



# DTCC

## 数 / 造 / 未 / 来

### 第十二届中国数据库技术大会

DATABASE TECHNOLOGY CONFERENCE CHINA 2021

2021年10月18日 - 20日 | 北京国际会议中心



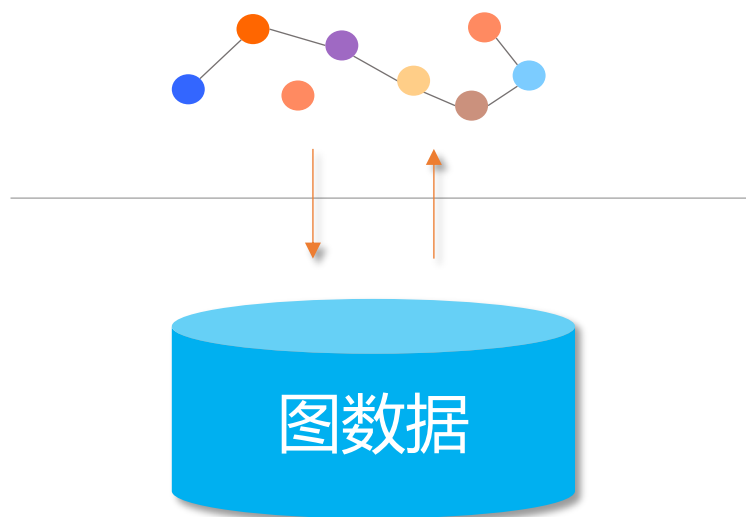
# 目录

1. 图计算是什么，能解决什么业务问题
2. 图计算与传统计算的区别，为什么需要图计算
3. 业界现状如何，大规模图计算所面临的挑战是什么
4. 如何应对挑战 -- 自研图计算
5. 图计算的架构与优势
6. 图算法开发与应用简介



