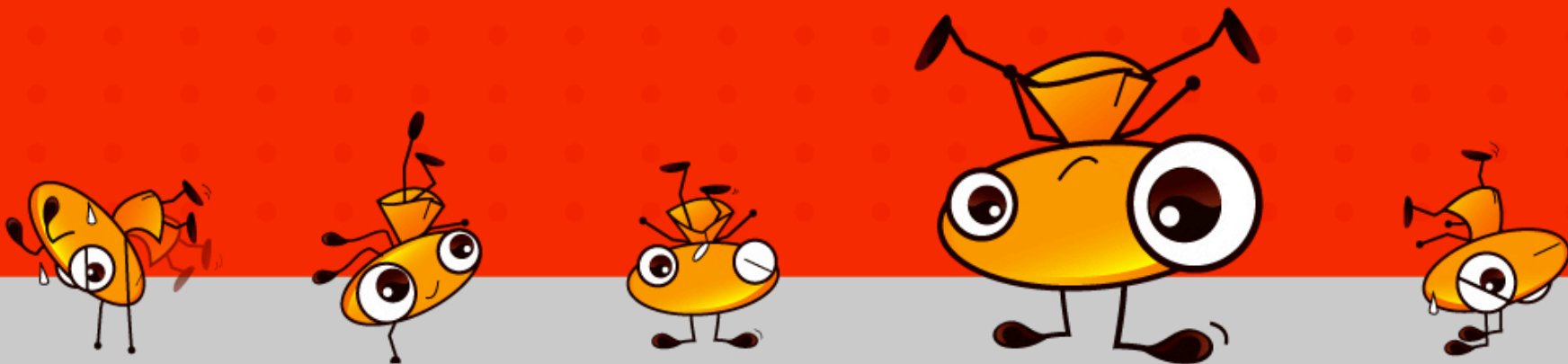


# 淘宝云梯分布式计算平台整体架构

淘宝数据平台与产品部  
云铮



## 目录

系统架构

数据同步方案

调度系统

元数据应用

## 目录

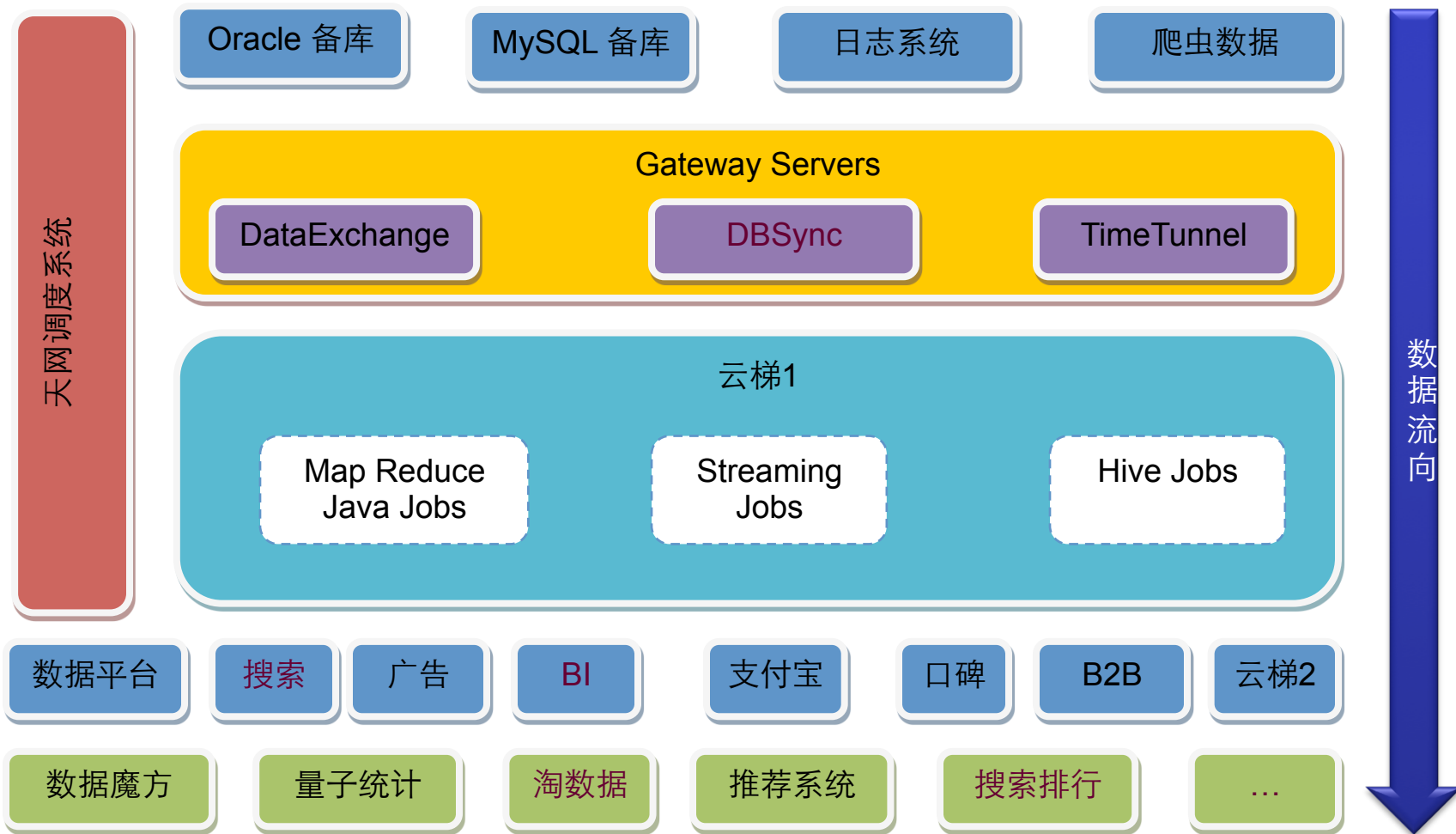
系统架构

数据同步方案

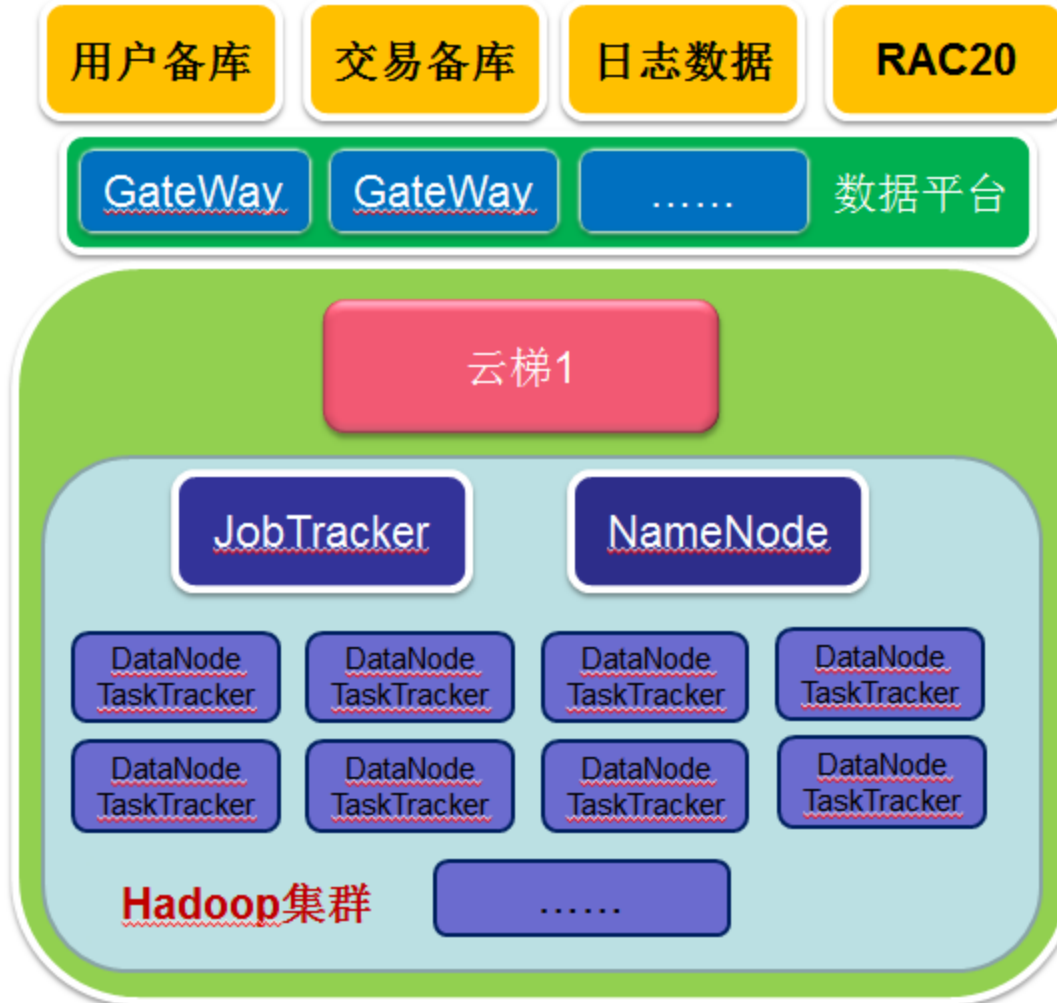
调度系统

元数据应用

系统整体架构



## 淘宝云计算介绍



## 目录

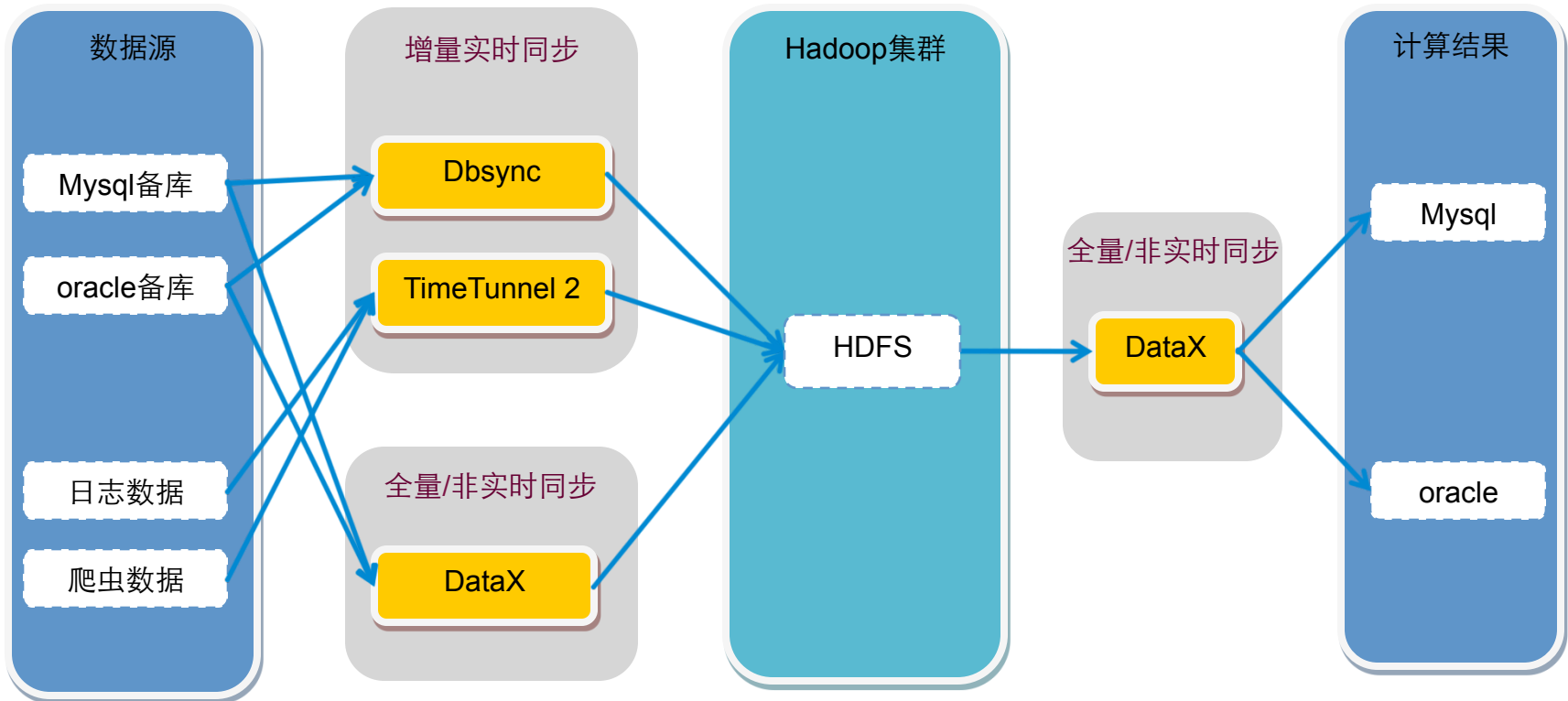
系统架构

数据同步方案

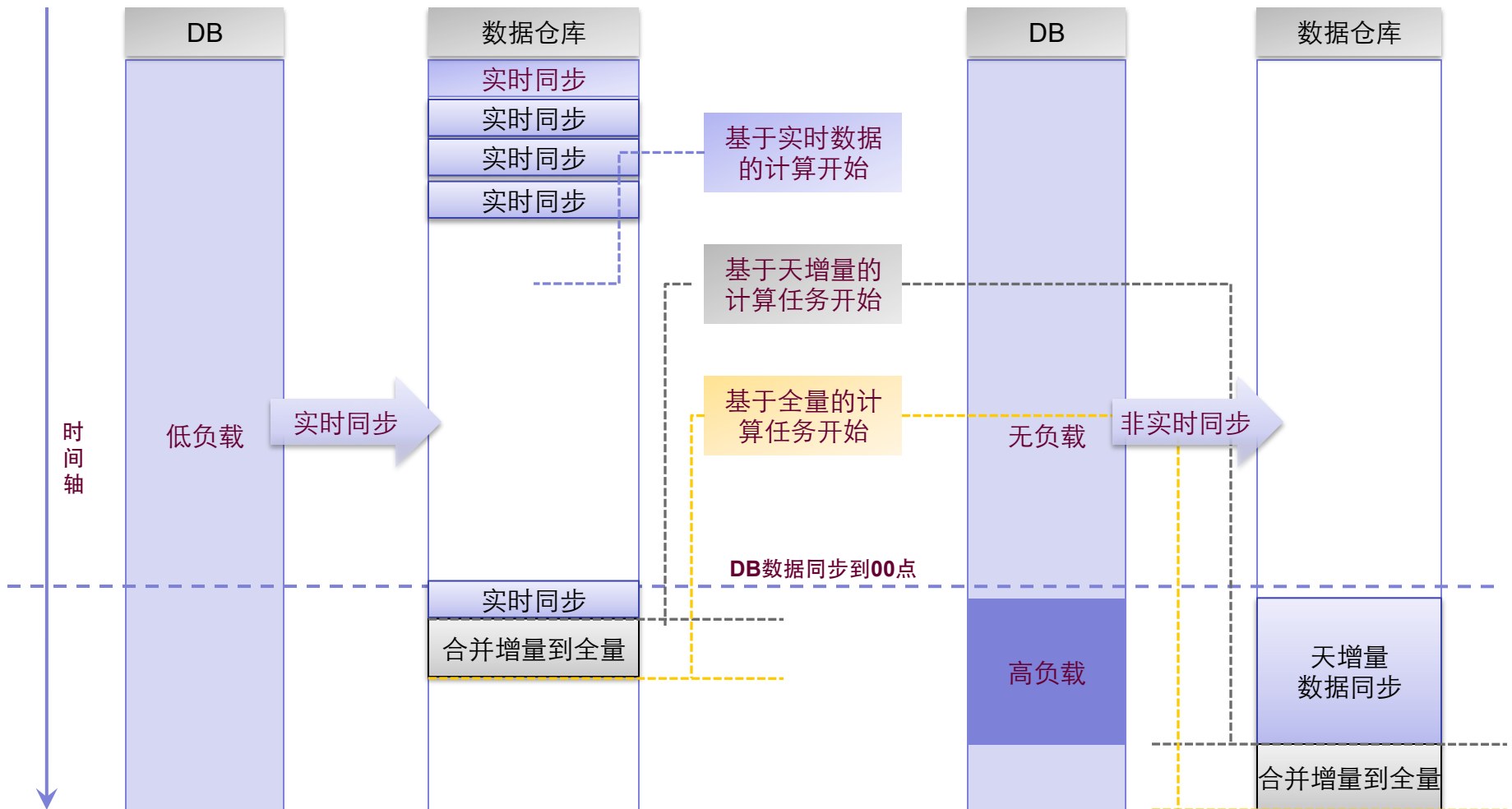
调度系统

元数据应用

## 数据同步方案——概览

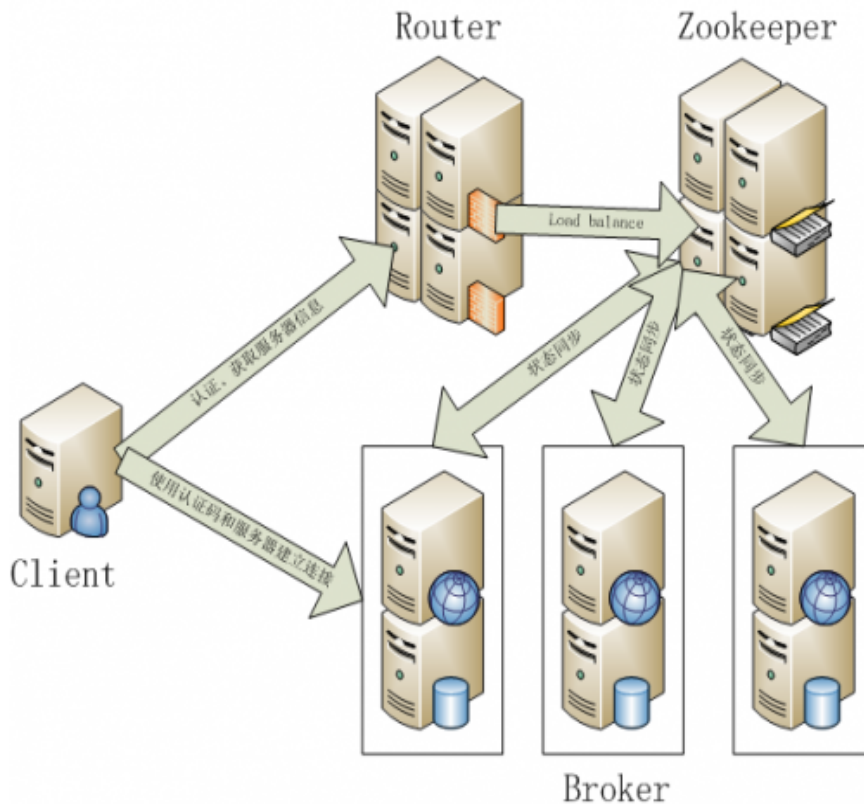


# 数据同步方案——实时同步VS非实时同步





## 数据同步方案—— TimeTunnel2 介绍



TimeTunnel是一个实时数据传输平台，TimeTunnel的主要功能就是实时完成海量数据的交换，因此TimeTunnel的业务逻辑主要也就有两个：一个是发布数据，将数据发送到TimeTunnel；一个是订阅数据，从TimeTunnel读取自己关心的数据。TimeTunnel作为一个实时数据传输平台具有以下特点：

