

GOTC

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

OPEN SOURCE , OPEN WORLD

「AI、大数据与数字经济论坛」专场

本期议题：AI融合大数据，助力产业数字化升级

韩炳涛 2021年07月10日

AI融合大数据，技术中台赋能数字化转型

1 数据统一



数字生活

N个
数字化场景



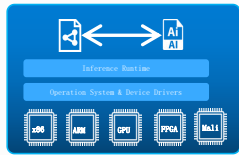
数字经济



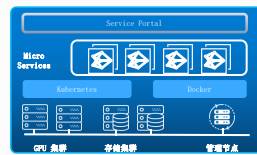
数字治理

2 高效开发

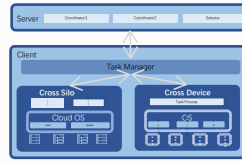
3种
应用部署模式



端



云



联邦

3 云端协同



开箱即用

自动分布式训练

推理加速



数据治理



低代码建模

1个

大数据+AI 中台

场景闭环

非侵入式

开箱即用

以分布式云为基础，实现大数据+AI全场景部署

1 全场景算力

工业控制 | 极速

- 5ms超低时延
- 裸容器轻量化部署

视频监控 | 极宽

- 420G大带宽
- 一体化部署，分钟级交付

AI训练 | 极强

- 硬件加速
- 2000节点规模部署

2 统一双核云底座



3 积木式行业云套件



4 集中化统一管理

AnyService

AnyScale

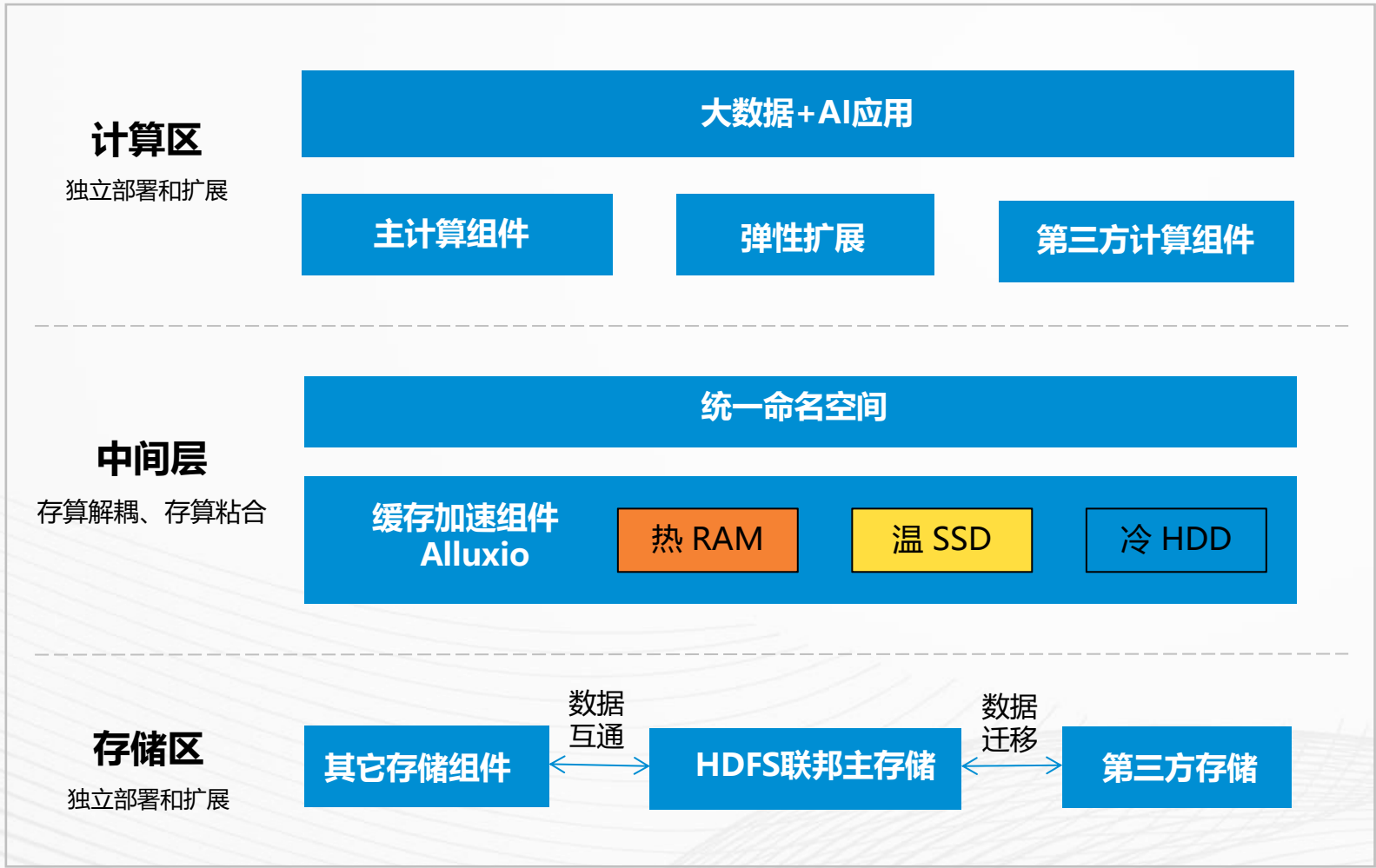
Anywhere



• 技术挑战

- 高效的存算分离架构
- 批流一体的计算架构
- 基于K8S的统一调度器
- 数据分析和深度学习统一建模
- 自动化模型部署

存算分离架构优势与挑战

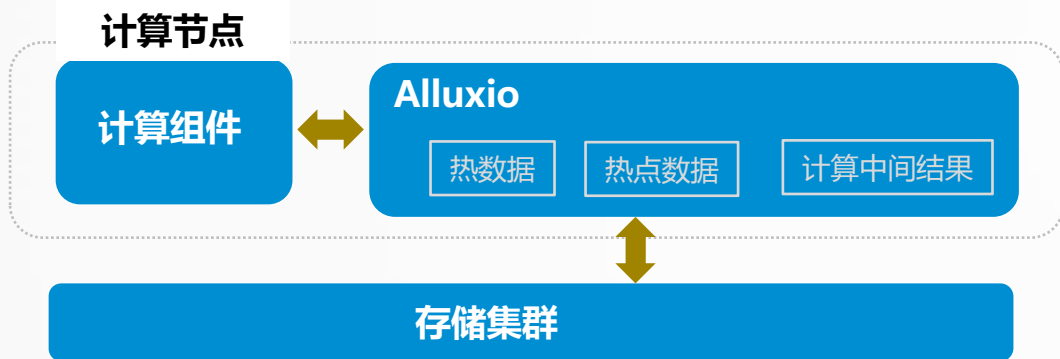


- **存储、计算解耦，各自独立集群**
 - 各自弹性扩缩容，减少浪费、提高资源利用率
 - 专用的存储集群可实现跨文件系统的数据融合
 - 计算集群可以更为灵活的部署算法
- **分级缓存加速数据读取**
- **基于开源接口实现，上层应用无感知**

