

基于超算云的高性能计算服务化平台

湖南大学 国家超级计算长沙中心 唐卓

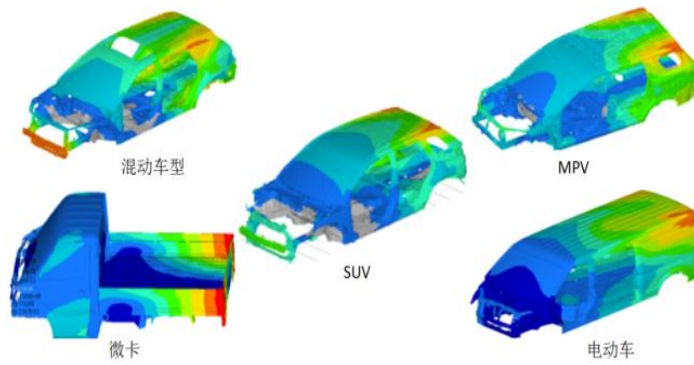
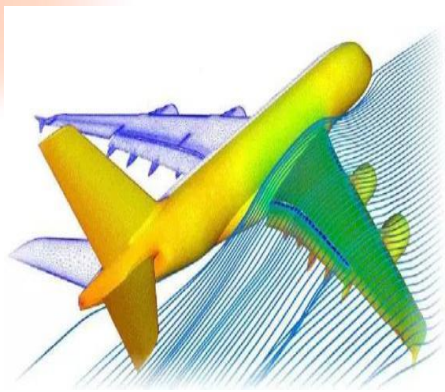
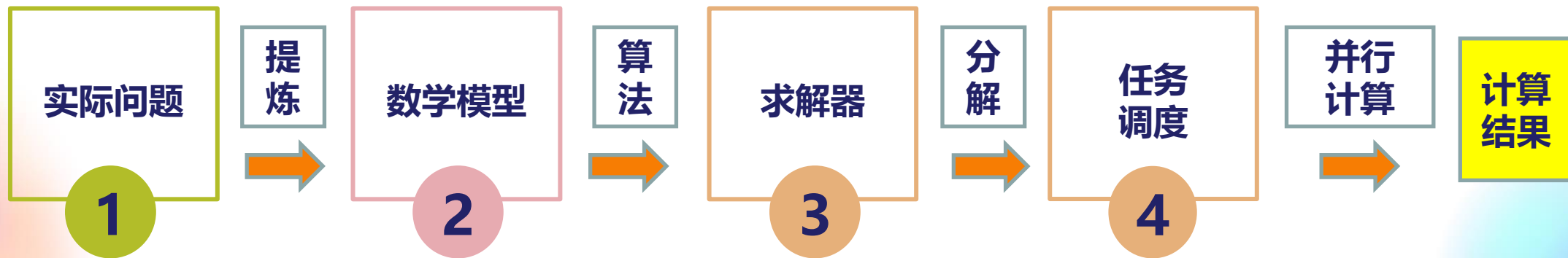


2021 可信云大会
2021 TRUSTED CLOUD SUMMIT
数字裂变 可信发展

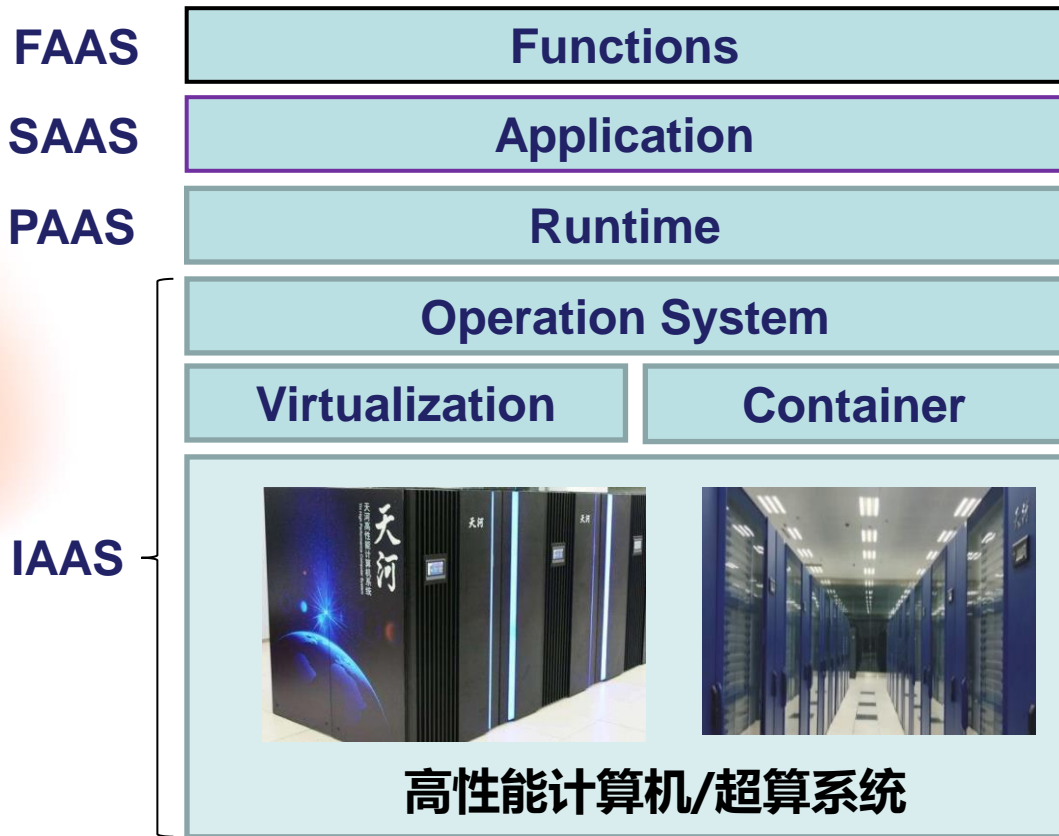
- 1、背景和挑战
- 2、高性能计算服务化关键技术
- 3、超算云服务平台及整体架构

背景：高性能计算服务化是面向业务需求领域的全流程服务化

TRUCS



背景：高性能计算服务化=高性能计算集群+云原生+函数计算



基础数学函数库（高性能计算算子）

归并
比对

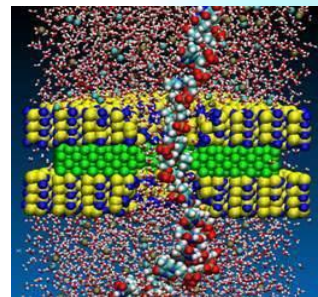
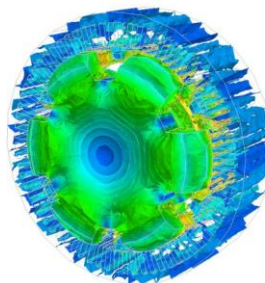
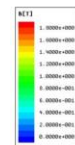
N体
方法

