

携程大数据平台与服务化实践

携程-----杨晓青

xq.yang@Ctrip.com

自我介绍

- 2014.8 ~ 携程-大数据平台，负责携程数据基础架构。
- 2012.4~2014.8 阿里巴巴-CDO-数据平台部，负责从无到有搭建galaxy流计算平台，集群规模到2K台。
- 2009.2~2012.4 腾讯-数据平台部，腾讯TDW平台从无到有建设。

目录

大数据架构

数据采集与同步

数据调度平台

数据分析平台

目录

大数据架构

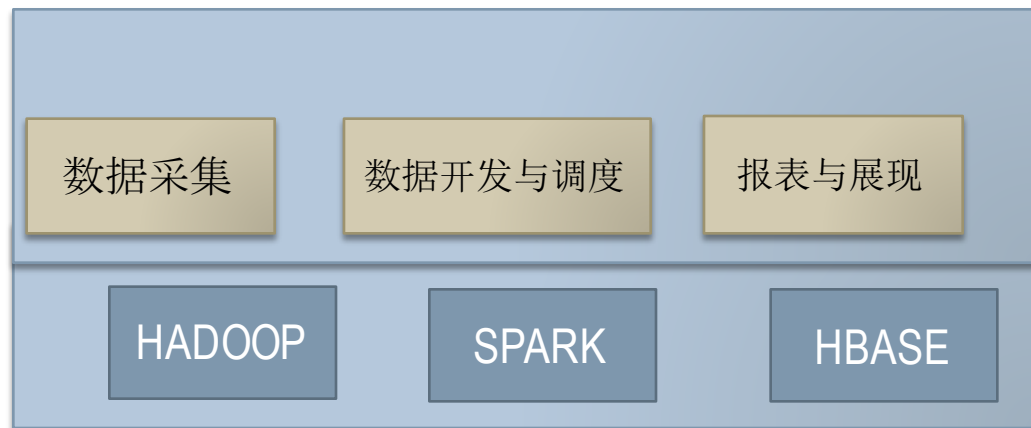
数据采集与同步

数据调度平台

数据分析平台

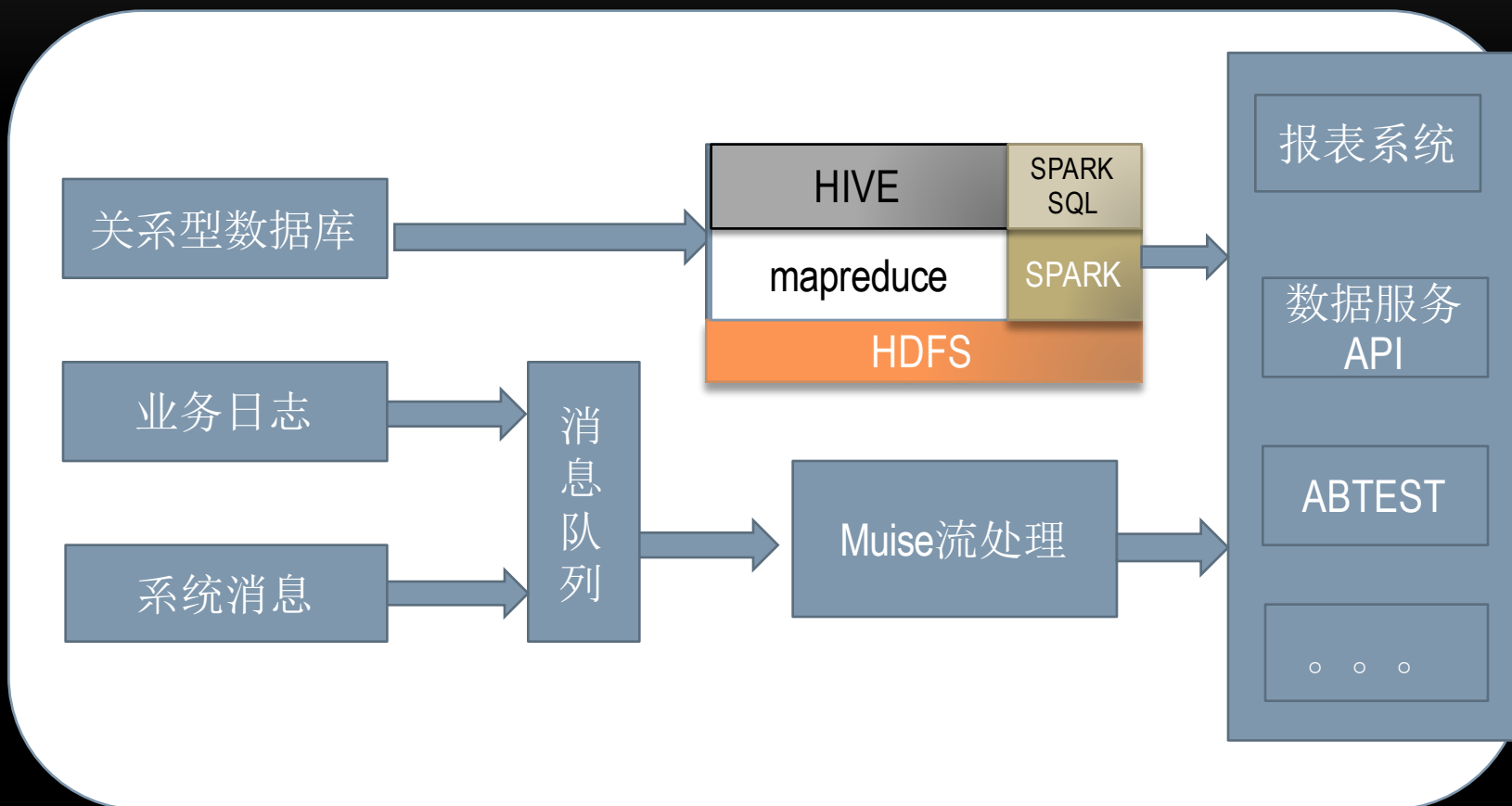
数据平台

数据产品



数据平台

数据处理流程



大数据架构

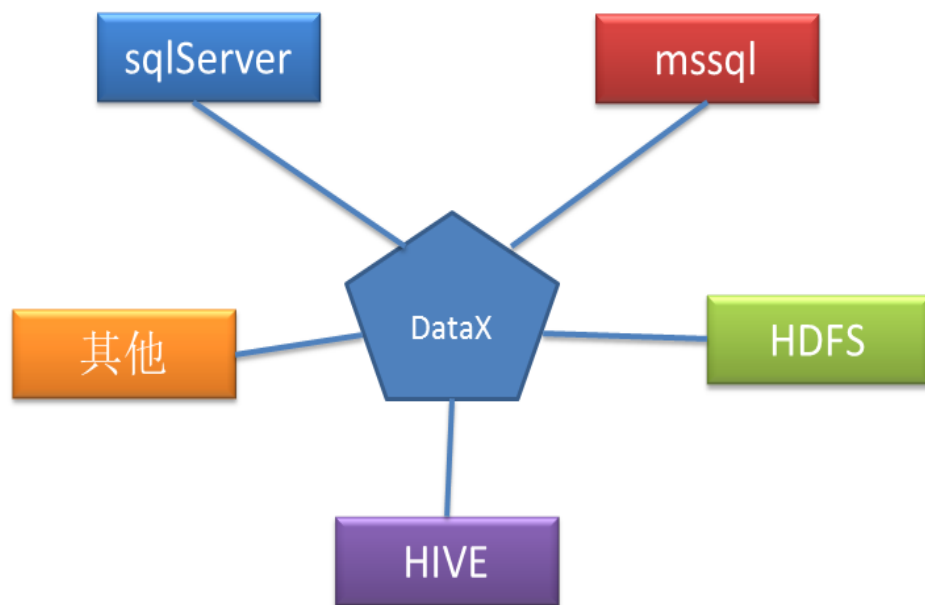
数据采集与同步

数据调度平台

数据分析平台

数据同步服务

- 1.支持SQL Server, My SQL 和Hive等携程主流数据库间任意同步（解决多种同步工具混用，效率难以保证，维护负责问题）
- 2.支持全量与增量数据同步模式
- 3.可线性扩展
- 4.以SASS（SOFTWARE as a service）提供服务



交易数据同步配置

datax配置工具

源

源 (【sqlserver|mysql】)

