



CCF BDCI CCF BIG DATA & COMPUTING
INTELLIGENCE CONTEST

2020^{8th}

CCF 大数据与计算智能大赛

赛题名：基于大数据的互联网虚拟身份归一处理性能优化

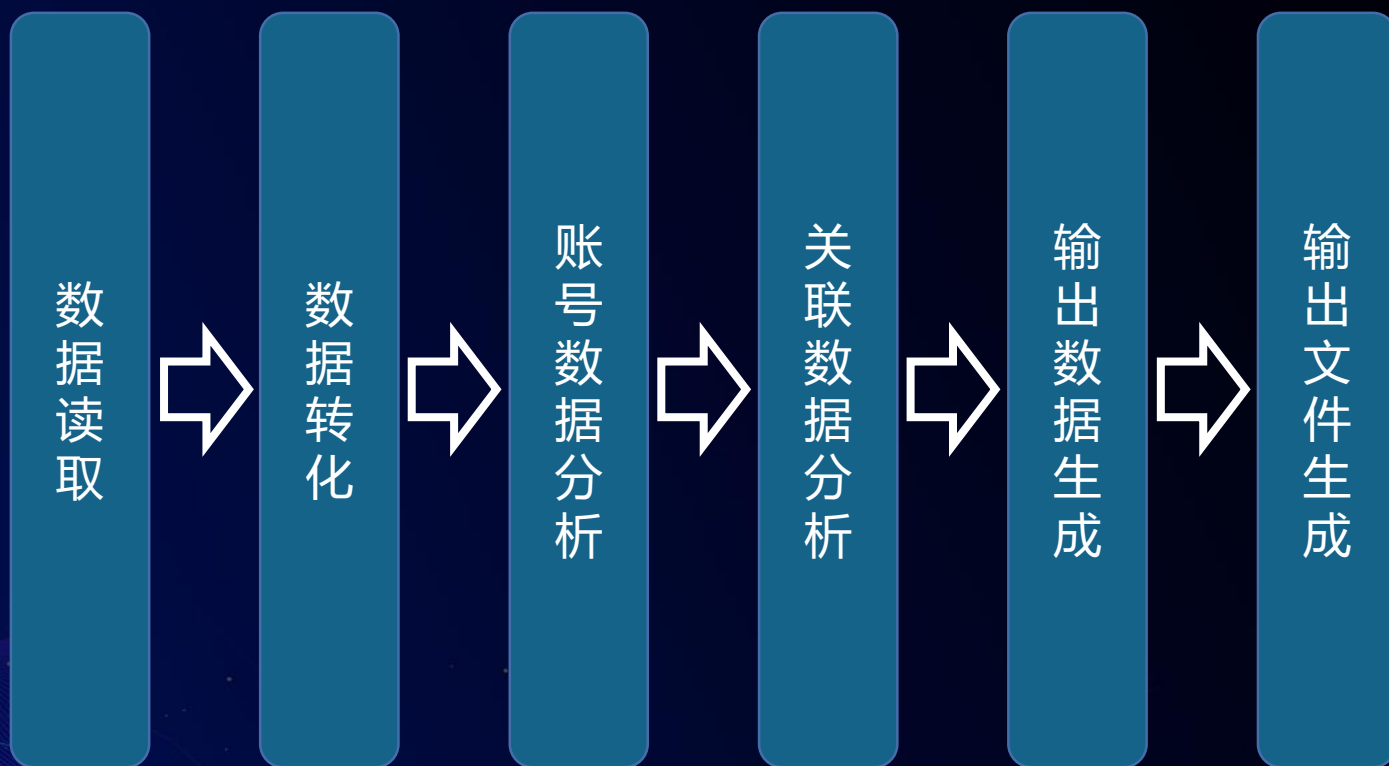
队伍名称：冀数所普拉斯

基本方案 (1.0)



CCF BDCI CCF BIG DATA & COMPUTING
INTELLIGENCE CONTEST

程序进程





超前预处理：预估账号数量和关系数量

数据映射：对文本文件进行映射读写



关系计算



CCF BDCI

CCF BIG DATA & COMPUTING
INTELLIGENCE CONTEST

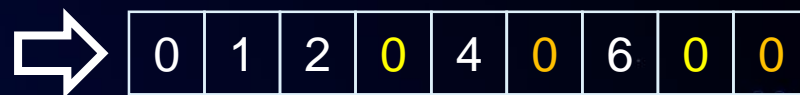
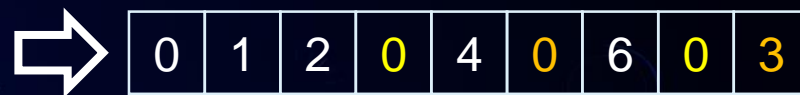
初始化



