

# 浪潮信息AI解决方案分享

郭磊 浪潮AI产品经理

■ AI&HPC应用软件产品部 ● 2023年7月21日

# 目录 / Contents

**PART 1**

**金融企业智能化转型**

01

**PART 2**

**浪潮全栈AI能力概述**

02

**PART 3**

**CV 场景解决方案**

03

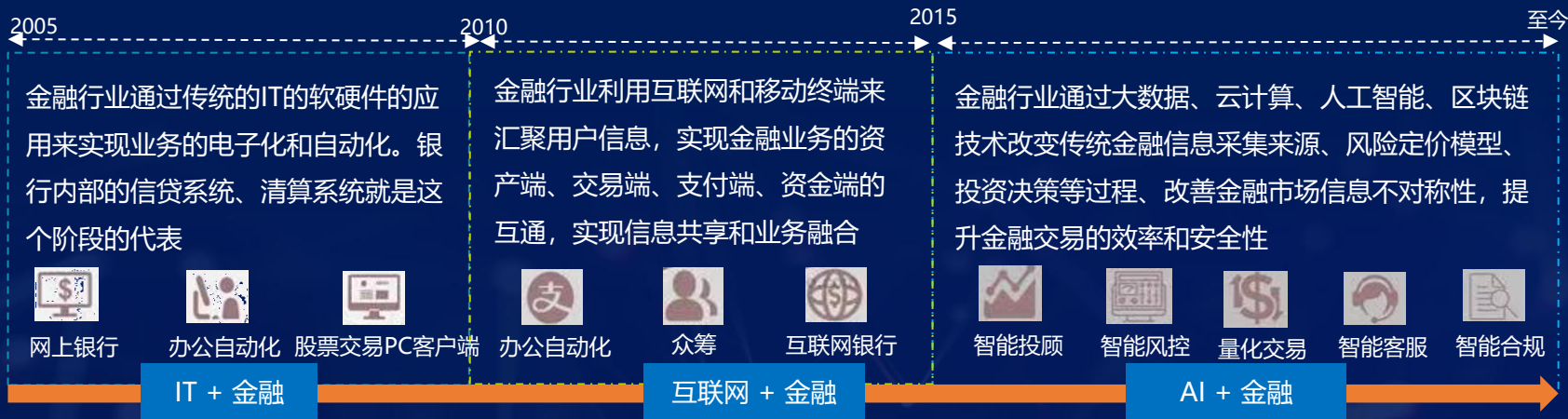
**PART 4**

**NLP 场景解决方案**

04

# PART 1

# 金融企业智能化转型



### 传统金融行业痛点

#### ◆ 人工成本高

传统金融行业存在业务流程复杂，处理周期较长，提供差异化服务能力较弱等问题，同时高度依赖人工审批及处理，导致人工成本居高不下。

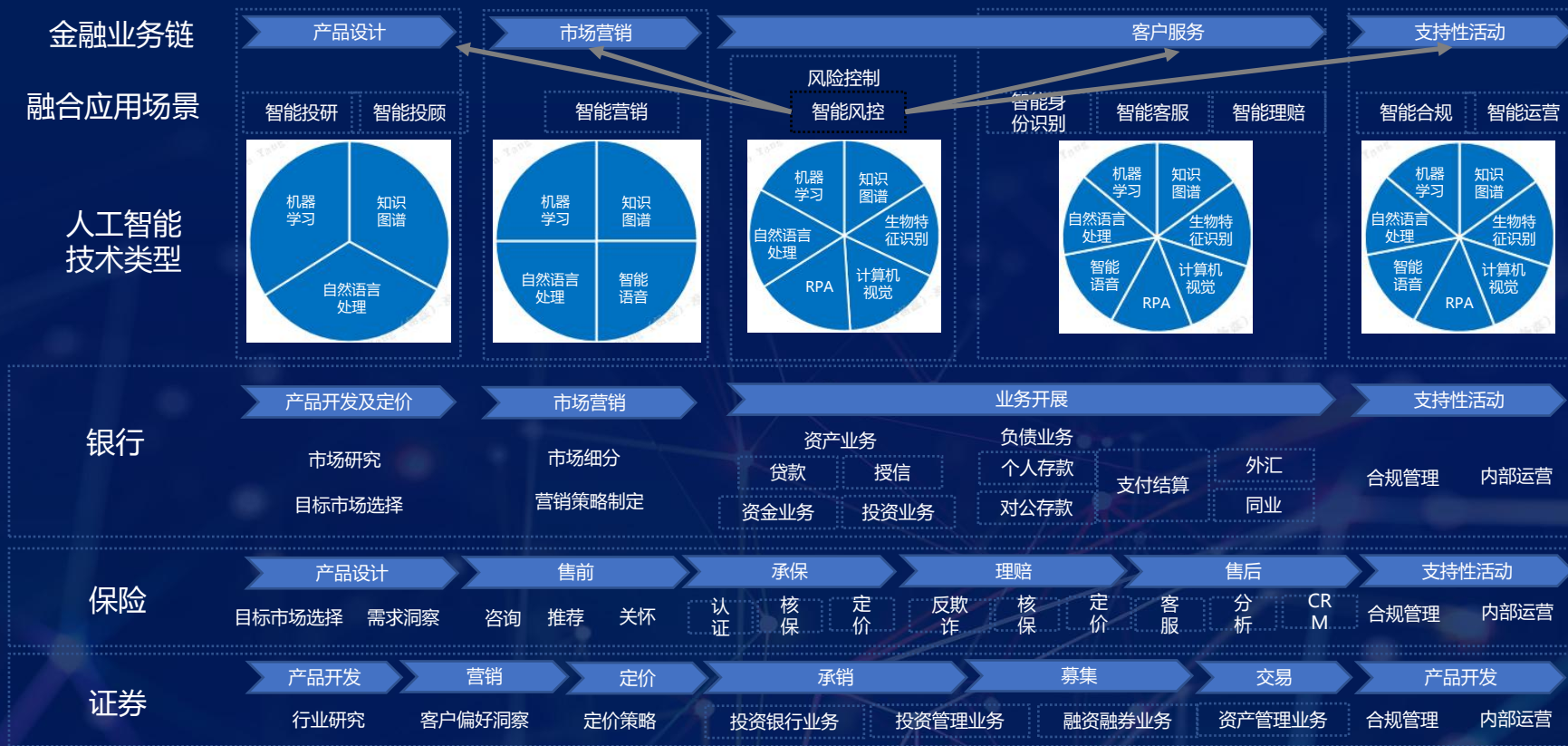
#### ◆ 金融市场信息不对称

金融各大业务场景数据量大，数据复杂、有效信息难以挖掘等情况导致了数据孤岛和数据分析效率低下等问题，增加了金融市场的信息不对称性。

#### ◆ 客户转化率低

面对金融消费者上，传统金融业务存在办理流程繁琐、差异化服务能力较弱等问题，导致部分业务场景下客户满意度较低，客户转化率难以提升。

AI + 金融，可提升金融行业数字化水平，通过实现金融业务流程自动化，降低人工成本，同时提升数据处理效率，减小金融市场信息不对称性，并且提供智能营销等差异化服务，提升客户满意度。