sqlserver

源数据库域名 *

请填写域名，不要使用IP

源数据库端口 *

源数据库名 *

源数据库用户名 *

源数据库密码 *

源数据库表名 (逗号分隔) *

源sql查询

querysplit 

目标

目标

hive

目标数据库名 *

目标数据库表名 (逗号分隔) *

hive导入模式

全量不分区

hive文件类型

TXT

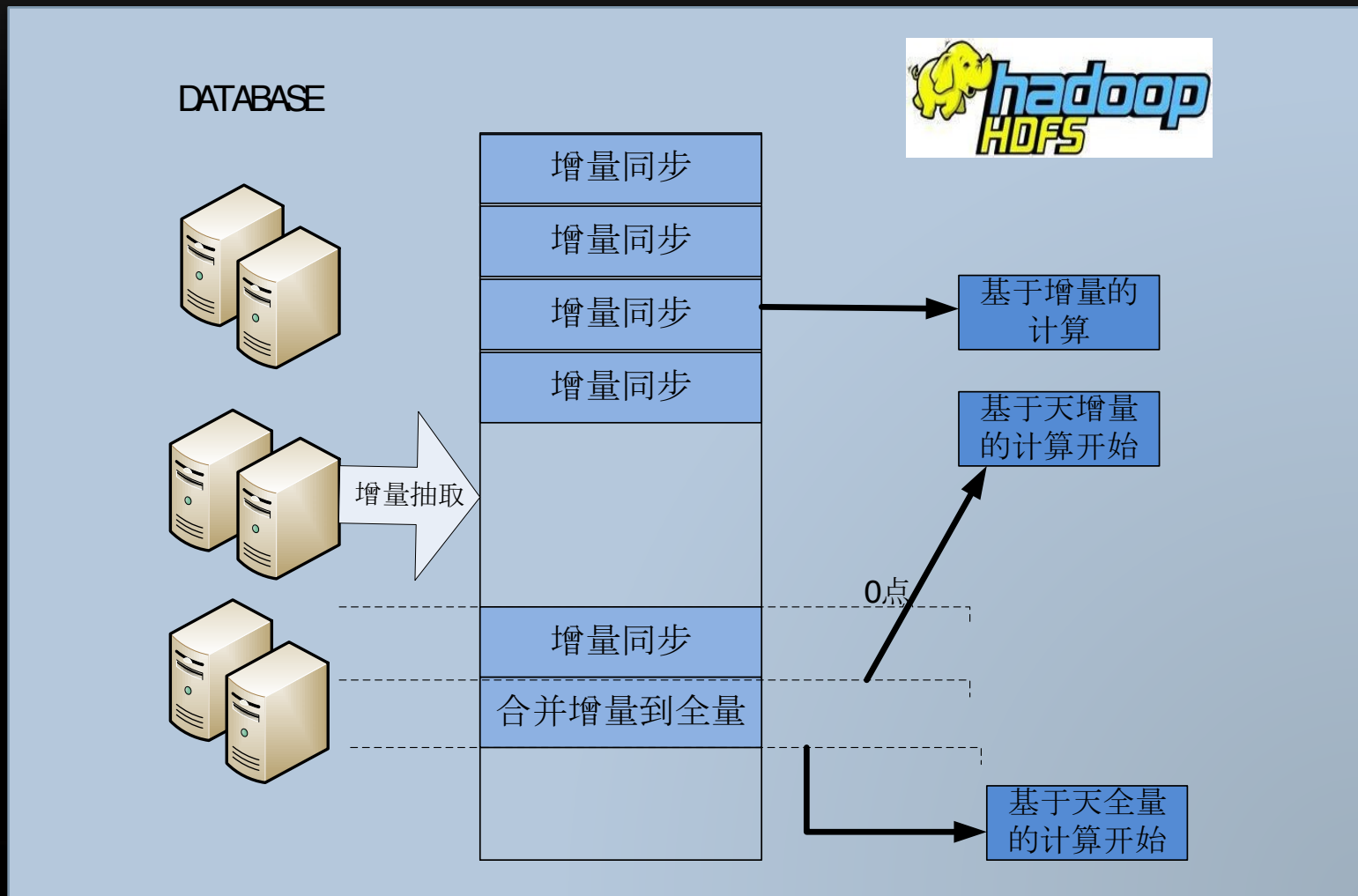
任务并发数

3

生成脚本

提示：需要拷贝生成的代码到编辑器中

准实时数据采集



运营数据