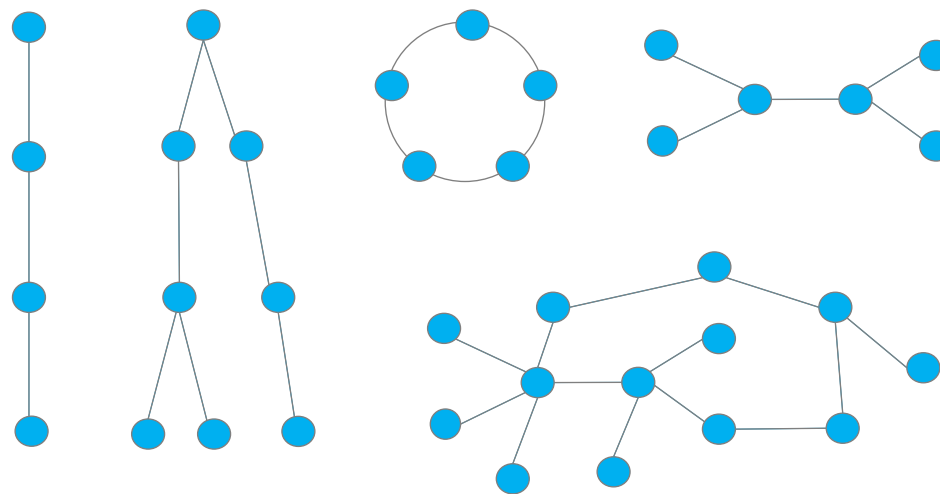
## 图数据库是什么？

- 存储、查询/分析图数据
- 灵活、高效处理图数据



## 图数据是什么样的数据？

- 由实体和关系组成的线、树、
- 环、叉、网等各种形状的数据



## HugeGraph是什么？

百度自研图数据库

国内首家开源图数据库

HugeGraph

Apache 2 License

Apache Tinkerpop Enabled

大规模分布式图存储与图计算



## HugeGraph应用分类

OLTP，即时查询，以图遍历为主



图查询

图计算

HugeGraph



## HugeGraph应用分类

OLAP，离线分析，全图多轮迭代计算



图查询

图计算

HugeGraph

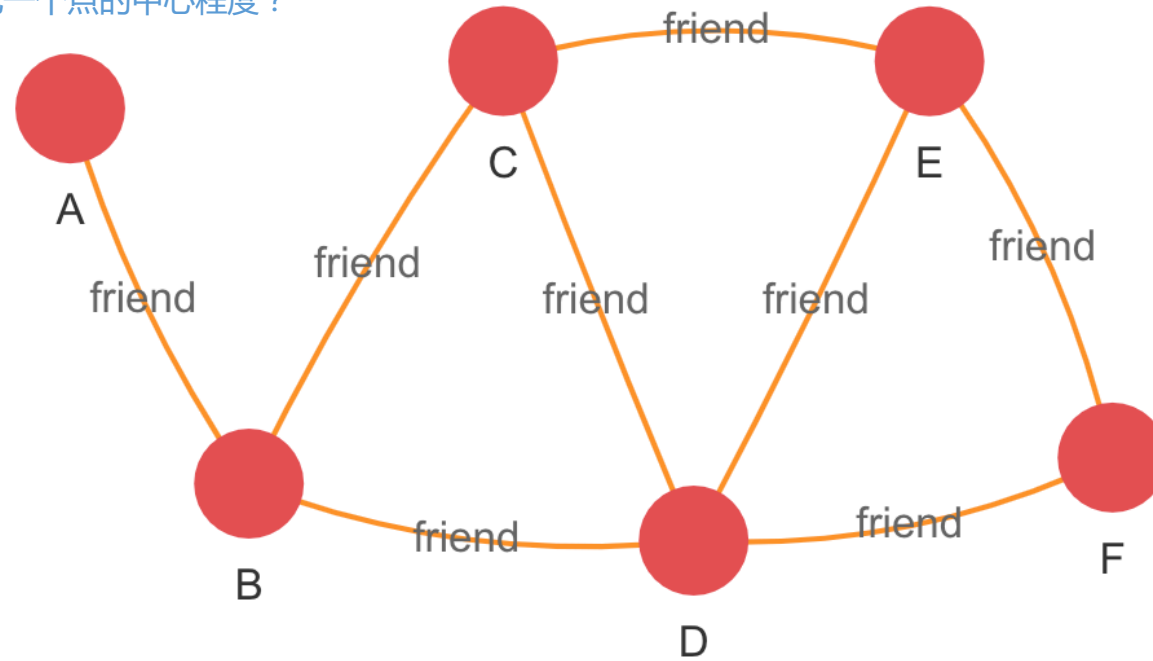




## 图计算是什么？

- 分析与抽象出复杂关联关系数据的特征，如实体/子图/路径等维度特征

如何量化一个点的中心程度？

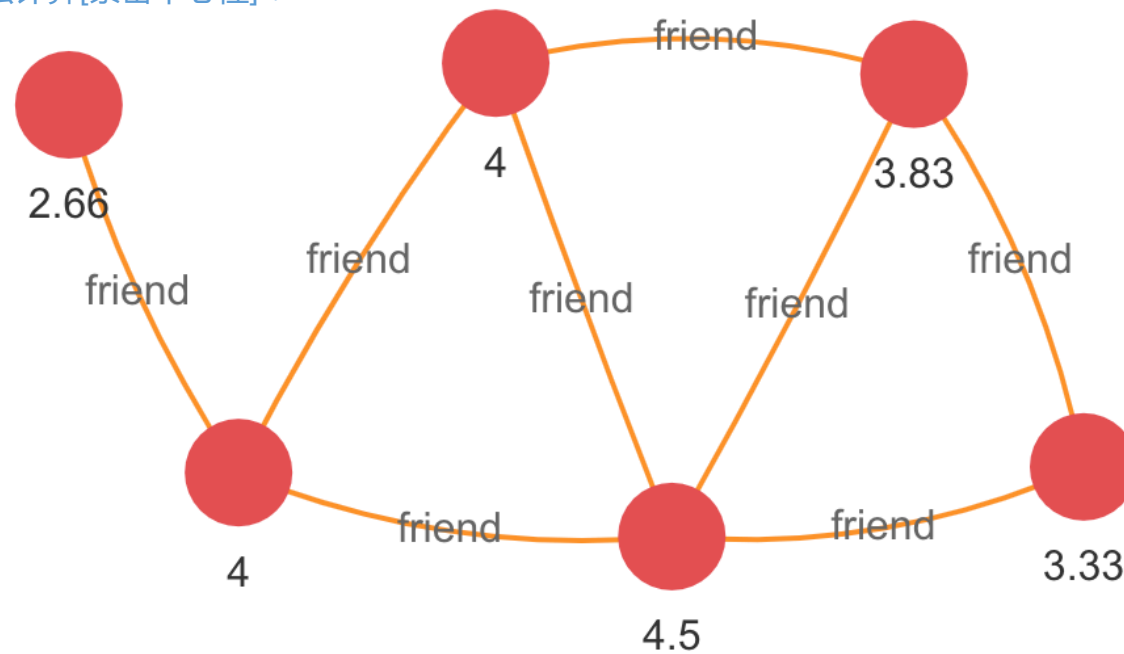




## 图计算是什么？

- 分析与抽象出复杂关联关系数据的特征，如实体/子图/路径等维度特征

通过图算法计算[紧密中心性]：







## 图计算是什么？

- 分析与抽象出复杂关联关系数据的特征，如实体/子图/路径等维度特征

## 图计算能解决什么业务问题？

- 金融风控、网络安全、情报关系、智能营销、智能推荐、智能运维等

## 图算法应用场景：

- 环路检测：循环担保
- 中介中心性：反洗钱
- 社区发现：团伙反欺诈



# 图计算与传统计算的区别



DTCC 2021

第十二届中国数据库技术大会  
DATABASE TECHNOLOGY CONFERENCE CHINA 2021

处理的数据维度多、关联关系错综复杂

非简单的单轮计算，而是具有依赖关系的多轮迭代

数据分布，逻辑上关联紧密的图无法简单按照单维度被划分到一起

分区不均问题，受现实世界中幂律分布而被放大，导致计算任务负载不均衡

多轮迭代决定上下游任务无法割裂，需要相互配合同步

每迭代中间结果量巨大，图算法复杂度 $O(n * (\text{degree}^{\text{depth}}))$ 随深度指数增长

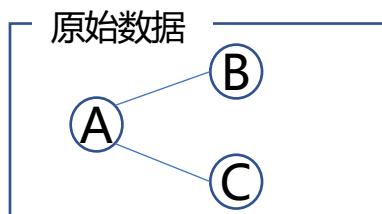
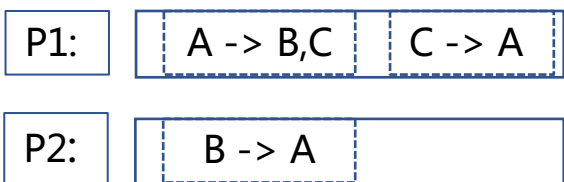


数，造，未，来



## 图的切分

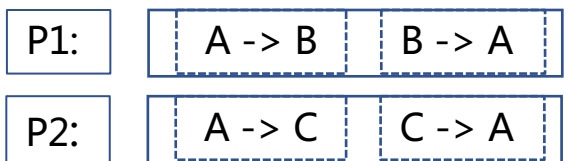
边切法 (顶点和边存储在一起)



采用的框架  
Pregel、Giraph

优点: 计算方便      缺点: 数据不平衡

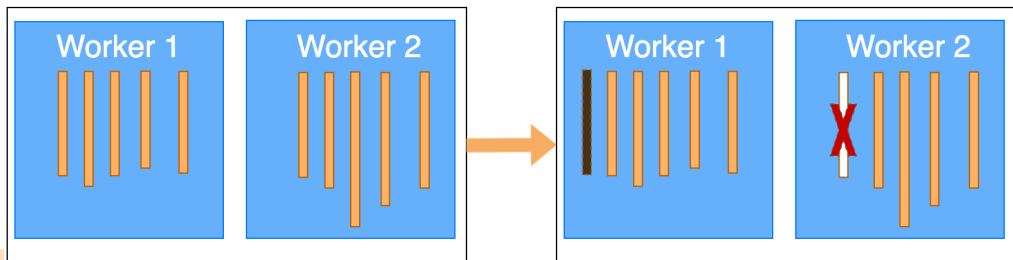
点切法 (顶点的边存在多个分区)