✓ **企业在传统的业务模式下，想要引入更多的AI业务，您需要面临多个方面的考虑...**

- GPU、FPGA、AISC等异构芯片如何统一管理，降低运维成本；
- 国产化背景下、如何统一纳管国产硬件集群与已有GPU集群，国产芯片兼容性问题；
- 基础AI业务流的构建，如何进行高效和规范化开发训练作业，解决数据孤岛，全局协同开发；
- 如何高效管理 自研/外采 第三方算法应用，降低运维成本；
- 如何解决软硬件绑定问题，应用无法共享资源池，资源利用率低；
- 中小银行客户群体端到端整体AI解决方案诉求；
- ...

→ **随着AI利用规模的扩大，基础计算环境的管控，成为制约企业高效部署AI的障碍**

✓ **我们能为您提供的解决方案包括：**

- 异构资源统一池化、GPU 资源细粒度划分、多维资源调度管理策略；
- 支持 28种AI 芯片、主流国产操作系统、混合集群统一纳管、国产化场景支持；
- 提供数据管理、模型开发、训练、测试、部署全栈AI 业务流作业；数据、算法等共享协同，提高效率；
- 多元应用灵活适配，快速集成、一键部署；
- 提供端到端 经过测试、验证联合解决方案，一站式服务；
- ...

## PART 2

# 浪潮全栈AI能力概述



左手伙伴

具备AI功能开发核心能力的科技公司  
(技术型伙伴)

MetaBrain

右手伙伴

具备实施行业AI整体方案交付的SI、ISV  
(方案型伙伴)