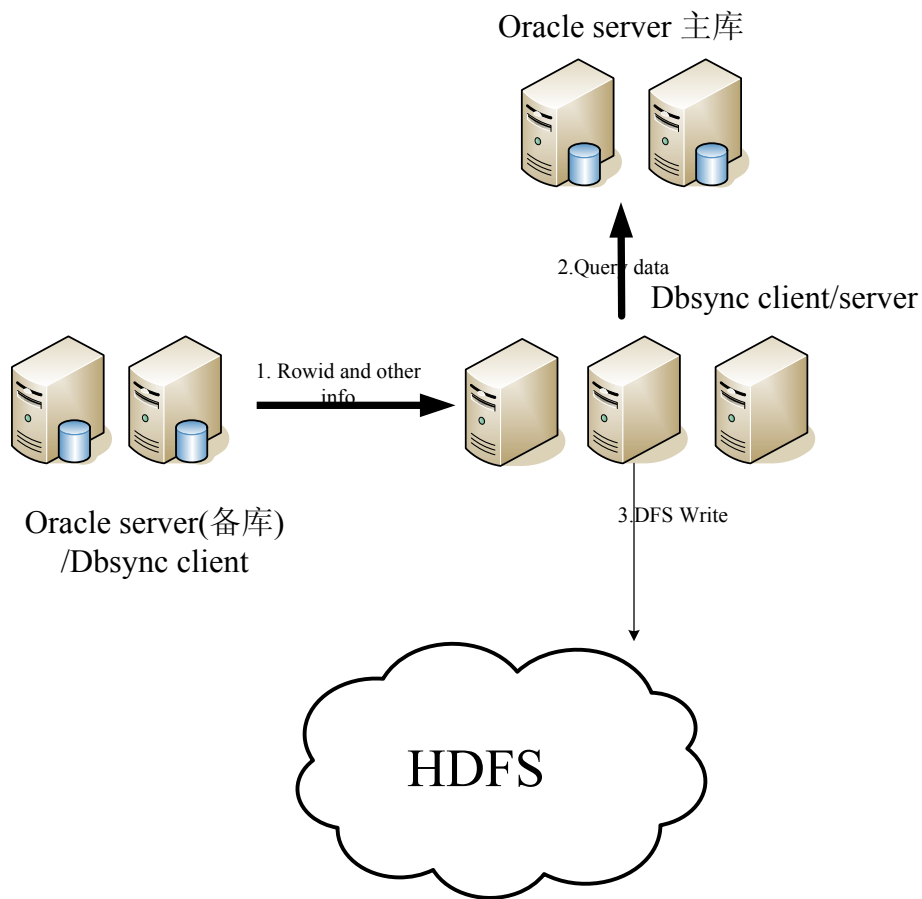
**高效性：**单点1k数据可以到4万TPS

**高可靠性：**M-S模式时保证数据不丢失

**高可用性：**单点故障不影响整个集群服务

**顺序性：**当没有故障发生时，保证所有传输都是顺序的，或者说一次连接内的传输是顺序的。

## 数据同步方案—— Dbsync 介绍



dbsync是一个用于同步服务库数据到HDFS的产品，通过分析数据库服务器的log文件来提取相应的数据库动作，进而达到数据库到HADOOP的数据同步，供相关部门提取增量数据。

### Dbsync实时同步性能

记录大小	速度
2K	4M/s
9K	10M/s

应用场景

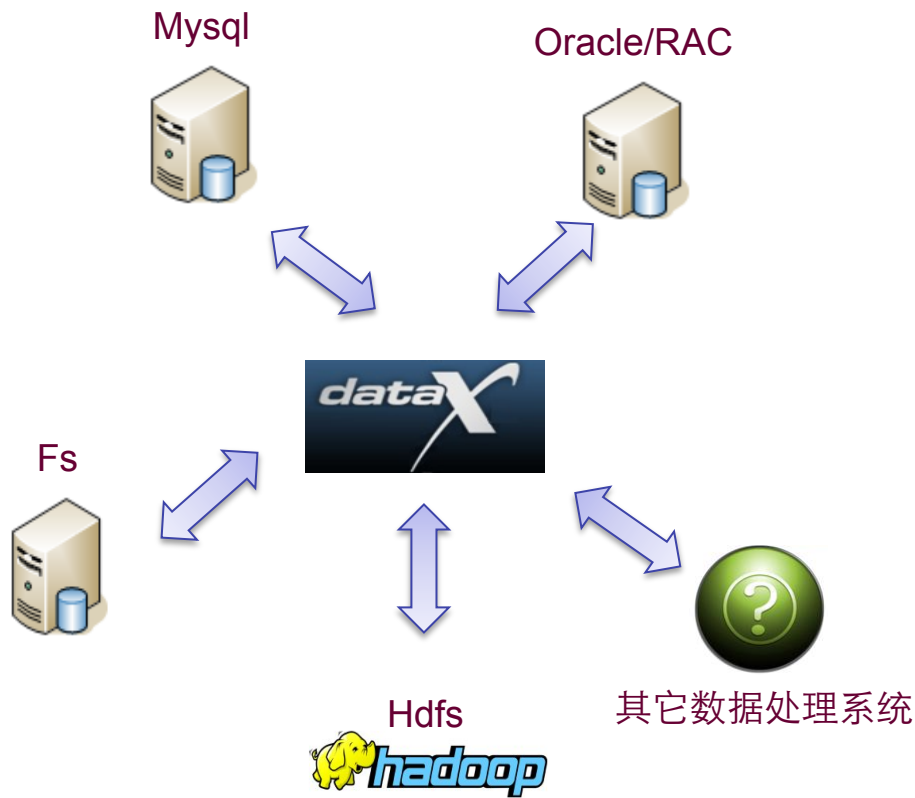
数据量 800G

00:10分备库打开:

非实时同步完成时间0:55

实时同步完成时间0:25

## 数据同步方案—— DataX 介绍



- DataX是一个在异构的数据容器之间交换数据的工具。用于在任意的数据处理系统 (RDBMS/NoSql/FS) 之间交换数据。

- Framework+plugin, Framework处理了高速数据交换的大部分问题, 插件提供对数据处理系统的访问。

- 运行模式: stand-alone / on hadoop

- Webui + cui 基于元数据的高效配置,例子: 表A sharding为32个库, 1024张表, 配置时间<1 min

DataX部分性能数据:

数据输出方向	Speed (M/s)
mysql->hdfs	29.9
oracleloader->hdfs	31.9
hdfs->oracle	18.1

## 目录

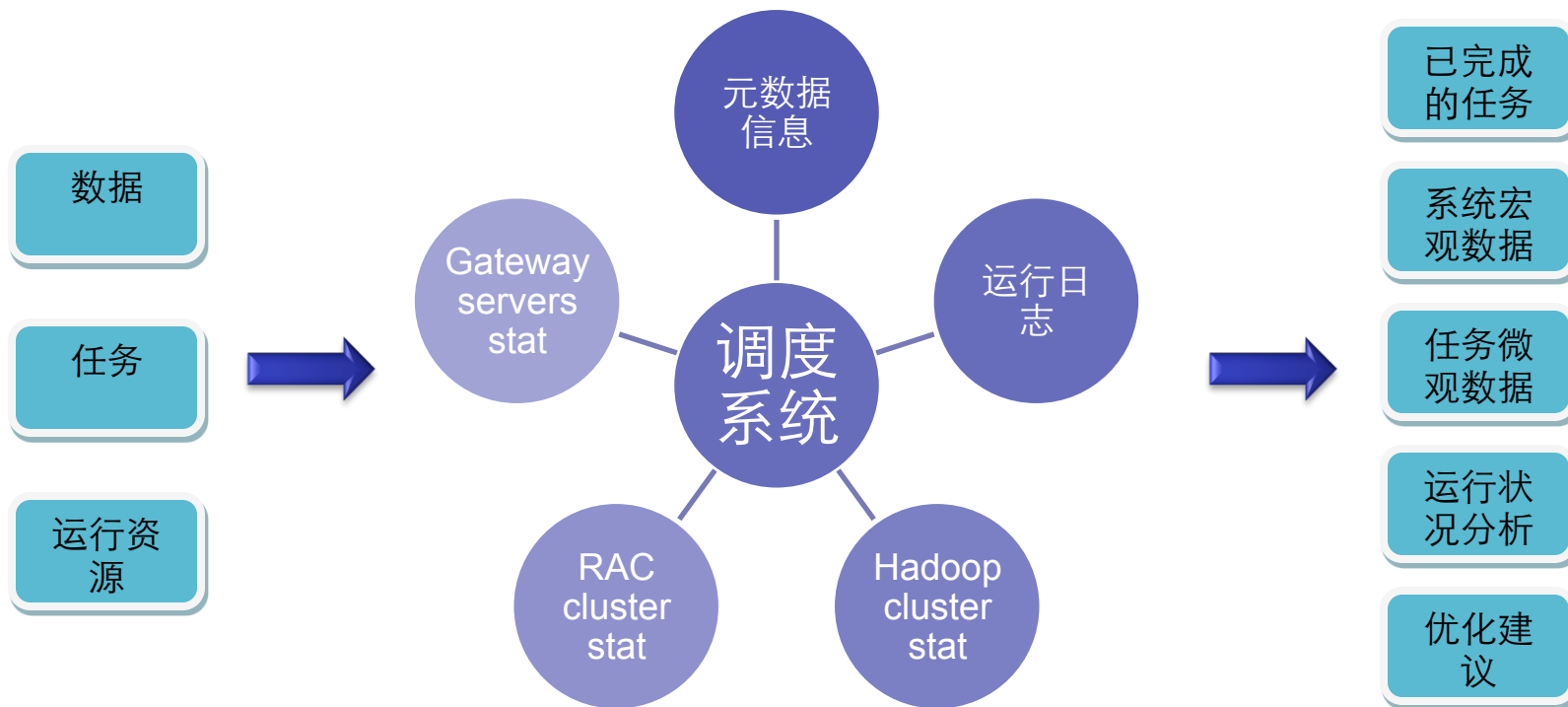
系统架构

数据同步方案

调度系统

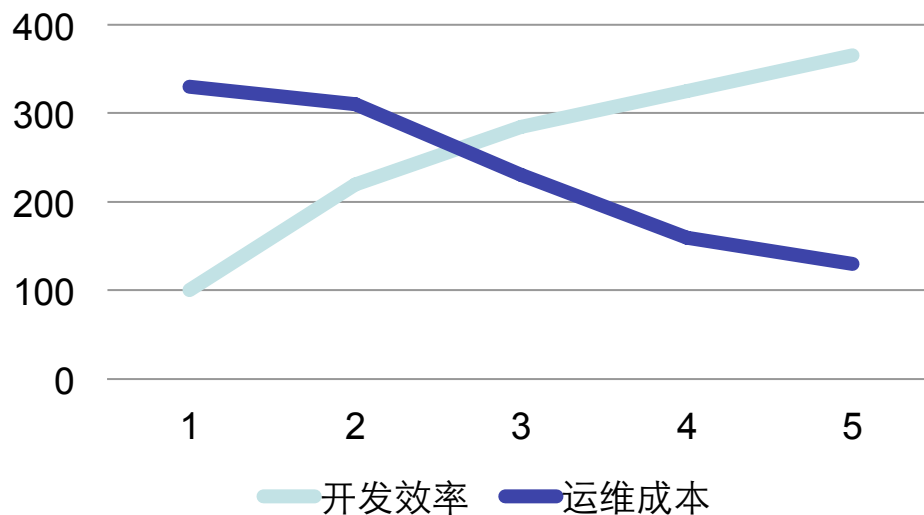
元数据应用

# 调度系统



## 调度系统——生产率银弹

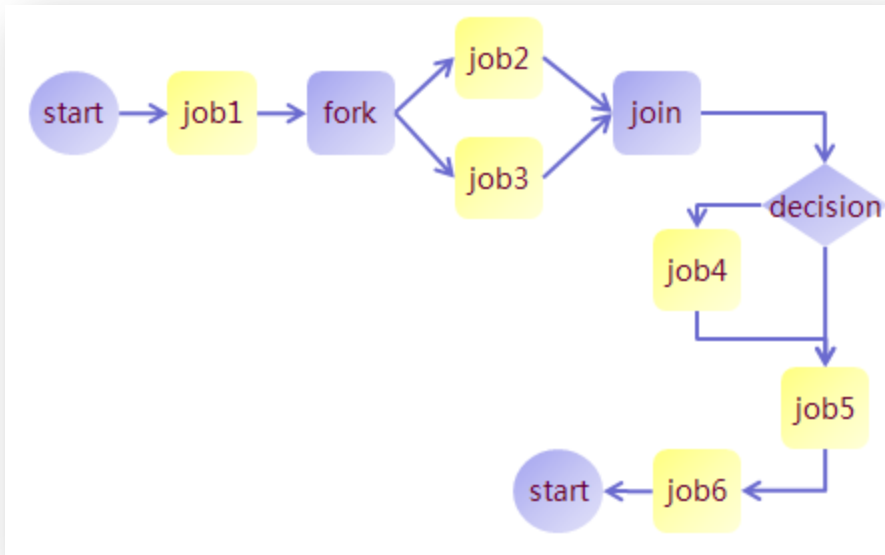
- 自动部署
- 智能调度
- 运维平台
- 监控告警
- Web UI
- 异构平台支持



## 调度系统——模块/子系统



# 调度系统——任务触发方式



Flow control/Data Trigger

Time Trigger

定时调度: 无

选择日期: 无

具体时间: 无 : 无 (小时: 分钟)

分钟调度设置

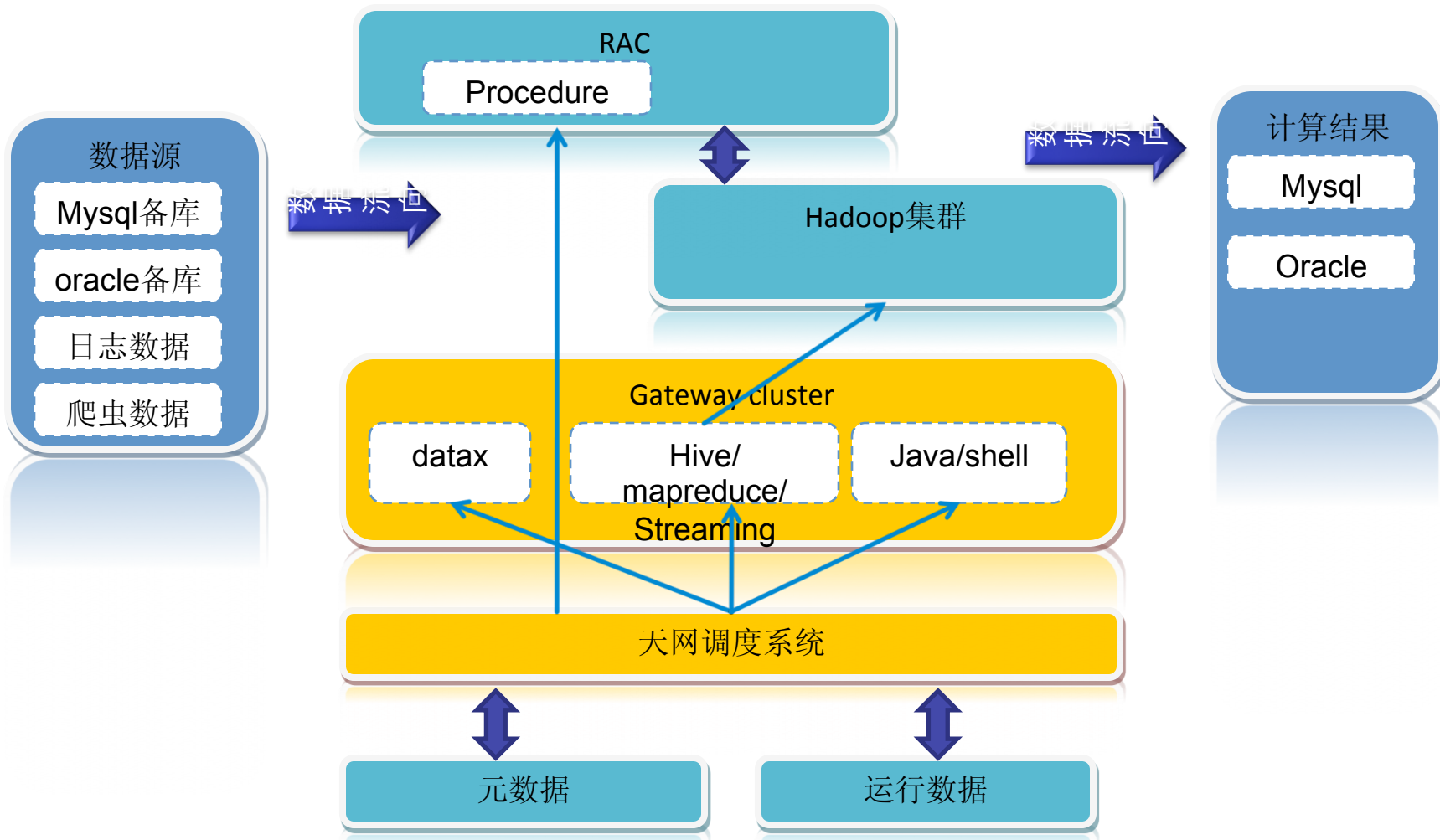
时间范围: [ ] - [ ] 间隔: [ ] 分钟

小时调度设置

自定义设置

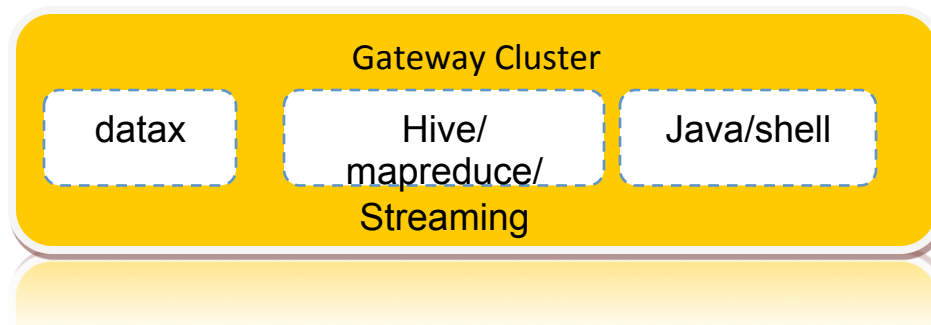


# 调度系统——调度方式



## 调度系统——什么是Gateway?

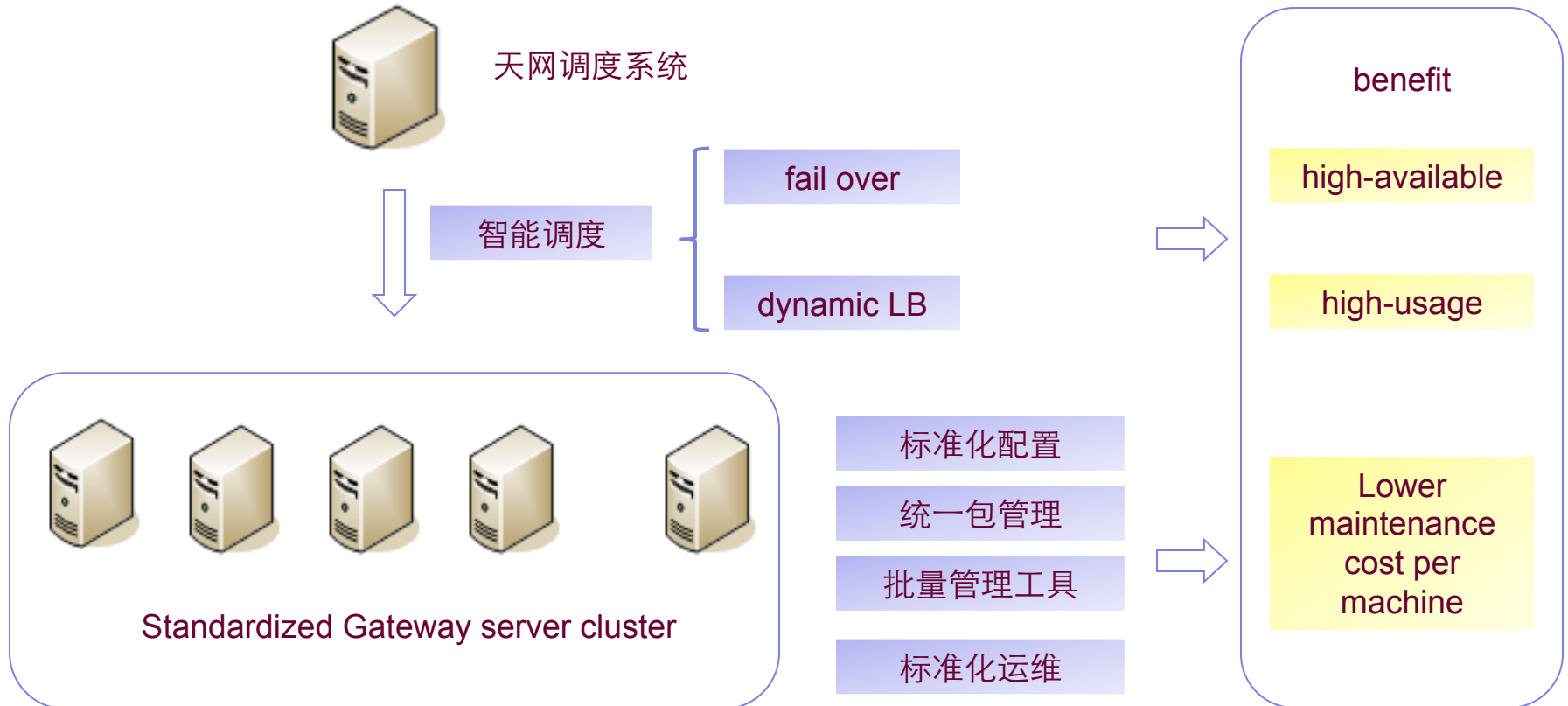
- Gateway:参与天网调度的资源
- 功能:
  - 数据同步(dataX, DBSync, TimeTunnel2...)
  - 数据上传/下载(hadoop fs -put/get/getmerge)
  - 日志收集
  - Hive sql语句提交运行
  - MapReduce程序提交运行
  - 集群间数据同步(hadoop distcp)