采用的框架  
PowerGraph、GraphX、Gemini

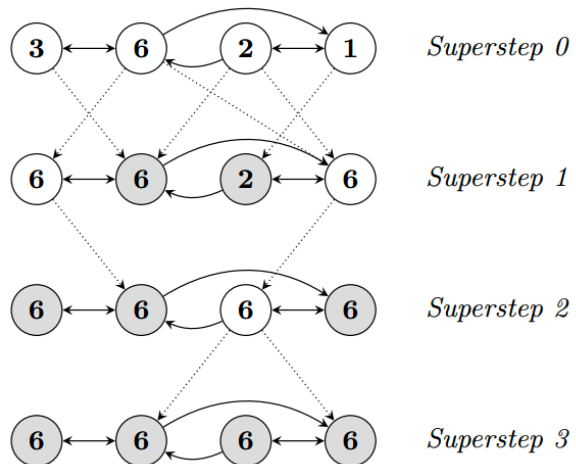
优点: 边分布均衡      缺点: 镜像顶点值

HugeGraph: 边切法 + Partition调整

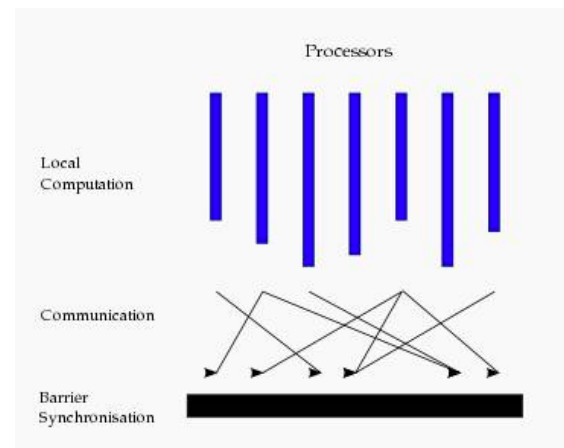


## 计算方式: 迭代计算和发送消息

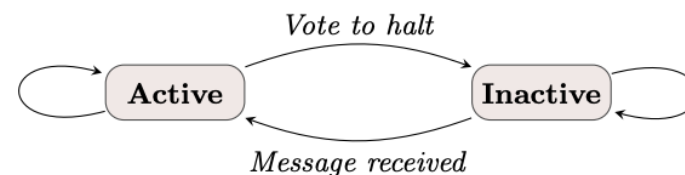
迭代计算执行流程



超步执行流程



顶点的状态转换



现有图计算框架

- 百亿以上图难以计算
- 面向开发者

HugeGraph图计算

- 大规模图稳定高效计算
- 面向普通用户



图随机访问多，邻接边幂律分布导致任务划分难以均衡负载

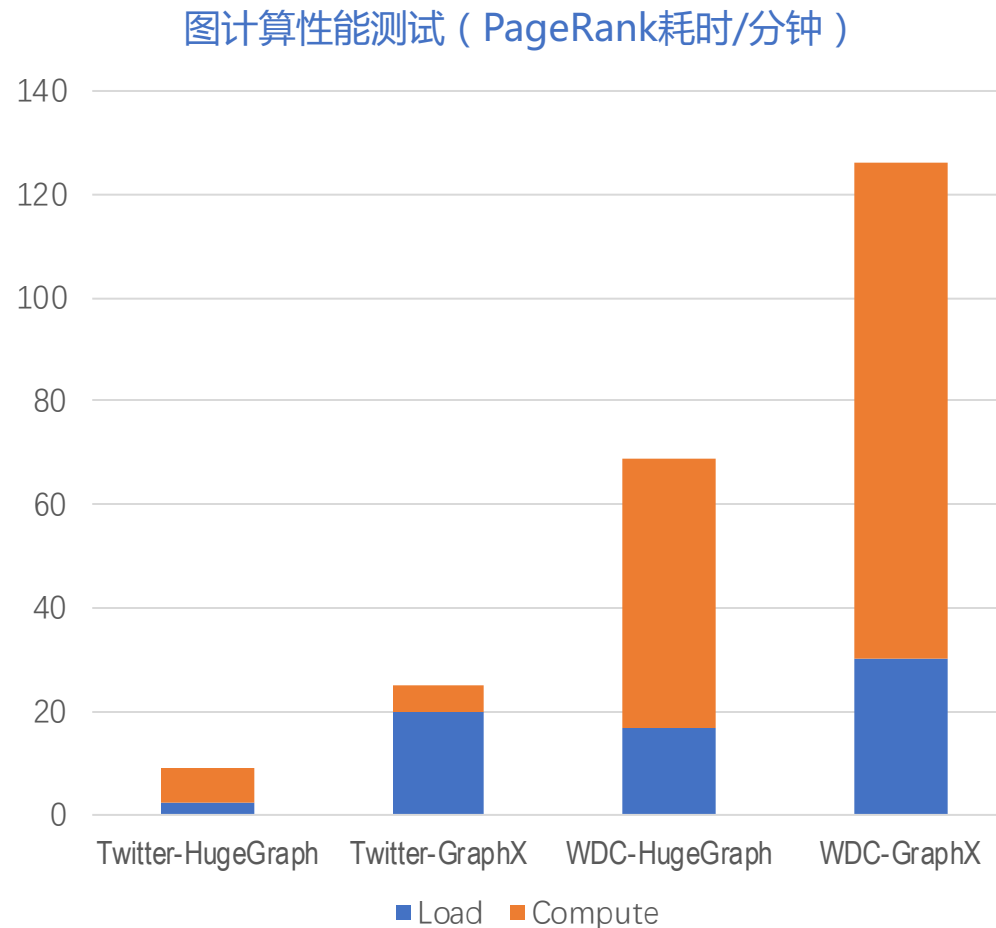
图算法复杂度高，消息与中间结果指数级增长

超过十亿规模图时，计算性能差、OOM严重

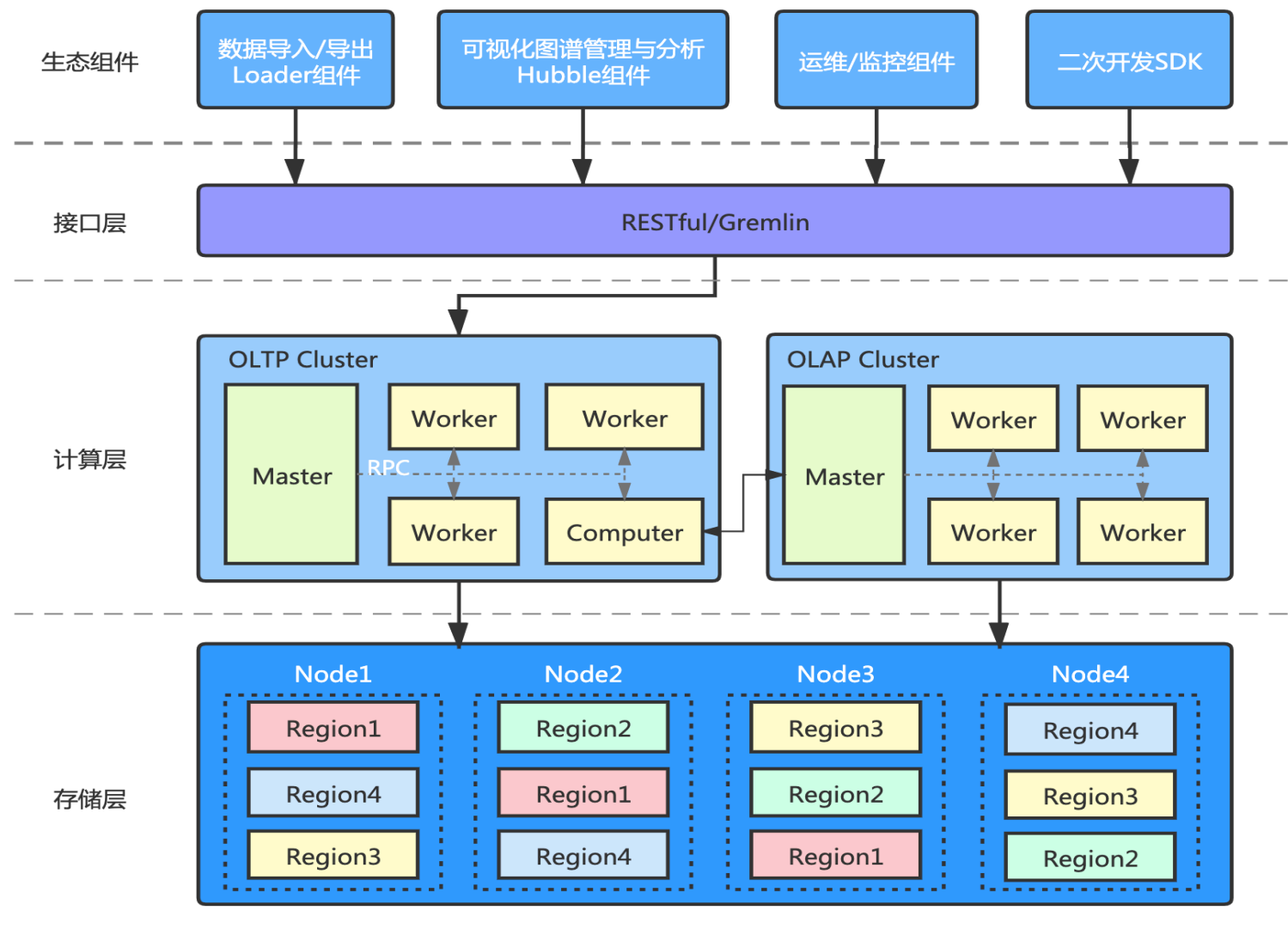


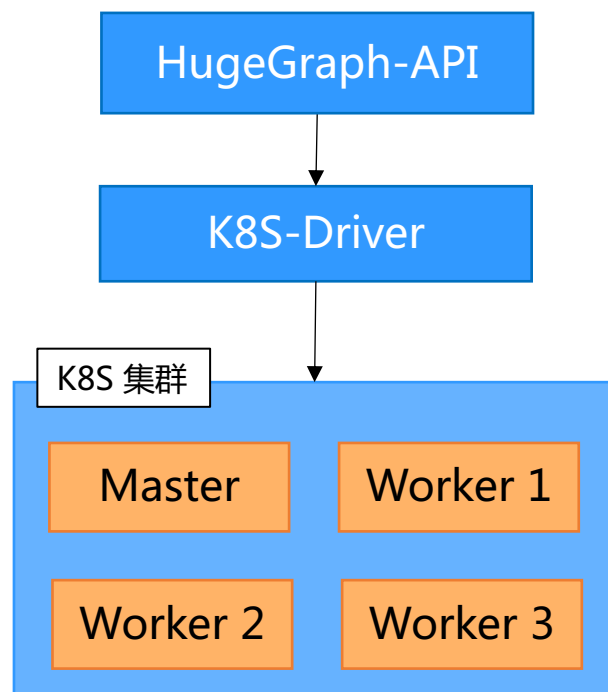
自研图计算效果：1) 千亿大规模图并行计算、2) 性能高效、3) 无需依赖内存容量、4) 高稳定性运行、5) 资源消耗低。

PageRank性能对比	Twitter - HugeGraph	Twitter - GraphX	WDC - HugeGraph	WDC - GraphX
数据集概览	14亿规模图	文本23GB	644亿规模图	文本1.2TB
Worker数	40	50	60	1500
CPU数	160	200	240	1500
内存	84 GB	1000 GB	1.8 TB	21 TB
磁盘	108 GB	-	2.4 TB	-
耗时	9 分钟	25 分钟	69 分钟	126 分钟
说明	每个Worker平均占内存2GB (0.1x)	使用200GB内存测试, 耗时133分	消耗资源低一个数量级, 稳定运行, 无需反复尝试执行	使用10TB内存测试~50台机器, OOM 失败