AI算法平台

AutoML Suite

云端&本地  
双部署

自动建模 自动调参 自动剪裁

TensorFlow-opt

512 GPU  
扩展效率高达90%

Caffee-MPI

首个并行版Caffee之一

TF2

无精度损失 加快FPGA开发

浪潮-源1.0 LMS

自研AI大模型计算框架

AI资源平台

AIStation 训练平台

模型开发  
模型构建 模型训练  
模型调优 模型导出

模型部署  
模型加载 模型部署 API服务

应用开发  
API调用 应用开发测试

AIStation 推理平台

兼容多种深度学习框架

AI模型在线测试与评估

多模型部署  
加权计算

T-Eye 天眼

AI应用与框架特征分析

AI算力平台

AI服务器

训练



NF5488A5



NF5468A5



NF5498A5



NF5688M6

推理



NF5468M6



NF5280M6

边缘



NE5260M5



F10A



F37X



N20X



M10A

AI加速卡

AI基础设施

MDC微模块数据中心



SDC2000  
AI单机柜数据中心



RDC2000  
AI单排数据中  
心



MDC2000  
AI微模块数据中  
心

CDC集装箱数据中心



CDC2000  
AI集装箱数据中  
心

LC液冷数据中心



LC2000  
AI液冷数据中  
心





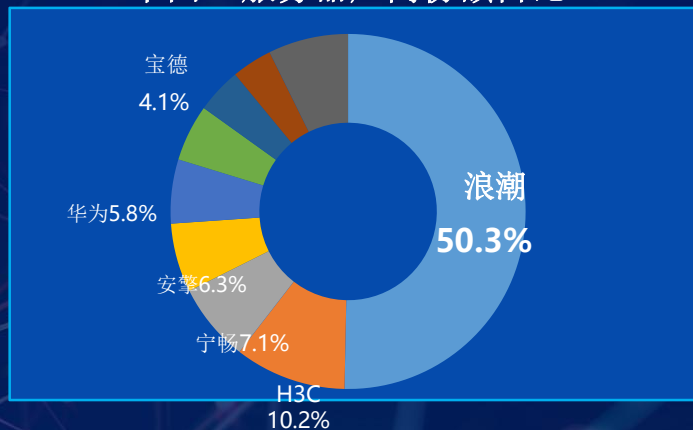
- 浪潮拥有业界最全的AI服务器产品线
- 浪潮AI服务器市场份额全球第一，超20%
- 浪潮AI服务器中国市场份额连续6年，超50%

全球AI服务器厂商份额占比

排名	厂商	市占率
1	浪潮	20.9% ↑
2	Dell	13% ↓
3	HPE	9.2% ↓

来源：IDC 2021H2 Semiannual Artificial Intelligence Tracker

中国AI服务器厂商份额占比



来源：IDC 2022H1 China Semiannual Accelerated Server Report

企业数据

模型开发与训练

模型部署与推理

AI服务

智能营销

智能风控

身份验证

智能投顾

智能保险

智能投研

智能推荐

## AIStation开发平台

数据管理 AI算力 AI开发环境 模型构建

数据导入

按需获取

构建环境

模型训练

数据预处理

自动调度

开发工具

模型可视化

数据加速

性能优化

开发流程化

超参搜索

PyTorch

TensorFlow

Caffe

PaddlePaddle

MxNet

Kubernetes 平台



AI计算资源池

AI存储资源池

AI网络资源池



2天 → 4小时

模型训练时间



40% → 80%

计算资源利用率

## AIStation推理平台

模型服务

部署发布

数据处理

应用场景

TensorFlow Serving

按需获取

构建环境

推荐系统

TensorRT Inference Server

自动调度

开发工具

CV场景

PyTorch Inference Server

性能优化

开发流程化

NLP场景



本地部署



私有云



公有云



2天 → 5分钟

模型一键部署



应用负载均衡

资源弹性伸缩

OCR

生物特征

虚拟数字人  
AIGC

双录质检

RPA

## 元脑升级

IPF 2021 智慧·创新

## 爆发

生态升级实现业务井喷发展

IPF 2020 筑基智慧世界

## 加速

汇集30家左手+20家右手  
开启元脑生态的磨合

AIC 2019 人工智能计算大会

## 发力

元脑三大基础平台

IPF 2019 智慧凝聚

## 启航

浪潮“元脑”顺势而生

## AIStore (爱士多) 发布

200

左手伙伴

400

2倍

200

右手伙伴

4000

20倍

## 认证伙伴指数级增长

## 200个行业应用创新案例



金融



智慧城市



铁路



交通



社区管理



石油



电力



水利



教育



医疗



新基建



泛行业

## PART 3

# CV 场景解决方案

# 3.1 计算机视觉可解决金融业务痛点分析

- ◆ 金融机构众多营业网点所产生的数据量巨大，数据处理工作量众多，对人力造成工作负担的同时，也提高了管理成本。金融机构的安防与风控场景急需改进工作方式，解放人力与降低运营管理成本。计算机视觉产品的引入不仅能提高金融机构内控管理效率，加速工作方式向智能化、标准化方向转变，而且能降低运营管理成本，以低成本换取高效益。
- ◆ 典型CV 场景包括：**OCR 票证识别、生物特征识别、RPA 数字员工、智能双录**等场景；

## 计算机视觉应用改善金融机构业务痛点



- 智能考勤系统、实时监控等计算机视觉产品可助力金融机构提高员工上下班管理、日常工作管理效率，降低管理成本。
- 实时监控、身份核验等计算机视觉产品助力金融机构安防与风控智能化升级，提升金融机构运营服务效率，降低运营成本。
- **OCR**产品、实时监控、人机交互等助力金融机构流程自动化与网点智能化探索，促进金融机构业务运营的发展。

- 图像识别技术应用在金融行业应用十分成熟且广泛，包括影像分类，手写识别，表格识别等，整体应用识别率达到96%以上；

<h2>图像识别</h2> <p>包括图像分类、内容识别、合同比对等多场景实现技术应用，并将持续积累图像识别应用领域知识</p>	身份证	驾驶证	银行卡	行驶证	名片
			车辆合格证	车辆登记证	营业执照
			护照	税务登记证	户口本
			军官证	对公开户申请书	
	结婚证	业务结算申请书		台胞证	士兵证
资产负债表	保单	合同	增值税发票	.....	