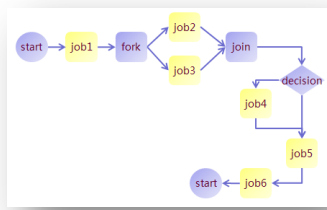
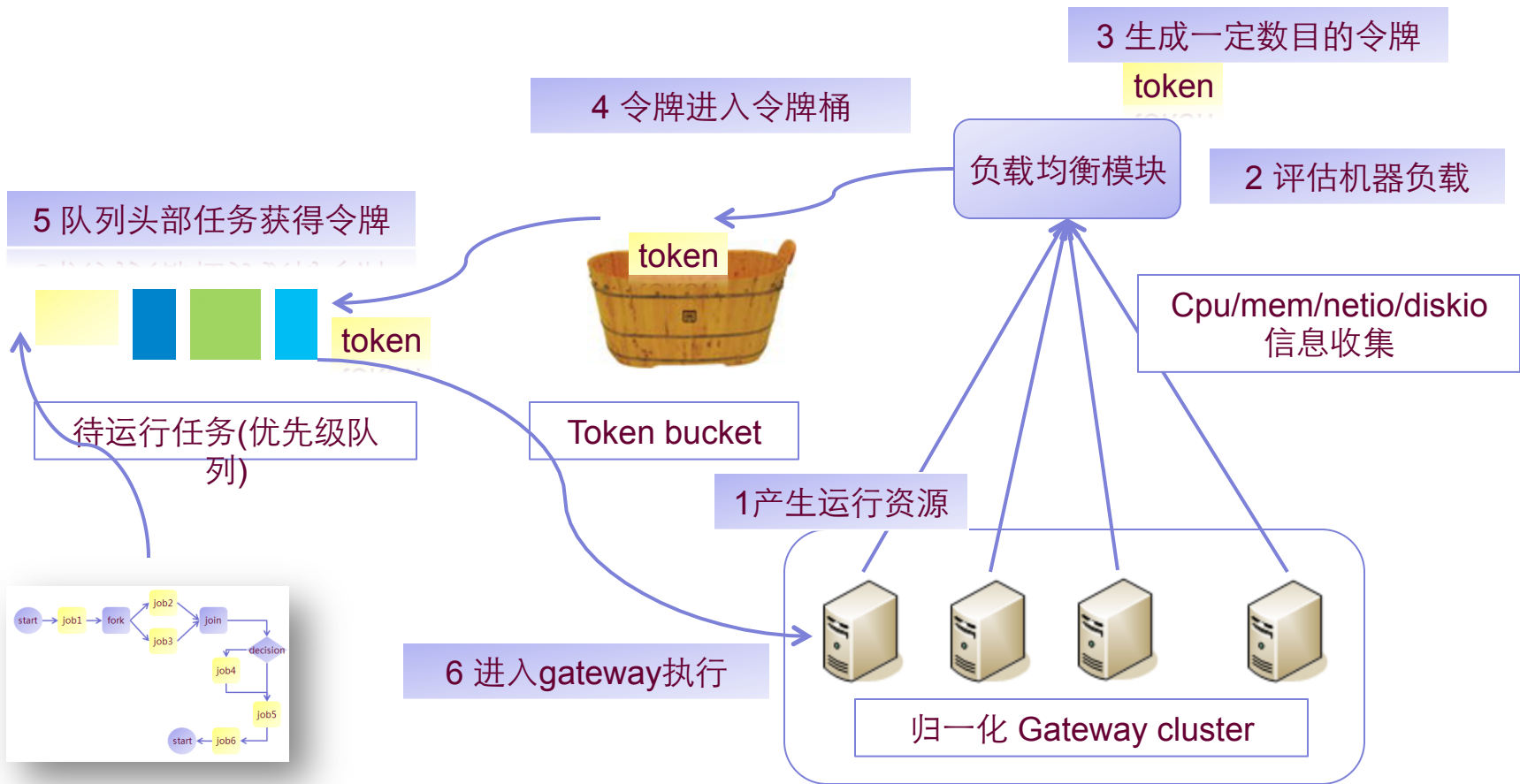
## 调度系统—— Gateway规模及规划

- 用于生产的Gateway约30台，由天网调度统一进行任务分发，并行控制。
- 数据同步(dataX, DBSync, TimeTunnel2...)
- 数据上传/下载(hadoop fs -put/get/getmerge)
- 日志收集
- Hive sql语句提交运行
- MapReduce程序提交运行
- 集群间数据同步(hadoop distcp)