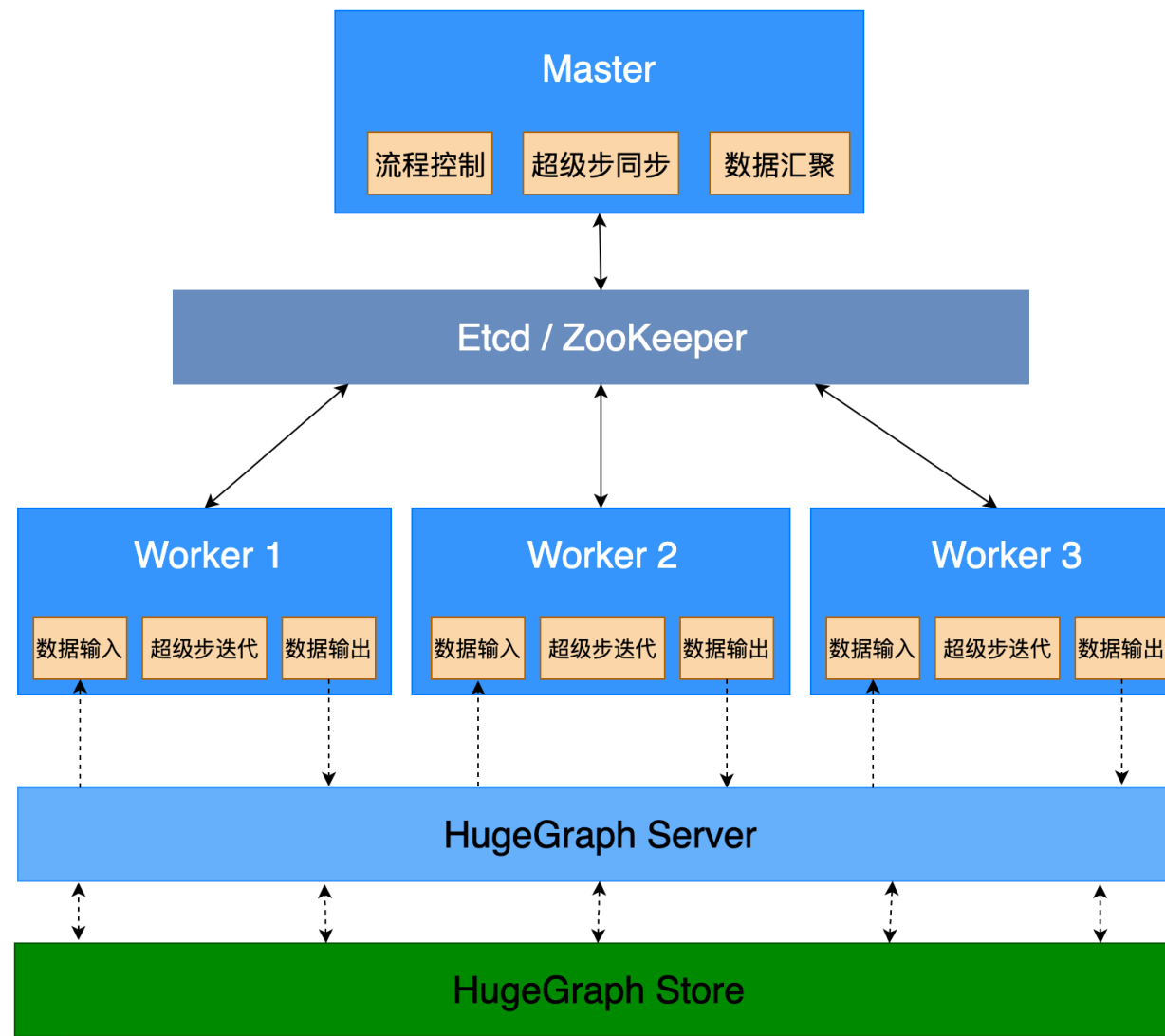


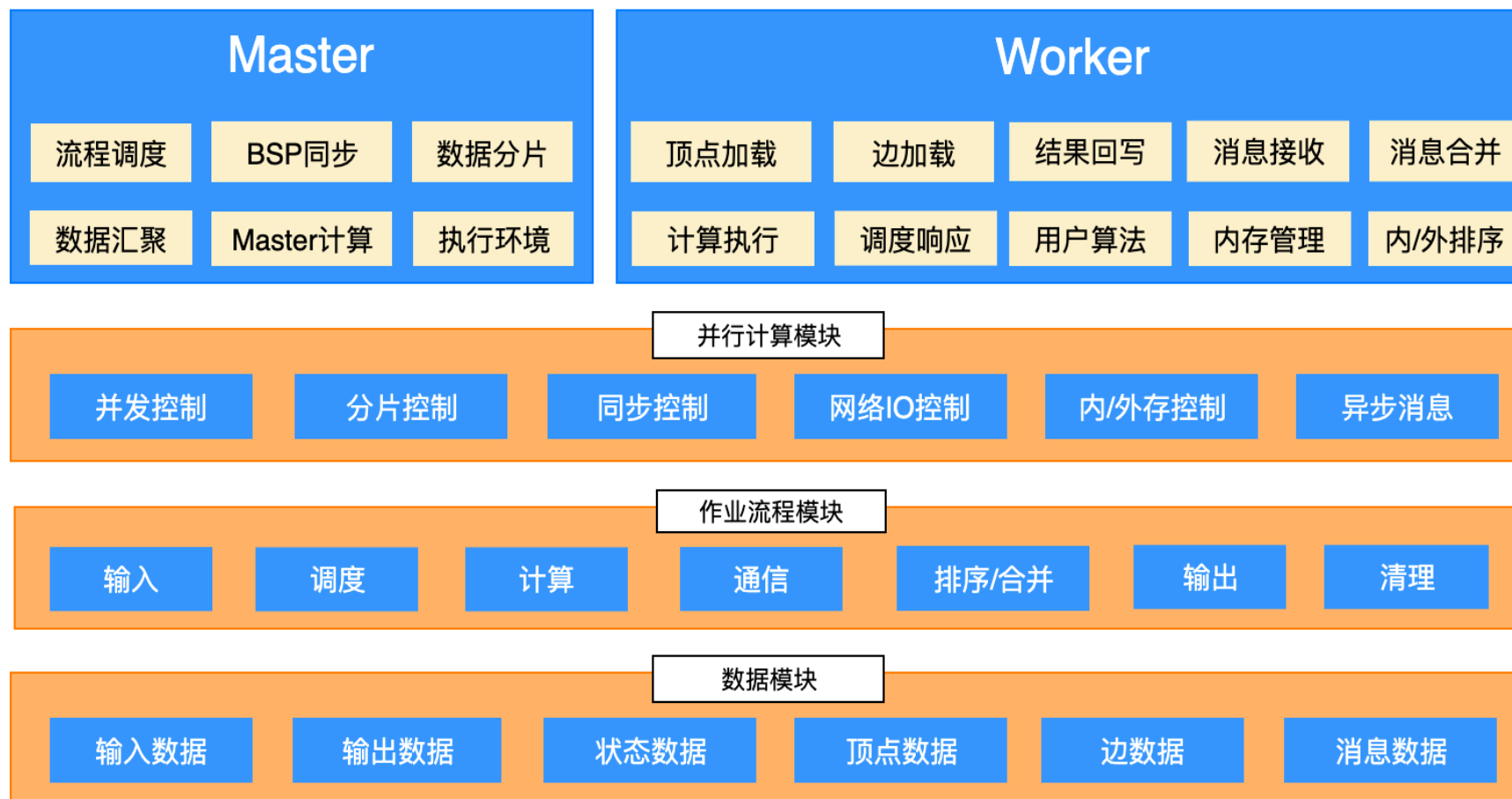
# 图查询&图计算 整体架构





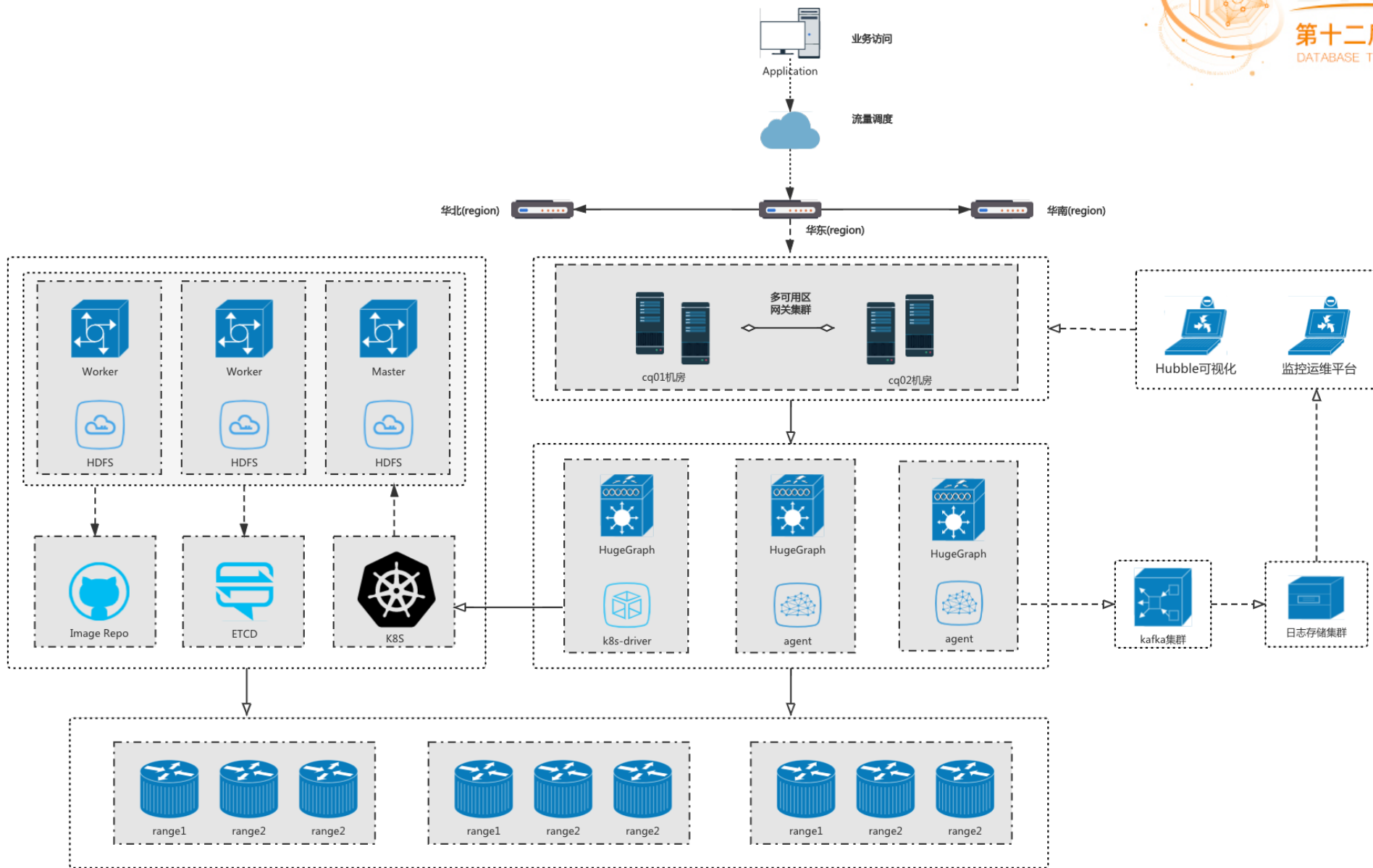
- 分布式并行计算
- 支撑千亿~万亿规模图
- 可横向扩展
- 打通图数据库







# 图计算架构 - 部署架构



规模大：千亿级规模图，自研图计算框架，横向线性扩展

性能好：并行计算，自研异步框架，多轮迭代BSP块同步

高稳定性：架构上批式计算，防止消息/超级点导致的OOM问题

简单易用：打通图数据库与图计算，以顶点为中心开发算法自然





# 图计算架构面临的关键技术挑战



并行计算  
负载不均衡

消息传输量  
巨大

内存OOM  
严重

分布式节点  
部分失败



数 / 造 / 未 / 来



并行计算  
负载不均衡

消息传输量  
