

# 苏宁云商的大数据平台架构

王志强

苏宁云商.大数据中心.平台开发部

2015-07-25

# 关于我和我的小伙伴们

## 大数据中心.平台开发部

### 职责:

提供集团各个业务所需要的存储和计算能力。  
保证平台的稳定、高效运行。  
提高平台易用性。

### 目标:

打造稳定、易用、高效的平台，提高数据分析效率，实现人人都是数据分析师。



团队

我

大数据攻城狮

5年的大数据工作经验

# 内容

- 发展概况
- 总体框架
- 平台介绍
- TO DO

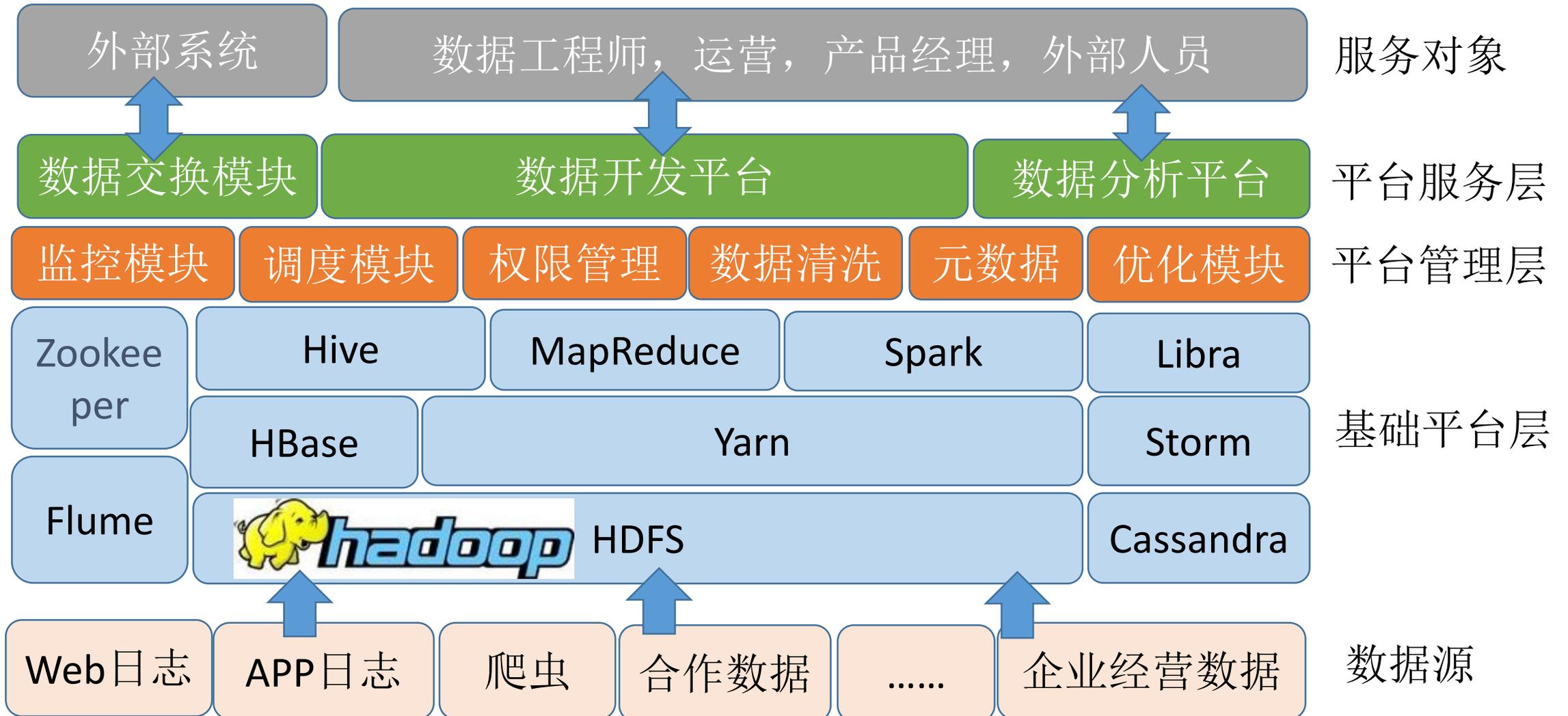
# 发展概况-特点

- 发展快
  - 600+节点，5w+ job/天
  - 作为苏宁基础数据平台，在公司内广泛使用
  - 30+中心或部门，200+开发者。
- 开源为主，辅以自研
  - 基础平台以开源为主
  - 部分组件自研