波尔兹曼磁
流体力学

快速傅里叶
变换

稀疏矩阵
向量乘

前后处理模块
冲压成形仿真
体积成形仿真
结构力学分析仿真
热力学分析仿真
裂纹扩张仿真
工程优化工具箱
材料数据库



SAAS

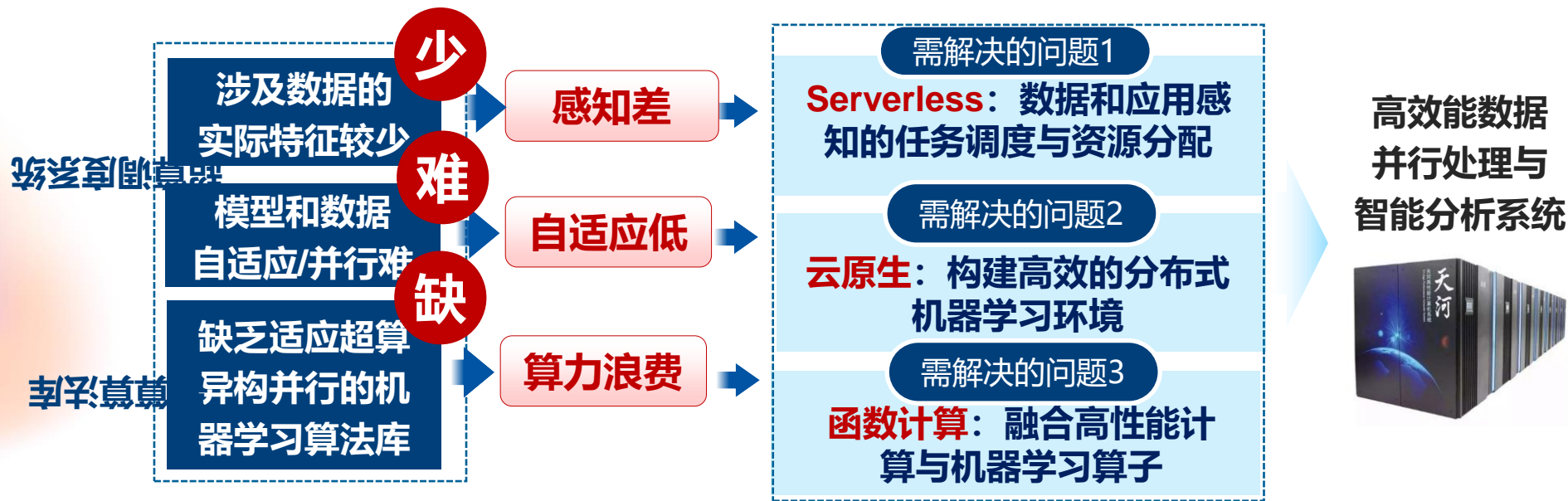
天河物理节点



背景：大规模异构是超级计算发展的主流



挑战：传统超算在操作使用上难以原生适应高效的计算服务化，智能应用场景、AI计算特性为高性能计算服务化提出了更高要求



如何基于现有主流超算系统的系统结构

构建高性能计算云服务基础设施，解决HPC for 大数据与AI计算的难题



- 1、背景和挑战
- 2、高性能计算服务化关键技术
- 3、超算云服务平台及整体架构

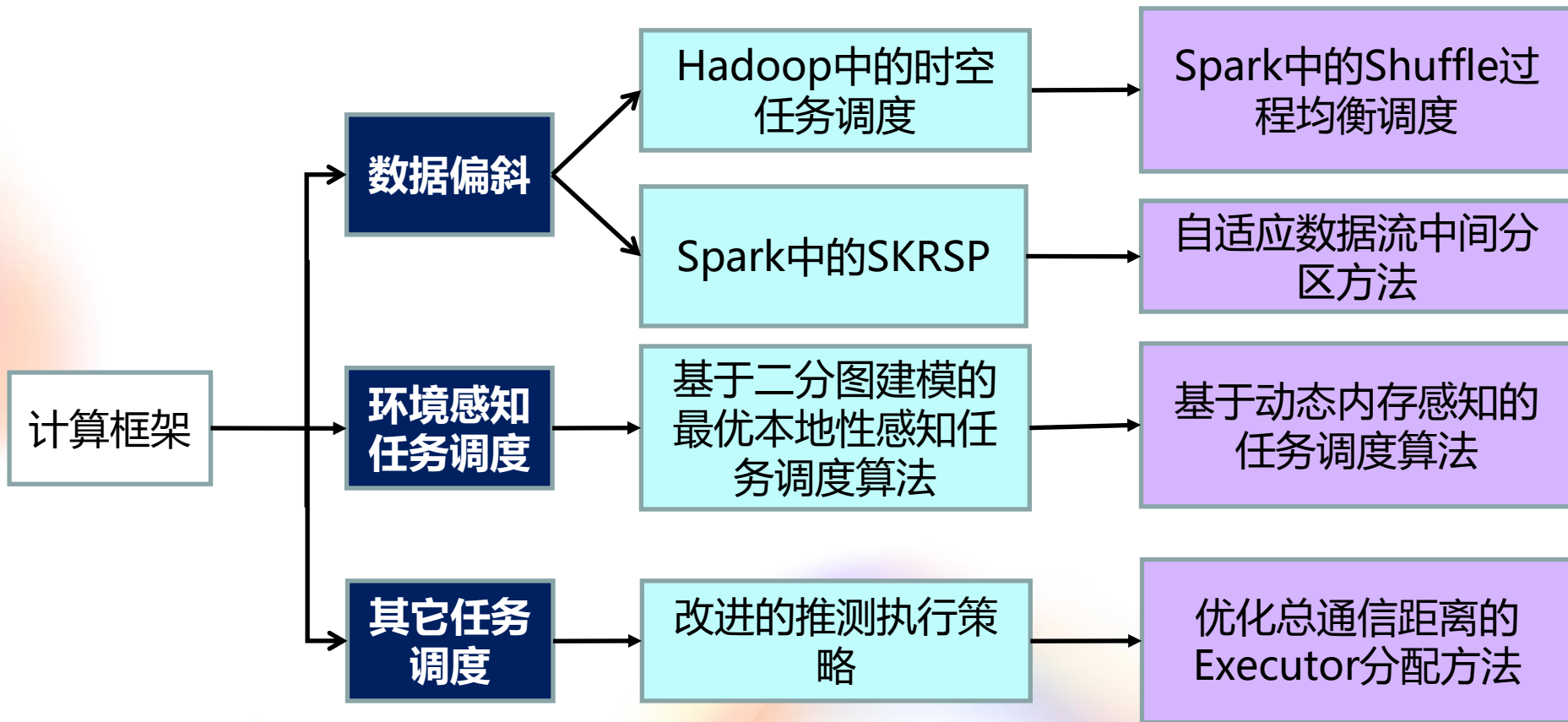
Serverless \neq FaaS

Serverless = FaaS + BaaS

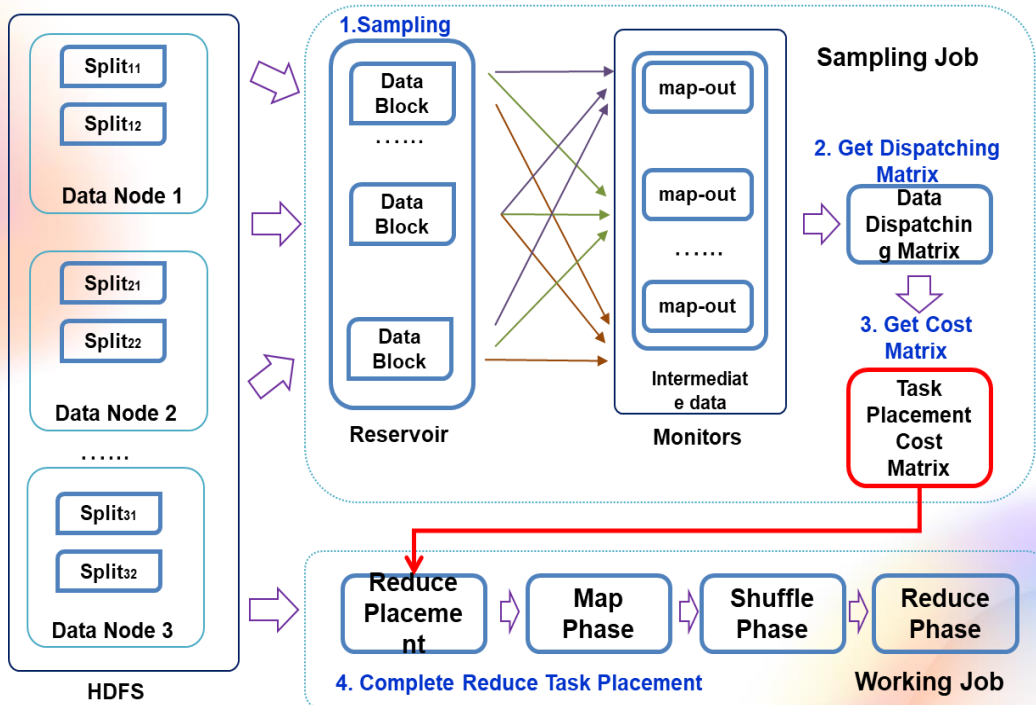
BaaS = DC + DaaS

构建数据和应用感知的分布式
计算和数据处理环境

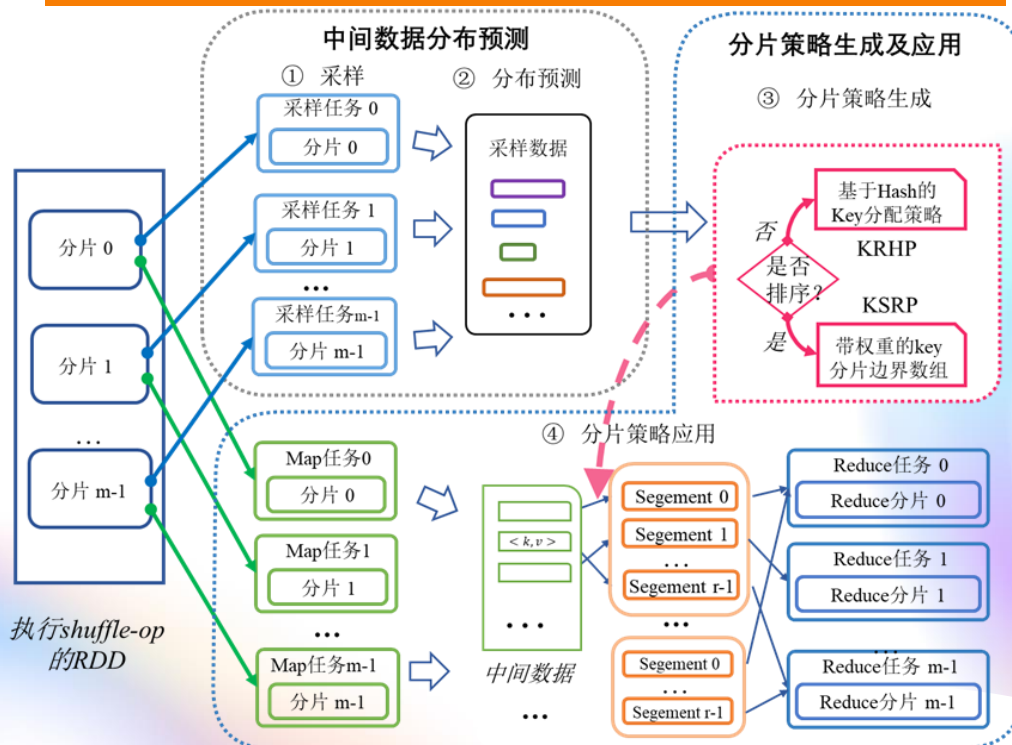
Serverless: 数据和应用感知的任务调度与资源分配



提出面向Hadoop架构的内部通信量优化的Shuffle过程任务放置策略



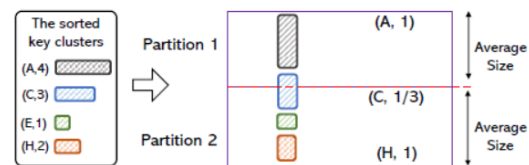
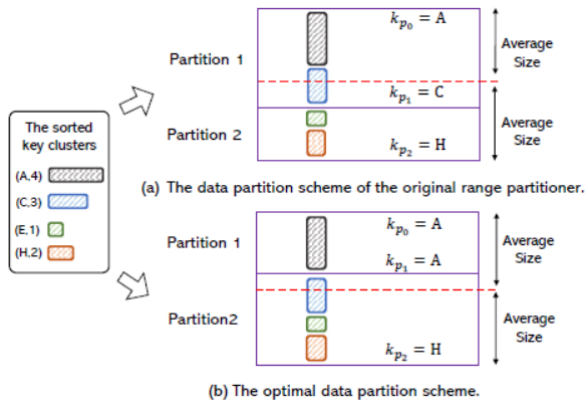
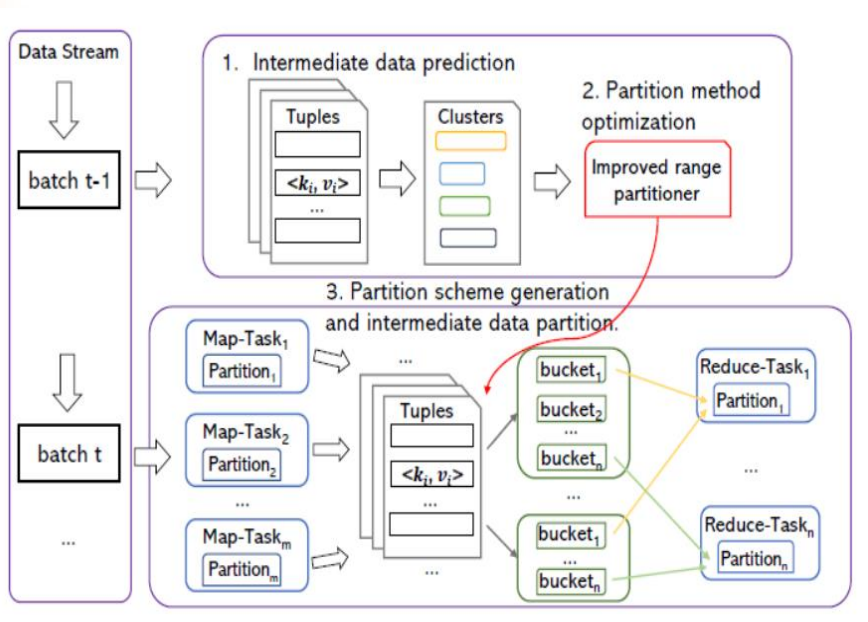
提出抗数据偏斜的Spark中间数据分片机制