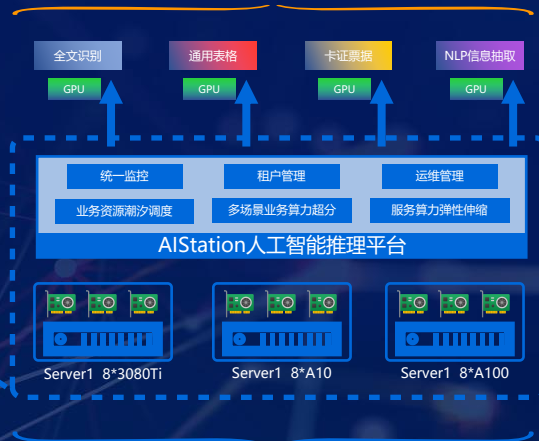
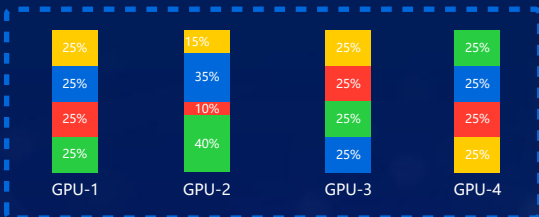
为有OCR自主开发需求的企业提供低代码、自动化的一站式OCR开发平台



Textin Studio文字识别训练平台

统一监控与管理，优化管理模式，简化运维操作

一键部署



AIStation人工智能推理平台

- 浪潮**AIStation** 推理服务平台目前已经同多家合作伙伴就生物识别场景进行合作和项目落地，包括眼神科技、**YC**、**STKJ**、百度、**KSKJ**等。



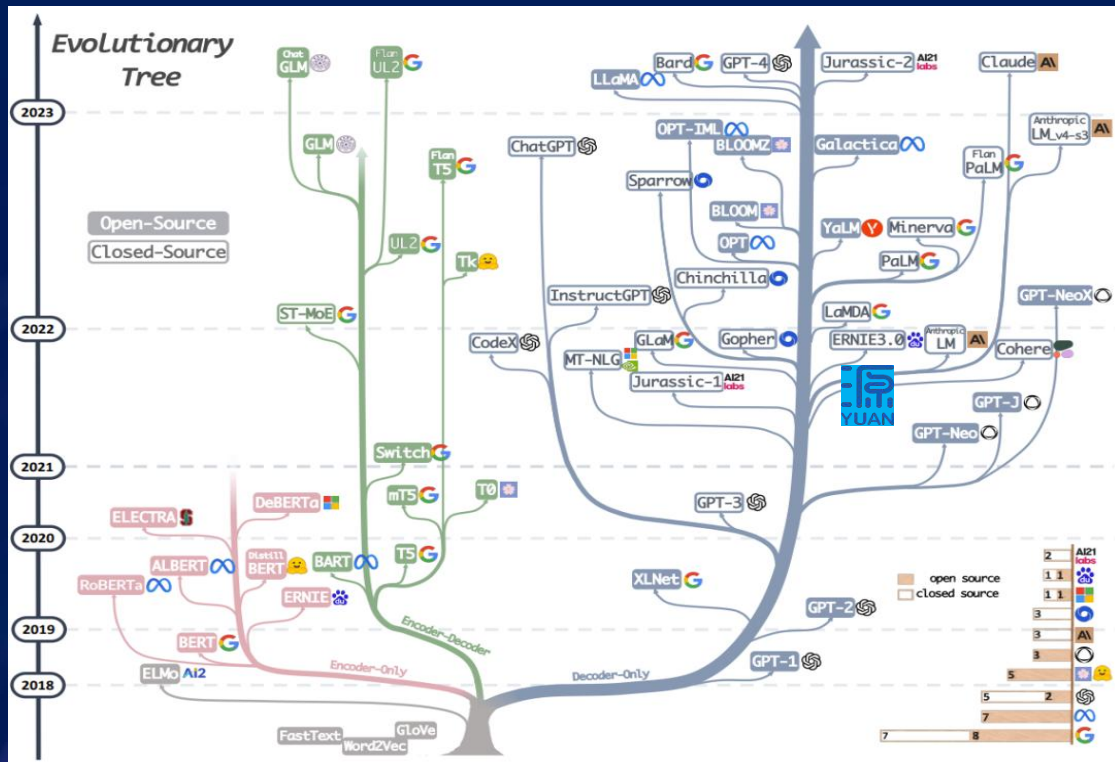


- ◆ 浪潮信息联合眼神科技共同打造生物特征识别平台，该平台基于浪潮NF5280M6服务器搭载英伟达NVIDIA L4芯片，预置人脸识别、虹膜识别、指纹识别、指静脉识别等多种核心算法及多模态融合算法，提供一站式“生物识别与身份认证”能力，搭配浪潮AIStation人工智能推理服务平台，支持应用一键部署及全生命周期管理能力，方便快速构建各类AI+应用系统和AI+解决方案，为各行业智慧化转型提供源源不断的AI动力，快速实现原有方案的AI+增值赋能。