业务指标（日）	8月
同步表数据	8000+张
数据量	5T+
可用度	100%
同步行数	80亿条

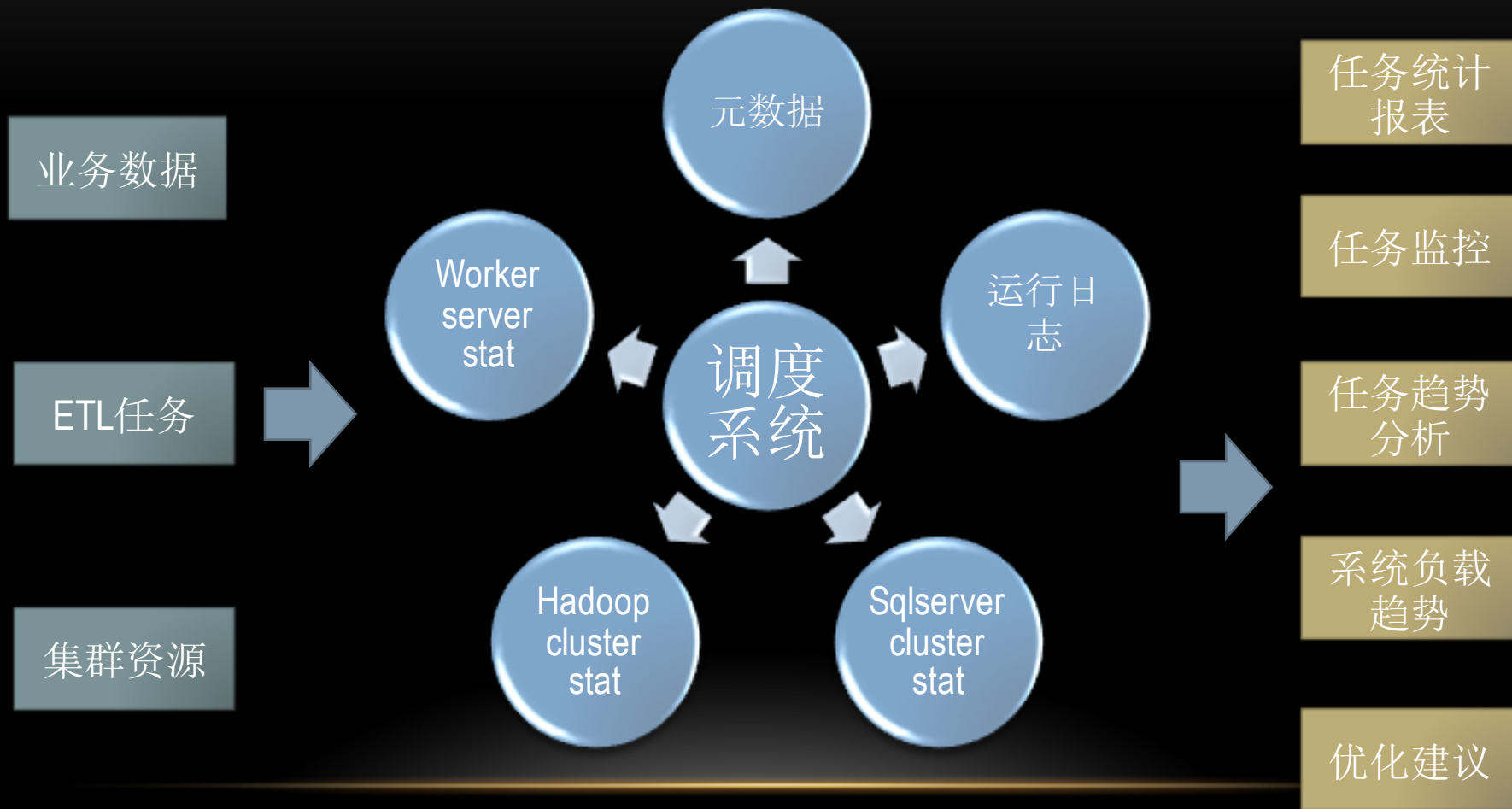
大数据架构

数据采集与同步

数据调度平台

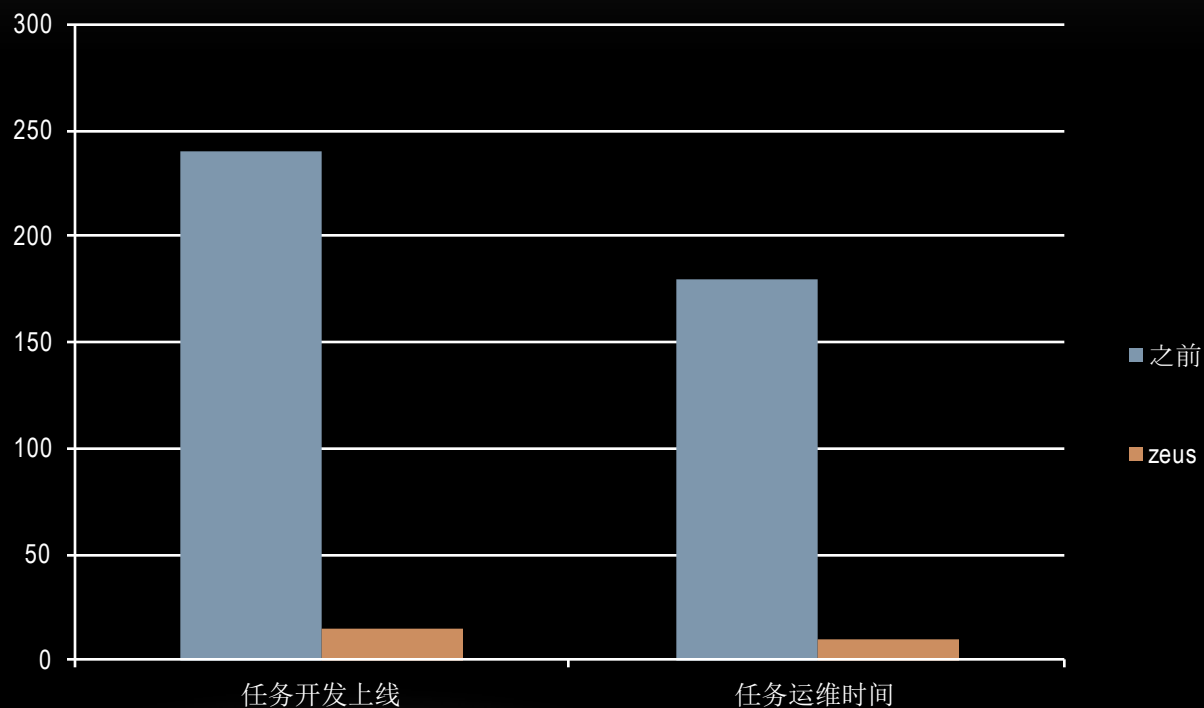
数据分析平台

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
- 异构平台支持
- Web UI



系统模块

开发IDE

- Shell, hivesql
语法高亮
- 调试运行