Serverless: 数据和应用感知的任务调度与资源分配

Spark-Streaming中抗数据偏斜方法

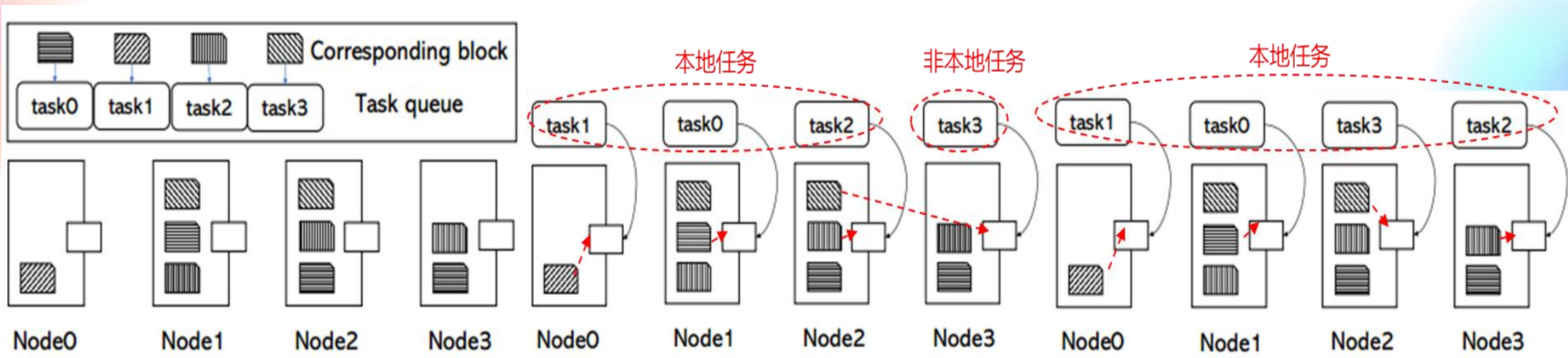
- ✓ 利用前一个批次已经处理产生的中间数据预测下一批次作业的中间数据key分布。然后针对中间数据分配不均，在范围分区方案基础上实现Shuffle操作前后的分区平衡



Serverless: 数据和应用感知的任务调度与资源分配

基于Spark平台二分图建模最优本地性感知任务调度算法

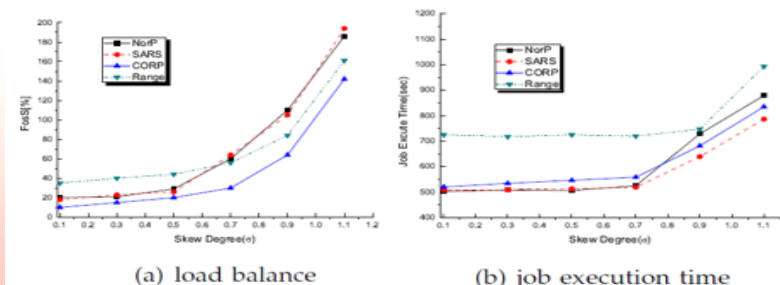
- ✓ Spark原生任务调度器采用贪婪的策略分配任务到executor上，会造成数据本地性局部最优
- ✓ 基于二分图建模提出了一种最优的数据本地性感知任务调度算法，最小化每阶段任务的总通信代价，以此减少map和reduce任务的跨节点/机架数据传输量，降低通信延迟



性能评估

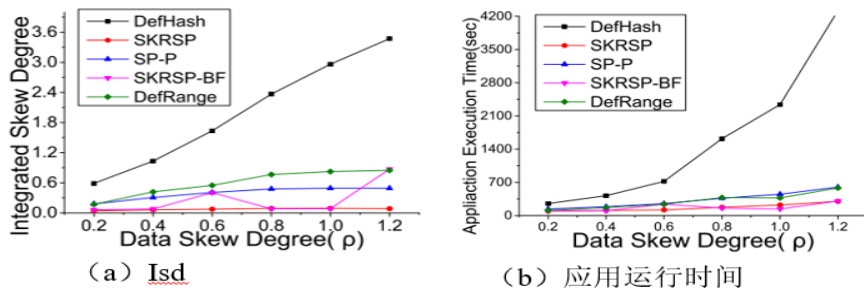


Hadoop Shuffle过程任务放置:



性能 vs. 数据偏斜

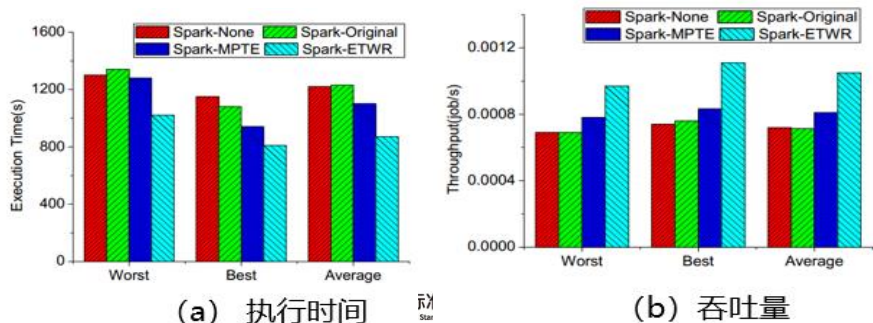
抗倾斜的Spark中间数据分片机制(SKRSP):



Join数据分片实验结果

Spark-Streaming抗数据偏斜框架(Spark-ETWR)

Spark本地性感知任务调度算法(OptLTS)