## PART 4

# NLP场景解决方案



语言模型的发展历程

什么样的模型架构更适合面向中文的巨量语言模型？

### 国内外发展趋势上：

三大类模型：编码器-解码器或仅编码器（encoder-only）语言模型和仅解码器（decoder-only）语言模型。自2022年开始，仅解码器模型正逐渐成为 LLM 发展的主导模型。

### 工程实践上：

测试结果表明，单向生成式结构拥有更好的Zero-shot和Few-shot学习能力。

### 源1.0大模型

采用了类GPT3的解码器结构

参数	源1.0
词表长度	53228
层数	76
隐向量长度	16384
前馈网络尺寸	65536
注意力层头数	128
输入序列长度	2048
参数量 (亿)	2457.3



- 大模型训练基础设施需要把大规模算力、海量数据、先进的AI算法有机结合，为长时间、多样化训练任务的稳定可靠运行，提供高效的生产力工具

## 大模型训练任务的特点

- 大规模算力，数千张GPU卡协同计算
- 海量数据的无监督/有监督计算
- 长时间持续计算
- 多样化任务



## 大模型训练基础设施须具备的特性

- 大规模、高性能计算资源
- 高效的数据标注和处理系统
- 先进的训练算法框架
- 稳定可靠的任务生命周期管理系统
- 大规模算力集群资源管理与调度系统
- 保障持续开发、持续训练的自动化工具

- 由高性能共享存储系统、高性能AI服务器、高性能网络系统构筑高性能算力集群，结合异构算力管理与调度系统、大模型训练作业管理平台，共同构建高性能、可扩展的大模型训练基础设施

- 完整的大模型开发训练工具链
- 多人协同
- 大规模异构算力池化管理
- **高性能、可扩展**的算力调度系统
- 容器化算力 按需供给

#### 计算系统可扩展

- AI服务器规模扩展：100台 -> 500台 -> 1000台 -> ...
- 异构算力扩展：NV GPU + XPU + ...
- 异构网络扩展：以太网 + Infiniband/RoCE

**存储系统可扩展** TB -> PB -> ...



- 大模型训练涉及海量数据集和超大规模参数量更新，例如GPT-3拥有1750亿参数，需要数百乃至数千台GPU服务器协同并行计算，训练过程的数据通信要求决定了算力集群的整体设计

### 大模型训练的常用并行策略：突破算力墙和存储墙

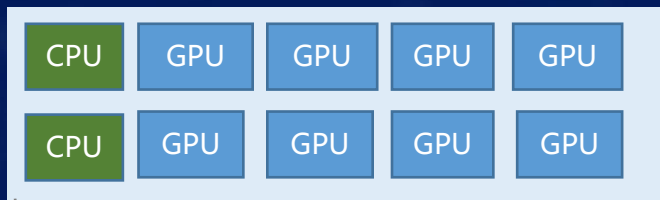
- 模型并行**：单台服务器内部多GPU之间采用模型并行，对Transformer网络进行均匀切分
- 流水线pipeline并行**：GPU服务器组内机间采用流水线并行，每组放置一个完整的大模型，每台服务器放1/N模型
- 数据并行**：GPU服务器组与组之间采用数据并行

并行策略	通信操作	对集群通信的要求
模型并行	all-reduce	节点内部高速互联通信
流水线并行	send/receive	节点间p2p低延迟
数据并行	all-reduce	节点间高吞吐

节点内 (优选)：NVLINK高速互联

节点间 (优选)：高速RDMA (Infiniband/RoCE)

GPU服务器



GPU服务器组



GPU集群



- 大模型训练集群可达数千节点，训练作业可达数千GPU卡并行计算，须兼顾低延迟和高吞吐

### 算力集群设计：

- GPU服务器配置 多IB/RoCE卡：提高机间互联带宽
- 集群网络分层：管理网络（业务面）、计算网络（参数面）、存储网络（存储面）分离
- 流水线并行的数据通信最小跳数可达
- 易于扩展到更大集群规模