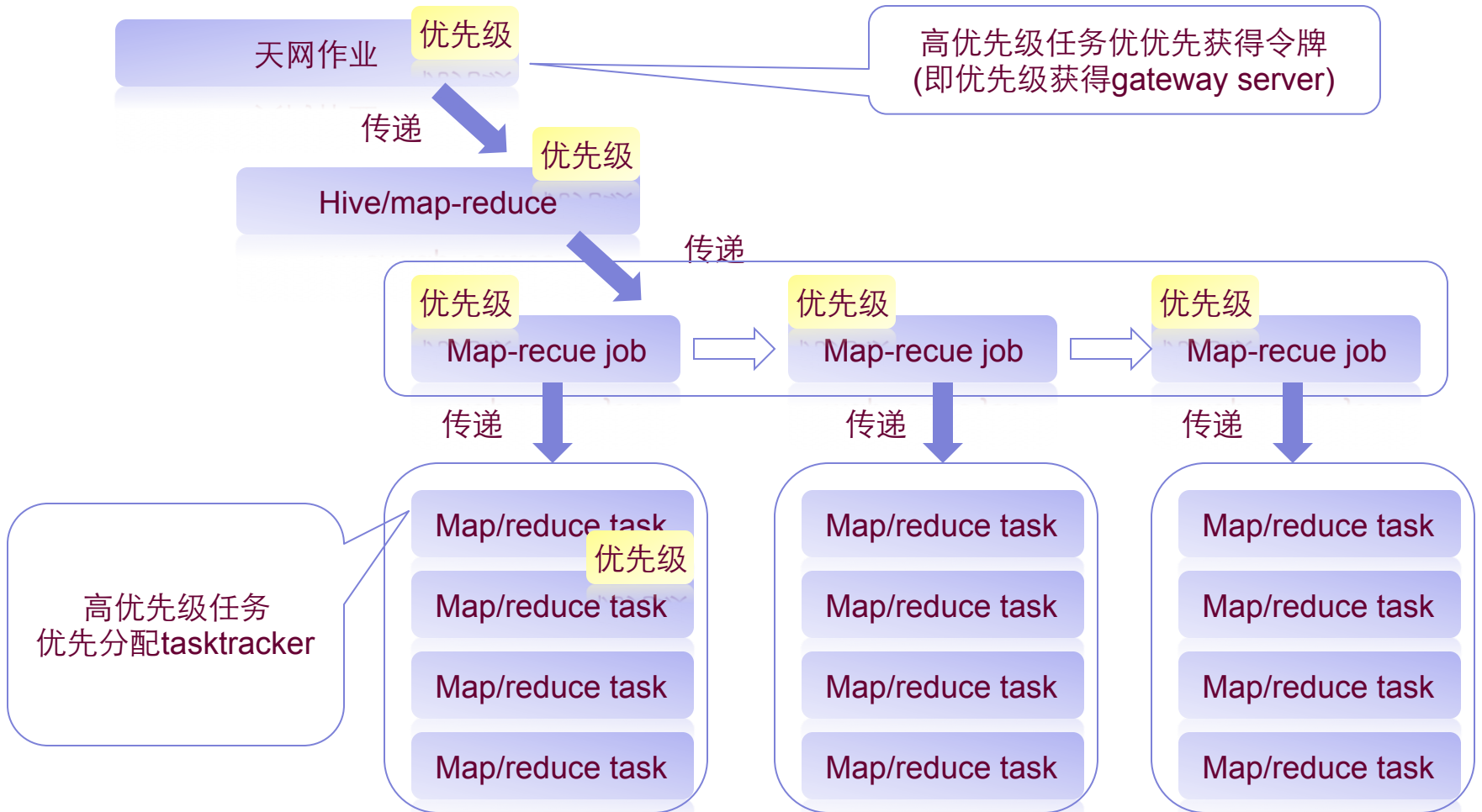
# 调度系统——gateway standardization



# 调度系统——Dynamic LB实现



# 调度系统——优先级策略(实现)

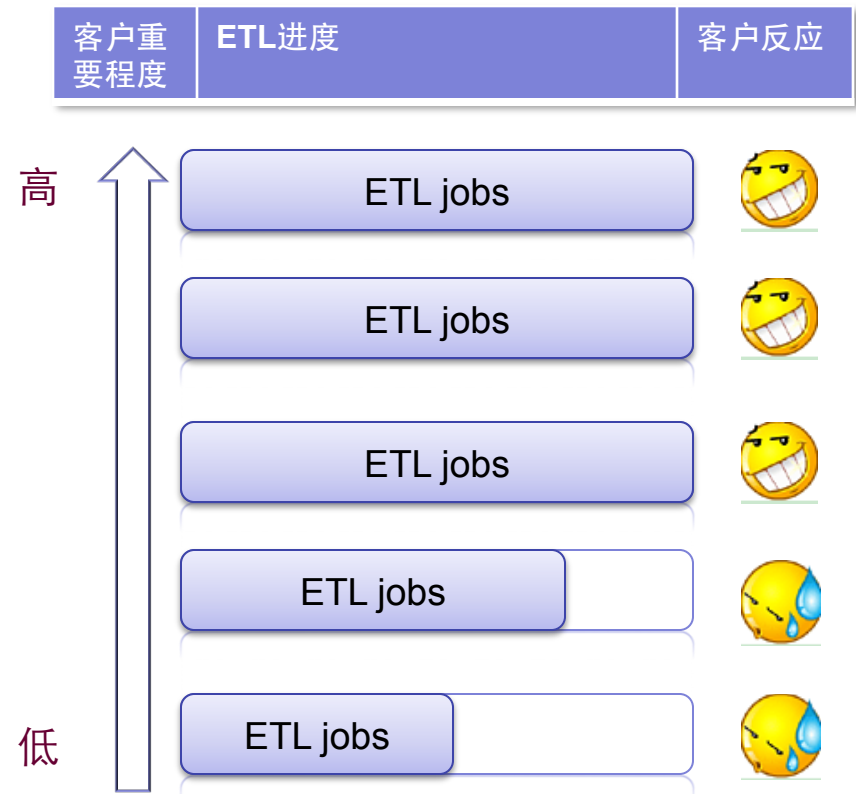


## 调度系统——优先级策略(意义)

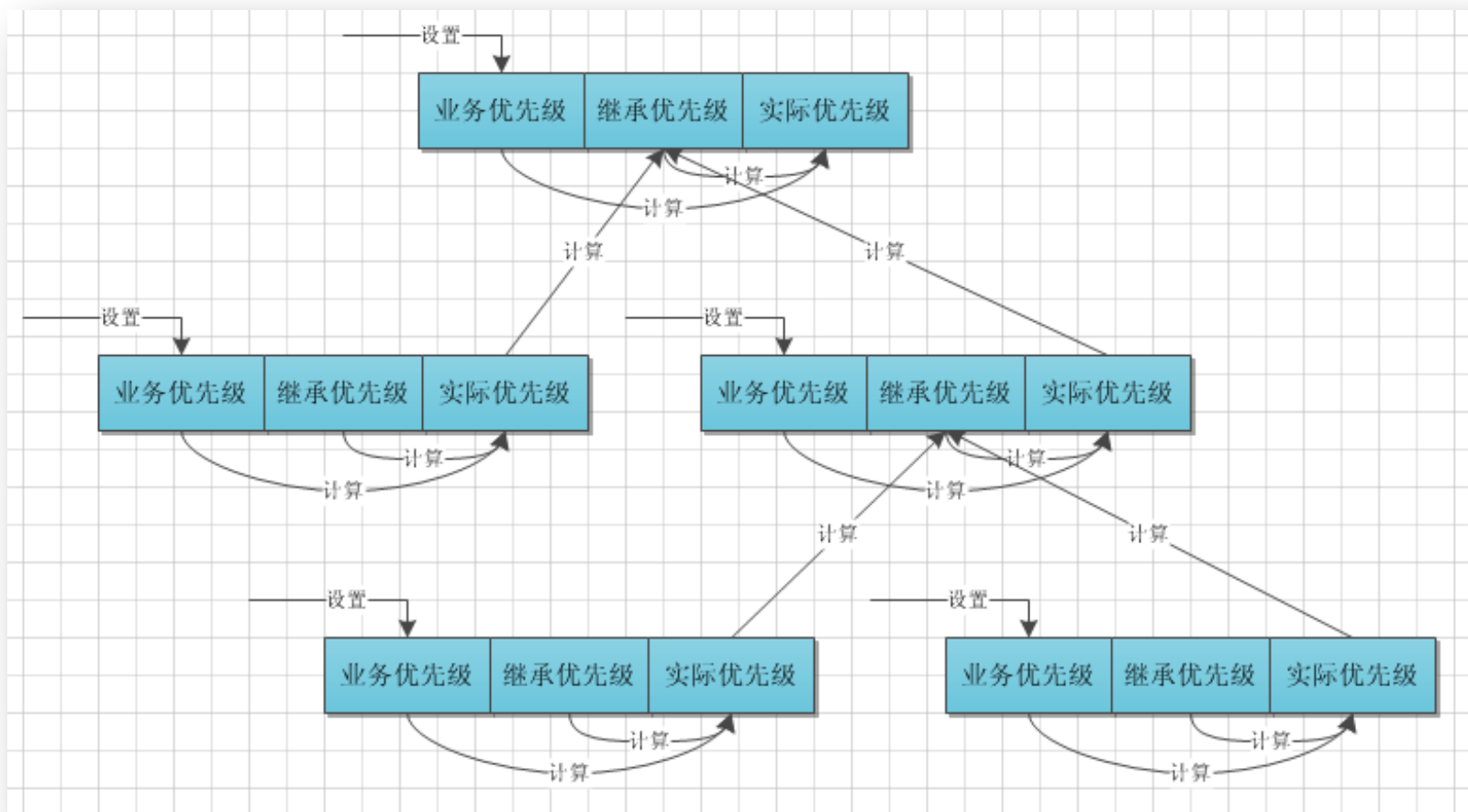
没有优先级



有优先级

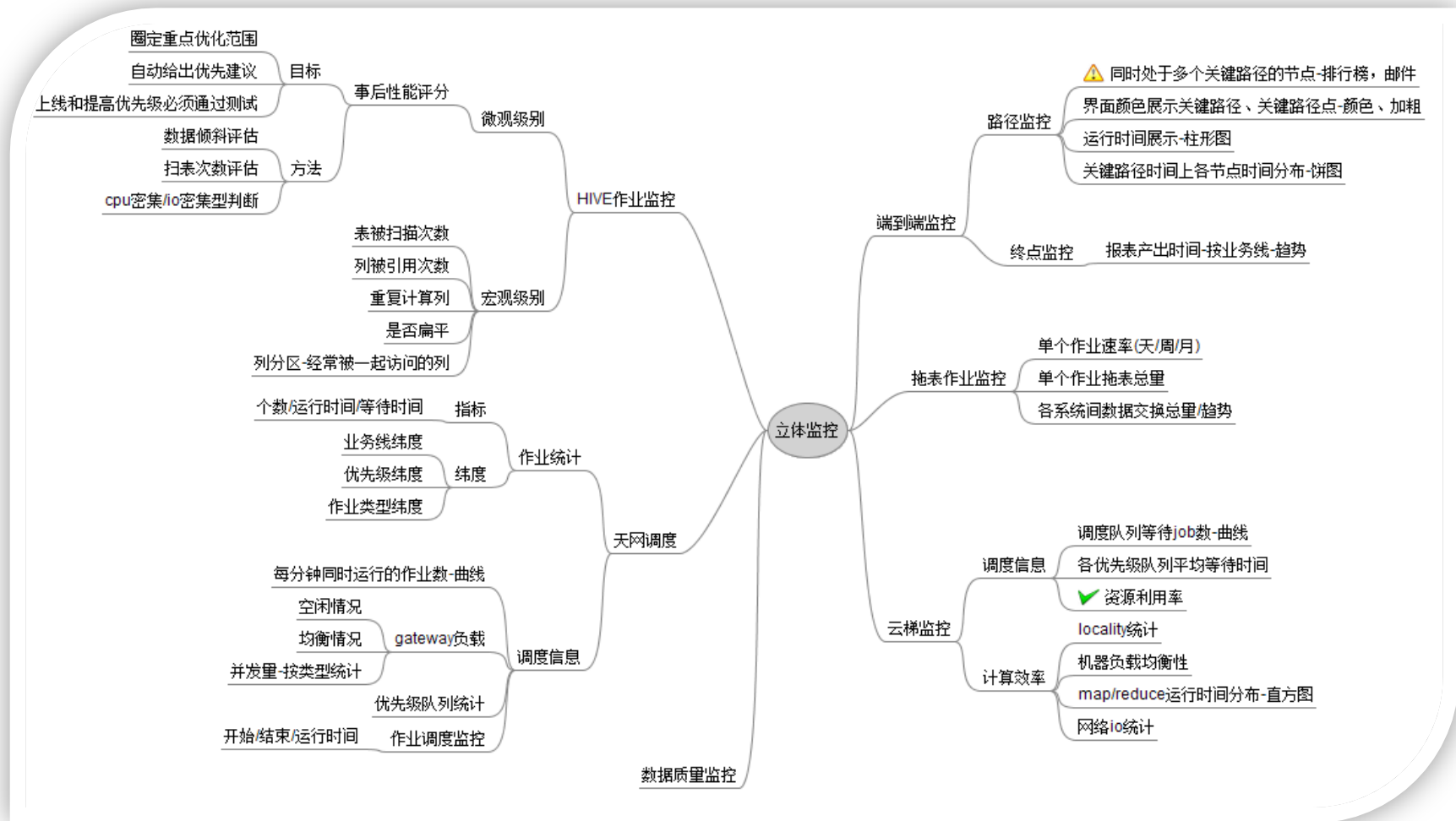