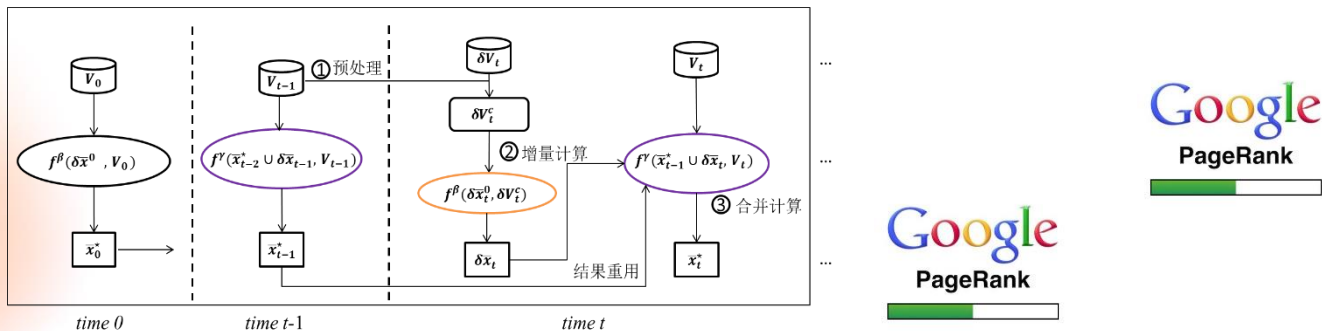


Sort下不同策略的性能比较

集群宏基准下不同算法的性能比较

构建面向机器学习的分布式并行计算体系结构

问题提出：对于迭代图算法，给定的迭代结果和，如何有效地快速地计算得到的迭代结果



(1) 预处理步。此步骤为增量计算做好准备，以便在图更新后识别已更改顶点。

(2) 增量计算。此步骤仅计算已更改顶点的状态。表示已更改顶点的计算结果。是通过在已更改顶点上以初始状态迭代轮所获得的。

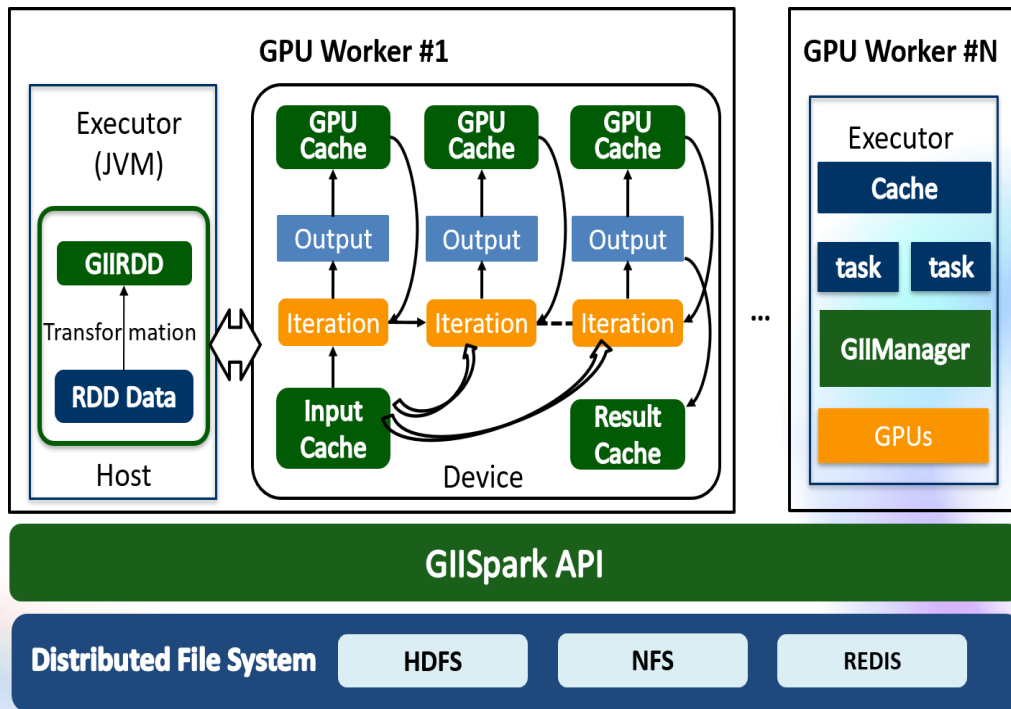
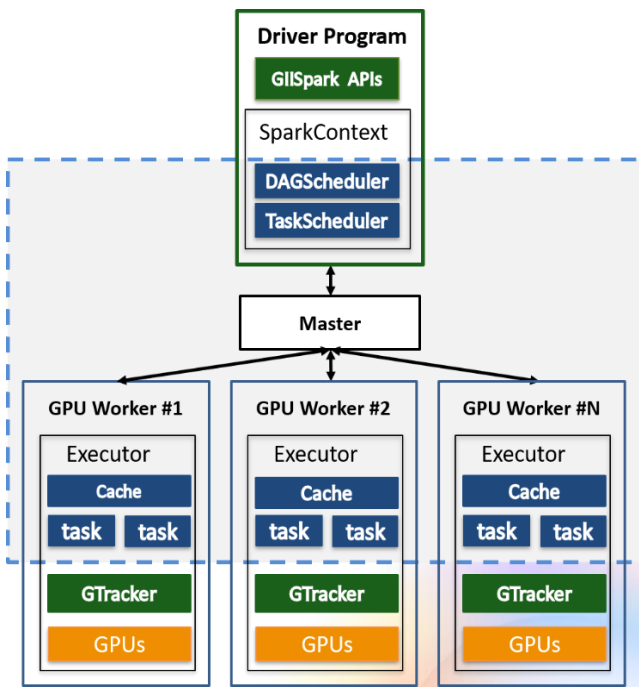
(3) 合并计算。此步骤是在整个顶点集上执行的，顶点初始状态为，其中为上一个图的迭代结果。最终的收敛顶点状态集是在轮迭代后所获得的。

云原生：构建高效的分布式机器学习环境



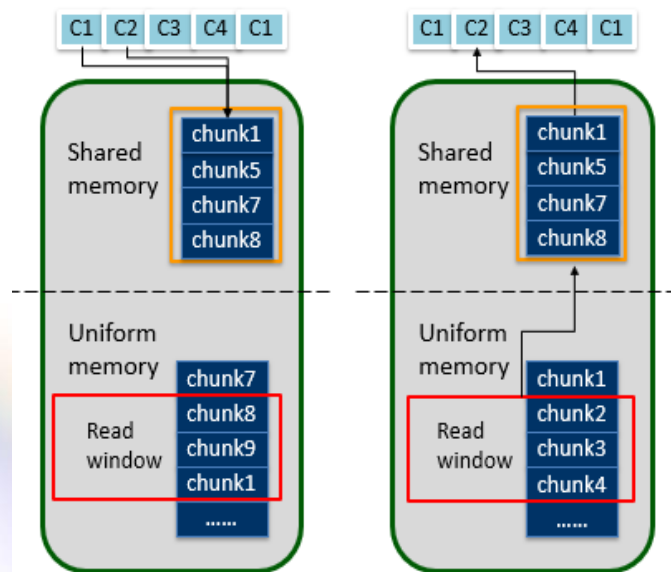
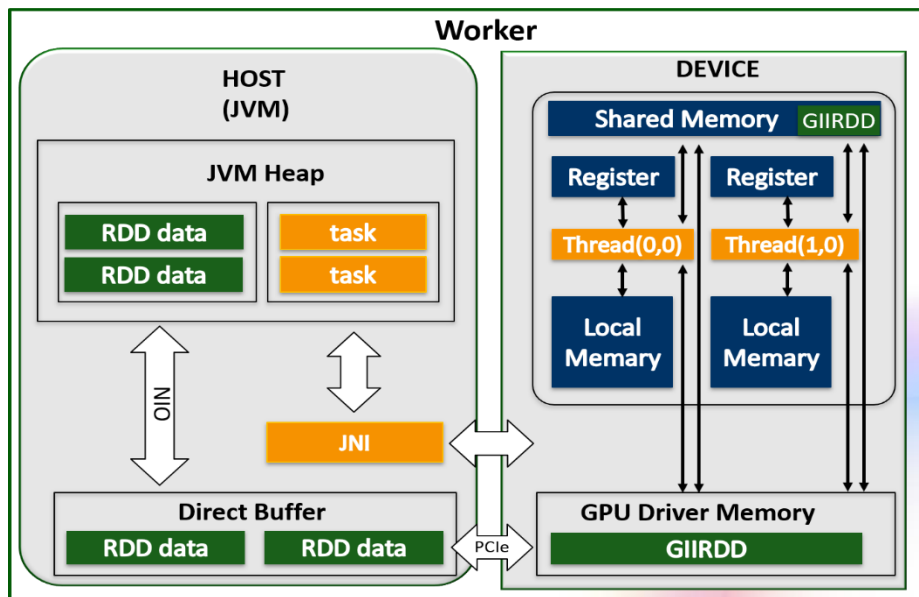
构建面向机器学习的分布式并行计算体系结构

- 扩展分布式并行框架 Spark
- 实现GPU资源的自动申请和释放，迭代计算任务的管理
- 扩展Spark原生的抽象数据类型RDD，符合GPU基于块的合并访问。
- 提供GIIRDD的缓存策略，实现迭代计算的数据复用，加速计算收敛



JVM-GPU通信缓存与滑动窗口机制

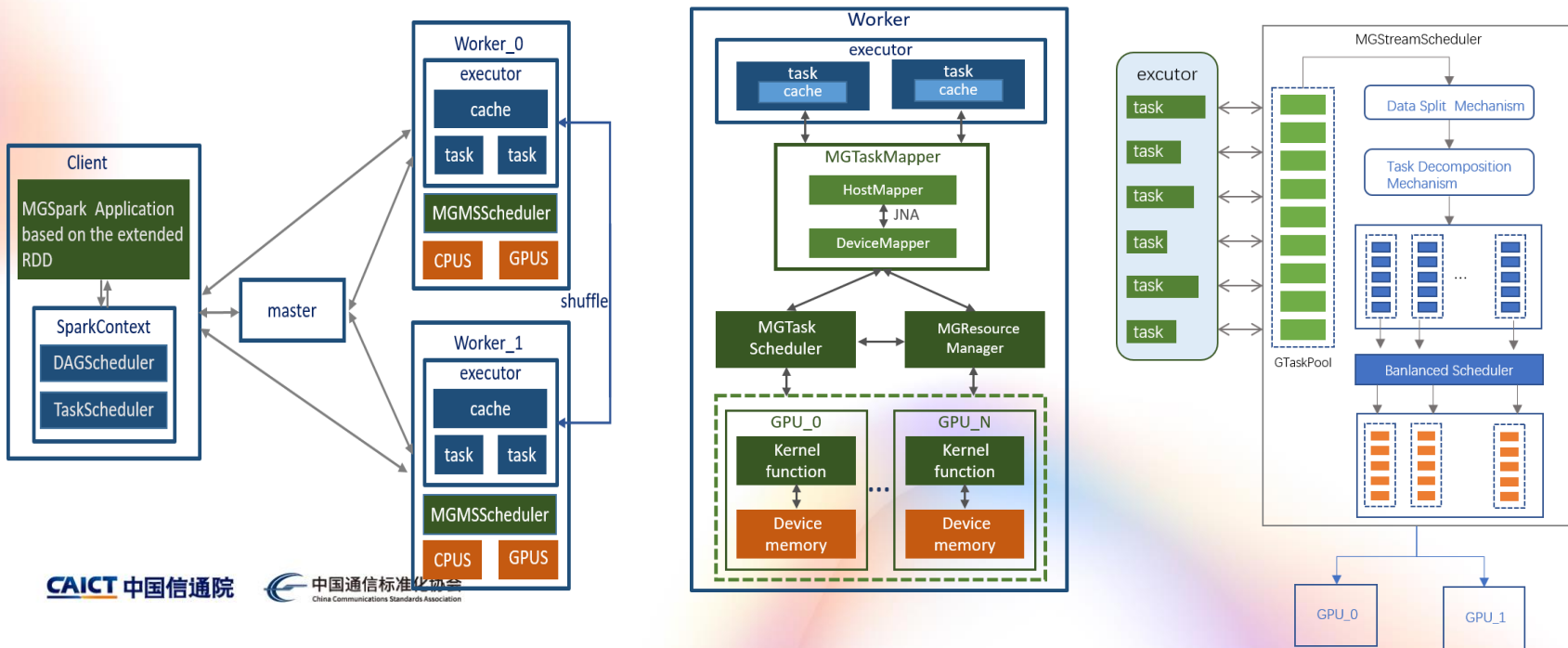
- 通过PCIe传输到设备端内存，并通过Gspark提供的接口，封装任务逻辑，由JNI将指令发送到设备端kernel执行；
- 使用一级缓存 (Shared Memory)开辟中间结果缓存区，加速迭代计算中共享数据的访问；使用基于block的滑动窗口机制并利用全局内存的合并访存机制来实现数据置换。



云原生：构建高效的分布式机器学习环境



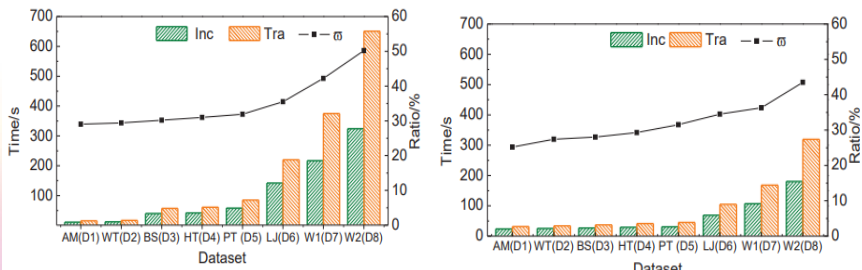
- AEML：用于分布式异构环境中的多GPU负载均衡的加速引擎，以插件形式整合到Spark框架中。其包含多个子组件，分别实现任务映射，GPU负载均衡调度和设备资源管理,可以有效地将GPU集成到分布式处理框架中，并在多个异构GPU之间实现负载均衡。
- MGMS：基于多个GPU和多个流的异构任务执行模型，有效地平衡多个GPU的工作量



