



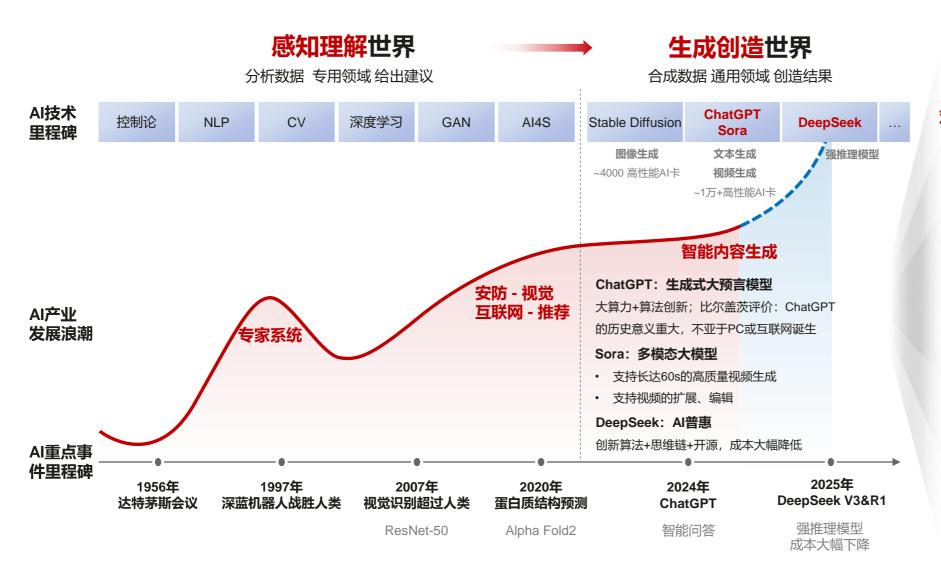
01 AI大模型发展趋势和挑战

02 DCS AI解决方案

03 典型案例



AI趋势:从感知到内容生成,从简单预测到复杂逻辑推理,AI普惠化带来新的机遇



生成式人工智能带来的机遇

对客: 利用大模型提升客户体验, 实现精准营销

远程对话助手

优秀话术推荐 对话情绪质检 智能客服多轮优化

数字化营销助手

营销文案策划 营销话术生成 产品推荐

办公: 打通业务堵点连接断点, 提升办公效率

文档助手

投研报告生成 信贷报告撰写 会议纪要生成

智能问答

业务知识问答政策制度问答

<mark>开发:全面提升开发效率,降低开发成本</mark>

代码助手

代码生成 代码注释 代码检查

数据助手

NL2SQL直接生成查询语句 图表自动生成

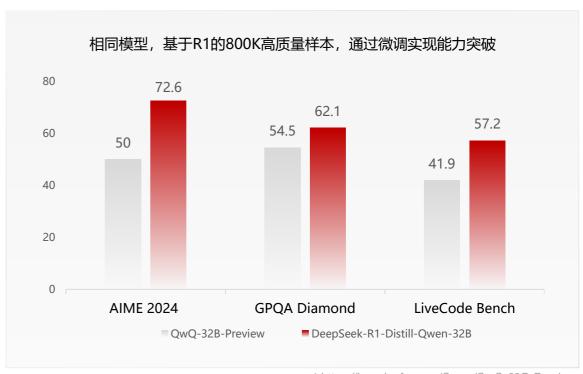


影响大模型效果的核心要素

数据是AI大模型的核心要素

数据 蓬勃发展 1000X↑ 数据 从单模到**多模** 从TB到PB 算力 逐渐趋同 以英伟达、昇腾为代表 算法 逐渐收敛 算力 算法 Transformer模型基础架构 Pytorch/TensorFlow/MindSpore开发框架

AI高度取决于数据规模和质量



^{*} https://huggingface.co/Qwen/QwQ-32B-Preview



训练&推理是大模型在行业落地必须经历的过程

大模型训练及推理基本概念

模型训练研发(开发态)从现有的数据中学习新能力,形成算法(模型) 未训练的 神经网络模型 训练框架 (开发环境/工具)

模型推理生产(部署态) 把模型部署到应用或服务中,完成具体任务



模型训练





预训练 微调/知识库 (义务教育) (专业课)

模型推理



推理/预测 (行业工作)



大模型在行业落地面临三大挑战

数据工程耗时长



数据工程占模型开发时长60%

数据孤岛、格式多样、资产繁多, 语料收集、清洗、标注等耗时长

模型训练和应用落地难



开发周期不可控

项目开发难度大,人员技术要求高

集群可用度低



集群可用度不足50%

因算力等待、任务潮汐、资源碎片化等原因



目录

01 AI大模型发展趋势和挑战

02 DCS AI解决方案

03 典型案例



DCS AI提供训推集群、训推一体机解决方案,满足不同场景用户AI使用需求



智慧医疗



智慧办公



智慧科研



智慧 金融



智慧城市

场景应用

行业大模型

伙伴大模型

开源大模型

商用大模型

AI全流程工具链

模型训推框架 | 数据处理套件 | 应用开发套件

ModelEngine

AI算力底座

高性能算力 | 超宽无损网络 | 高性能存储

企业训推场景





面向企业训推场景(百卡以内)