存储计算一体	存储计算分离
性能高	性能中
不够灵活（计算力、存储量、应用量）	组网灵活，按需增减
硬件采购成本高	异构硬件，降低成本

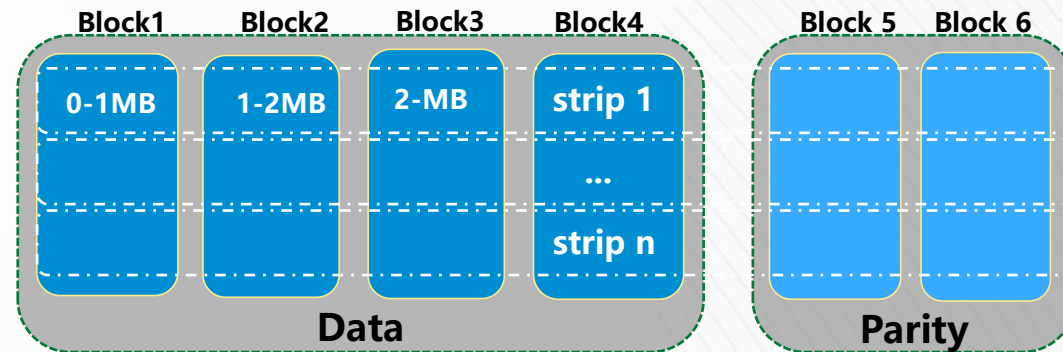
解决存算分离的性能降级问题

Alluxio缓存



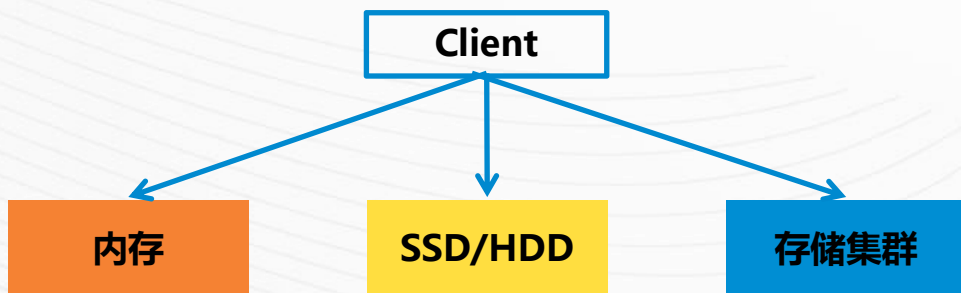
IO密集型场景性能最多可提升40%，总带宽节省10%-50%

纠删码



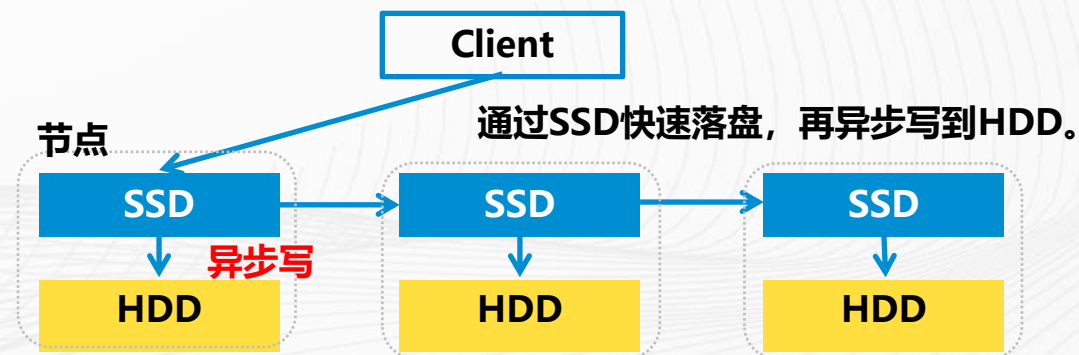
相较HDFS（三副本），写性能提升40%，存储节省50% (4+2)