巨大

内存OOM  
严重

分布式节点  
部分失败

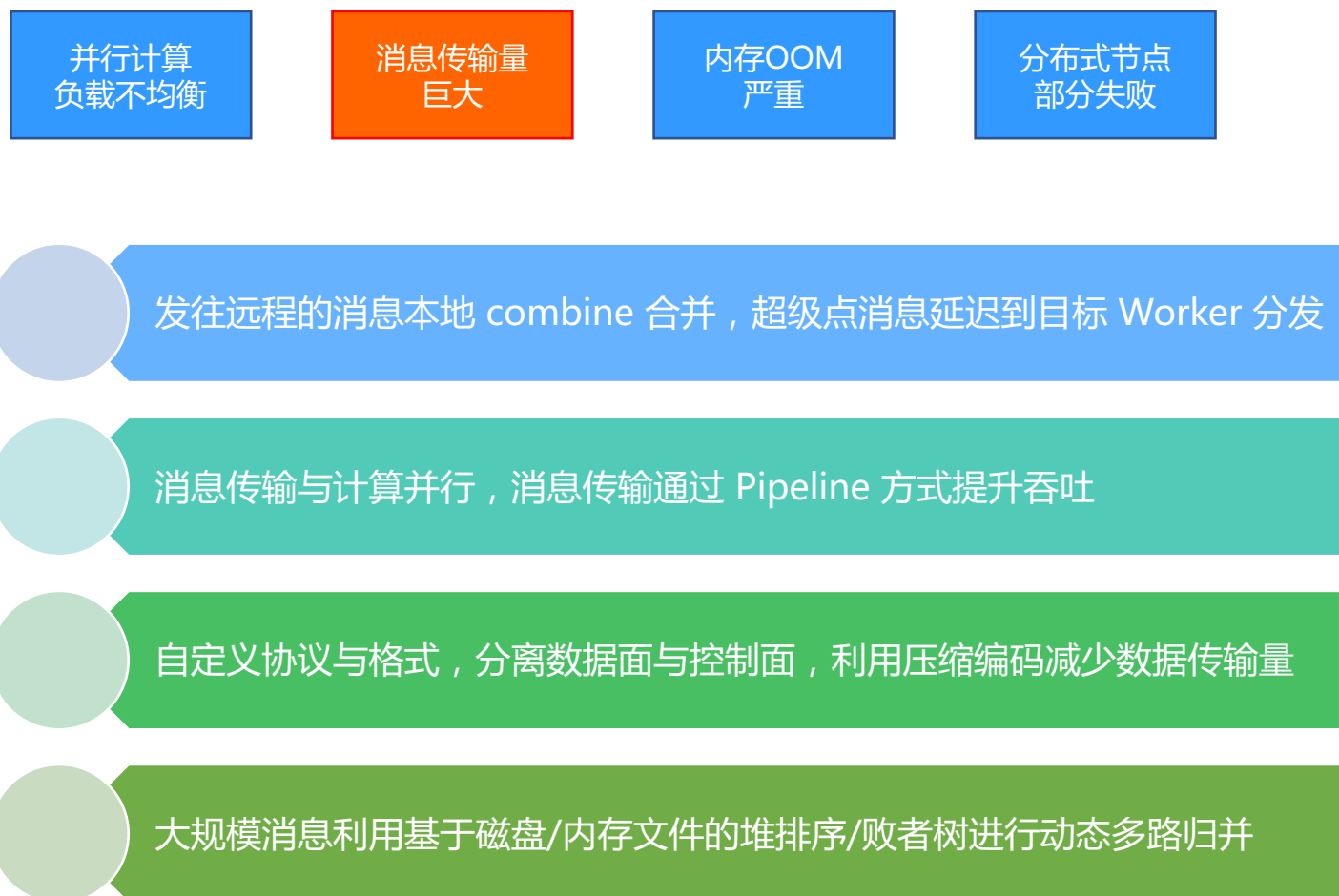
并行计算：存算分离架构，使用 K8S 极易扩展，多 Worker 并行计算

若干 Partition，每个 Worker 负责部分 Partition 数据，按照点粒度分区

一个 Worker 按边总数量分配 Partition 个数，确保性能与负载均衡

Worker 内部多 Partition 多线程并行计算，计算与消息异步并行





并行计算  
负载不均衡

消息传输量  
巨大

内存OOM  
严重

分布式节点  
部分失败

不依赖内存容量来存储图计算的输入和中间结果，根本上解决 OOM

内存-磁盘自适应，不依赖内存 Map，通过顺序 Join 内存/磁盘 列存文件

重用内存对象，自分配内存方式管理内存，使用JVM堆外内存降低GC压力

批式计算，通过流水方式进行大数据计算，矩阵向量方式平衡性能