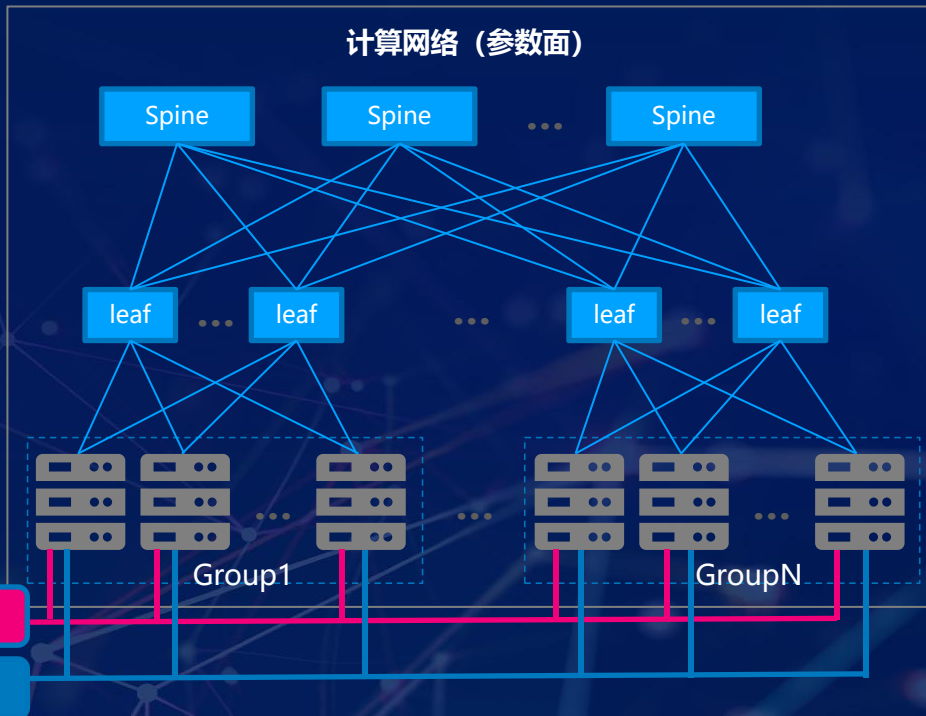
### 案例：

- 每台GPU服务器配 8 IB，带宽8x100Gb
- 集群网络：管理网络10Gb，计算网络8x100Gb，存储网络25Gb

集群网络设计影响算力管理与调度系统的设计，容器网络须与集群物理网络一致。

管理网络（业务面）

存储网络（存储面）



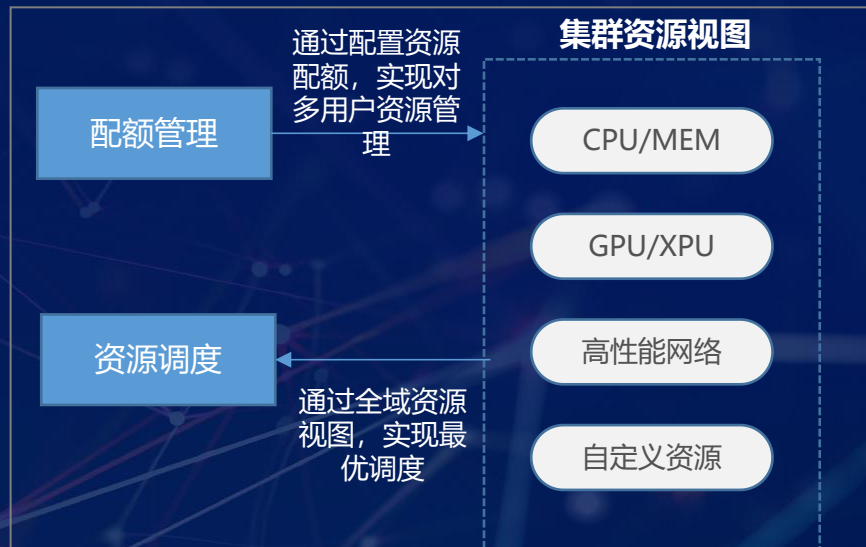
注：在实际场景中，集群网络需要根据项目具体场景来设计，设计思路相同



- 算力管理与调度系统对高性能算力集群进行统一池化管理，根据大模型训练作业对资源的需求，按需提供集群内最佳资源，并以容器化形式供给算力

#### 算力集群管理与调度系统设计要求：

- **支持团队协作**：多租户、配额
- **用户作业不感知底层资源的复杂性**：容器化、资源池化、全局资源视图、全局资源调度、亲和性策略、拓扑感知
- **大模型开发模式遵循“开发-调试-小规模短期预跑训练-调试-大规模长期训练”**：多尺度作业调度，包括小尺度资源调度，大尺度资源调度、高性能调度
- **容器网络与集群物理网络一致，保证容器互联性能**
- **支持数千上万张GPU卡**：高可扩展性
- **可统一管理异构算力集群**