多模shuffle



根据不同的配置客户端shuffle到不同的系统中

写缓存技术



比ALL_SSD性能低17%，比ALL_DISK性能提高108%

批流一体的计算架构

业务层

OLAP、报表取数服务

AD-Hoc即席查询、监控、在线类服务

数据湖层



数据
修补

Log Message

RDS



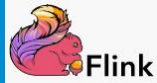
kafka

明细层加工

清洗

关联

转换



Flink



kafka

汇聚层加工

轻度
汇总

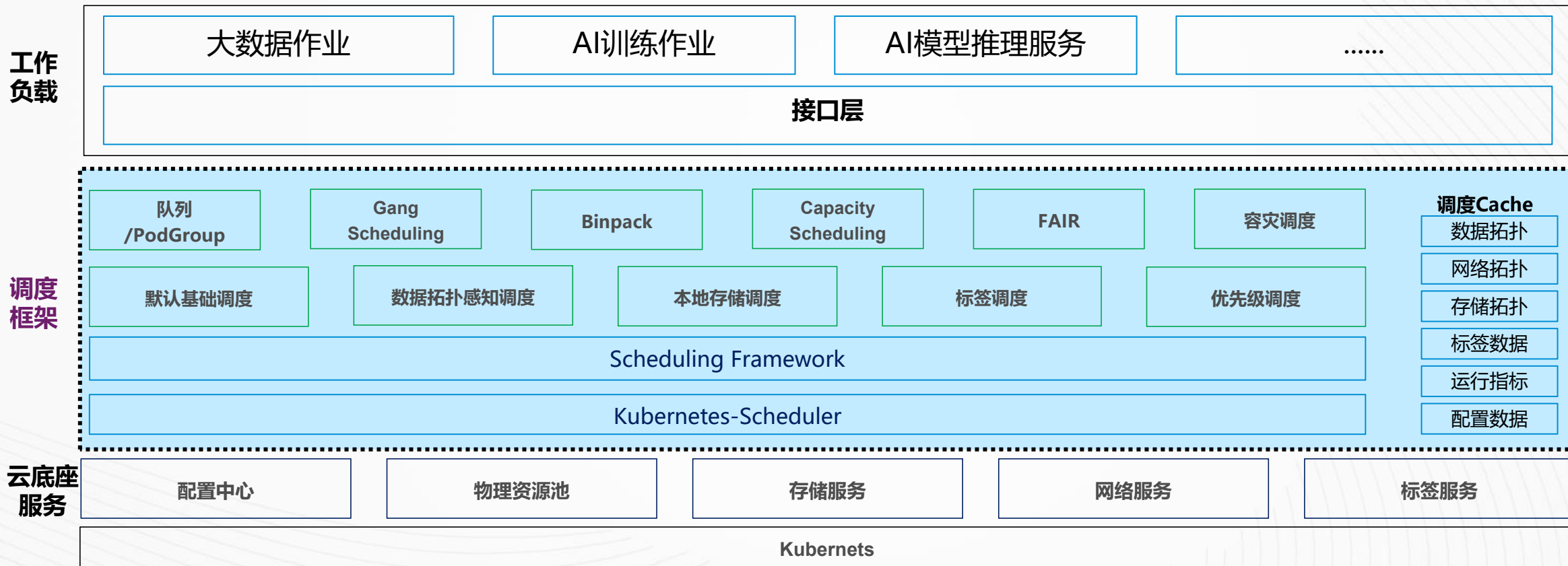
