类别, 生成作业级别 的故障现象

故障传播链分析技术

- 结合作业、进程级别的异常事件信息,构筑 基于大模型的故障传播链技术
- 异常信息提取技术
- 多种手段提取各节点的日志异常信息,生成 大模型的有效输入数据

异常日志识别技术

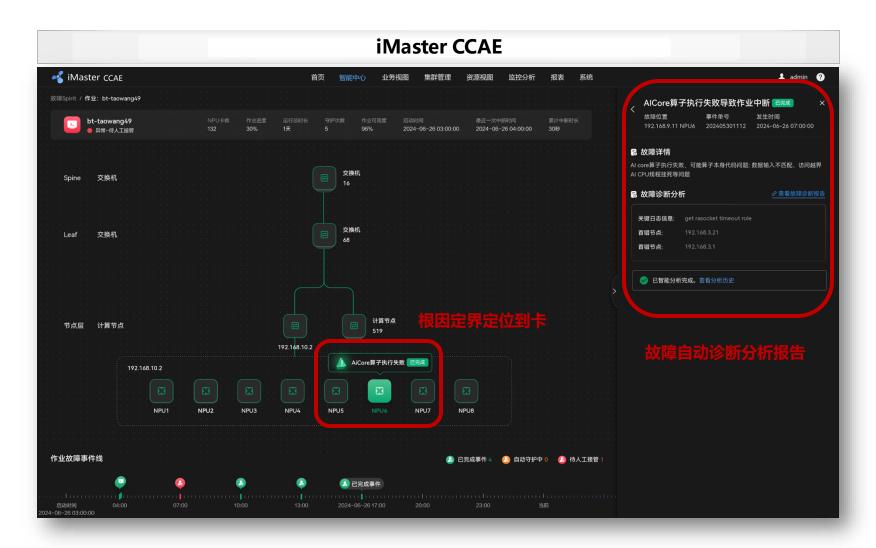
- 日志预解析&预标注,跳过正常日志分析
- 异常日志聚合,降低入库数据量

第8届 AI+研发数字峰会 | 拥抱 AI 重塑研发



▶ 客户视角目标态:故障自动感知定界,支撑端到端无感闭环 ◇IDD Sth





- 目标用户:客户运维、客户算法工程师、 华为驻场人员
- 使用场景:作业中断时,自动触发诊断 定界,10min内给出定界报告,一键下 发讲行故障隔离恢复
- 能力目标: 算网存全场景的精准定界
 - 计算侧定界到节点和NPU卡
 - 网络/存储侧定界到域



PART 04

技术实践: LogAnalyzer

LogAnalyzer基于大模型全面分析各域日志,显著提升故障 ∧jDD δth. 定位效率



模块Z日志

挑战

- 某业务团队在运行2千多卡的MoE长文本微调训练任务,每运行2个小时候,偶发中断,阻塞任务训练
- 37G运行日志,没有典型错误日志,在千卡 (测试) 环境上,多次跑任务验证,问题不复现

海量日志 37G 模块B日志 模块A日志 人工分析 专家合议 10+专家 ##诊断流程 问题干卡不复现 添加日志复现 - 训练作业因重复运行时错误 (`rtKernelLaunchWithHandleV2 failed: 507048`, `fftsplus task timeout`) 中断。 - 日志条目显示HCCL超时和通信操作失败。 **根因缩小**: - HCCL日志明确报告了`ReduceScatter`超时(`hcclfftsplus task timeout`) - RUNTIME日志显示任务失败错误\innerCode=0x715006b\mmyhh\ [fftsplus timeout]\, 确认EI0002。 预期不一,多次验证 - **网络不稳定**: DRV日志 (`drv<u>HdcSessionAccept error`, `socke</u>t accept failed`) 表明节点间通信故障