- 元数据查看
- 数据预览

workflow

- 依赖管理
- 时间触发
- 事件触发

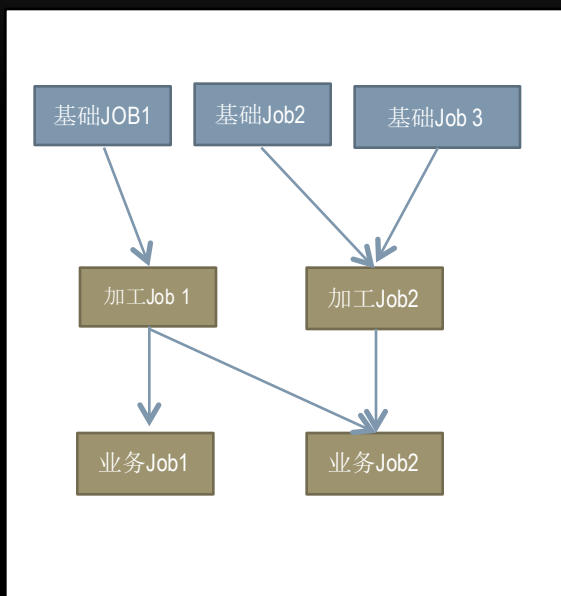
运维管理

- 手动重跑
- 依赖重跑
- 查看任务状
- 查看错误日
志
- 失败, 超时
告警
- 邮件, 短信,
电话告警
- 每日任务报
告

智能管理

- 优先级
- 负载均衡
- 过载保护
- 分组调度

调度系统---任务触发方式



依赖触发

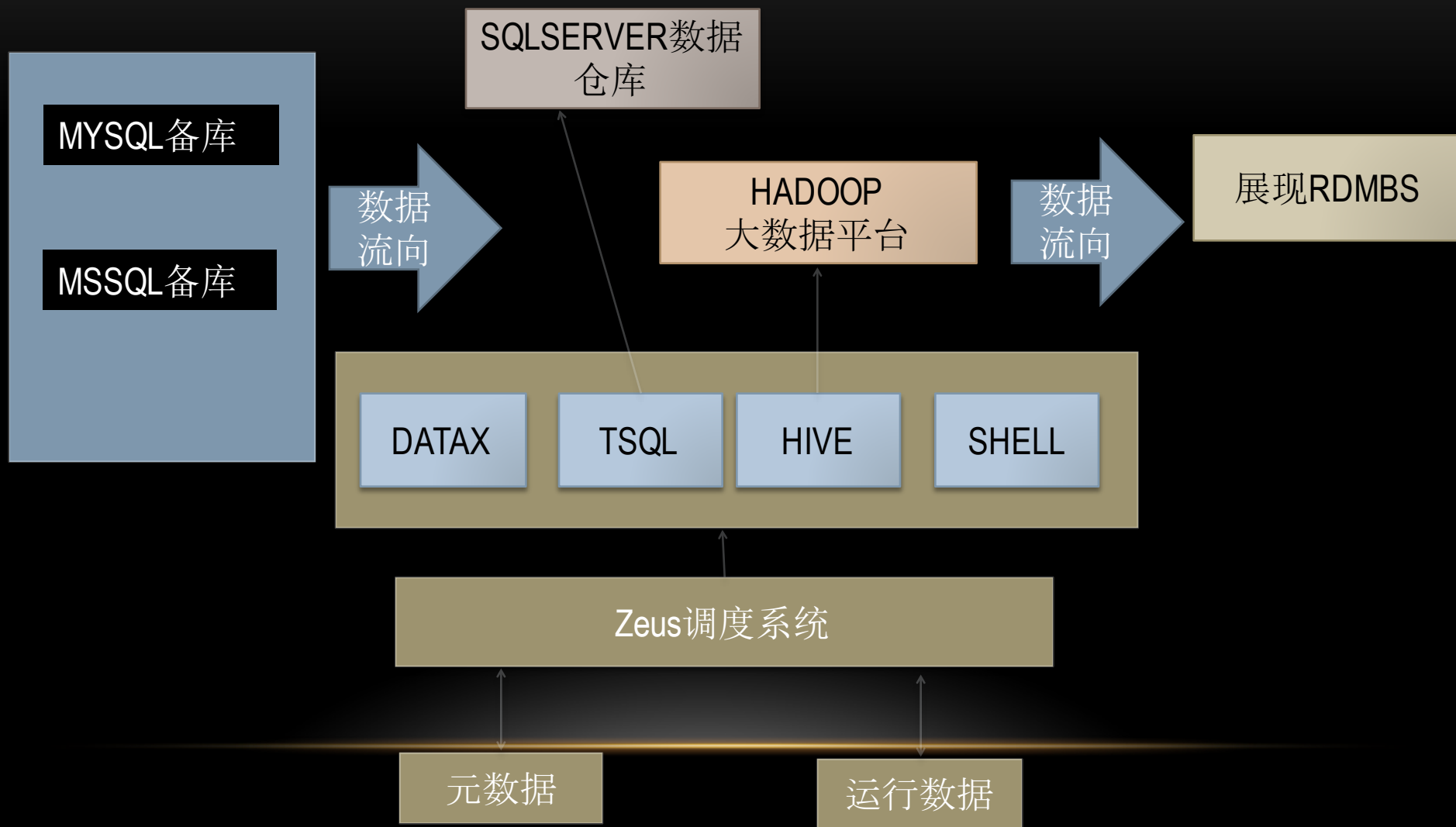


循环触发
(敬请期待)