高度
汇总



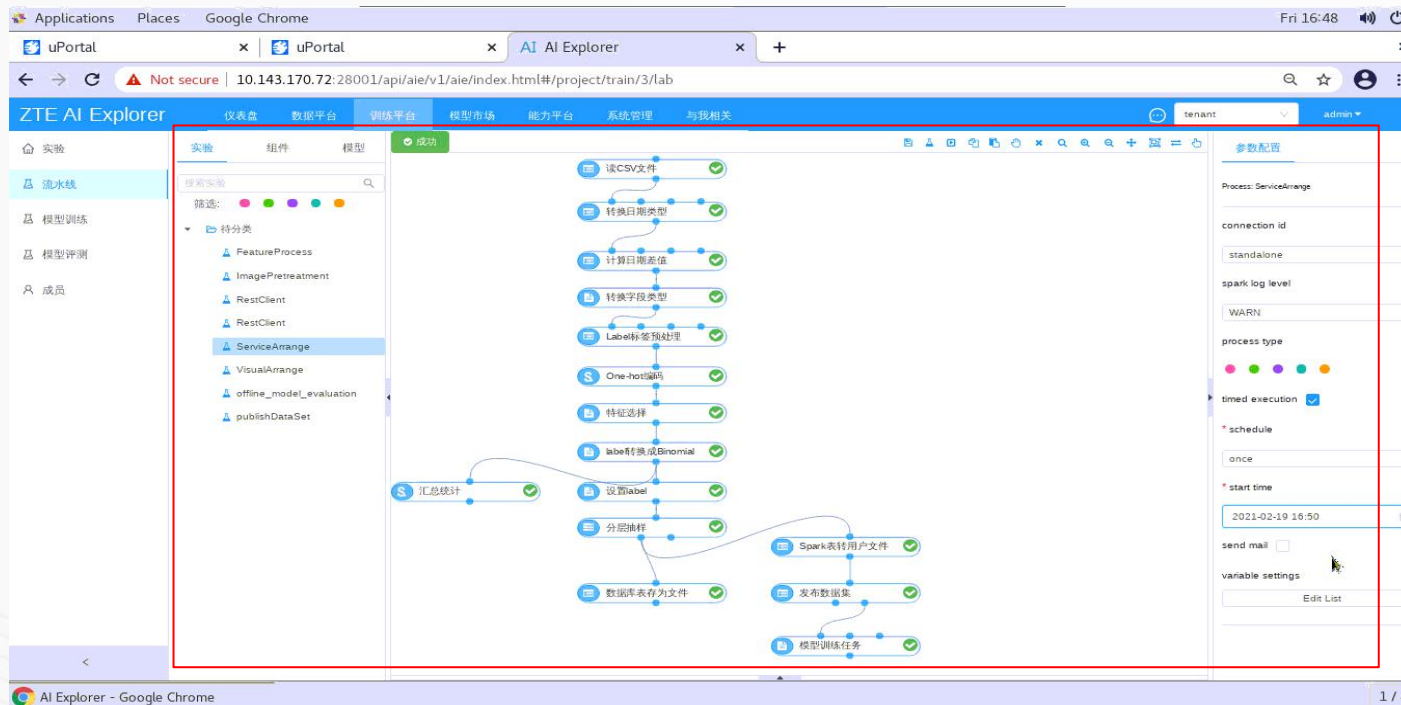
Flink

特点：流批一体，计算引擎统一；湖仓一体，统一存储，降低数据孤岛问题

基于K8S统一调度框架



- ◆ 面向批量计算：支持大数据、AI类型任务的调度，将数据计算类型中常用多Queue、Gang Scheduling、Capacity Scheduling等特性，融入到原生Kubernetes中，保证对社区原有调度能力完全兼容性，并与K8S解耦
- ◆ 支持资源队列，支持多租户场景下的资源调度
- ◆ 支持细粒度资源调度，保障资源共享与隔离