# 发展概况-平台

- 海量数据平台
  - HDFS/Yarn/Hive/MapReduce/Spark
- 流式计算平台
  - Storm/Libra
- KV存储平台
  - HBase/Cassandra
- 数据开发平台（CBT）

# 总体框架-架构



# 海量数据平台

- 包括：Hadoop、Hive、Spark
- 海量数据存储
- 离线计算
- 内存/迭代计算
- 算法平台



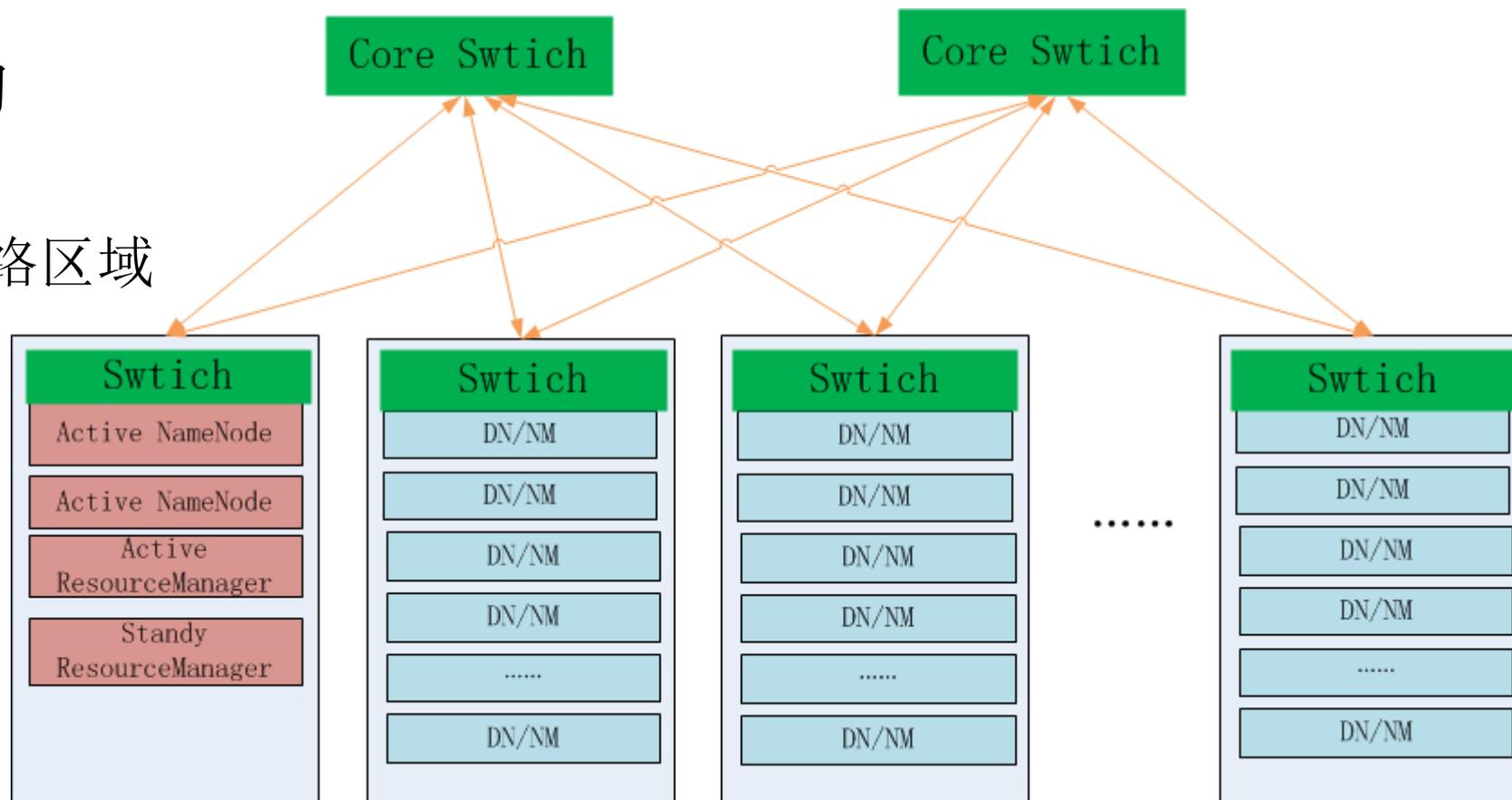
# 海量数据平台-Hadoop

- 部署架构
- 平台管理
- 监控体系
- 运行分析
- 分享
  - 升级过程
  - 问题及解决

# 海量数据平台-Hadoop

- 网络部署架构

- 单独划分网络区域
- 主节点隔离