性能评估



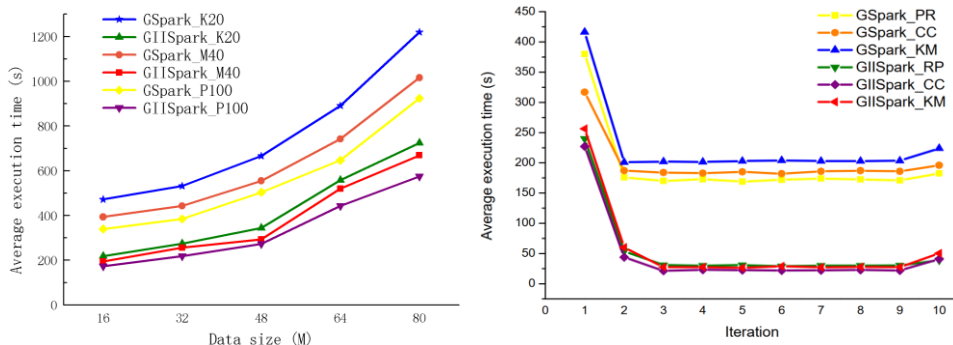
IncGraph-分布式增量图计算模型:



(a) PageRank

(b) SSSP

分布式异构增量迭代加速体系结构:

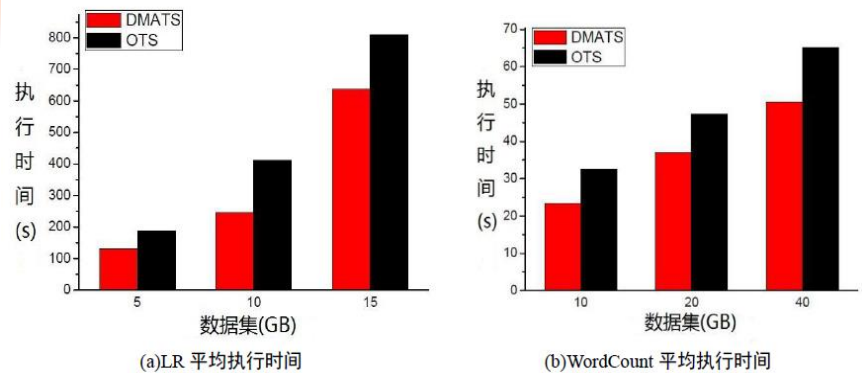


不同数据集执行时间 传统迭代 vs.增量迭代

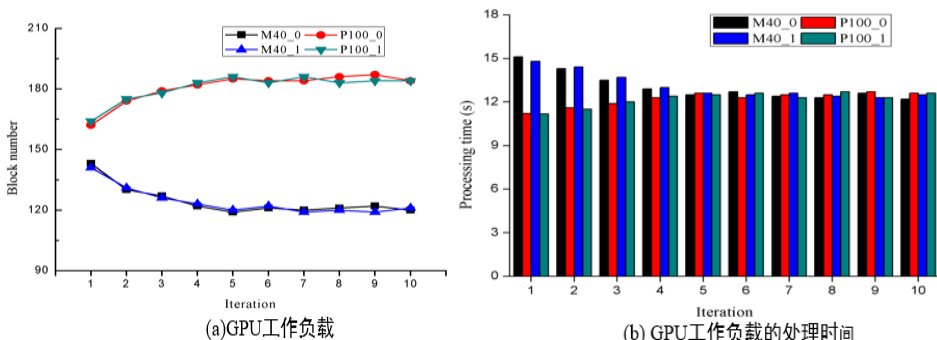
不同数据集大小和迭代次数的加速效果

基于Spark的动态内存感知调度算法:

AEML-分布式异构多GPU负载均衡加速引擎:

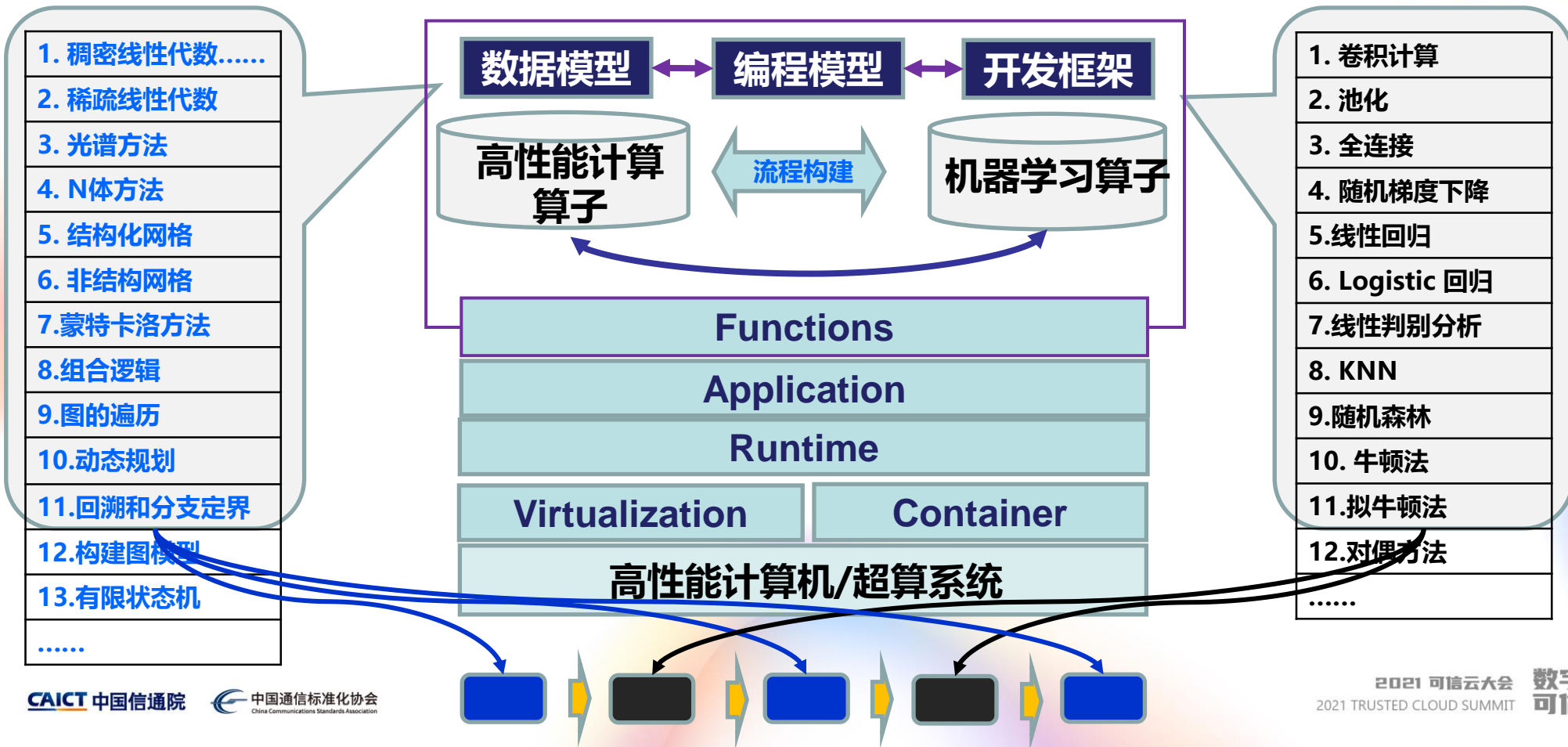


与原生spark任务调度策略对比 (执行时间)



不同GPU工作负载 (FSA) 下的任务执行时间
FSA(The feedback based streams adjustment scheme)

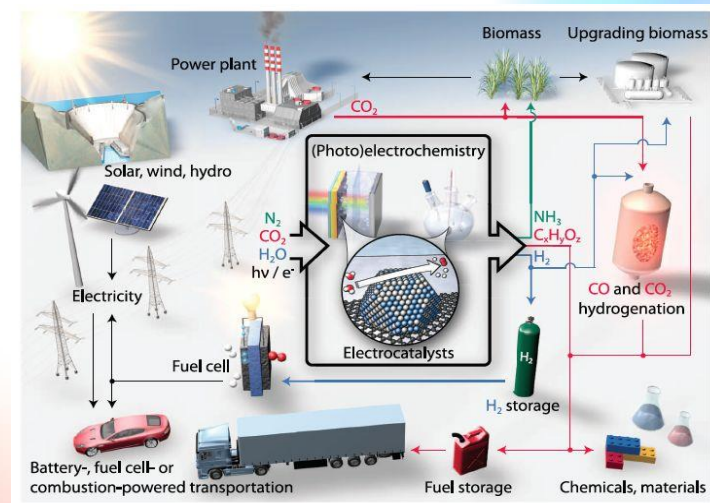
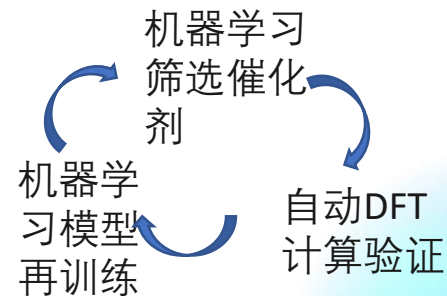
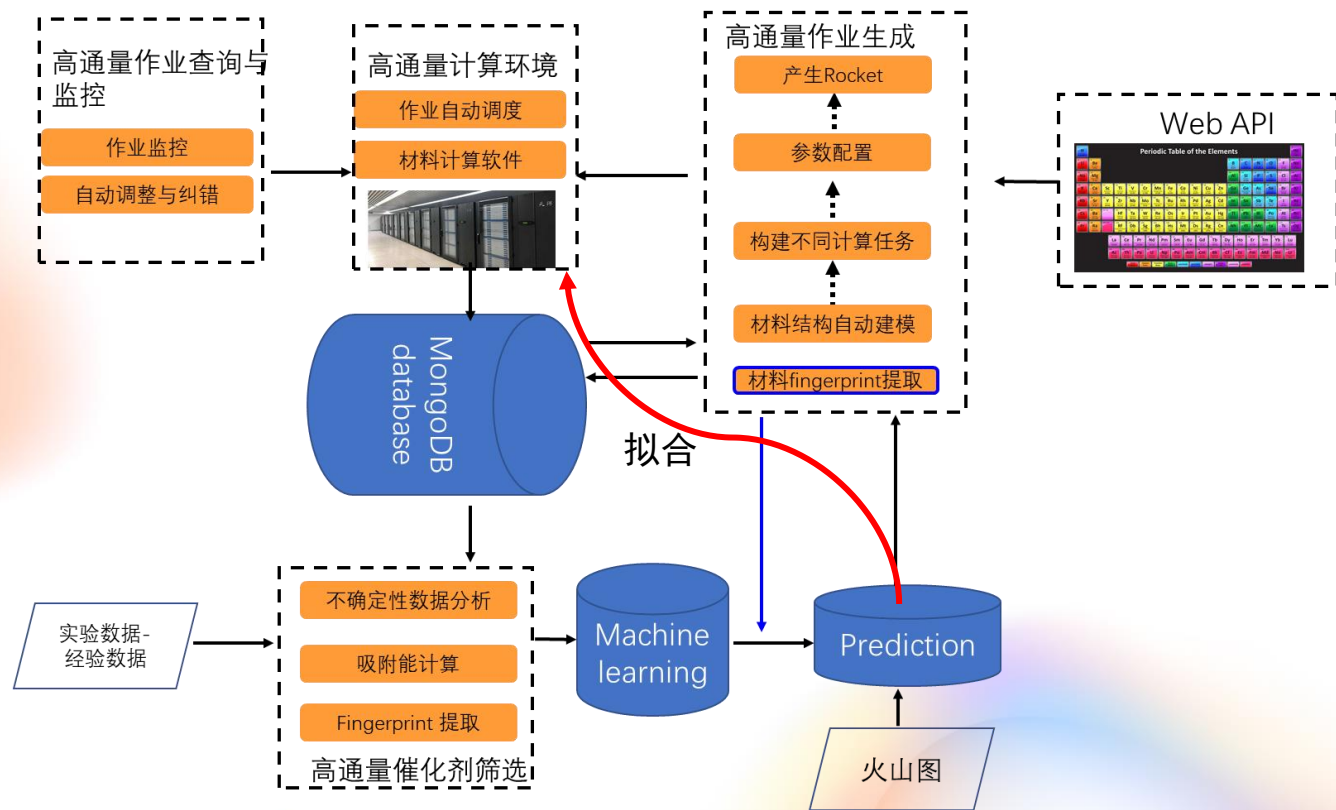
函数计算：融合高性能计算算子与机器学习算子的开放式框架