经典并查集



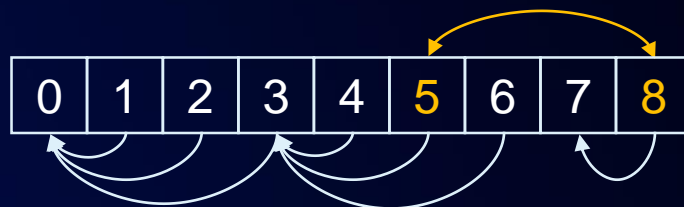
示例



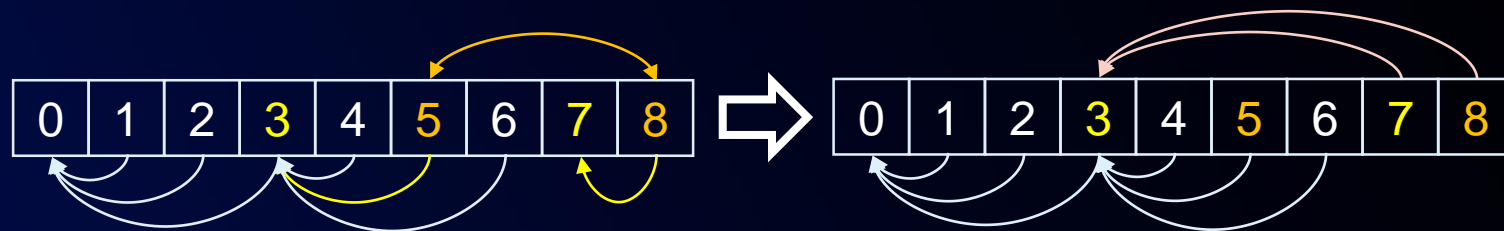
核心：通过递归进行深度并查。

问题：关系越复杂，路径越深，从底部找到根节点会越来越难

初始化



局部并查集



示例



特点：通过判断进行有限深度并查

优势：路径压缩，有限深度，局部并查

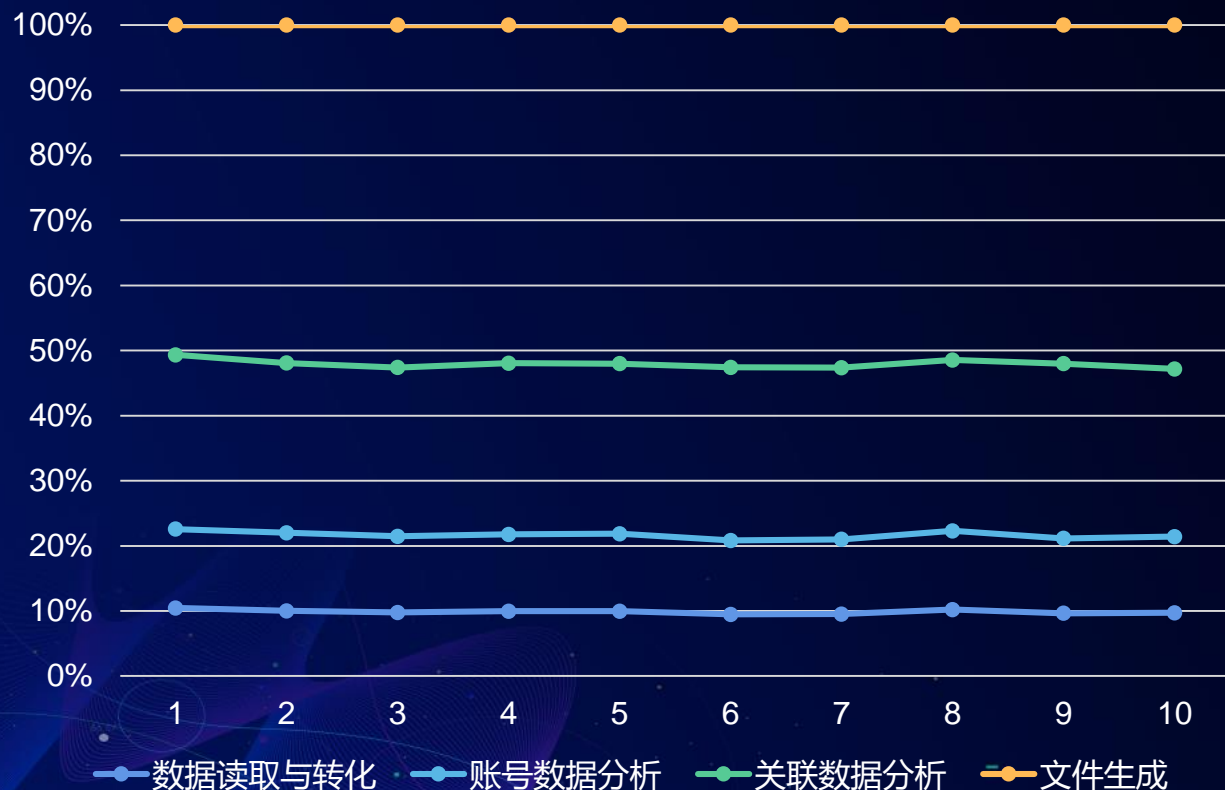
线程分配



CCF BDCI

CCF BIG DATA & COMPUTING
INTELLIGENCE CONTEST

数据处理各部分耗时一览表



特点：数据读写、关联关系计算耗时较长

处理：分批次递增式多线程读取;依赖关系分析任务即时调用；分批次数据生成

关联数据提取 0

关联数据提取 1--4

关联数据提取 5

关联关系计算

输出文件生成

输出数组生成0

输出数组生成2

输出数组生成1

最终方案 (2.0)