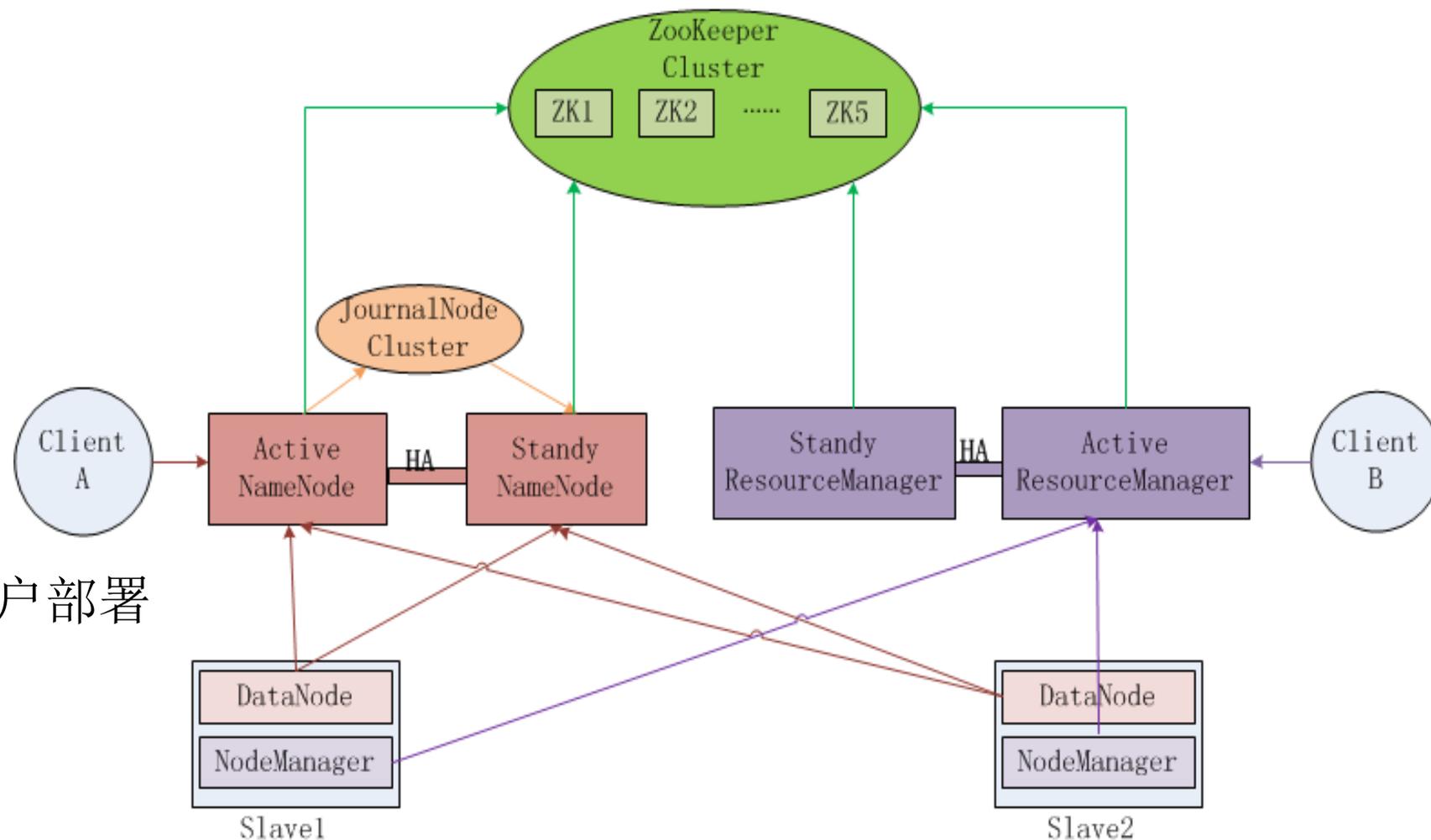


# 海量数据平台-Hadoop

- 逻辑部署架构

- HDFS HA
- Yarn HA

- 不同组件不同用户部署
- 客户端统一管理



# 海量数据平台-Hadoop

- 平台管理

- 自动化扩容
- 统一配置管理



- 用户管理

- 按部门/业务分配hadoop账户

- 权限控制

- 机器接入权限
- 目录文件访问权限



# 海量数据平台-Hadoop



- 监控体系

- 主机层监控

- CPU, Memory(Free,Swap), Network (ping, usage, mode), Disk (inode, free, iowait), Linux kernel (dmesg)