传统方案

问题复现成本高,耗时长,小规模集群没有开启 DP 并行,因而流数量小于

方案

- **跨领域错误检查无遗漏,效率高**,全面分析日志,不会因目的性受限范围查 **跨领域故障诊断需要协调多个模块专家投入分析,耗时长**,37G日志难全覆 找,导致关键错误遗漏。
 - 基于一次作业故障数据进行深层分析,成本低,展示故障传播关系,为专家 进一步验证和判断提供明确依据。

- **资源耗尽**: RUNTIME日志警告 hccl streams greater than 8! 表明HCCL资源过度分配

LogAnalyzer方案

LogAnalyzer

分析、合议、诊断、报告

成果

涉及多个模块的疑难故障,LogAnalyzer能够全面分析,显著降低定位时长:天级 → 分钟级

8,导致问题无法复现。

验证设想

盖, 跨模块间的故障影响需要跨领域知识。



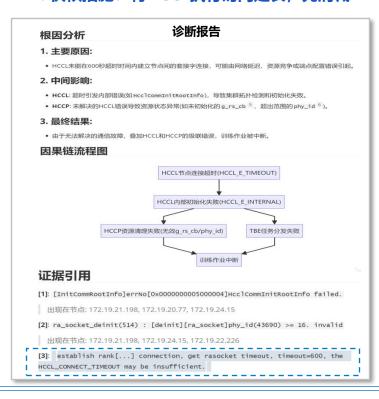
▶ LogAnalyzer构建集合通信域透视图,提升通信域初始化超时\JDD ᠍ 问题诊断效率

挑战

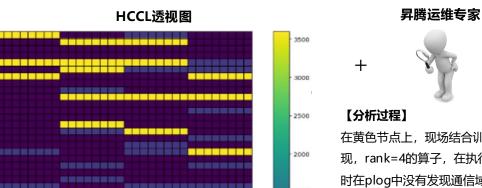
生产环境大规模集群,模型训练多个任务发生偶现卡死,影响模型正常训练,降低集群资源利用率,高度关注

方案

1. 快恢措施:将HCCL执行时间延长,先消减



2. 根治措施: HCCL透视图 + 运维专家, 找到根因



在黄色节点上, 现场结合训练脚本的日志发 现, rank=4的算子, 在执行时卡住了, 同 时在plog中没有发现通信域创建日志,此 故障为通信域创建未下发,导致超时。

【问题根因】

通信域创建指令未能成功下发至计算卡导致 超时中断

【修复措施】

1、Rank4卡隔离; 2、减少数据量

成果

针对集合通信库业界难题,LogAnalyzer通过HCCL透视图+大模型,HCCL相关问题,诊断效率提升2倍

黄色: 未下发通信域初始化

蓝色:通信域初始化超时(默认600s)



▶ LogAnalyzer上线使用4个月,未知故障诊断达到预期效果



4月~7月期间, LogAnalyzer有效故障诊断任务数72个, 离线诊断覆盖率100% (有日志数据即可诊断覆盖);

诊断正确63个,整体准确率: 87.5%; 其中未知故障51个,诊断正确43个;未知故障准确率84.3%

问题分类	问题类别	数量	典型故障
未知故障诊断	1、跨栈定位类	合计: 30 准确率: 80%	异常表现在HCCL,但根因可能是用户脚本,PTA、CANN、RTS等多模块耦合导致,典型故障包括: CANN软件Bug导致通信/建链超时 硬件问题导致算子超时
	2、环境配置类	合计: 13 准确率: 92.3%	训练环境和软件依赖导致作业无法拉起,根因发散无法详尽穷举,典型故障如: 训练通信环境配置错误软件版本依赖残余进程占用…
	3、硬件类故障	合计: 8 准确率: 87.5%	训练任务因传输链路故障,NPU 故障导致的任务中断或者超时, 典型故障包括: > AICORE TIMEOUT导致算子执行超时 > Kernel无效内存访问
已知故障优化	4、用户脚本及任 务初始化类	合计: 15 准确率: 93.3%	训练任务初始化失败,根因与用户代码和依赖库相关: IndexError、配置参数错误等用户脚本类问题 GE初始化失败,TBE依赖项缺失 R条进程占用导致任务拉起失败
	5、资源占用类	合计: 6 准确率: 100%	由于存储、内存、计算资源冲突或分配问题导致,包含: > 磁盘配额满问题 > OOM问题,指向片上内存内存分配