定时调度

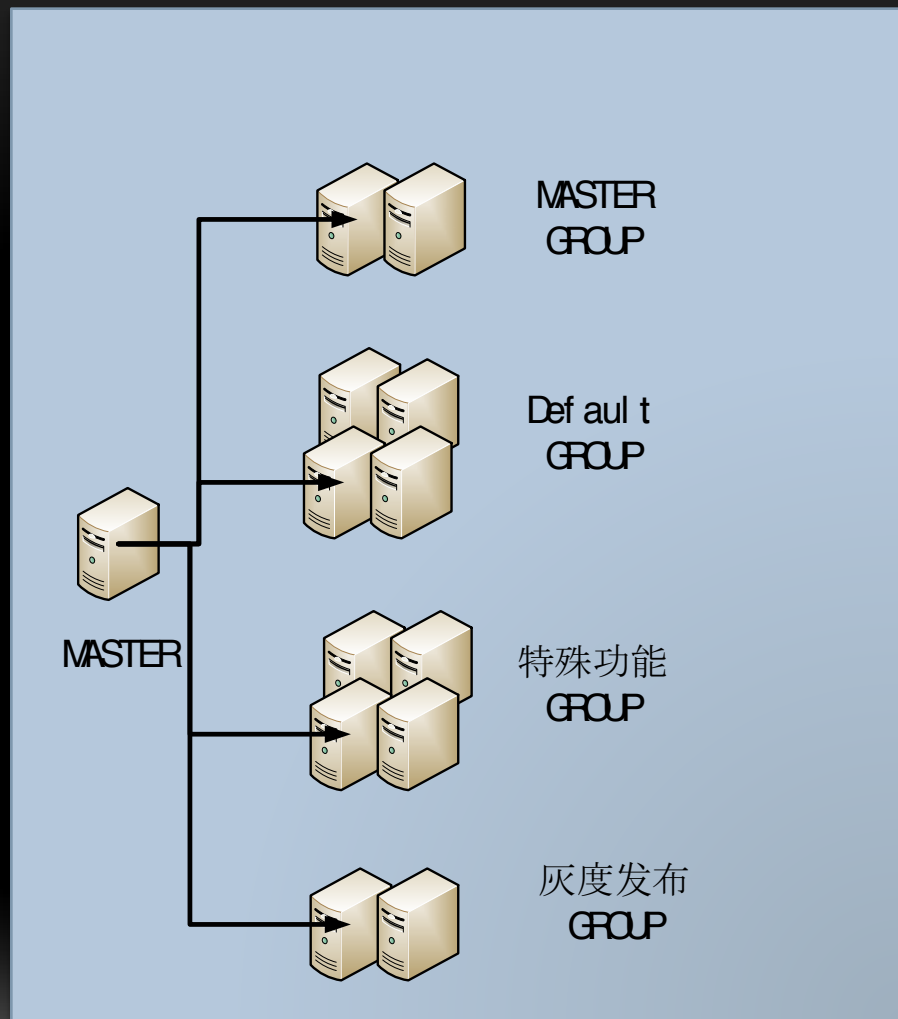
The screenshot shows a software interface for scheduling. A modal dialog box titled '构造定时表达式' (Construct cron expression) is open. It contains input fields for time components: '分' (Minute) with value '0', '时' (Hour) with value '3', '天' (Day) with value '*', '月' (Month) with value '*', and '周' (Week) with value '?'. There is a '确认' (Confirm) button at the bottom right. In the background, a form is partially visible with fields for '调度类型' (Scheduled type) set to '定时调度' (Cron scheduling), '定时表达式' (Cron expression) set to '0 0 3 * * ?', 'host组' (Host group), '脚本是' (Script is), and '预计时' (Estimated time).