- 平台层监控

- 端口监控, 进程监控, **JVM GC**, 丢块监控, 长job监控, 容量监控, HA切换监控、**pending app**监控, dead节点数, 失败原因分析

- 应用层监控

- Hdfs应用层监控: put/get
    - Yarn/MR应用层监控: wordcount, 单独资源池, 防止误报
    - Hive应用层监控: drop/create table

# 海量数据平台-Hadoop



- 运行分析

- HDFS

- 总体情况：总容量，已使用容量，剩余容量，Live Nodes，DeadNodes，目录数，文件数，新增数据量，新增目录数，新增文件数
    - 分用户统计：总数据量，新增数据量，平均文件大小，空间使用率

- Yarn:

- 总体情况：总资源量，按小时统计CPU/Memory使用量，资源使用率。
    - 分pool统计：资源量，资源使用率，每个时间点的资源使用率。

# 海量数据平台-Hadoop



- 运行分析

- MapReduce统计

- 统计信息包括：提交任务总数，成功/失败任务数，map/reduce tasks数，datalocal比例，rack local比例，总处理数据量，map平均处理数据量/处理时间，map处理效率，reduce平均处理时间等。

- Job评分机制

- 通过分析Job Counters，计算Job的得分，督促提交者优化job。

- 评判标准：尽量减少overhead和IO。**

- 权重因子：map overhead，map spill数据量大小，中间结果是否压缩，是否长尾，reduce数是否合理等。

# 海量数据平台-Hadoop

- Hadoop升级过程

- 背景

原有的Hadoop 1.2.1版本存在问题：可靠性差、资源利用率低、不支持多种计算架构。

升级到2.4.0版本。

- 那么，问题来了.....

这么多业务，兼容性？

升级不成功咋办？？

回滚失败咋办？？？

数据丢了咋办？？？？



# 海量数据平台-Hadoop

## • 方案对比

### 方案一

Hadoop upgrade硬升级

优点:

简单  
代价最小

缺点:

风险大  
不可控

### 方案二

Upgrade硬升级 + 数据备份

优点:

风险较小

缺点:

非常复杂  
代价大

### 方案三

数据迁移, 切业务

优点:

风险小

缺点:

复杂  
代价大



# 海量数据平台-Hadoop

- 技术点：
  - 全量拷贝+增量拷贝：通过hdfs audit日志分析更改的目录或文件。
  - DistCp改造：跨集群同路径拷贝。
  - 回滚方案
- 实施过程：
  - 1. 启动audit定时分析
  - 2. 数据全量拷贝
  - 3. 数据增量拷贝
  - **4. 集群进入safemode模式**
  - 5. 数据增量拷贝
  - 6. 数据准确性校验
  - **7. 切客户端软链接**
  - 8. 验证