DCS AI训推超融合一体机

单节点起步,快速上线

• 灵活轻量: 单节点起步, 一柜式全集成, 提供计算存储网络

• 统一管理: 全栈管理, 硬件+软件平台 4小时开局

工具链:集成数据处理套件、模型微调套件、应用开发套件,加速大模型场景应用快速落地

面向智算中心场景(百卡以上)

DCS AI训推集群方案

一站式预集成、高性能训推集群

• 高效算力调度: 统一纳管异构算力资源, 提高算力利用率50%

• 超宽无损网络: 200G RoCE网络, 负载均衡有效带宽98%

• 高性能存储: 百万IOPS, 百G大带宽, 智能分级存储

• 一站式AI开发工具链, 高效AI训推开发、部署、算法迁移



企业训推场景: DCS AI训推超融合一体机助力大模型快速行业落地

客户挑战

设备部署运维难

"七国八制"部署繁杂, 设备运维、故障定界难

数据工程耗时长

数据格式复杂,清洗难度大; 海量数据集标注, 耗时耗人

模型训练和应用 落地难

大量AI组件, 手动部署适配繁杂; 应用开发周期长

AI集群可用度低

资源碎片,造成资源浪费;集群 业务存在明显的潮汐现象



典型 科研、高端企业服务 场景 复杂逻辑推理

问答问数场景 成本与性能平衡

智能办公场景

低成本和低时延

DME 资源

检查

数据处理强

- 数据清洗效率提升60%
- 问答对自动生成,留用率可达60%

方案价值

• 一键式安装部署,业务上线从天->小时级

部署运维易

• 存算网全集成,5合1全融合管理

应用上线快

- 3步完成DeepSeek模型部署,快速拉起
- 应用上线从天到小时

系统效能局

- 系统任务并发提升30%
- 基于负载的资源弹性伸缩



运维管理 模型层 DeepSeek | ChatGLM | LLaMA | Qwen | ... AI全流程工具链 | ModelEngine 监控 数据使能 应用使能 告警 AI使能 数据归集 数据 评估 日志 模块化RAG 应用编排 平台 模型使能 性能 **(/**) 监控 模型蒸馏 模型评测 模型管理 模型部署 资源 大屏 裸金属容器平台 | eContainer 报表 平台 AI开发框架 拓扑 分析 基础 软硬件 健康 CE交换机.

昇腾算力 OceanStor存储

DCS AI训推超融合一体机

部署运维易:一键式安装部署,5合1融合管理,降低部署和运维门槛

SmartKit 整柜开局 DME统一全栈管理 安装/部署/初始化 设备管理 + 可视监控 + 设备自动发现 硬件组网可视 资源一眼可视 SIRKE BIRKE 集群统一管理 向导式安装 设备可视监控&故障一眼定界







数据处理强:数据清洗效率提升60%,问答对留存率达60%



爺 统一数据归集

- 支持NFS v3、OBS协议和第三方存储
- Dorado与Pacific异构存储数据可归集

清洗效率提升60%

- 内置50+算子, 多维数据处理
- 1000+小文件处理时长1小时

留存率可达60%

• 原始语料自动生成QA对, 问答对留用率 可达60%, 领先业界50%

评估精度80%以上

• 多维自动化评估: 完整性, 一致性, 有效性等



应用上线快:三步完成模型部署,应用上线天级->小时级





模型部署复杂

100+ Al组件, 部署流程复杂, 人工适配, 效率低



技术门槛高

需要跨领域的技能(数据, 算法等),技术门槛高



业务上线慢

需解决数据格式、通信协议 等兼容性问题,工作量大

方案价值

三步完成模型部署

(以DeepSeek为例)