调度---驱动数据系统方式



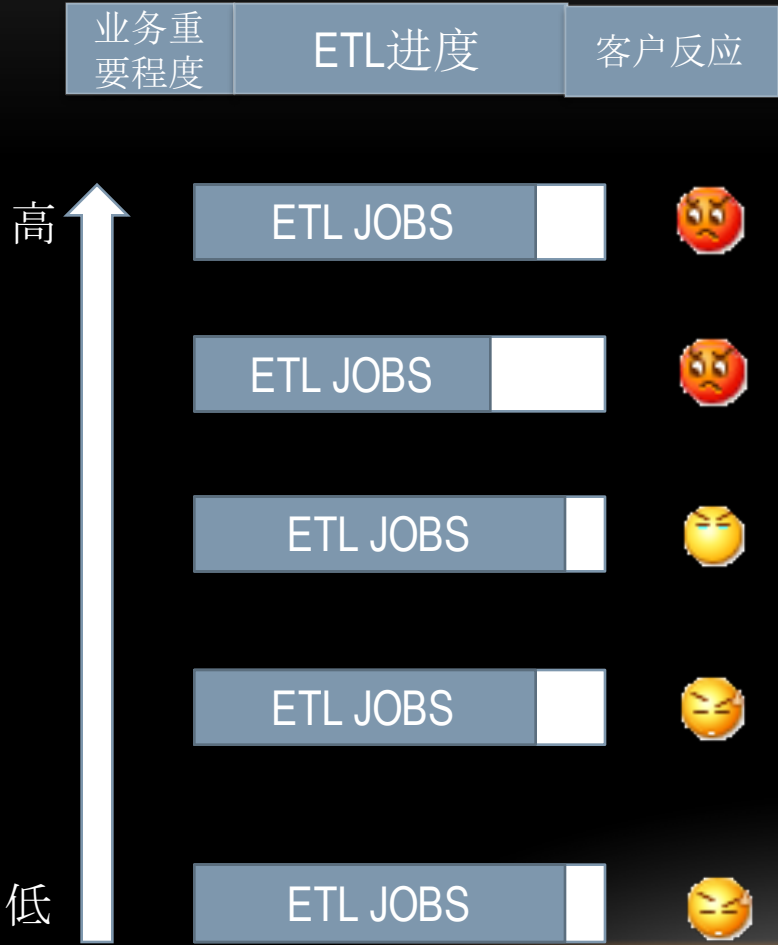
调度---资源分组

- 基础功能隔离
- R集群化
- Python集群化
- Zeus HA
- Zeus灰度发布
- 特殊业务线隔离



ZEUS---优先级

没有优先级



有优先级



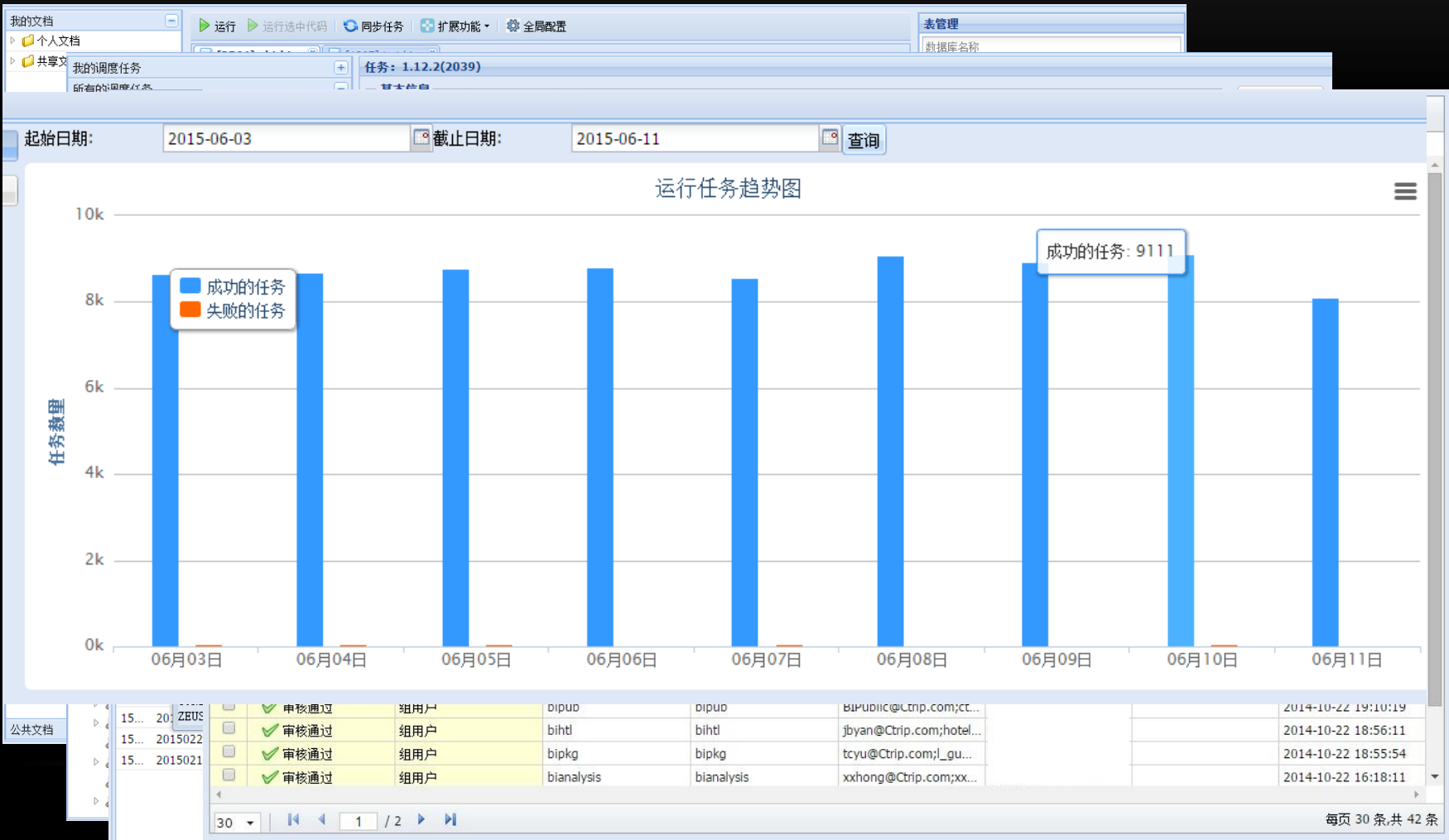
应用情况

指标	
Zeus JOB数（日）	2w
集群规模	33台

开源地址：

<https://github.com/ctripcorp/dataworks-zeus>

功能演示



系

host组信息:

组id	名称	描述	host	CurrentPosition
1	default	默认组	8	
2	master	master抢占组	0	
3		单个节点	0	
4		单个节点	0	