# 调度系统——优先级策略(DAG继承算法)

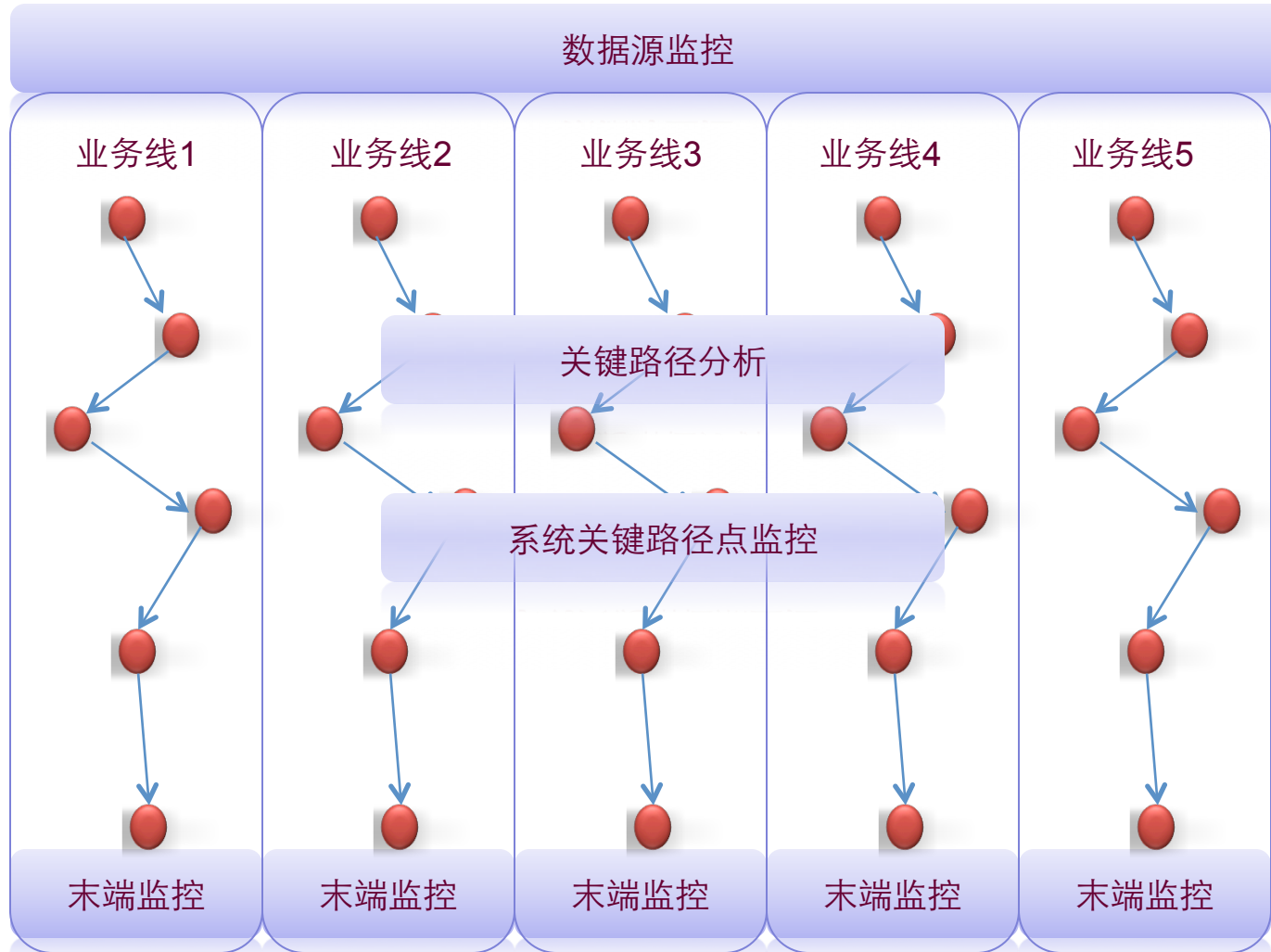




# 调度系统——监控全景



# 监控重点——系统边界和主干



## 目录

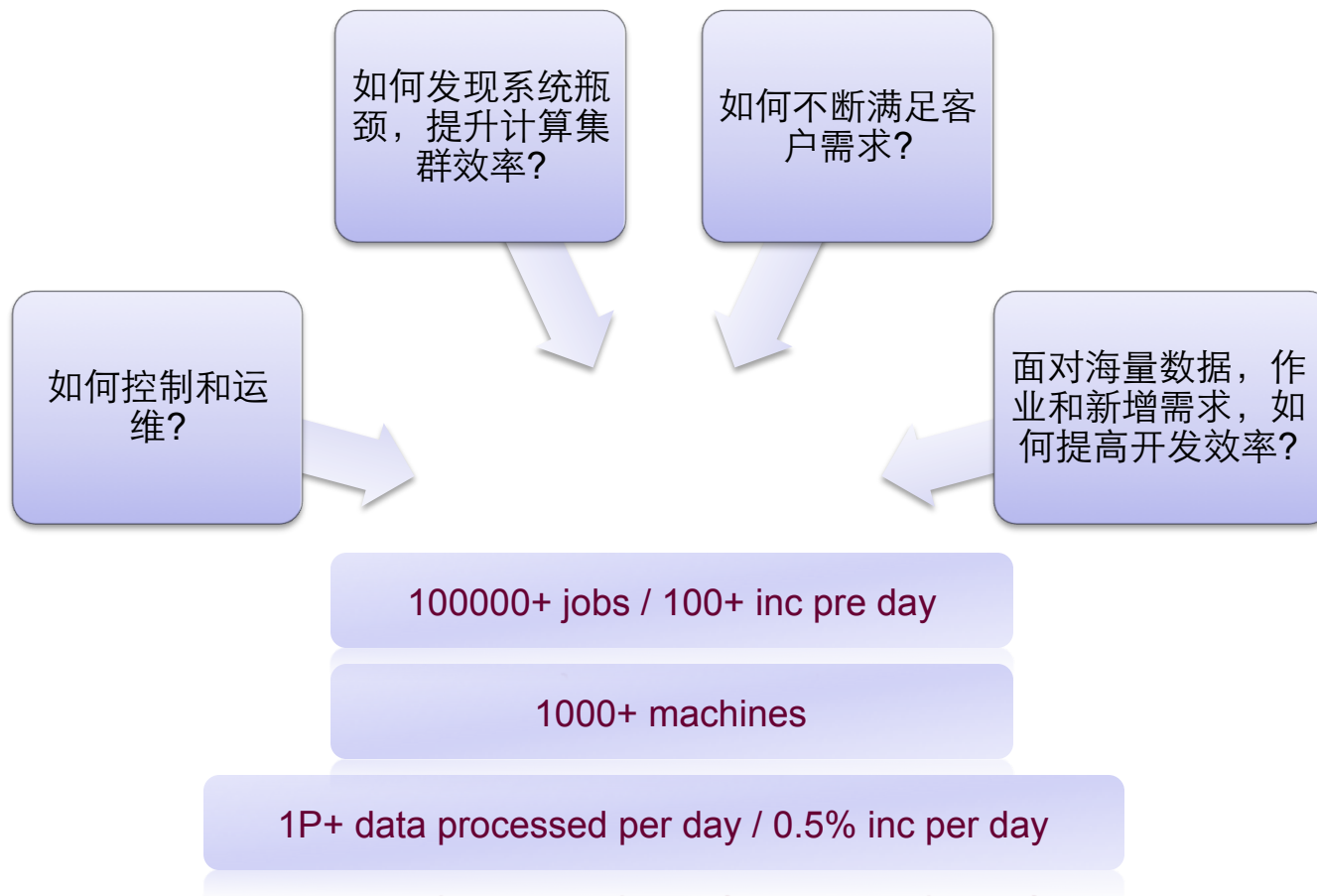
系统架构

数据同步方案

调度系统

元数据应用

## 问题

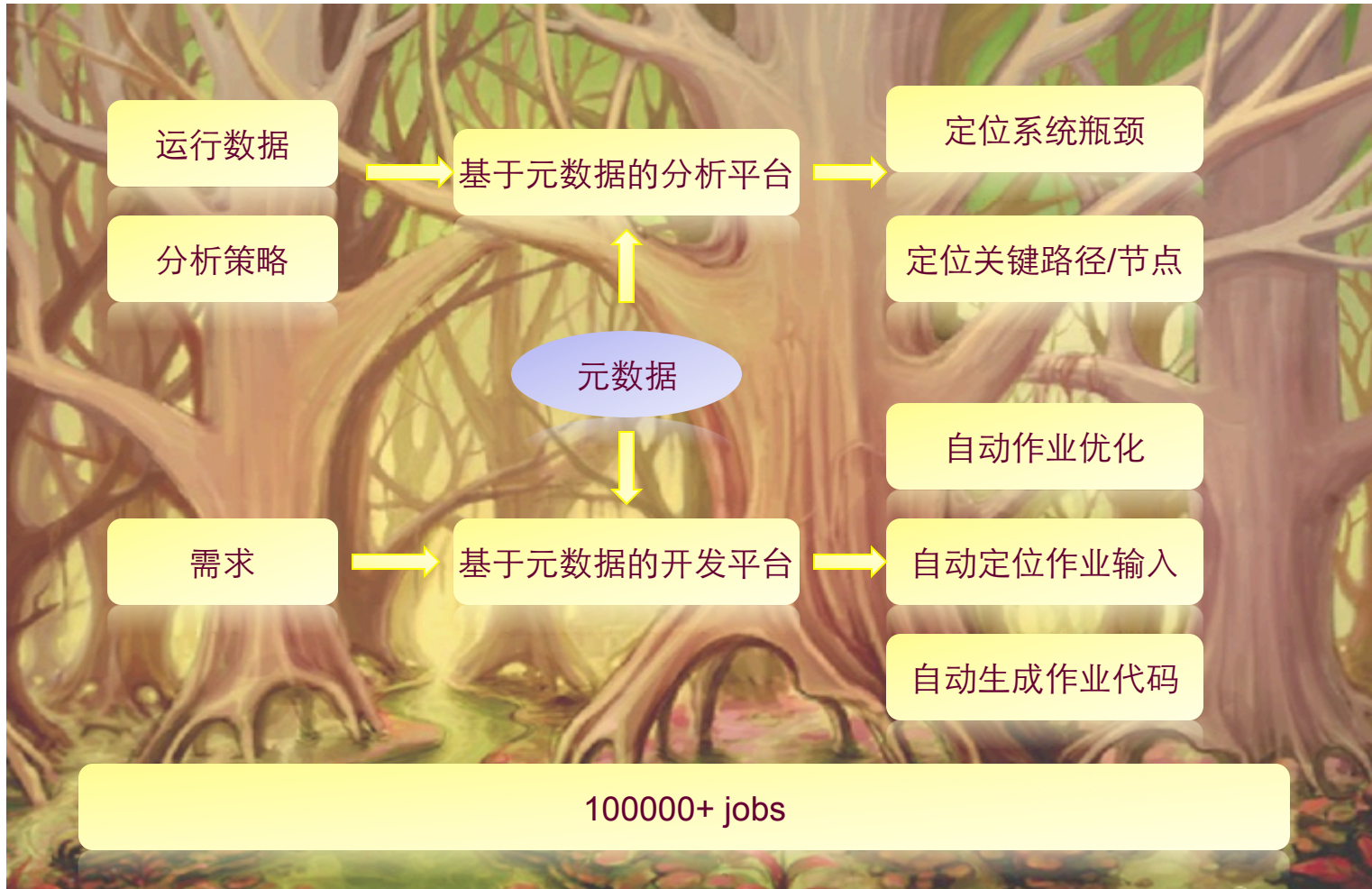


## 问题

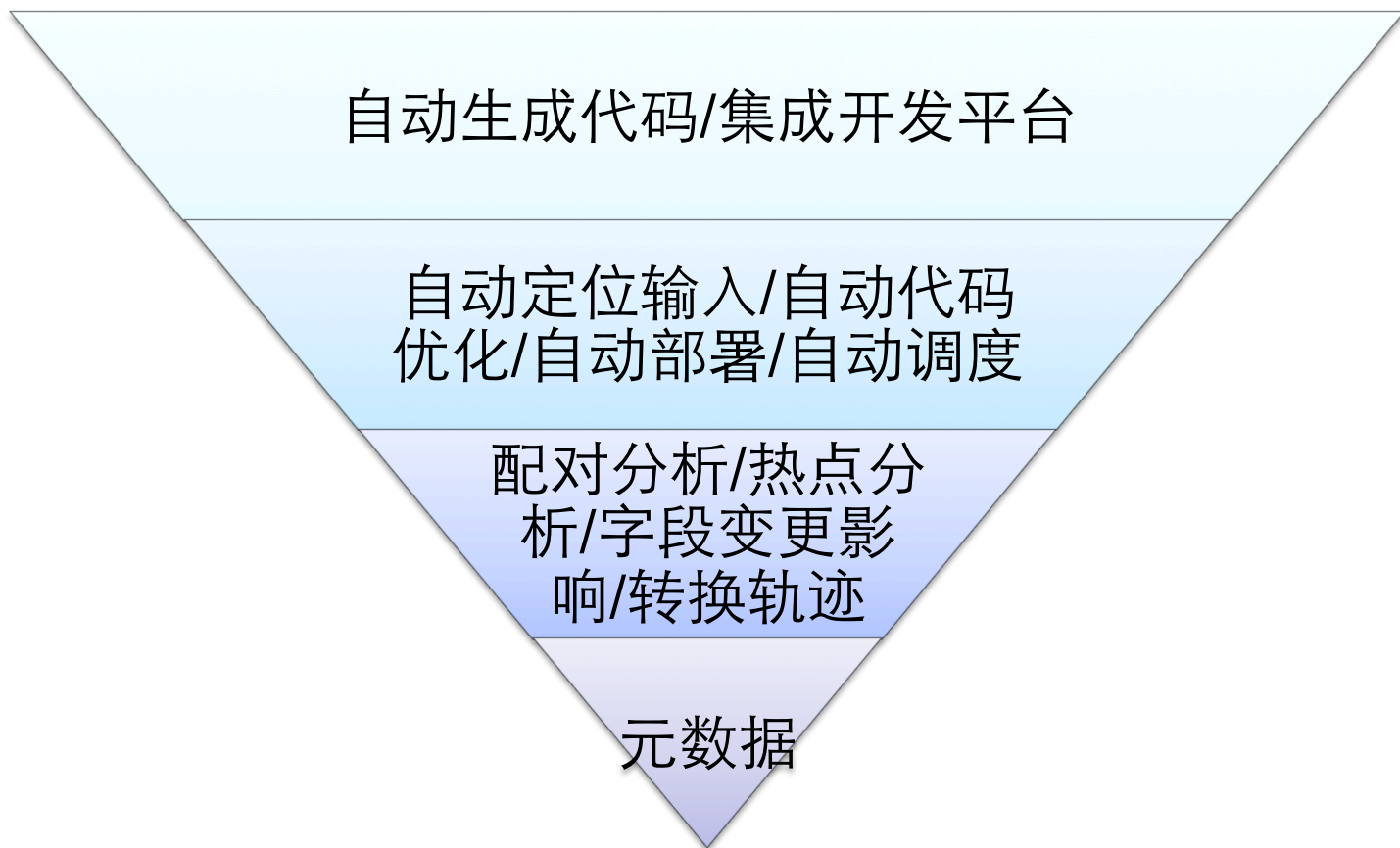
面对上面的问题  
靠经验丰富的架构师?  
还是靠智能的分析系统?