并行计算  
负载不均衡

消息传输量  
巨大

内存OOM  
严重

分布式节点  
部分失败

部分节点失败，将导致整个集群计算中断，造成时间、计算资源的极大浪费

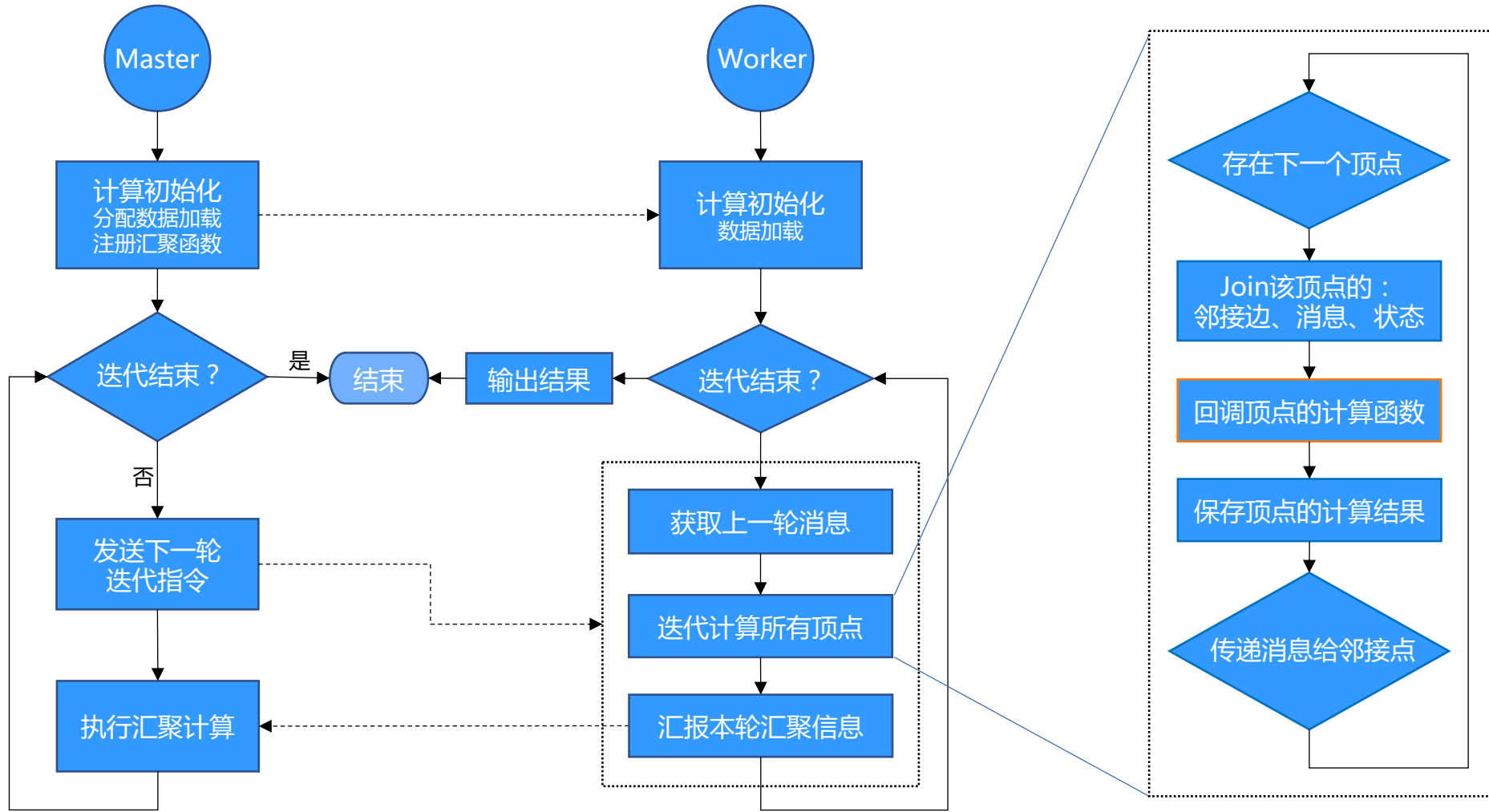
建立全链路负载反馈机制，尽可能避免部分节点因瞬时负载过高导致的失败

建立故障恢复机制，每一轮计算的中间结果快照到分布式文件系统上

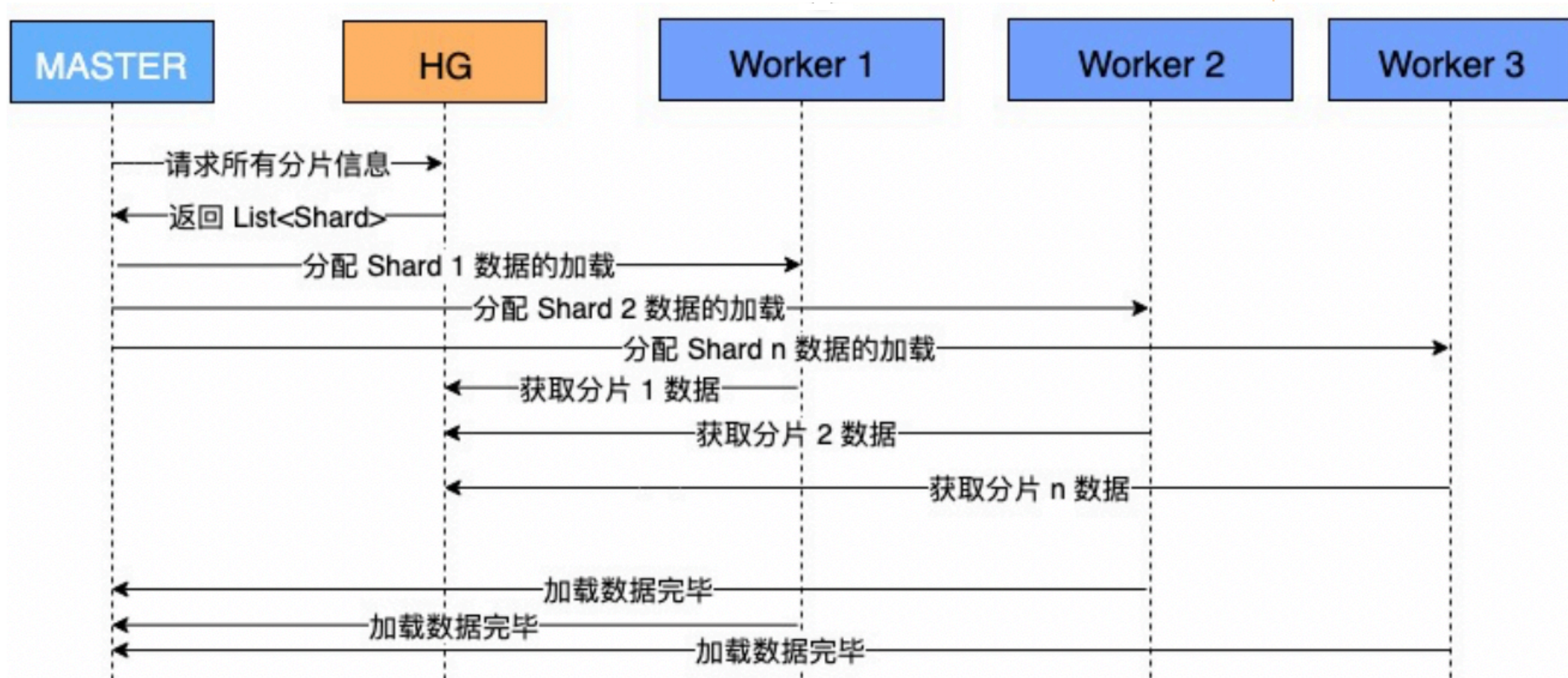
部分节点失效时，无需从头计算，通过启动新的节点替代故障节点即可



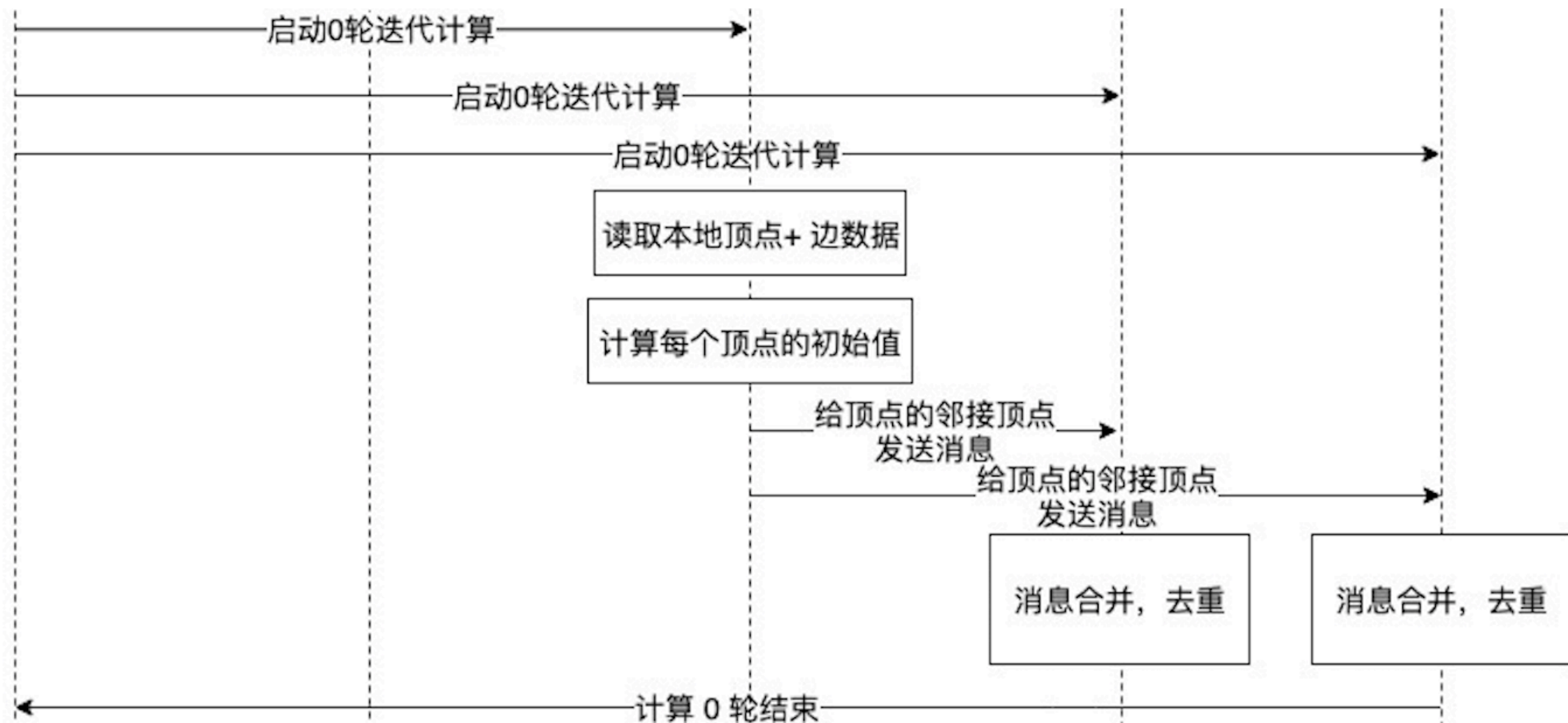




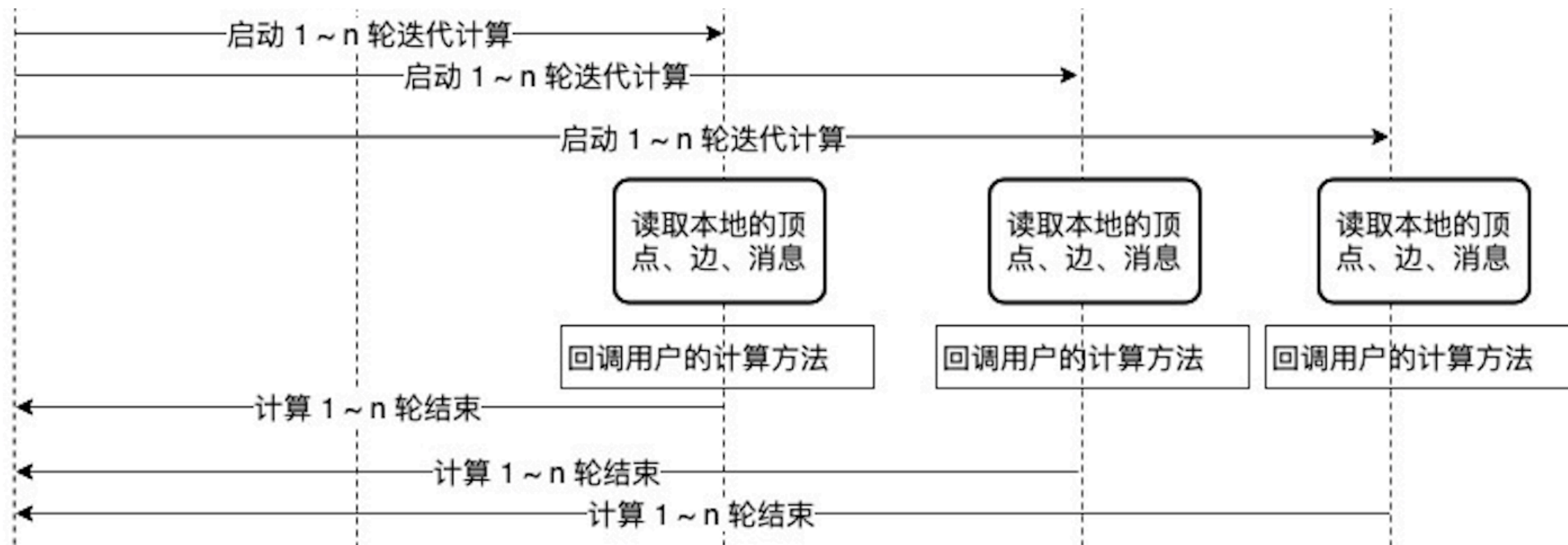
# 图计算流程 - 并行加载



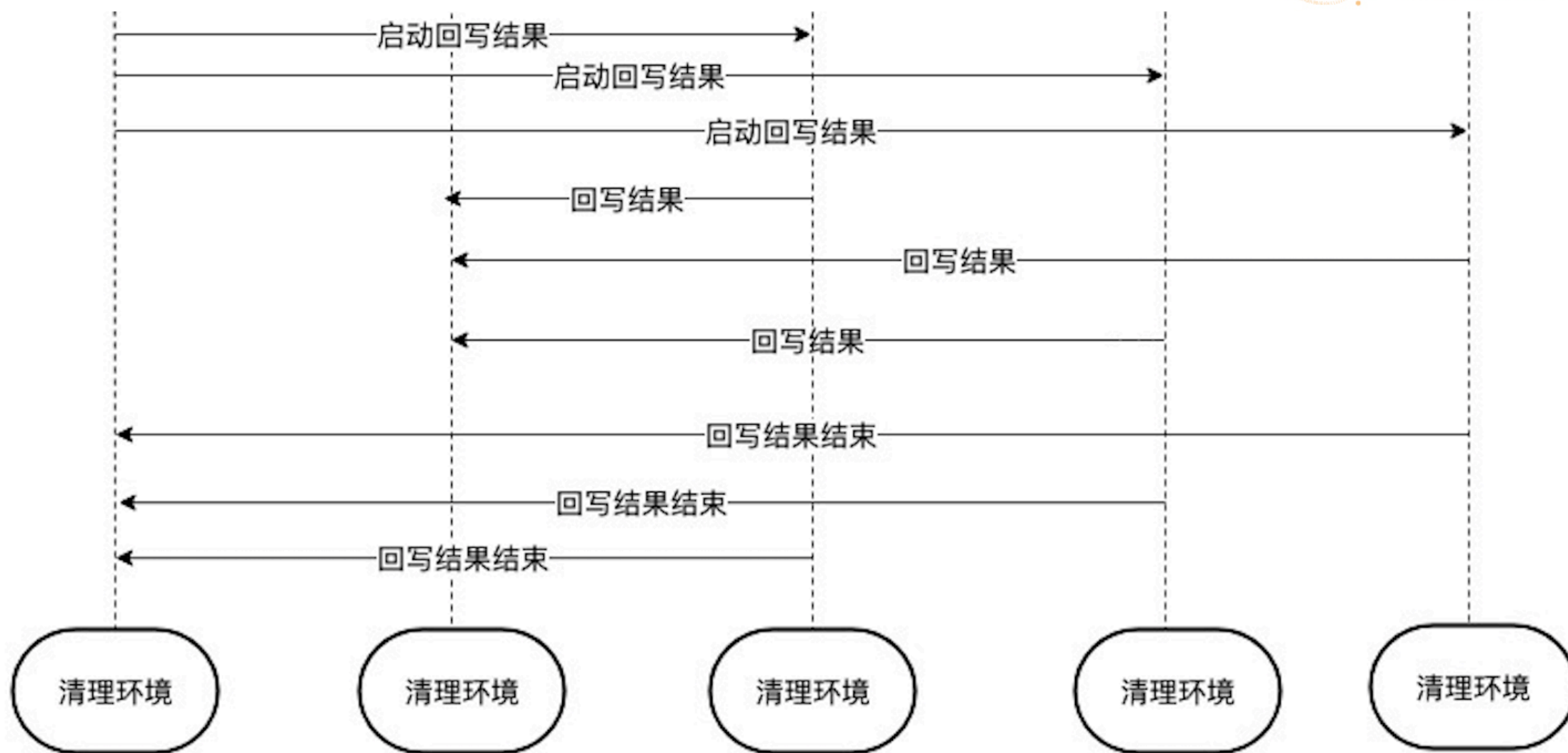
