

服务不确定性高

系统吞吐效率低

资源抽象能力差

科学问题



目标3：异构计算系统高吞吐

改善缓存成本

提高资源分配效率

高吞吐、低延时的面向AI的多样化算力需求场景

服务质量优化

研究内容3：面向用户服务的长冷启动时延优化技术

多级缓存问题
抽象及建模

冷启动率约束下的
L1热池缓存决策

服务性能驱动的
L2冷池缓存决策

成本效益提升

研究内容2：基于中心化控制的低成本函数实例缓存技术

控制平面分离的
中心化缓存架构

效益最优的热点
函数可区分缓存方法

热点感知的实例
缓存负载均衡调度

算力扩展增强

研究内容1：高吞吐Serverless异构计算系统

CPU/GPU异构
计算资源统一抽象

QoS保障下的
端到端推理时延管理

吞吐最优的异构
资源协同分配方法

目标1：异构算力高扩展性

优化启动性能

优化执行性能

目标2 异构计算系统高性能

研究内容2：基于中心化控制的低成本函数缓存技术

研究内容1：高吞吐Serverless异构计算系统