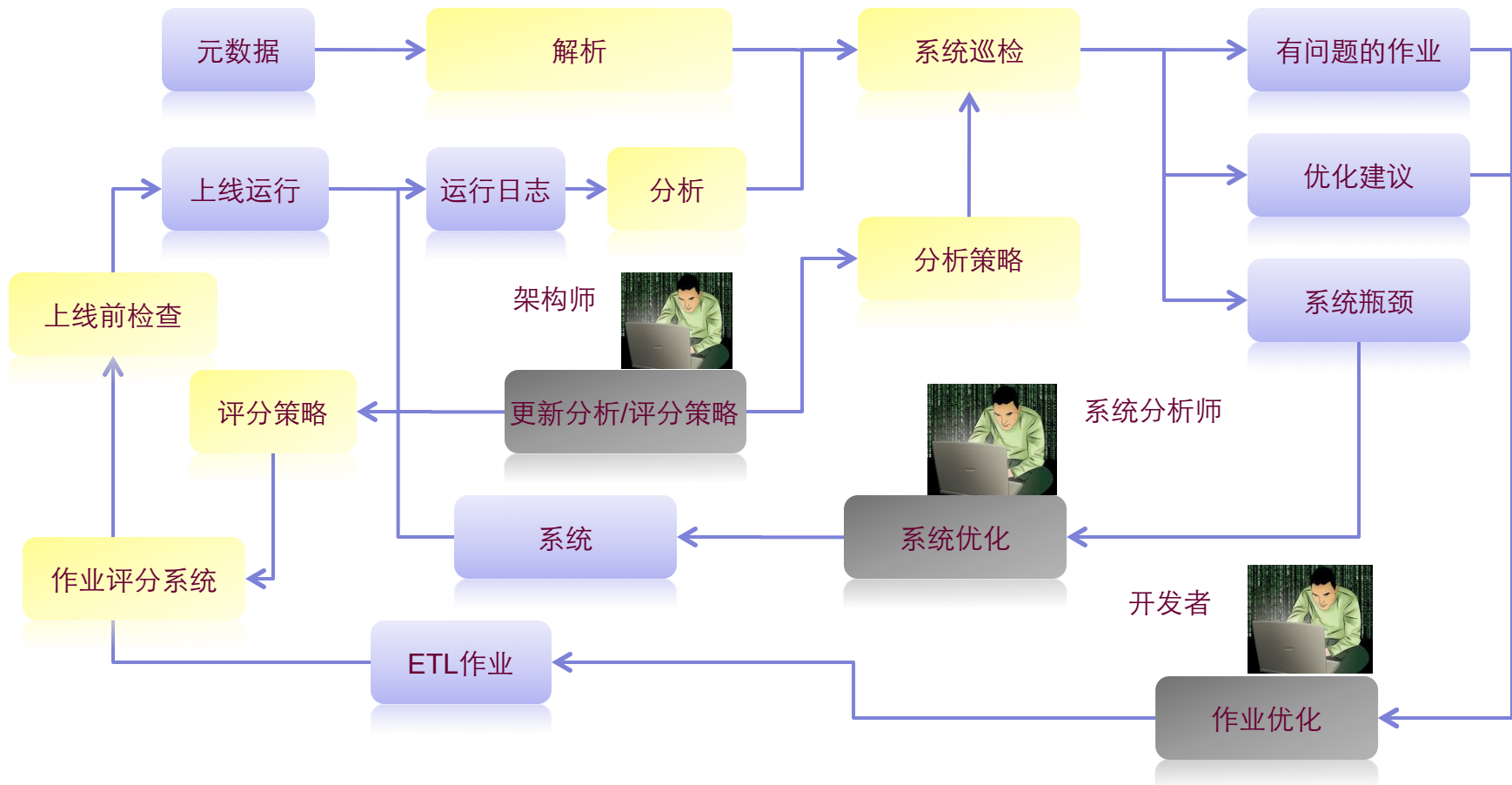
## 挖掘元数据金矿



## 基于元数据的开发平台

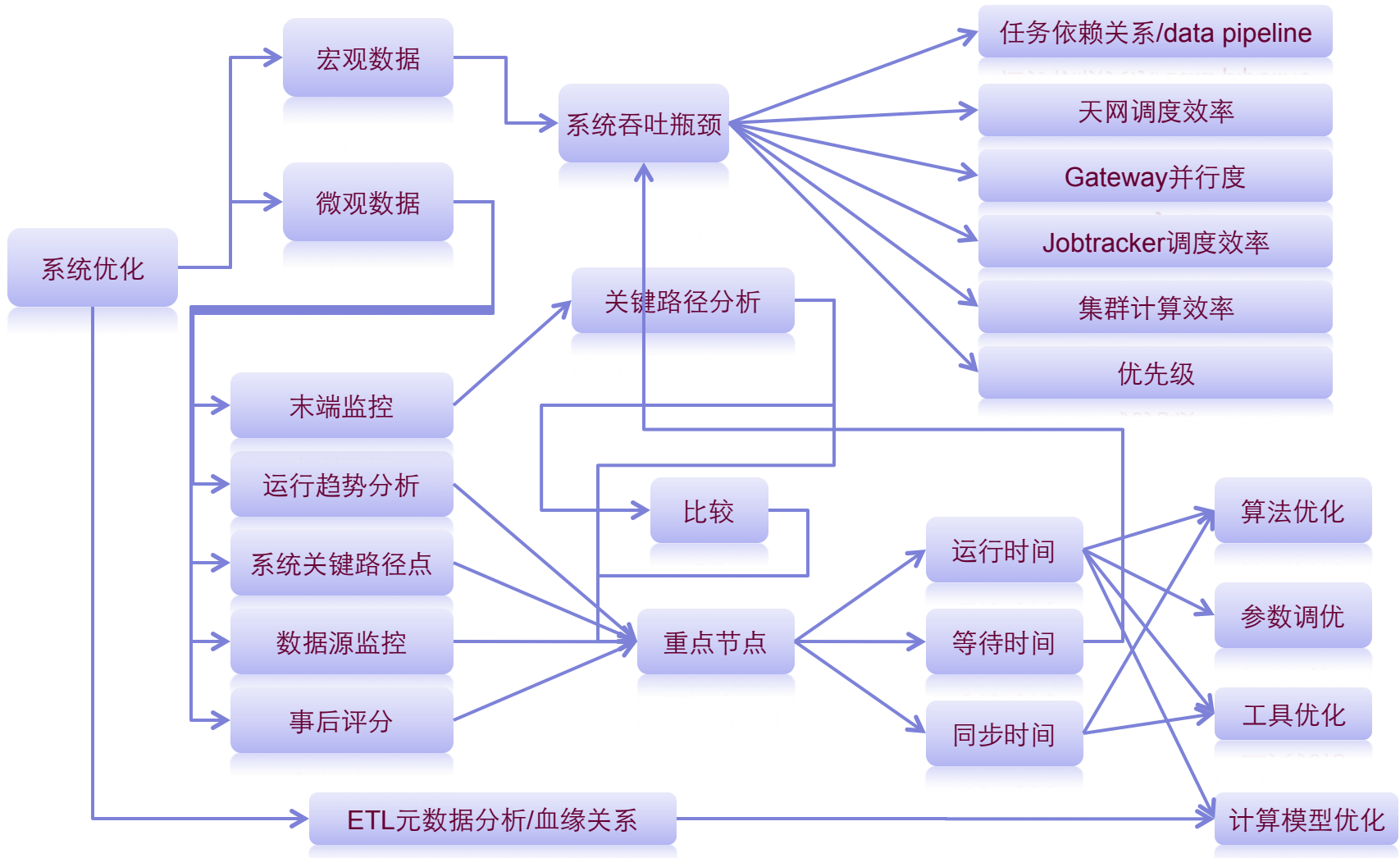


# 基于元数据的分析平台——运行分析系统

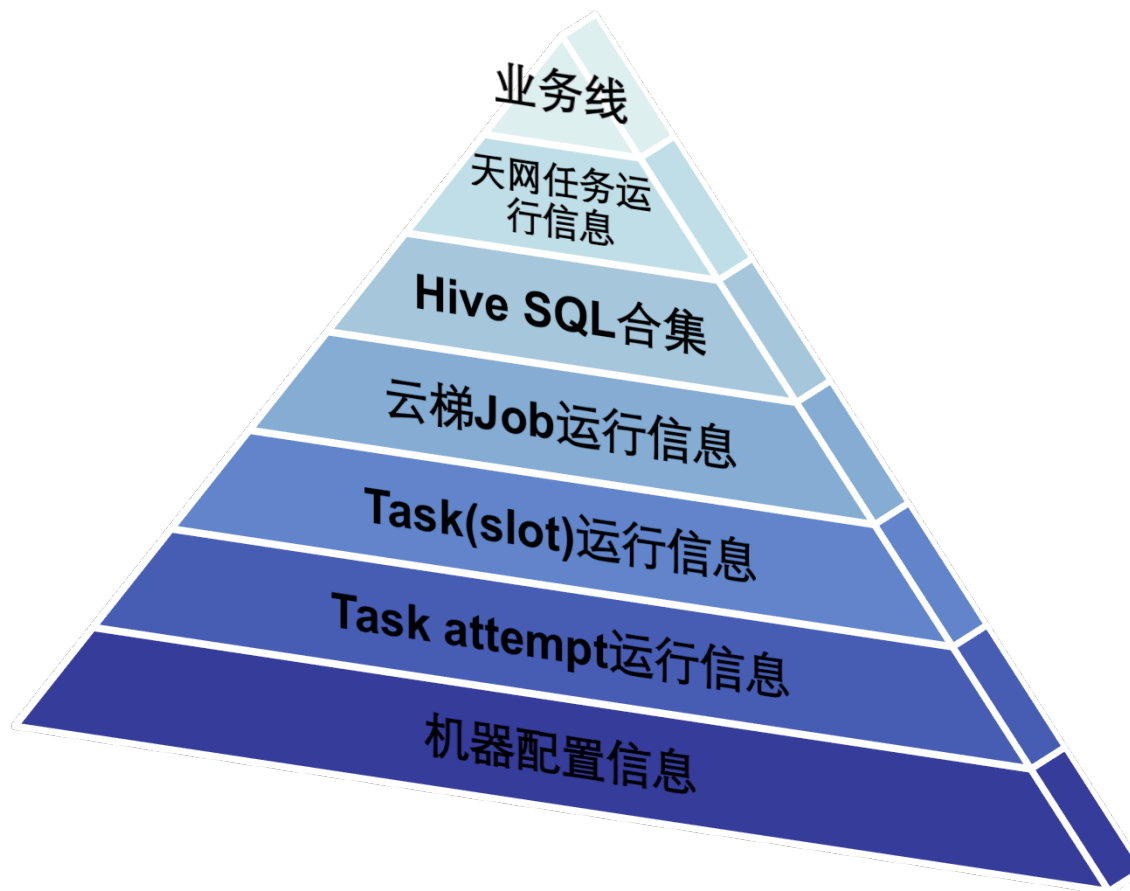




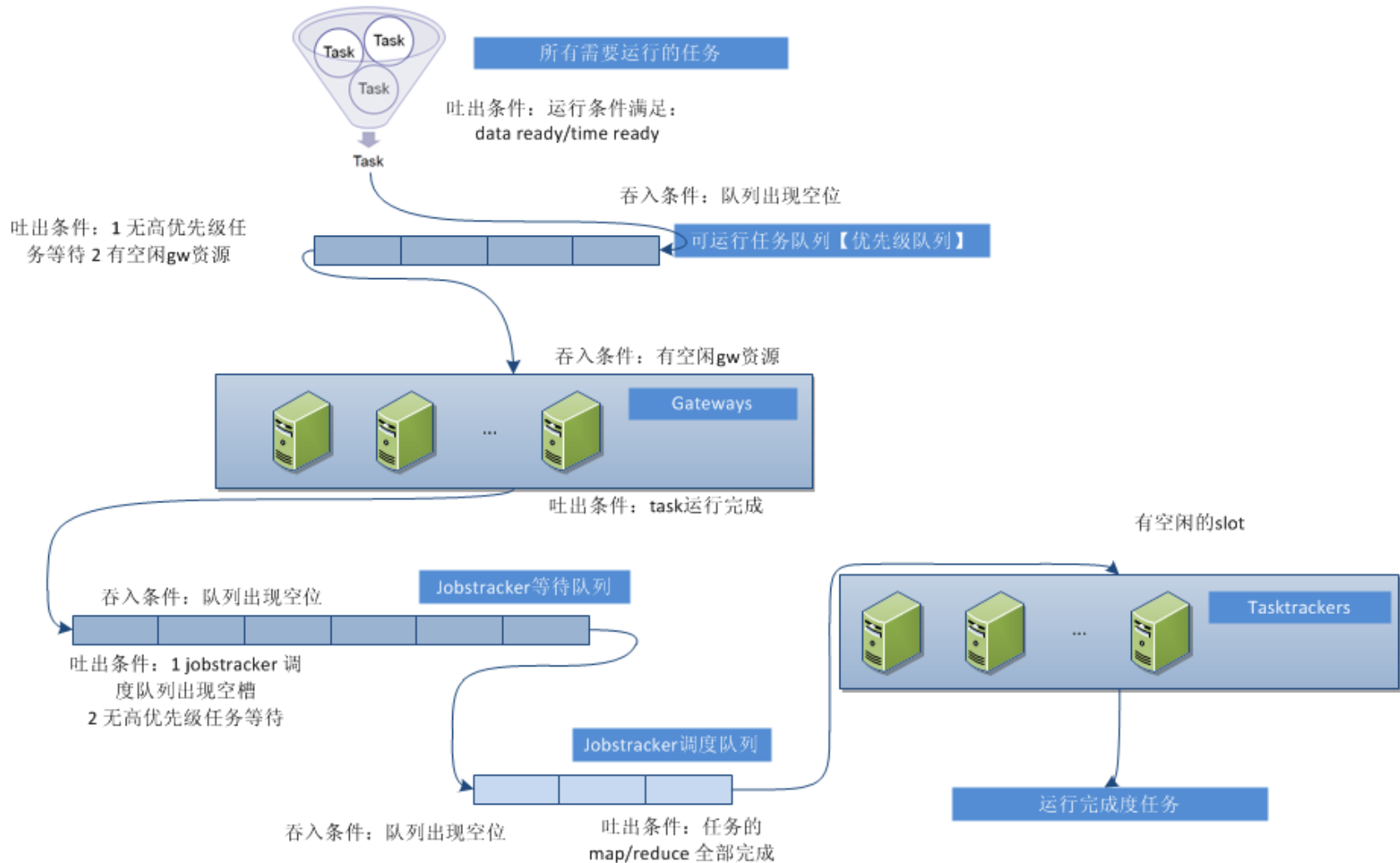
# 基于元数据的分析平台——分析策略概览



## 基于元数据的分析平台——运行数据收集



# 基于元数据的分析平台——宏观分析策略

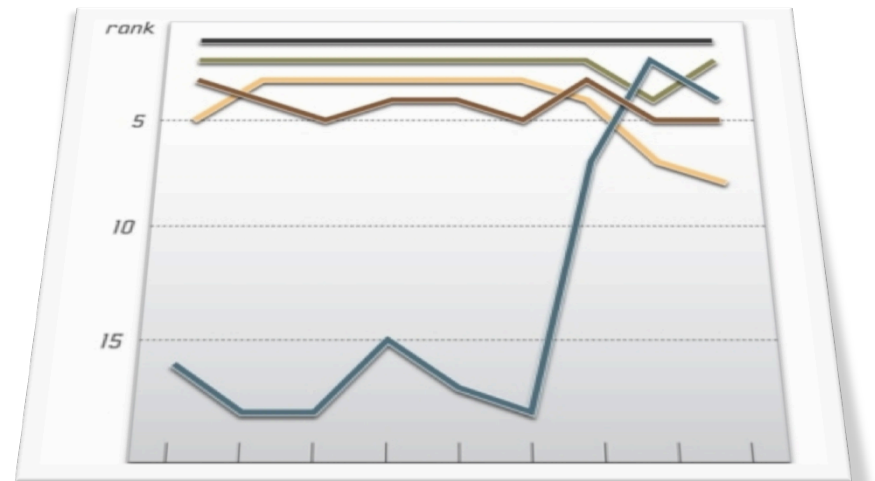


## 基于元数据的分析平台——定位系统瓶颈

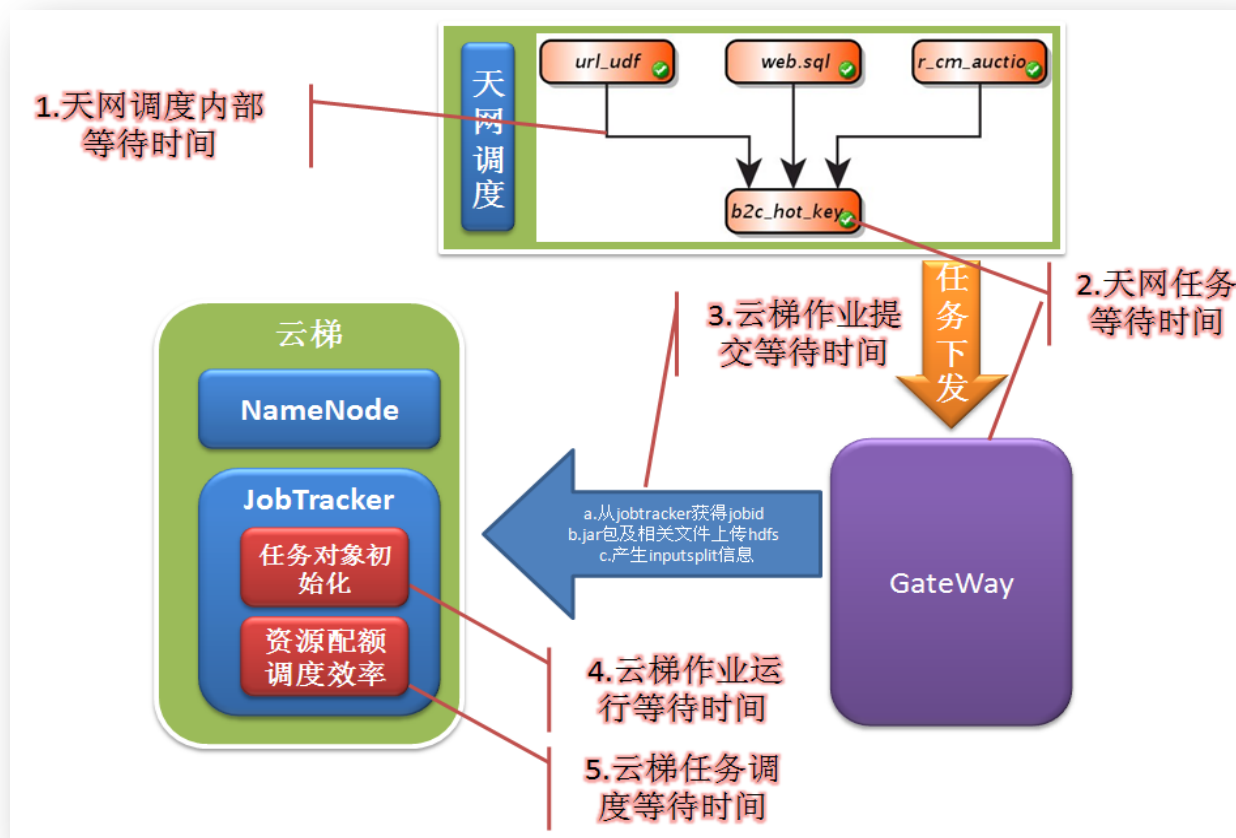
- 每个环节的吞吐能力都是动态变化的。
- 在某个特定时间区间内，整个系统的吞吐能力由吞吐能力最小的一个环节决定。
- 如果需要发现系统的短板，需要对每个环节的吞吐曲线绘制出来。
- 针对系统的短板进行重点优化。
- 对于吞吐能力抖动比较大的环节，需要在前面设置队列进行缓冲。