全流程可视化

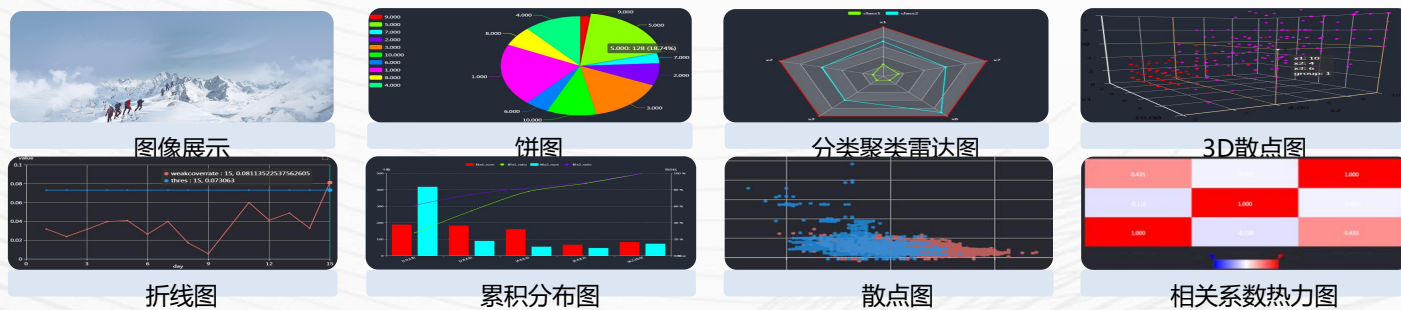
- “数据预处理、模型训练算法、模型效果、模型发布、能力部署”全流程可视化编排

机器学习、深度学习、强化学习可视化算子

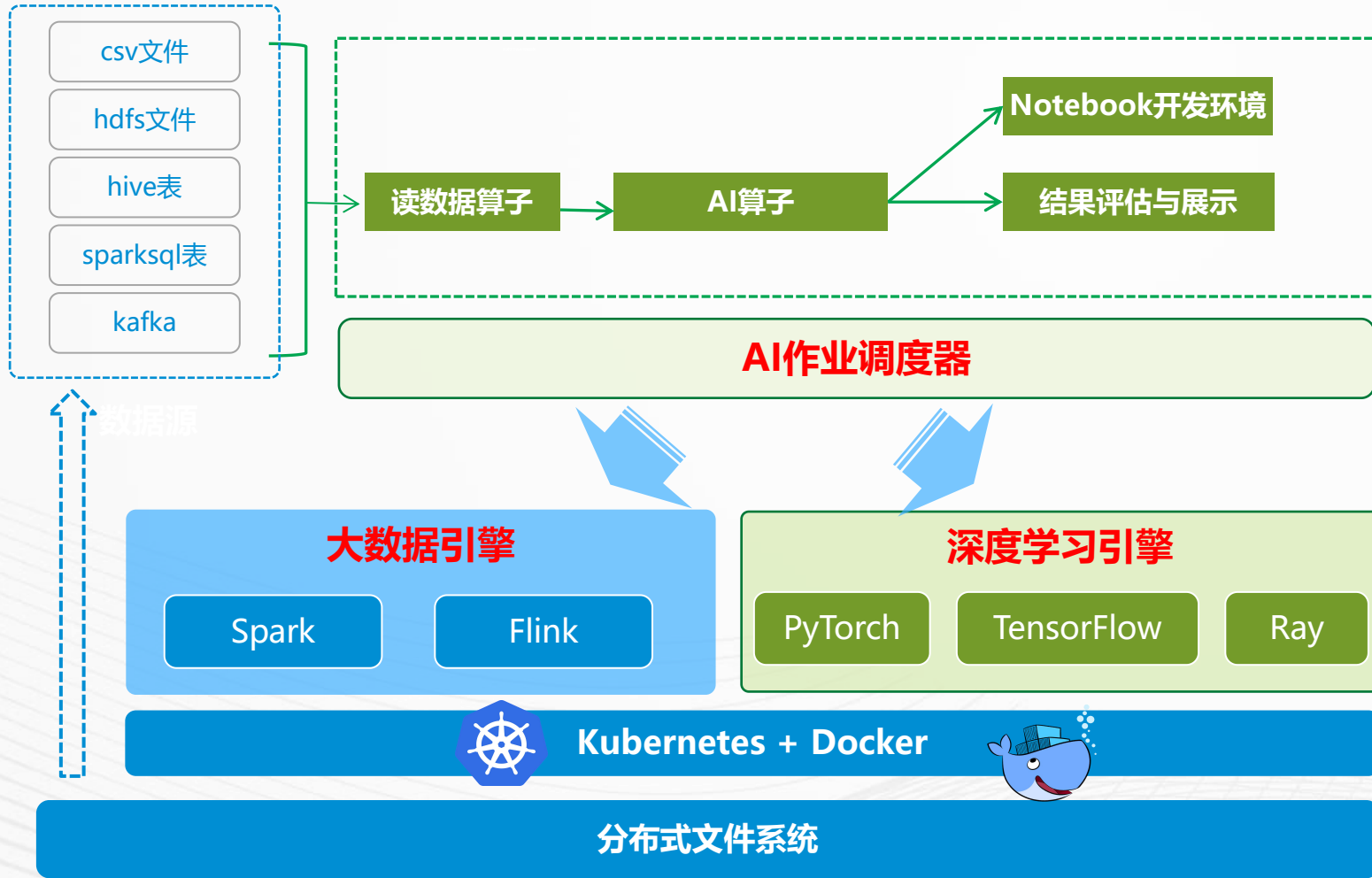
- 分类、回归、聚类、推荐等机器学习算子
- CNN、DNN、RNN、GAN、BERT等深度学习算子及模型
- 深度DQN、DDPG等强化学习算子