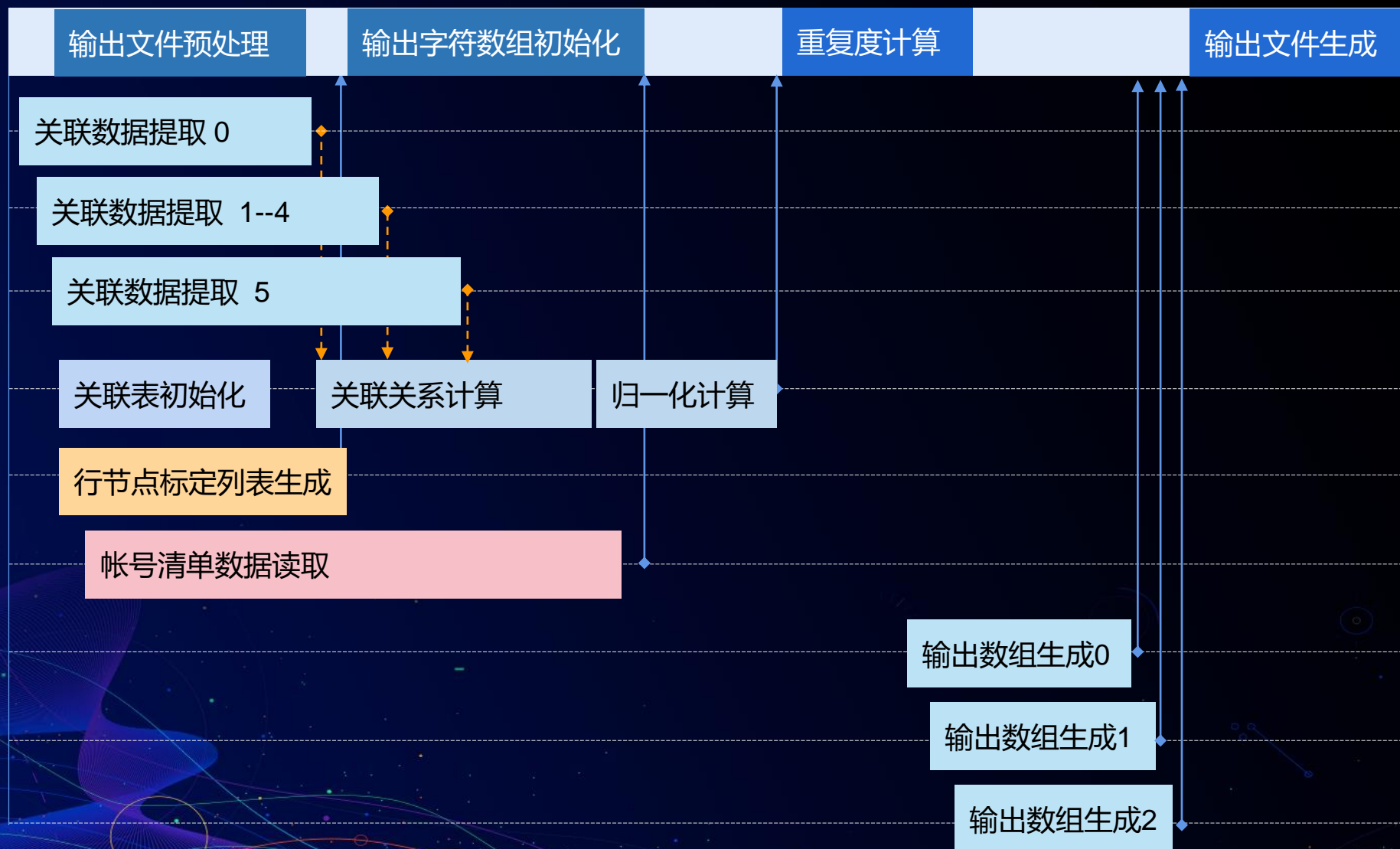


CCF BDCI

CCF BIG DATA & COMPUTING
INTELLIGENCE CONTEST

主线程

线程池