函数计算：融合高性能计算算子与机器学习算子的开放式框架



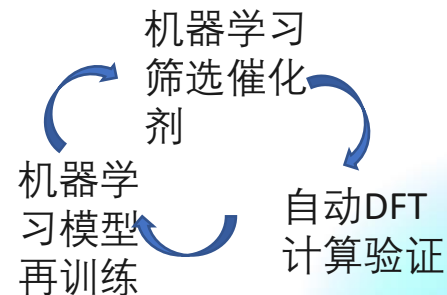
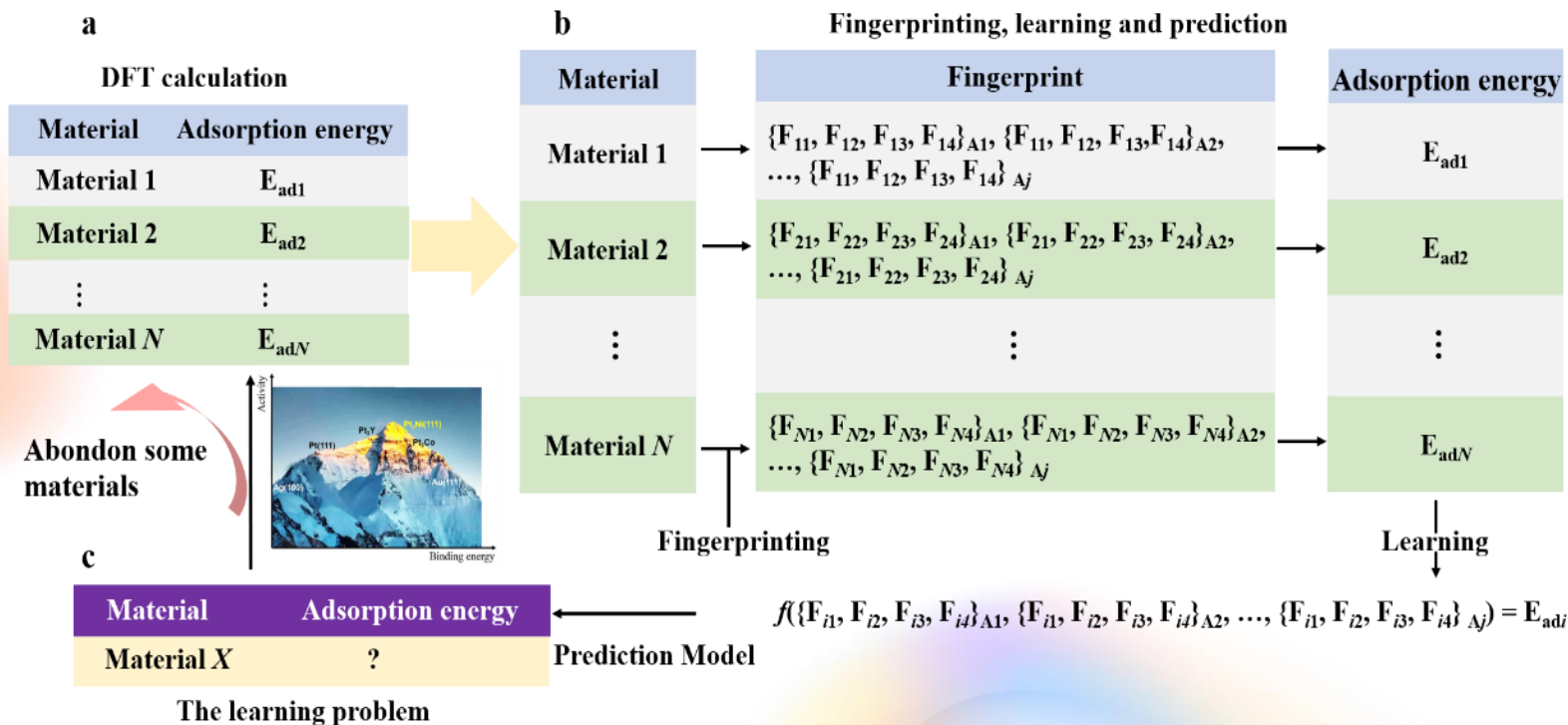
计算化学：机器学习+高性能计算助力高通量催化剂筛选