瓶颈定位方法：

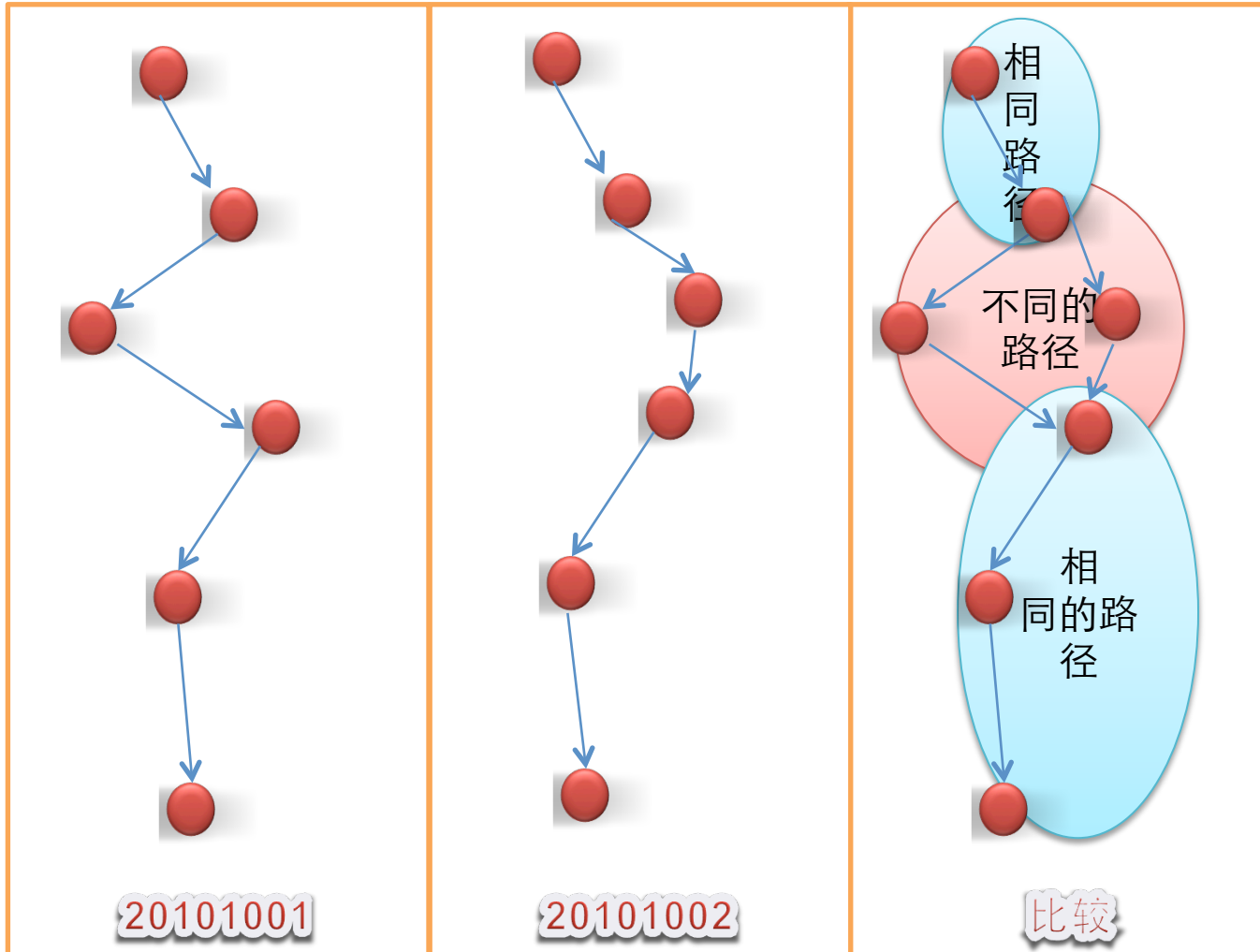
- 1 每个环节的吞吐曲线
- 2 两个环节之前缓冲队列的状态曲线
- 3 统一单位到task级别



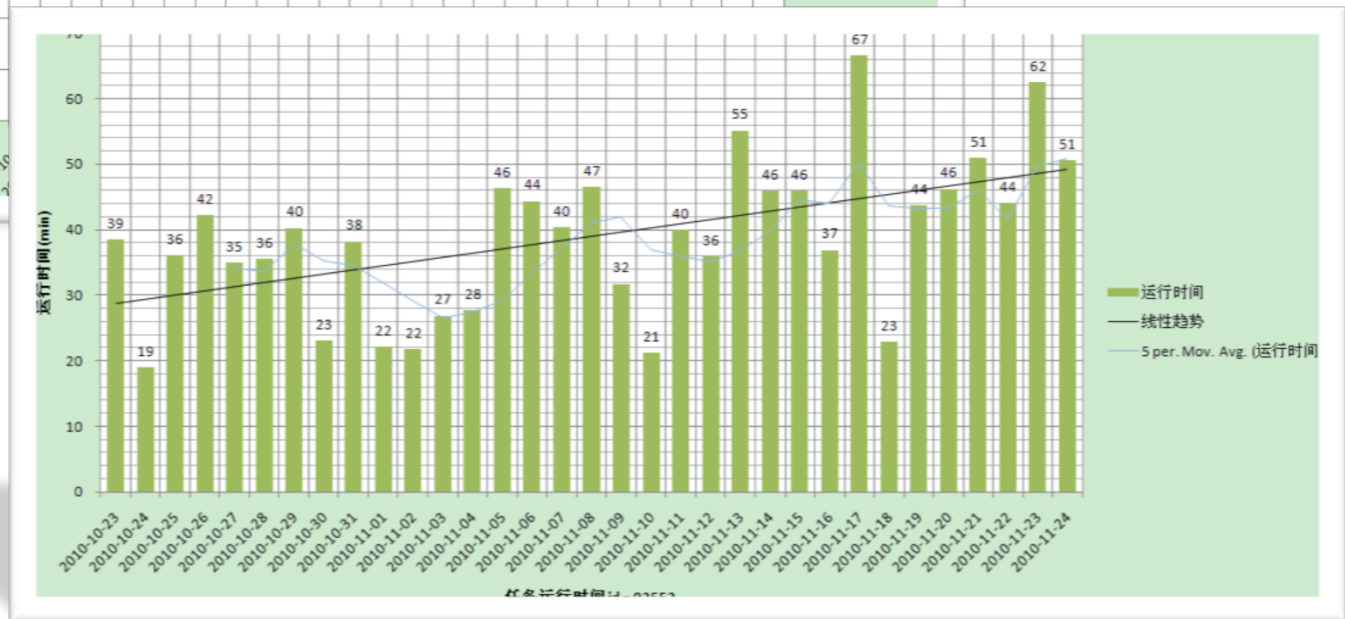
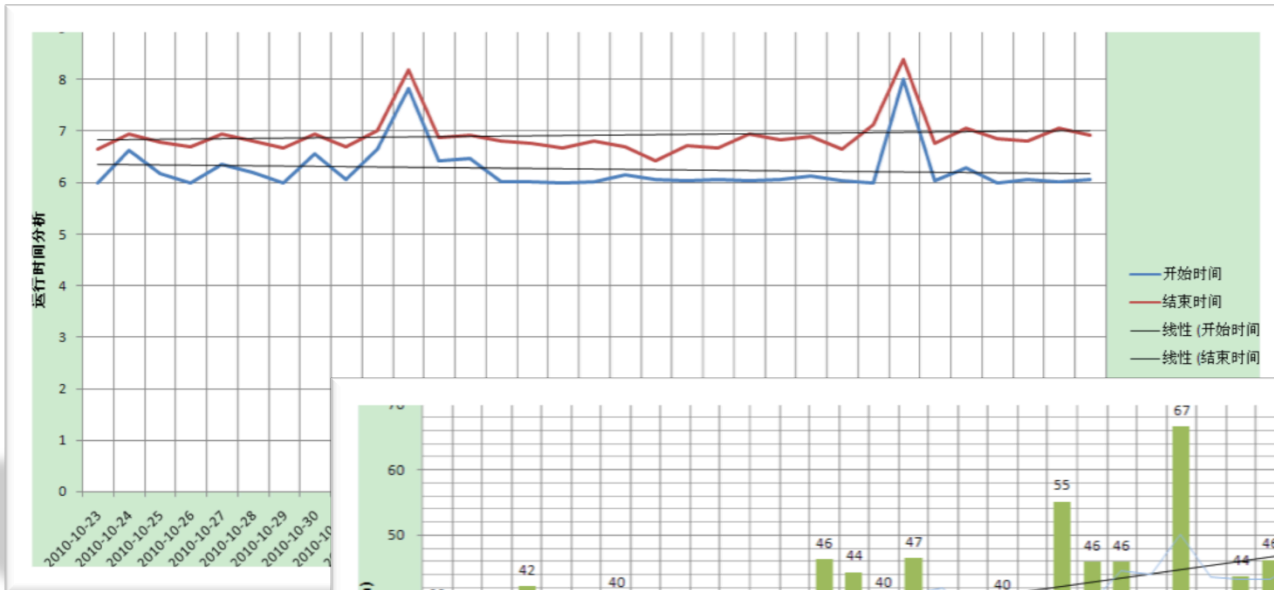
# 基于元数据的分析平台——任务等待时间



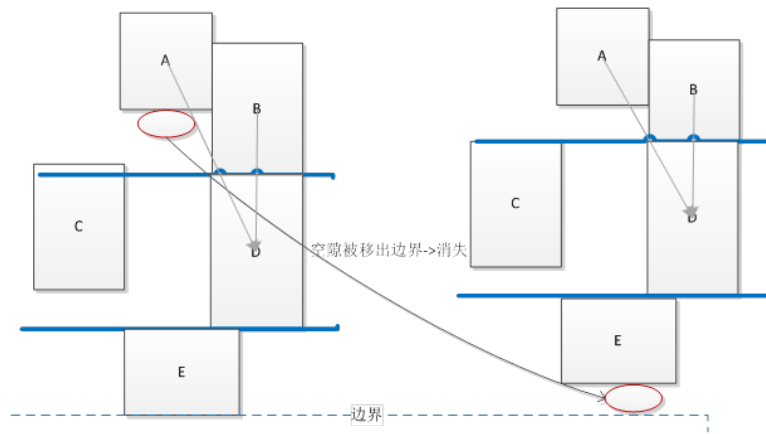
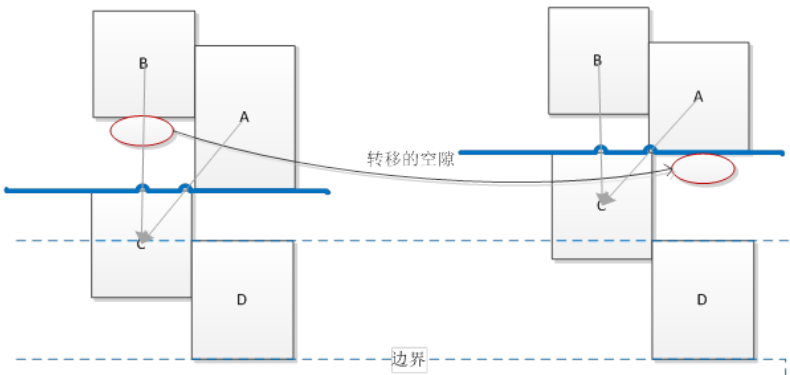
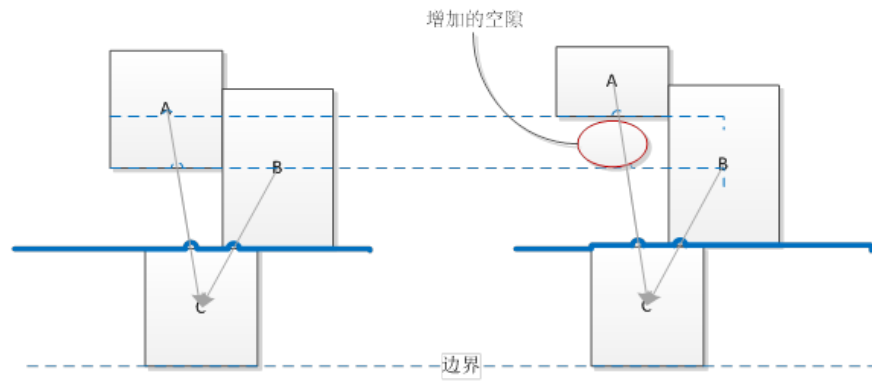
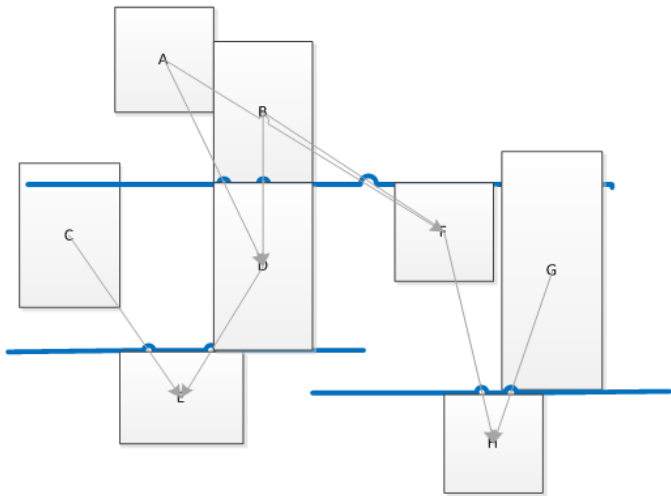
# 基于元数据的分析平台——关键路径分析



# 基于元数据的分析平台——任务运行趋势



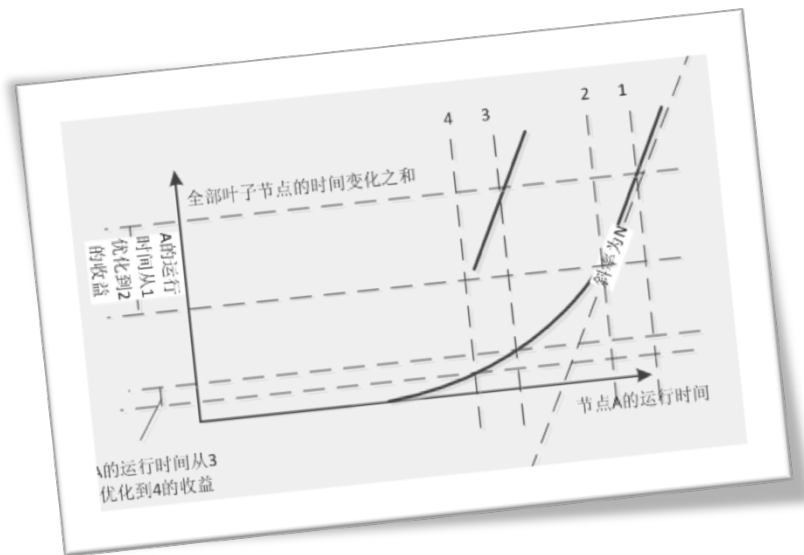
# 基于元数据的分析平台——调度紧凑度分析





# 基于元数据的分析平台——最值得优化的任务

从关键路径的角度考虑，任务A运行时间变化，对系统影响可以用下图中阴影的面积来计算，它取决于下面几个因素：



A 任务的当前运行时间——决定了当前位置的曲线斜率。

B 任务在几个叶子节点的关键路径上——决定了当前位置的曲线斜率。

C 其它关联任务的运行时间——决定A何时会离开某个关键路径，也就是决定了斜率的斜率。

最值得优化的任务：

以下三项评分中综合评分最高

- 1 运行时间长
- 2 同时处于多个关键路径
- 3 孔隙度大

谢谢大家

Q&A