# 海量数据平台-Hadoop

## • 问题及解决

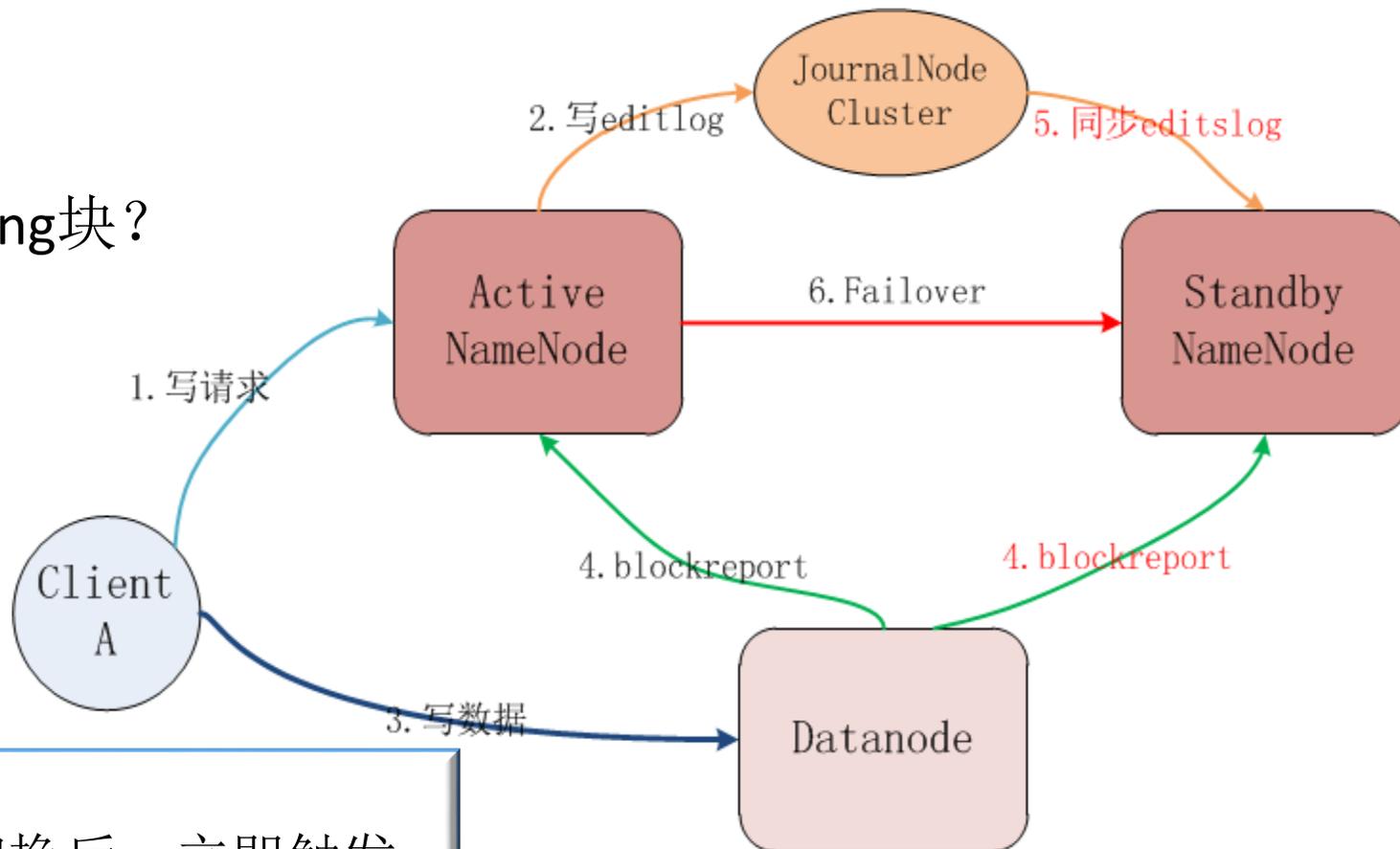
- HDFS HA切换后，Missing块？

原因：

standby namenode 回放 edits log的操作落后于 datanode blockreport

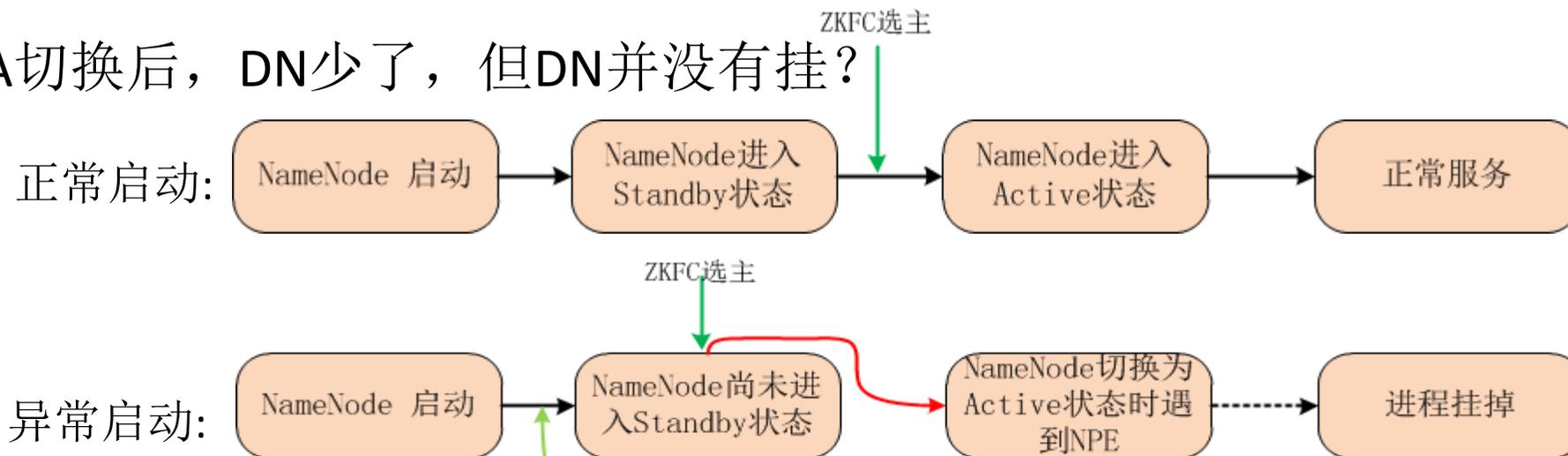
解决：

datanode发现namenode切换后，立即触发一次全量块汇报。



# 海量数据平台-Hadoop

- HDFS HA切换后，DN少了，但DN并没有挂？



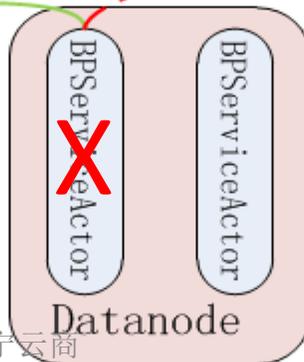
解决:

1. 启动hdfs时，待namenode都 standby状态后，再启动ZKFC。

方案2和3，参考

<http://weibo.com/p/1001603838201184837423>

连接 注册



# 海量数据平台-Hadoop

- 问题及解决
  - 数据误删问题？
    - 屏蔽skipTrash命令，删除数据必须先进入Trash。
  - 数据量/Job一直在增长，用户无成本意识，无优化动力？
    - Job评分机制，排名，优化建议
    - 计费
  - 空间问题：自动压缩，冷数据归档，自动合并小文件。
  - ResourceManager Crash退出、Hive Client CPU 100%的问题、。。。。。