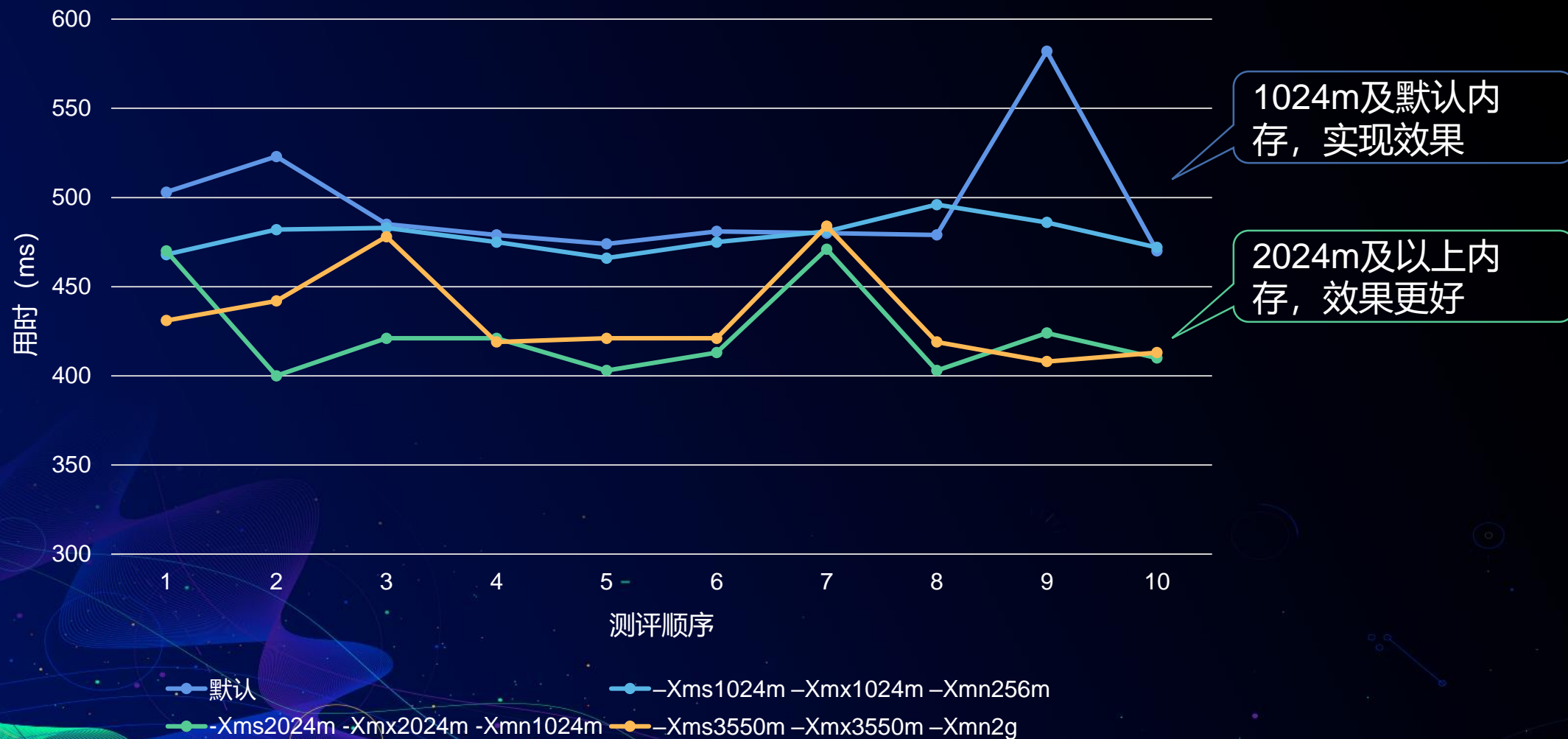
结果分析



CCF BDCI

CCF BIG DATA & COMPUTING
INTELLIGENCE CONTEST

不同内存限制下程序运行效果比较



优化方案 (3.0)