技术途径:基于专业模型压缩故障上下文,融合算网存上下文/JDD Sth 自主诊断

理念: 大模型作为诊断的"大脑"(负责推理、分析),调用各种专业工具和专用模型的"双手" (执

,从而解决复杂故障定位问题。 行特定任务)

智算故障诊断专业算法方案设计 诊断报告

> 生成 报告

融合算网存上下文的自主诊断

推理 模型

① 大模型友好的故障上下文

压缩故障上下文、结构化故障上下文, 构建推理模型友好的故障上下文

日志 专业 模型

时序 专业 榵型

结构 化告

错误

还原 拓扑 结构

治理 历史 案例

系统可观测性















历史故障案例



根因推导:多种报错类型间建立传播关系

- 根因识别: 通过Trainlog信息模 型分析为资源问题
- 根因传播:磁盘满引发关键资源竞 争,导致All-Reduce操作期间任
- 最终体现:通信延迟持续过长,流 同步失败, 任务无法拉起



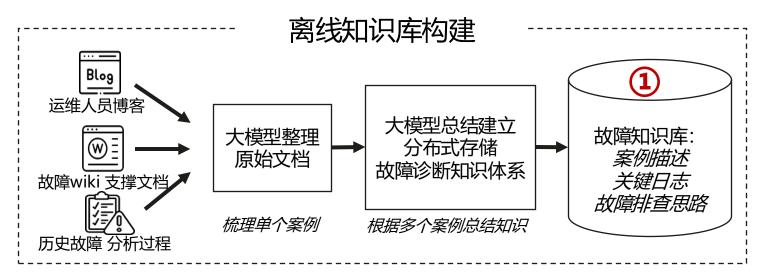




▶ 关键技术1: 大模型自驱动构建故障诊断完备的背景知识



①知识库构建成本太高





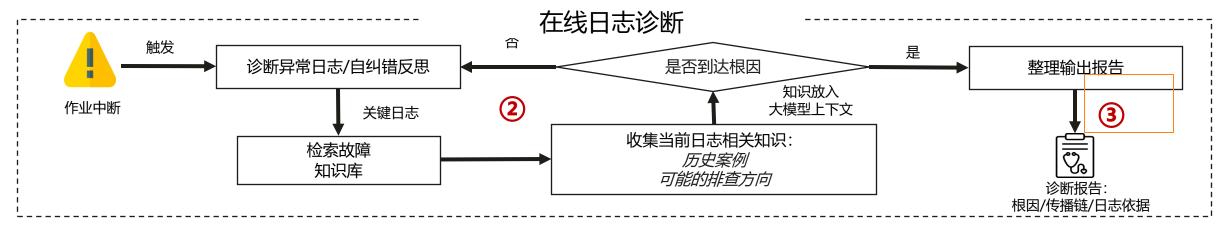
大模型自驱动知识库构建,高 效赋能在线诊断

通过大模型将内部开发文档,用户 支撑文档,排障案例,一线人员博 客中的海量运维知识自总结为大模 型易理解,易检索的文档,高效构 建分布式存储领域故障知识库,克 服人工构建耗时耗力的瓶颈。



▶ 关键技术2: 融合算网存上下文自主诊断

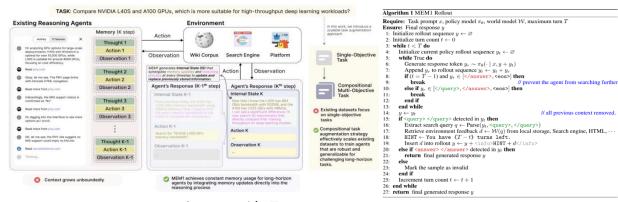




② 背景知识太多了,一次放不下

DeepResearch诊断流程, 小步获取知识

借鉴DeepResearch方法提高大模型诊断准确性,在诊断过程中遇到大模型无法理解 的信息时,通过检索知识库在上下文中自适应补充领域知识;获取新知识后自纠错反 及时修正故障诊断方向, 避免一错再错。