5		单个节点	0	
6		单个节点	0	
7		R节点	0	
8		Python节点（地面专用）	0	

大数据架构

数据采集与同步

数据调度平台

数据分析平台

分析数据粒度层次

High Level
Aggregation

非常高层次，比
如：全网交易额

Analysis
Query

中等层次，比如：某
业务线的app交易额

Drill Down to
Detail

订单级别

Low Level
Aggregation

Uid 级别

Transaction
Level

事务级别

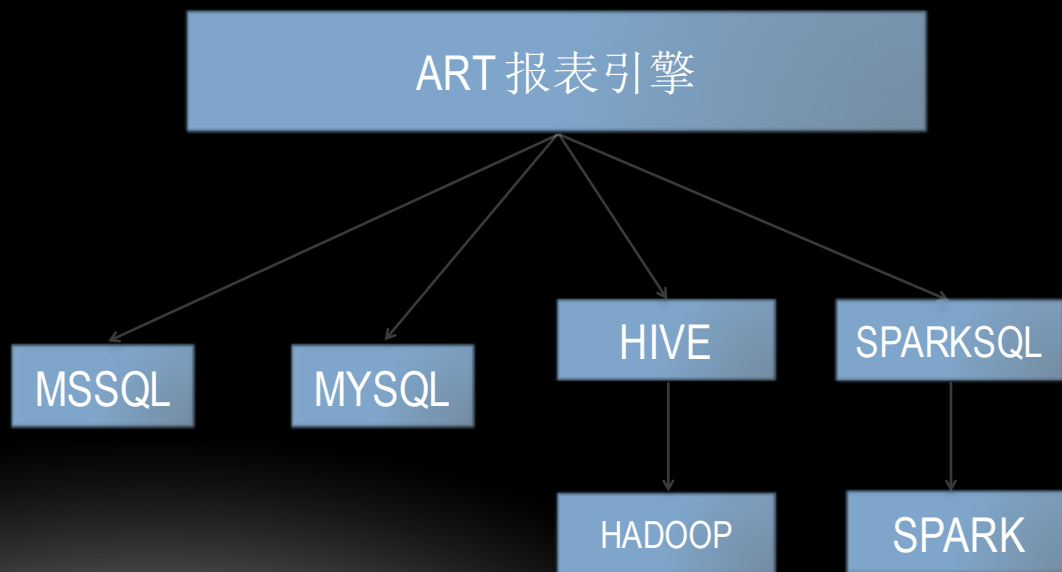
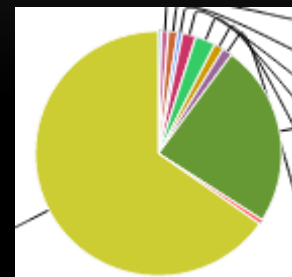
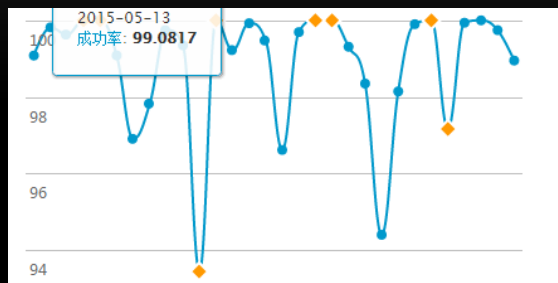
Active Reports 3.0