数据、过程、结果可视化

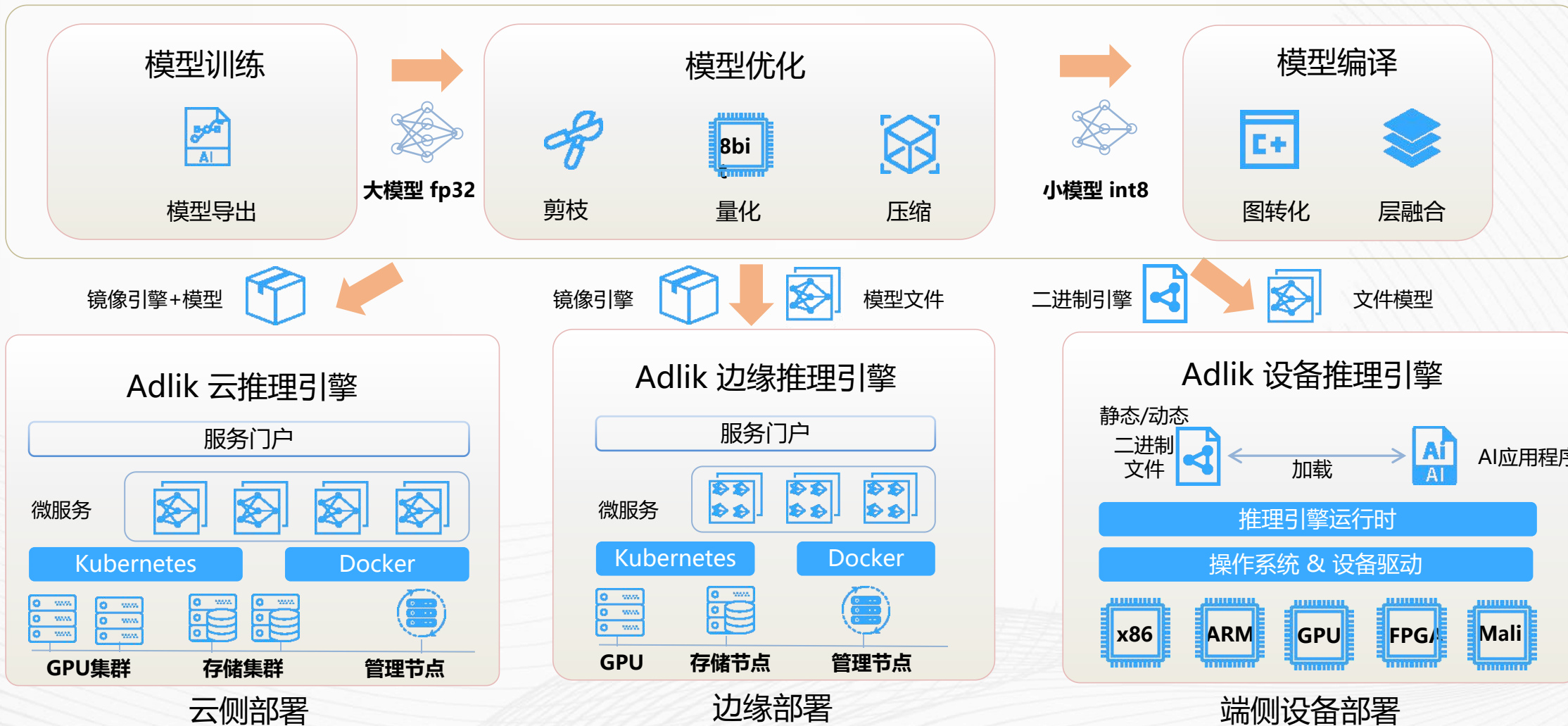
- 柱状图、折线图、散点图等数据可视化工具
- 训练过程 LOSS 曲线可视化
- 训练结果可视化评估工具



根据算子类型派发至对应的计算引擎



- **与现有大数据集群共部署**
 - 集中数据存储，节省存储空间，也避免大规模数据迁移造成的效率降低
 - 计算资源得到充分利用，减少浪费降低总体硬件投入
- **多引擎统一调度**
 - AI算子可选Spark、TensorFlow等不同计算引擎，调度器派发任务到相应引擎执行
 - 基于相同的分布式存储，实现多引擎间数据交换
- **统一编排**
 - 将数据和算子编排在一起，通过DAG描述算子间依赖关系



大数据+AI融合技术趋势，进一步提升全场景部署能力

GOTC



全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE

GOTC

THANKS

全球开源技术峰会

THE GLOBAL OPENSOURCE TECHNOLOGY CONFERENCE