# 图计算流程 - 并行计算



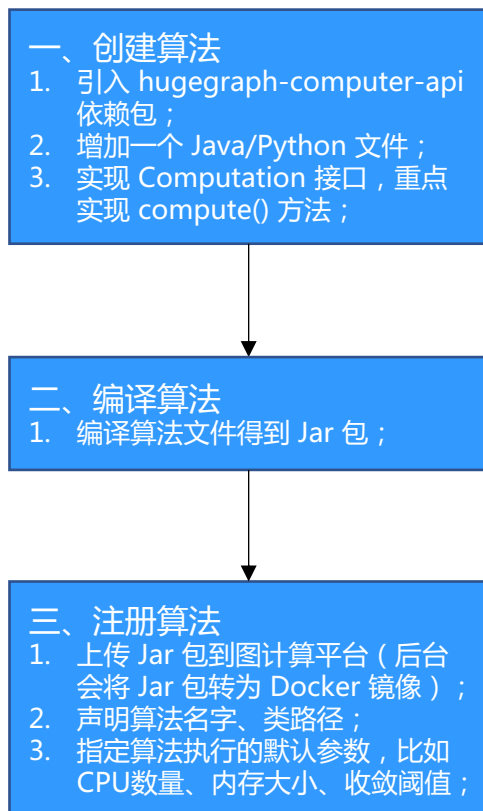
# 图计算流程 - 并行计算



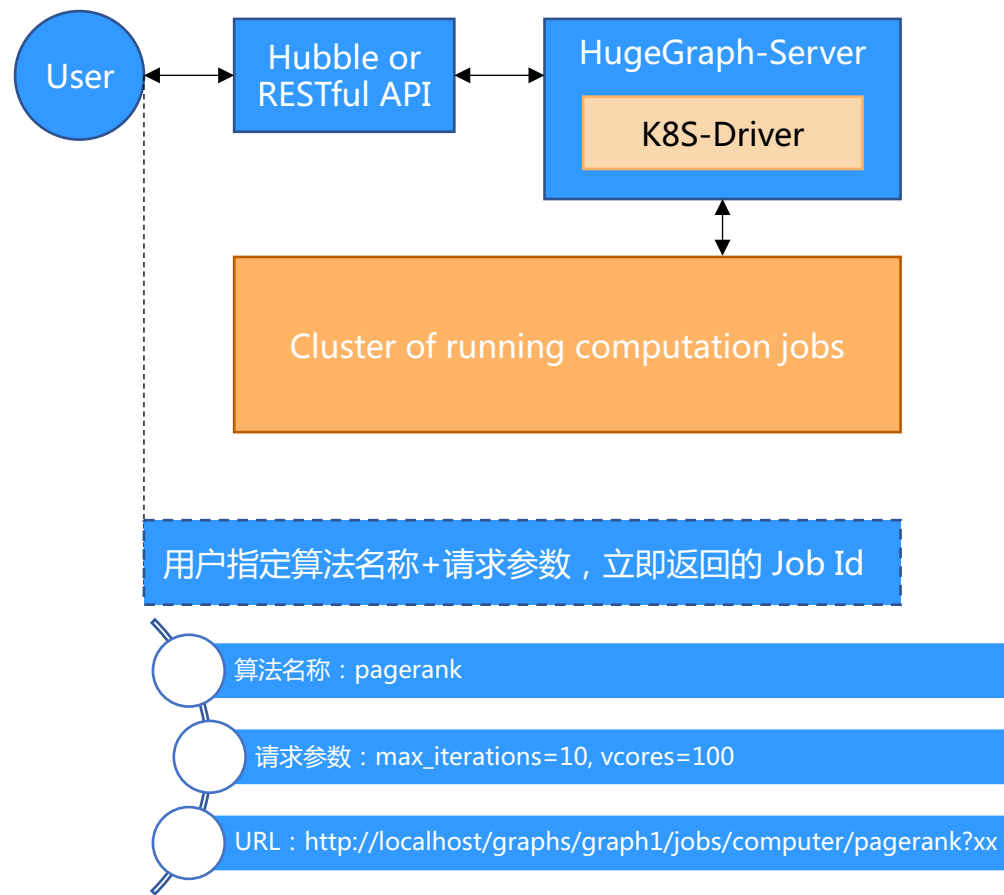
# 图计算流程 - 并行输出



## 算法开发流程



## 算法使用/执行流程





## HugeGraph图算法开发的特点：

- 算法开发简单，以顶点为中心编程直观；支持Java、Python等开发语言，轻松动态扩展算法。
- 无需关注框架流程与算法的执行细节，仅关注算法逻辑本身，专注每个顶点收到的消息及其处理。
- 计算的输入及输出格式均与算法本身解耦，采取HugeGraph标准的OLTP接口访问计算结果。

```
public void compute(ComputationContext context, Vertex vertex,
                    Iterator<DoubleValue> messages) {
    double rankRecved = Combiner.combineAll(context.combiner(), messages);
    double rank = rankRecved * (1.0 - this.alpha) +
                  this.alpha / context.totalVertexCount();
    vertex.value(rank);

    int edgeCount = vertex.numEdges();
    double contributeValue = rank / edgeCount;
    context.sendMessageToAllEdges(vertex, contributeValue);
}
```



# THANKS

时序

开源

自研

分布式

文档

Spark

SQL

HANA

OLAP

Oracle

Hadoop

MPP

DB2