DeepResearch Agent流程

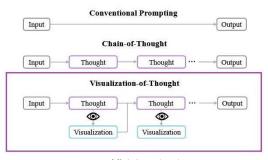
流程伪代码

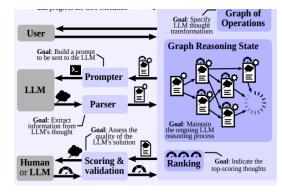
模型生成结果用户不相信

提取大模型思维过程,构建可信故障传播路径

从大模型诊断过程中的思考过程中提取关键定位步骤,反向构建故障传播链 从原始日志中反查关键定位日志证据,验证模型定位的结果,提供支撑诊断 结果的日志证据链,提升结果的可信度。







思维链可视化

将LLM的思维转为图结构

第8届 AI+研发数字峰会 | 拥抱 AI 重塑研发



PART 05

总结与展望: 智算DevOps

联合实践 & 持续创新,打造业界领先的智算集群智能化运维 / IDD Sth 解决方案 **Operation Al Agent** Ops **AI SRE Team** MR 代码提交 发现 7 * 24 代码管理 User 任务协同 Robot 调用 Al Developer **Coding AI Agent** 软件部署 信息收集 异常感知 故障分析 行动建议 Robot Robot Robot Robot 根据上下文各自诊断根因,并合议 Deploy M CP Schedule **KPI Analyzer** Log Analyzer **Alarm Analyzer** • C: 基于日志分析是否因 • C: 基于告警分析是否因 • C: 基于KPI分析是否因 **Deployment Al Agent** 软件问题导致的任务中 计算或存储硬件导致的任 网络拥塞、磁盘IO导致 断? 务中断? 的任务中断? R: xx进程异常退出 R: 任务期间无告警 • R: 未发现性能拥塞 监控运维 • A: 非计算/存储硬件问题 A: xx软件升级后, 致任 • A: 非网络拥堵, 建议排 务中断,建议查验软件 建议排查其他原因 查其他原因 Log 可观测数据 KPI

智算运维大模型(LLM) + 知识库

Operation Al Agent

科技生态圈峰会+深度研习



——1000+技术团队的共同选择





时间: 2026.05.22-23



时间: 2026.08.21-22



时间: 2026.11.20-21



AiDD峰会详情











产品峰会详情



EDEAI+ PRODUCT INNOVATION SUMMIT 01.16-17 · ShangHai AI+产品创新峰会



Track 1: AI 产品战略与创新设计

从0到1的AI原生产品构建

论坛1: AI时代的用户洞家与需求发现 论坛2: AI原生产品战路与商业模式重构

论坛3: AgenticAl产品创新与交互设计

2-hour Speech: 回归本质



用户洞察的第一性

--2小时思维与方法论工作坊

在数字爆炸、AI迅速发展的时代, 仍然考验"看见"的"同理心"

Track 2: AI 产品开发与工程实践

从1到10的工程化落地实践

论坛1: 面向Agent智能体的产品开发 论坛2: 具身智能与AI硬件产品

论坛3: AI产品出海与本地化开发

Panel 1: 出海前瞻



"出海避坑地图"圆桌对话

--不止于翻译: AI时代的出海新范式



Track 3: AI 产品运 AI 产品运营与智能演化

从10到100的AI产品运营

论坛1: AI赋能产品运营与增长黑客 论坛2: AI产品的数据飞轮与智能演化

论坛3: 行业爆款AI产品案例拆解

Panel 2: 失败复盘



为什么很多AI产品"叫好不叫座"?

--从伪需求到真价值: AI产品商业化落地的关键挑战

智能重构产品数据驱动增长



Reinventing Products with Intelligence, Driven by Data



感谢聆听!

扫码领取会议PPT资料