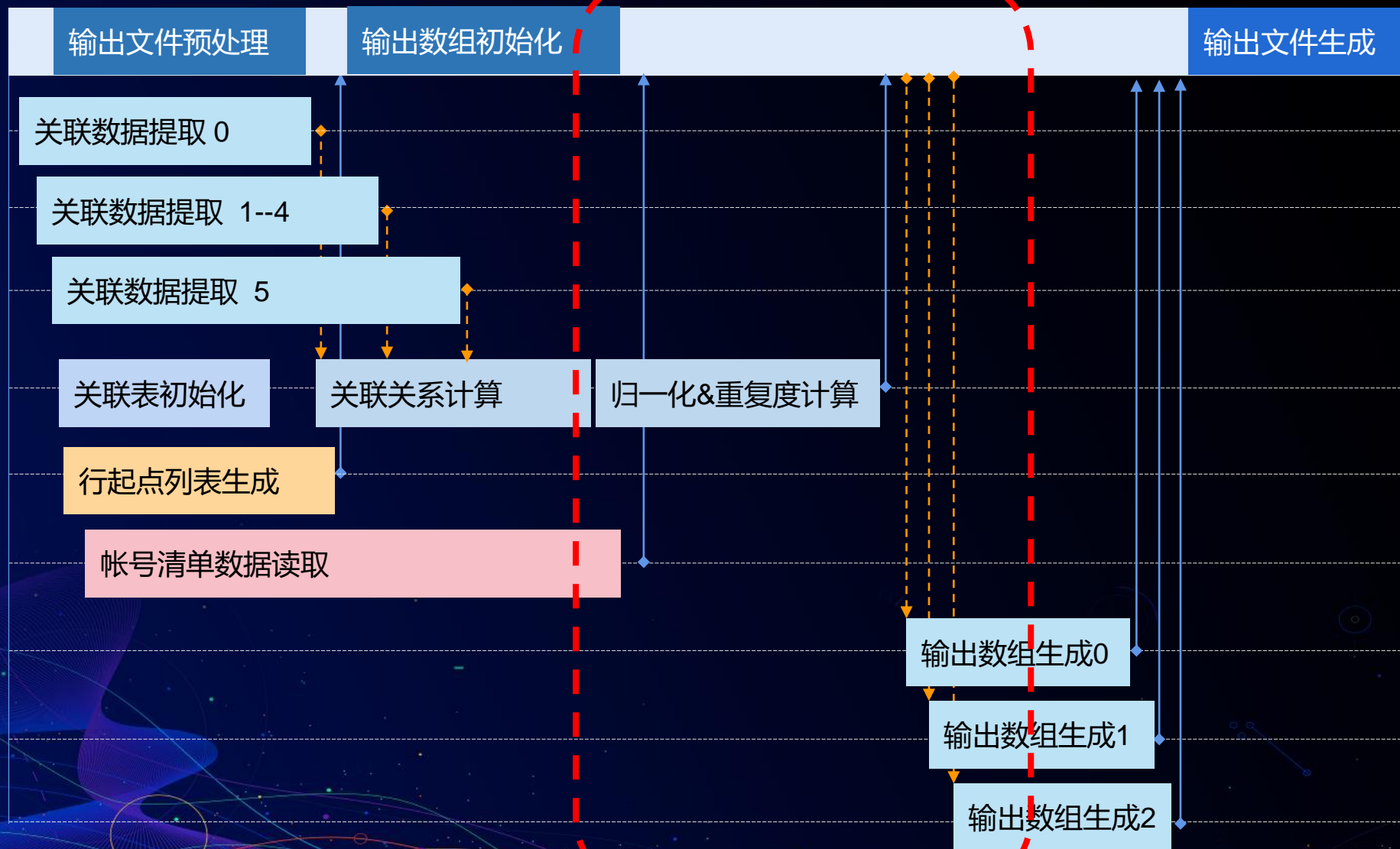


CCF BDCI

CCF BIG DATA & COMPUTING
INTELLIGENCE CONTEST

主线程

线程池



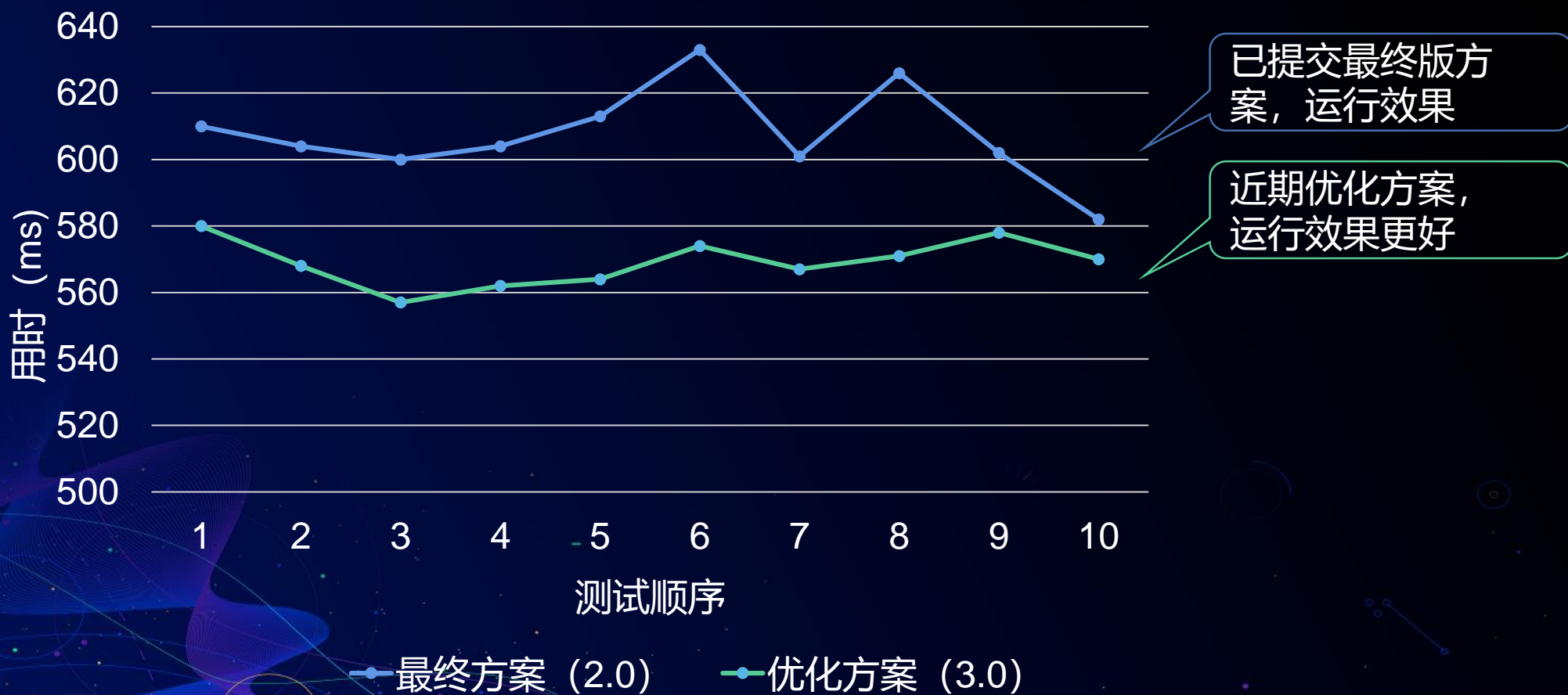
优化方案 (3.0)



CCF BDCI

CCF BIG DATA & COMPUTING
INTELLIGENCE CONTEST

已提交最终方案与优化方案程序运行效果比较



- ✓ 方案对数据集进行分批次处理，适用于更大规模数据集；
- ✓ 方案尽量避免过多的中间变量生成，能够合理复用已经分配的内存空间，适用更小内存；
- ✓ 方案采用的任务调度、索引优化和查询等算法都对程序整体运行效率产生重大影响，适用于多种大数据应用场景。



CCF BDCI CCF BIG DATA & COMPUTING
INTELLIGENCE CONTEST

谢谢!