SAIKU

CUTTING EDGE OPEN SOURCE ANALYTICS

简单报表引擎ART

- 多种数据源
- 多种图例（饼图，柱图，线图）
- 定时发送邮件，附带pdf/png附件
- 简单易用
- 分级权限管理



流程监控 >> 宙斯监控 >> Job运行分布柱状图

StartDate

2015-08-20

EndDate

2015-08-24

查看方式：

显示报表 (Table)

☐ 显示参数

☐ 显示SQL

查询

10,000

8,000

6,000

4,000

2,000

0

08-20

08-21

08-22

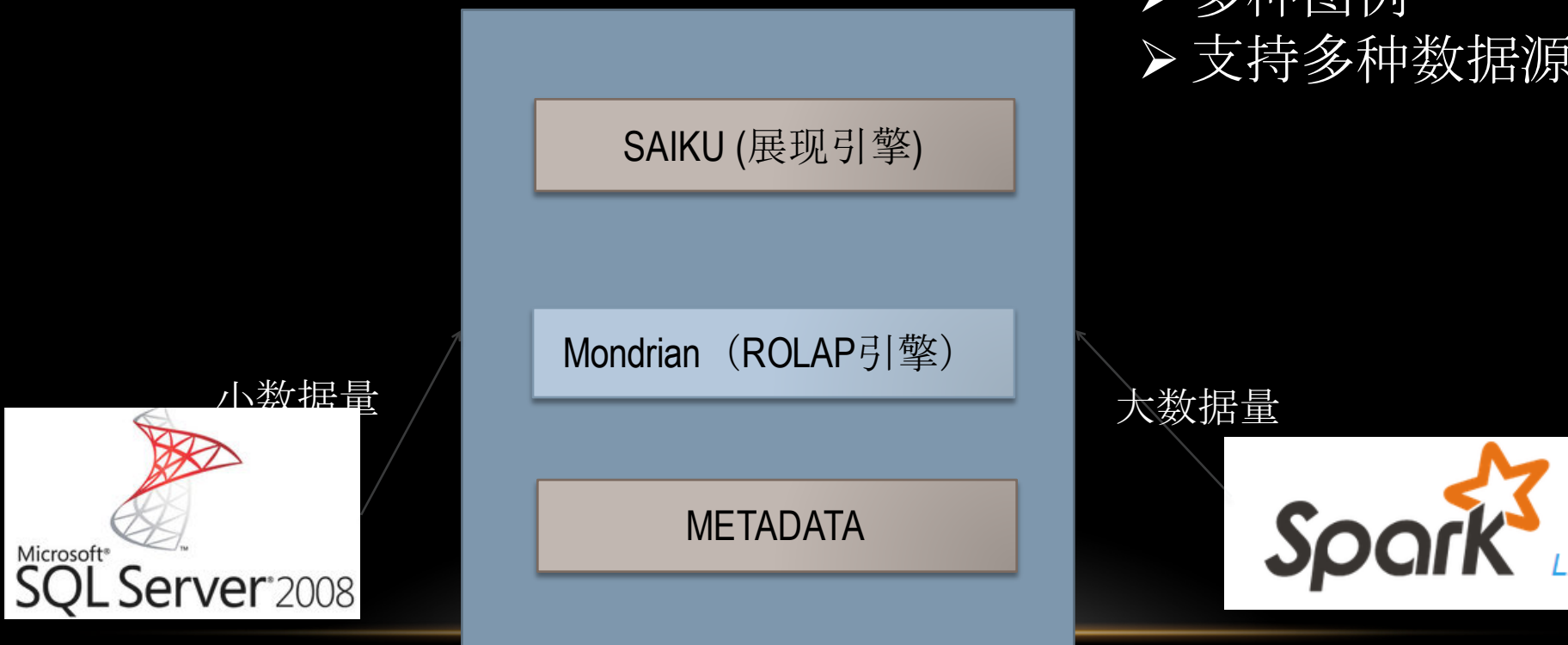
08-23

08-24

0_1m 1_5m 5_10m 10_15m 15_30m 30_45m 45_60m 60_90m 90_120m 2_3h 3_4h 4_5h greater_5h

多维分析引擎SAIKU

- 上钻，下钻
- 拖拽多维分析
- 多种图例
- 支持多种数据源



谢谢大家