函数计算：融合高性能计算算子与机器学习算子的开放式框架



计算化学：机器学习+高性能计算助力高通量催化剂筛选

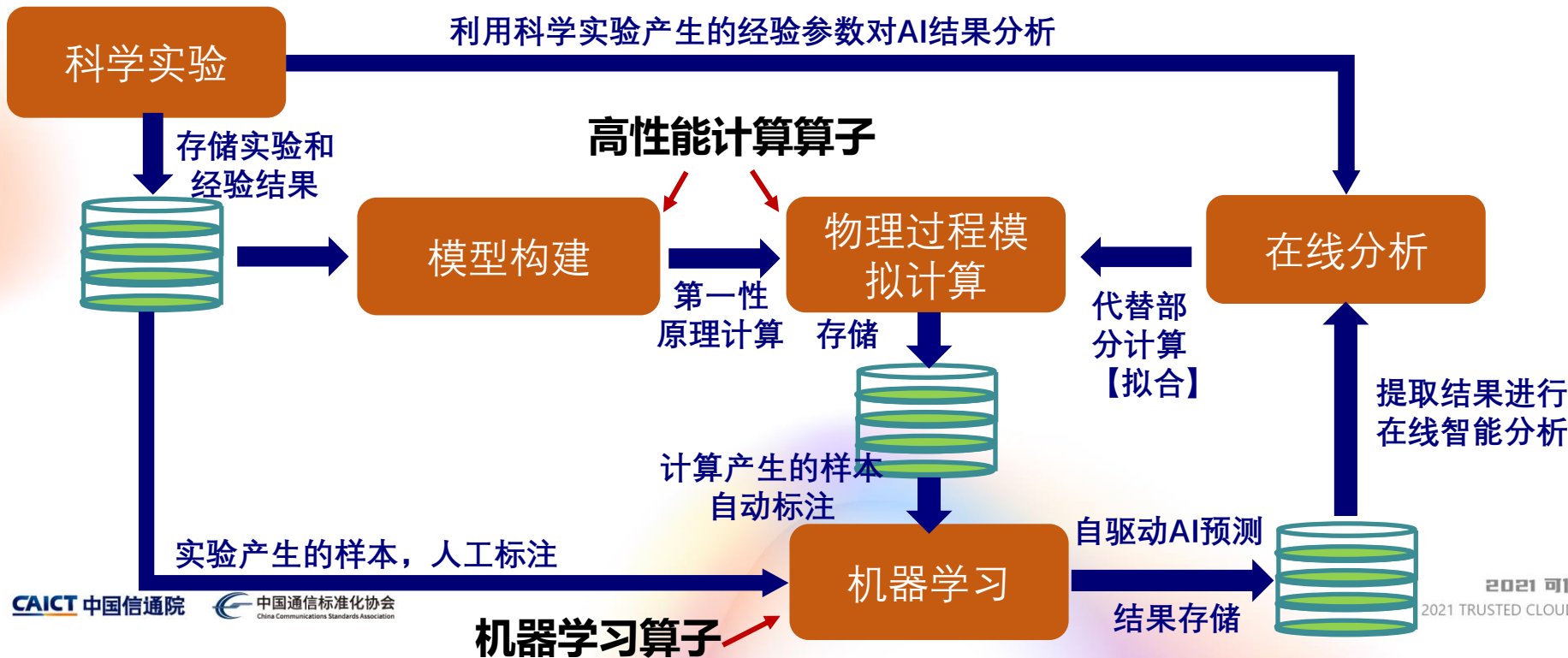


函数计算：融合高性能计算算子与机器学习算子的开放式框架



第五范式：双向互反馈的创新

- 1、机器学习过程加速数值计算，节省实际计算时间
- 2、数值计算为机器学习迭代提供可训练样本

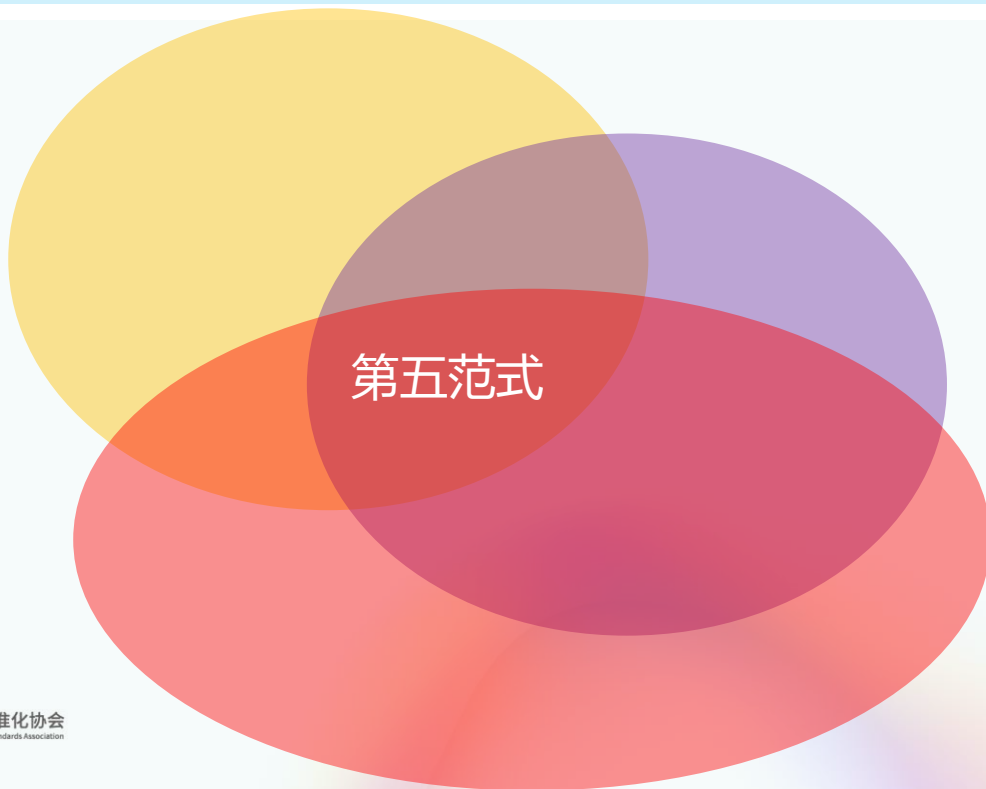


函数计算：融合高性能计算算子与机器学习算子的开放式框架



第五范式：双向互反馈的创新

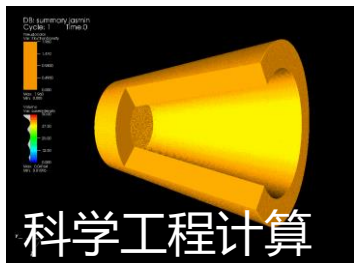
- 1、机器学习过程加速数值计算，节省实际计算时间
- 2、数值计算为机器学习迭代提供可训练样本



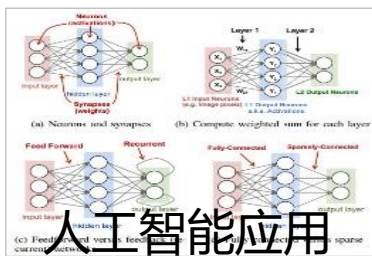
- 未知世界:A
- 科学实验:B
- 理论计算:C
- 机器学习:D

- 1、背景和挑战
- 2、高性能计算服务化关键技术
- 3、超算云服务平台及整体架构

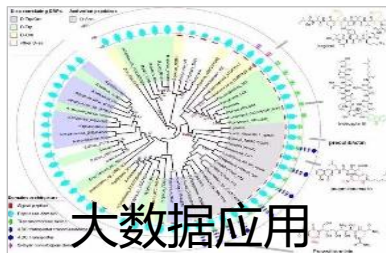
高效能数据并行处理与智能分析系统



科学与工程计算



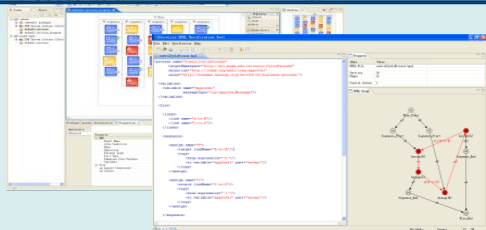
人工智能应用



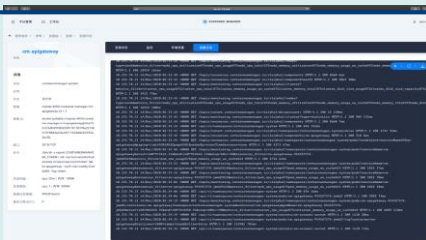
大数据应用



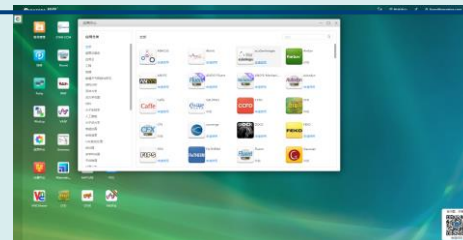
区块链应用



计算服务建模、验证与调试



虚拟机、容器与微服务管理



计费管理与APP工具

高效能数据并行处理与智能分析系统

超级计算机集群及其算力网

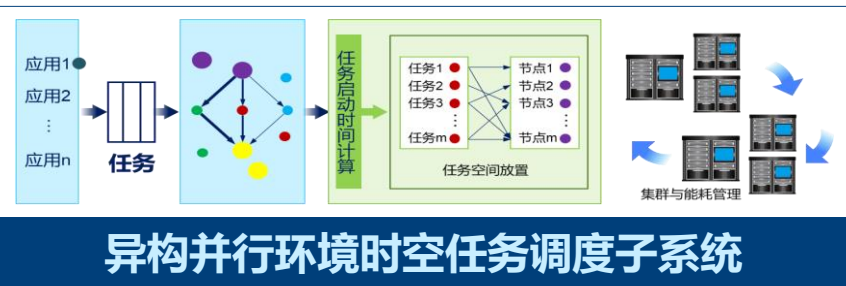
高效能数据并行处理与智能分析系统



分布式体系结构优化技术

异构计算融合关键技术

高效能数据并行处理与智能分析系统

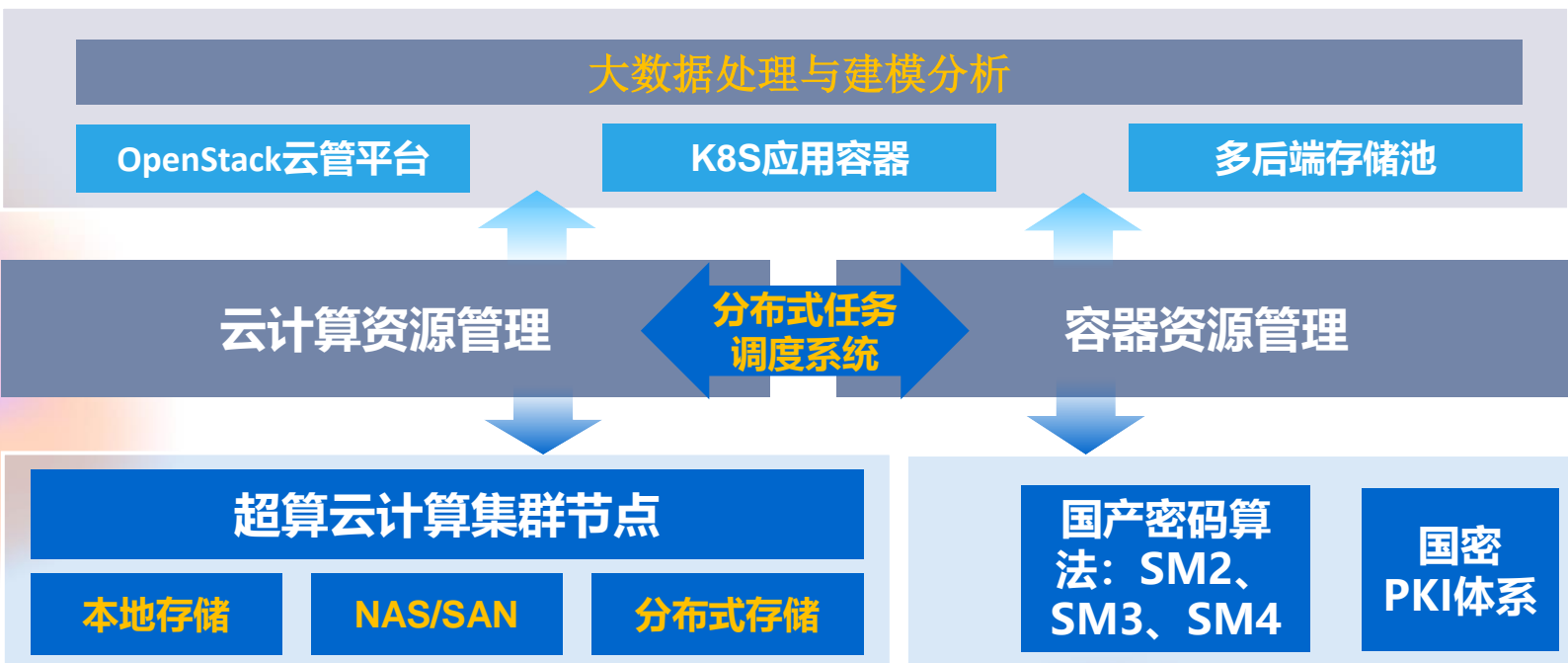


云服务器ECS 实例列表

实例名称	镜像	IP地址	状态	配置	创建时间	可用域
云盘	d5689e3-66be-4678-9ccb-5f4d1da...	centos-8-aar...	192.168.1.7 (内)	运行中	vCPU: 1核 内存: 1GB	2020-10-29 09:40 飞腾鲲鹏
云盘快照	5e4d079a-e494-447a-ba0a-100c3ad...	centos-8-aar...	192.168.122.3 (内)	运行中	vCPU: 1核 内存: 1GB	2020-10-12 12:22 飞腾鲲鹏
云盘备份	607d5065-62e6-4827-bd7e-6979fce...	centos-7-x8...	192.168.1.19 (内)	运行中	vCPU: 1核 内存: 1GB	2020-10-10 10:18 Intel-金蝶-6261
自动快照策略	d016d7c0-923-4944-8518-bd996f7...	debian-10-a...	192.168.1.3 (内)	运行中	vCPU: 1核 内存: 1GB	2020-09-25 17:57 飞腾鲲鹏
网络和安全	d9a0826-71d6-4036-5cda-1ad5adf...	debian-10-a...	192.168.1.20 (内)	运行中	vCPU: 1核 内存: 1GB	2020-09-25 17:23 飞腾鲲鹏
安全组	a1f639d8-bd9c-479b-8441-075b128...	windows-2...	192.168.1.8 (内)	运行中	vCPU: 1核 内存: 1GB	2020-08-22 16:23 Intel-金蝶-6261
密钥对	5b626f15-c163-4655-502c-475288...	ubuntu-16.0...	192.168.1.4 (内)	运行中	vCPU: 1核 内存: 1GB	2020-09-25 17:15 飞腾鲲鹏