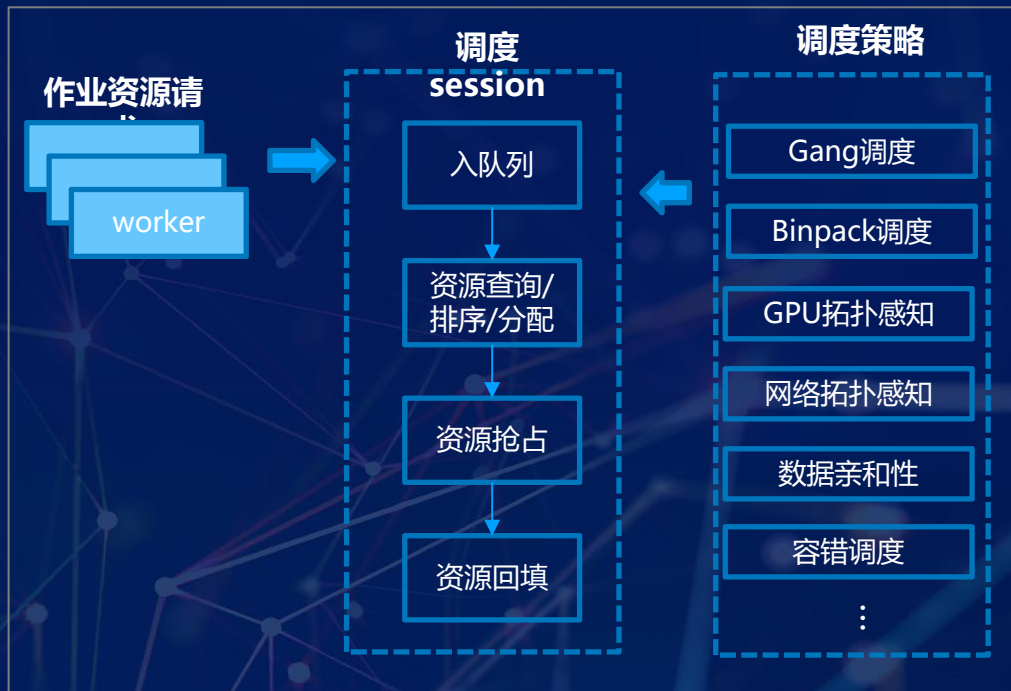
- 算力调度器通过动态、智能地管理和调配集群计算资源，制定合理的作业执行计划，以最大限度地利用资源，满足各类训练任务的时延和吞吐需求，保证作业高效稳定运行，实现算力平台高利用率、强扩展性、高容错性

### 算力调度特性：

- 全局资源池化调度，容器化供给
- 高扩展性：万卡集群、千卡作业
- 调度策略热插拔

### 调度策略：

- **Gang调度**：大规模并行作业
- **弹性调度**：弹性作业
- **Binpack调度**：减小资源碎片
- **GPU拓扑感知**：节点内资源分配
- **网络拓扑感知**：多节点资源分配
- **数据亲和性**：节点数据缓存感知分配
- **容错调度**：任务级资源亲和性感知
- ...

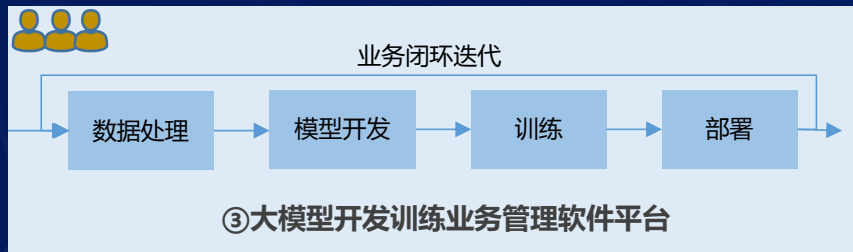


# 4.1 ③大模型开发训练业务流管理软件平台

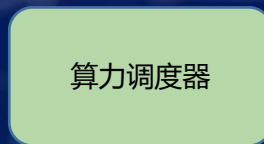
- 仅有高性能算力硬件集群和算力智能调度系统，还不足以高效地开展大模型训练业务，通过嵌入大模型开发训练业务流管理的软件平台，将高性能算力、智能调度、大模型开发业务有机结合，充分发挥算力价值

## 大模型开发训练业务管理软件平台具备：

- 从数据处理到模型部署全流程的工具链
- 统一的软件环境，兼容各种框架和工具
- 灵活的作业提交与监控，友好的UI界面和命令行工具
- 与调度系统深度集成,弹性管理资源，动态扩缩容
- 方便高效的模型开发、管理与部署
- 丰富的日志、监控与优化工具
- 实现作业高可用与容错能力
- 开放与兼容性
- 完备的数据与模型安全保障