Step1 DeepSeek社区原始权重下载

Step2 上传权重至ModelEngine

Step3 下发推理服务部署任务

应用上线天级->小时级







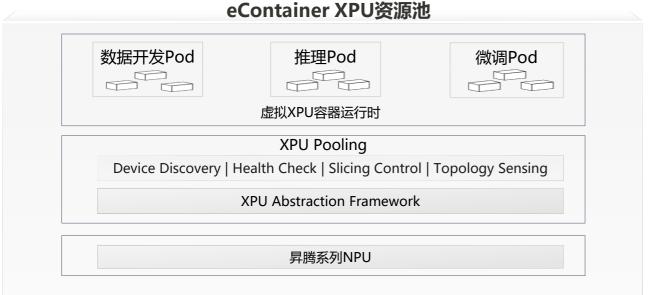
知识库构建 2小时

模型服务拉起 分钟级 应用编排 半小时



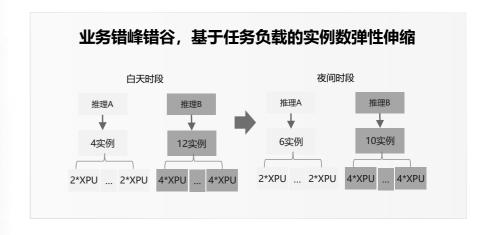
系统效能高: 资源随需伸缩, 系统任务并发提升30%





资源需求 随需伸缩

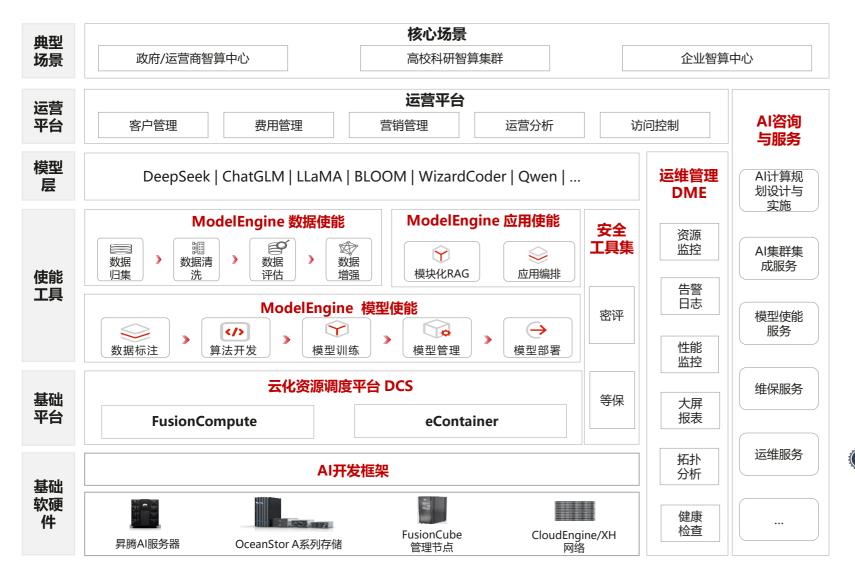
系统并发 提升30%







智算中心场景: DCS AI一站式方案助力模型高效训推开发



A BC 轻量起步,平滑扩展

> 算力8卡起步、 管理4节点起步,乐高式服务 叠加扩展

极简易用,全栈使能 Ai

> 一站式全流程AI工具链,支持数据使能,模型 使能,应用使能,零代码开发训练,应用落地 时间缩短3倍

存算协同, 极致性能

> **软硬协调优化**,集成DataTurbo 文件加速技 术,KVCache缓存推理加速,性能提升30%

智能调度,高利用率

XPU池化与高效调度,构筑业界最小1:20算 力切分粒度



AI专业服务:从建设到运维,助力客户建好用好算力平台,释放算力价值

辅助运营 建设 维保与运维 AI应用 AI应用 算力平台 (5) 维保服务 运维服务 模型使能服务 使能 AI模型 AI 工程服务 (可选) 4 运维与辅助运营 AI平台 算力平台集成服务 问题处理专享支 集群 持 基础软件 Hi-Care 子系统部署实施 集群系统集成 模型使能服务 风险预防消减 数据中心集成使能服务 ΑI 基础软件 集群系统集成服务 驻场服务 AI运力网络规 介质保留服务 AI计算规划 存储规划设计 划设计与实施 设计与实施 与实施 服务(含工程安 (含工程安装) (含工程安装) 基础 数据中心运维 硬件



目录

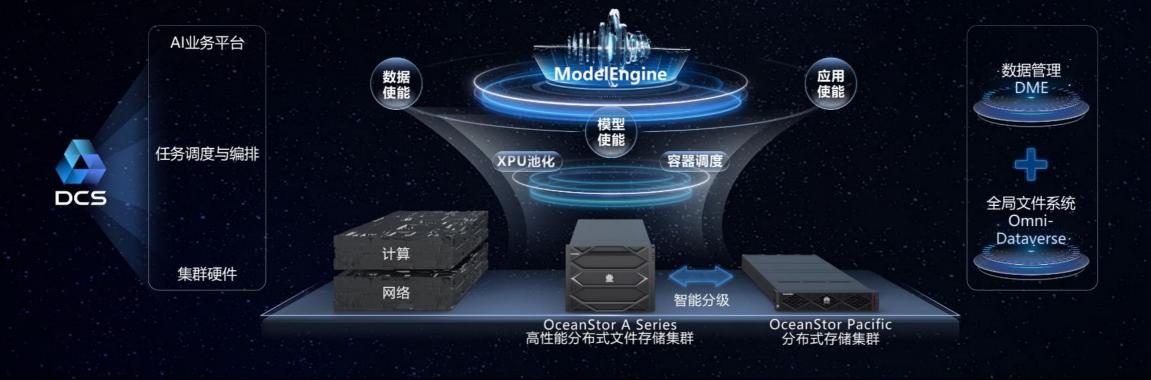
01 AI大模型发展趋势和挑战

02 DCS AI解决方案

03 典型案例



华为DCS AI解决方案



Thank you.

把数字世界带入每个人、每个家庭、每个组织,构建万物互联的智能世界。

Bring digital to every person, home and organization for a fully connected, intelligent world.

Copyright©2018 Huawei Technologies Co., Ltd. All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. Huawei may change the information at any time without notice.