飞腾、鲲鹏处理器 X86处理器

关键子系统1：高性能计算资源池及管理子系统



- ✓ 实现对天河超算及国产服务器等计算资源的弹性管理与按需分配
- ✓ 自主虚拟化管理平台、分布式云存储、云审计运维，以及应用容器
- ✓ 兼容飞腾、鲲鹏、海光等国产ARM芯片和X86芯片

异构资源池环境



国产CPU/GPU异构超级计算机



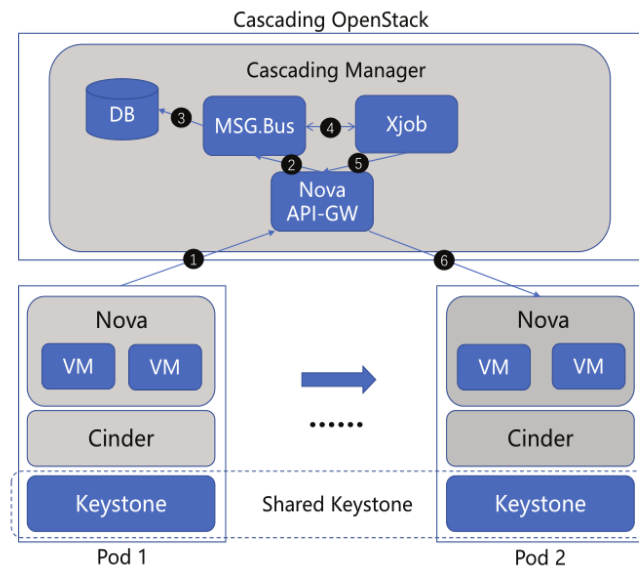
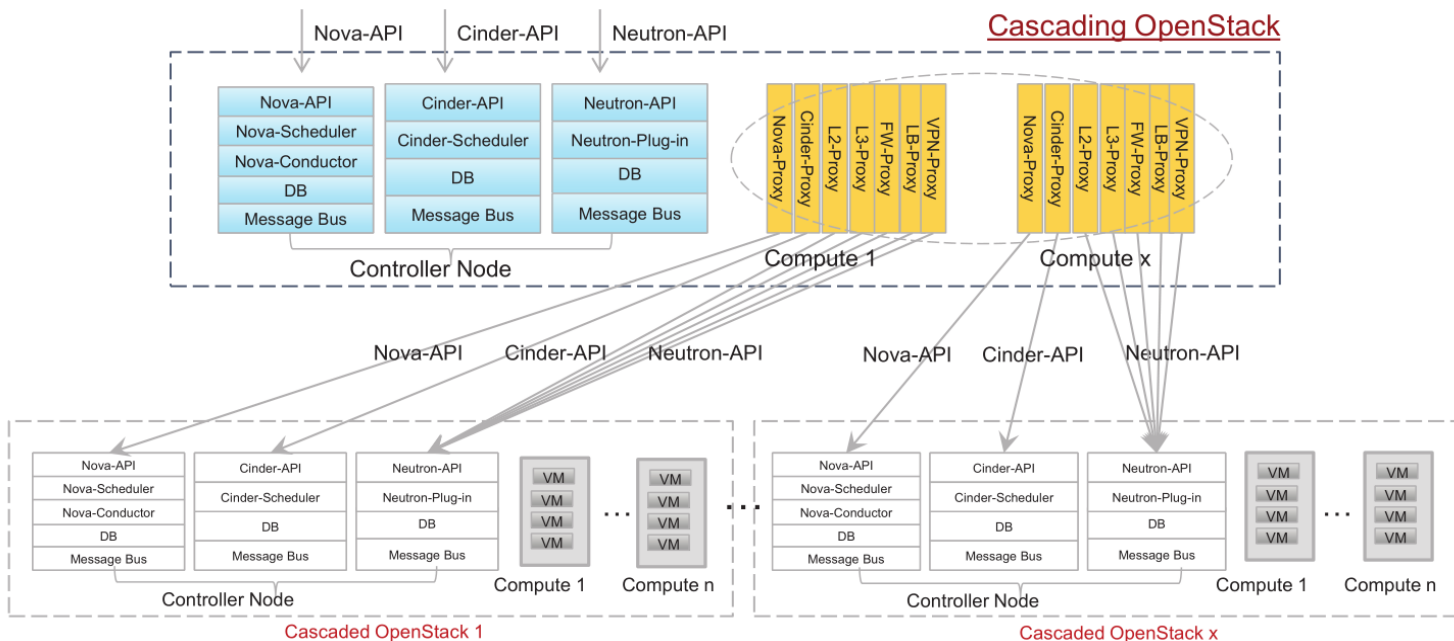
ARM架构：飞腾、鲲鹏 + 麒麟操作系统

兼容国产X86 CPU

关键子系统2：跨域资源管理与多云级联



- 优化对多云环境下各云实例上资源的调度途径，研发OpenStack原生组件网关接口和相应的资源管理策略
- 针对多云架构中顶层资源信息特征，提出VM镜像跨域冷迁移策略及内存数据跨域传输机制



大数据并行处理与建模分析子系统

智能应用

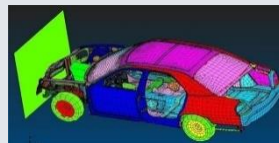
机器学习算法

并行处理环境

云计算资源池



设备状态监控



工业仿真



生产管理



智能制造

基于云的大数据并行处理与挖掘平台

算法库

高性能计算算子库

深度学习算子库

基于GPU集群的分布式深度学习并行处理框架

并行编程模型

面向大数据应用的多机多核并行编程模型

面向异构内存体系的数据存储特性描述

数据局部性编程表达

任务调度策略

异构资源和应用感知的负载均衡调度框架和策略

数据局部性驱动的自适应任务调度

基于非易失混合内存的任务迁移机制

并行支撑技术

基于大内存的并行计算运行框架与机制

高效的数据传输与通信优化机制

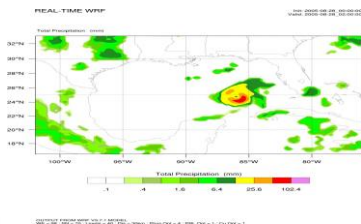
内存并行文件系统

云服务资源分配与资源池构建

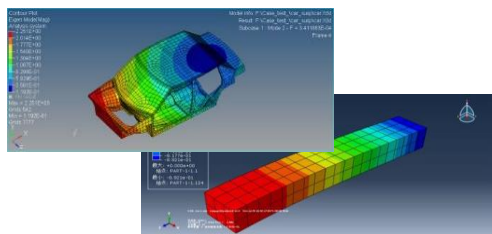
国家超算长沙中心：高效能数据并行处理与智能分析



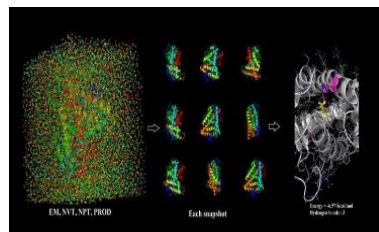
助力国家数字新基建战略，为超算中心的大数据与人工智能应用提供平台支撑



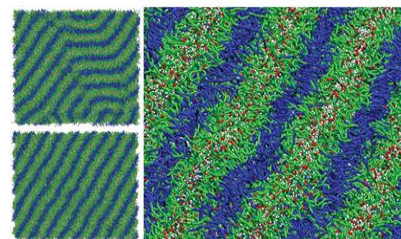
大气海洋环境模型 WRF



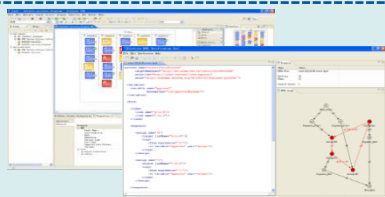
工业制造软件 ABAQUS



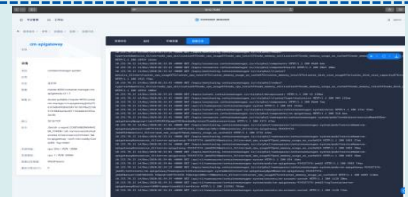
分子动力学Gromacs



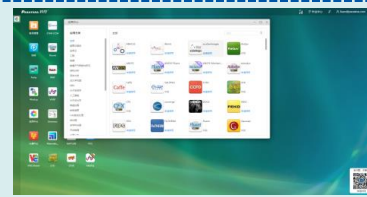
lammmps材料计算服务



服务建模、验证与调试系统



容器与微服务管理原型系统



计费管理系统与APP工具

高效能数据并行处理与智能分析系统

国家超算 广州中心



国家超算 无锡中心



国家超算 天津中心



数字裂变
可信发展

国家超算长沙中心：疫情大数据分析处理



支撑国家科学防治疫情，为湖南新冠疫情监控提供并行处理与分析平台支撑

对湖南省境内七千万人口的行为轨迹进行实时关联计算

基于MR信令数据，将个体行为轨迹精确到建筑物内

识别疑似和确诊病例轨迹1018例成功预警3起高风险外地人群输入事件



湖南卫视、科技日报、新华网等都进行了详尽报道



三大运营商信令数据



发现确诊患者后倒查活动轨迹碰撞计算
虚拟数据：红色—密集接触者
蓝色—健康

数字裂变更可信发展

THANKS!

2021
TRUSTED CLOUD
SUMMIT