①高性能算力集群



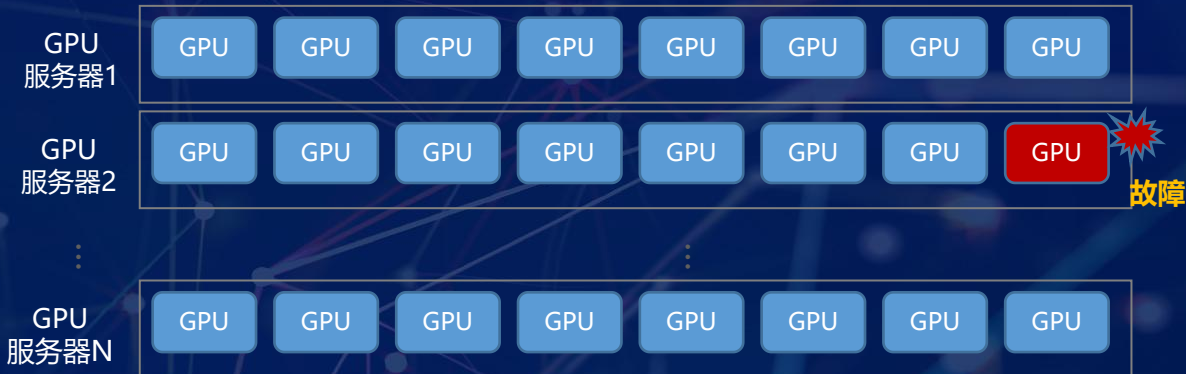
②智能算力调度

**目标：实现一个易于使用、强大灵活而可靠的大模型训练管理平台,使开发者可以专注于模型和算法创新**

- 大模型训练作业需要数百~数千卡并行计算数天至数十天，高容错性与完整的作业生命周期管理，是实现高稳定与高效率的训练过程并获得最佳结果的保障。



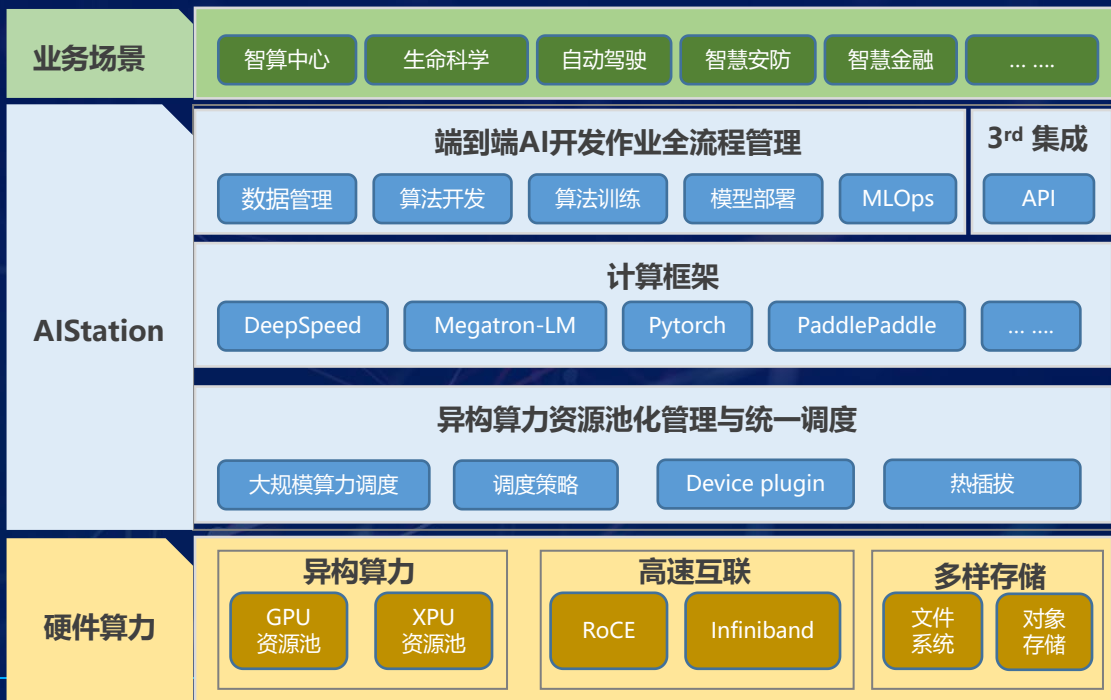
- 长期训练阶段需要大量GPU卡并行计算，任一GPU卡故障均会导致训练任务中止
- 单卡故障率 x 卡数导致 集群硬件故障率猛增
- 实时监控节点和GPU状态，检测到硬件故障后作业从最近checkpoint重启训练
- 作业整个生命周期托管运行，无需用户干预，释放低效时间



- 浪潮信息人工智能开发平台AIStation，聚焦企业级AI算法开发业务场景，提供异构算力集群统一池化管理和多策略调度、多样化建模、作业全生命周期管理、故障容错、集群监控运维等一体化能力，**为用户提供敏捷可靠的AI开发流程支持和可扩展的人工智能基础设施管理能力**

- 上述大模型开发训练所需的所有能力均已融合到AIStation中
- 支持异构算力集群、RoCE和IB集群统一管理
- 支持共享文件系统、对象存储、大数据平台对接
- 支持细粒度GPU作业
- 镜像分发加速、数据缓存加速
- Restful API可快速被集成到客户现有系统，扩展大模型开发训练能力

- ✓ **已成功训练「源1.0」2457亿参数大模型**
- ✓ **金融、能源、政府、教科研等百余家客户正在使用**





# THANKS

**浪潮信息**

浪潮电子信息产业股份有限公司 | 地址：中国·北京市海淀区信息路2号 | 电话：0086-400-860-6708 | 网站：[www.ieisystem.com](http://www.ieisystem.com)