# 海量数据平台-Hive

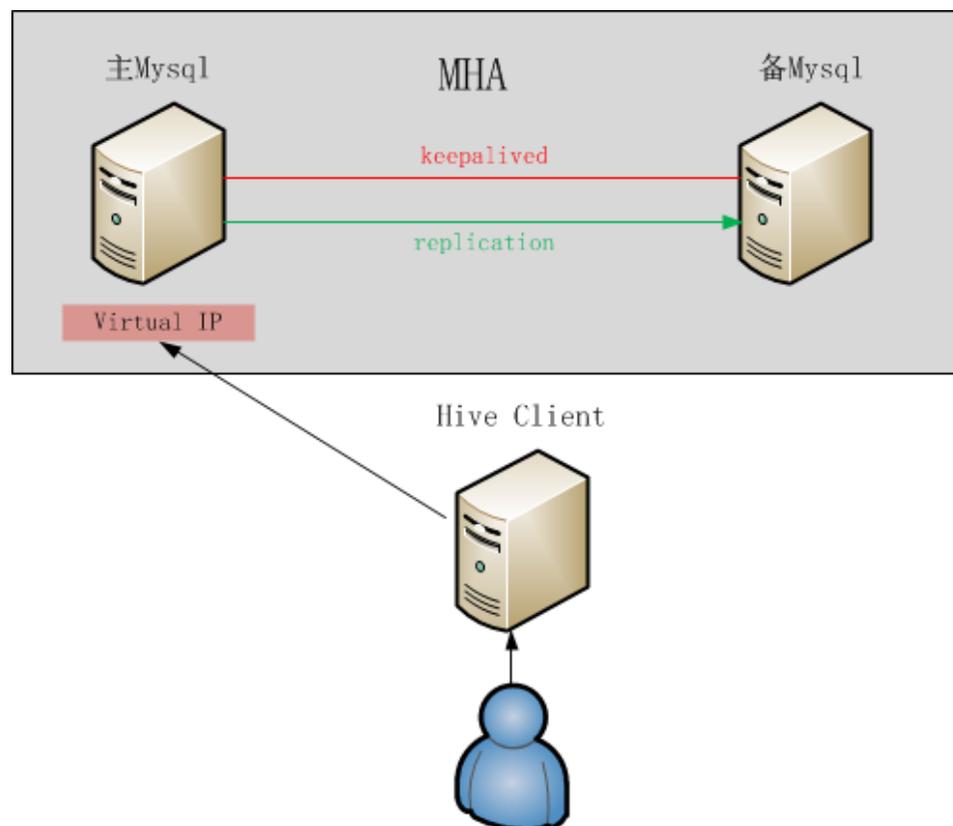
- 内容
  - 元数据存储
  - 安全
  - 权限管理

# 海量数据平台-Hive

- 元数据存储

- Mysql
- MHA
- VIP

- 高可靠
- 0事故



# 海量数据平台-Hive

- 安全

- 定时dump元数据并备份。
- javax.jdo.option.ConnectionPassword配置参数加密

```
<property>
  <name>javax.jdo.option.ConnectionPassword</name>
  <value>hivepassword</value>
</property>
```



- 原始数据和结果数据用外部表  
中间结果数据用内部表

# 海量数据平台-Hive

- 权限管理
  - 采用Hive自带的权限控制方案。
  - 实现管理员功能，只有管理员有赋权权限。
  - 跨用户访问数据，必须申请权限。
- HiveServer2的安全管理
  - 为每个Hadoop用户启动一个hiveserver2，不同机器不同端口。

# 海量数据平台-Spark

- 发展情况

- 刚起步
- 生产集群Ready
- 即将上线一个迭代式计算
- 正在对接三个中心的算法需求

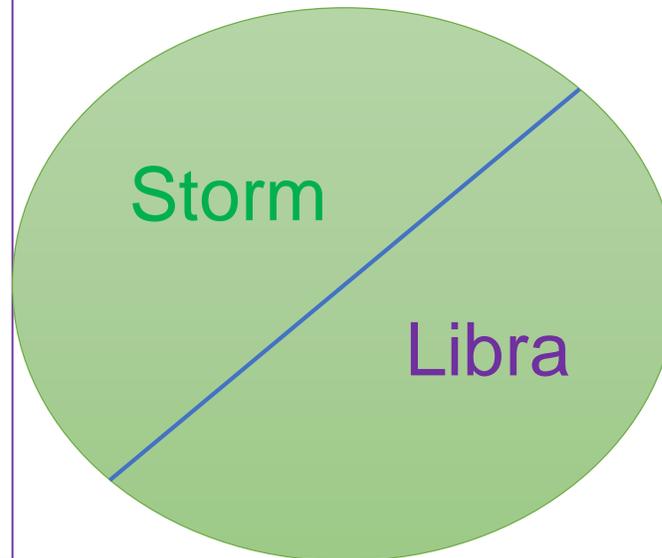


- 定位

- 快速计算
- 迭代计算
- 算法平台

# 流式计算平台

- 开源Storm 0.9.3
- 50+节点
- 12个业务
- 复杂算法
- 数据清洗、实时推荐、数据易道等



- 自研
- 50+节点
- 近300个sql
- 简单、sql可描述
- 流量分析、性能监控等

# 流式计算平台-Libra

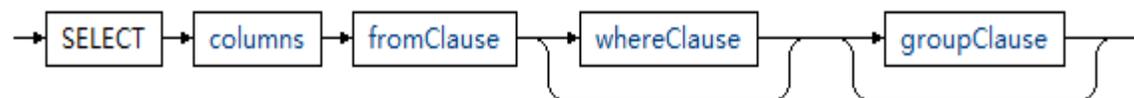
- 简化Storm应用开发，实现“Sql on Storm”



# 流式计算平台-Libra

## • 功能

- 支持的语法格式:



- 支持的运算窗口:

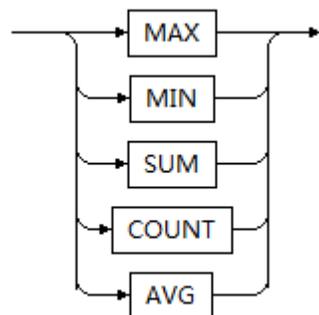
时间滑动窗口: 例如最近5分钟内某个IP的登陆次数。

时间批量窗口: 例如某个会员账号每天的付款总额。

长度滑动窗口: 例如某个会员最近5000次登陆的平均时间。

长度批量窗口: 例如某个会员每10次的平均付款金额。

- 支持的运算操作:



# 流式计算平台-Libra

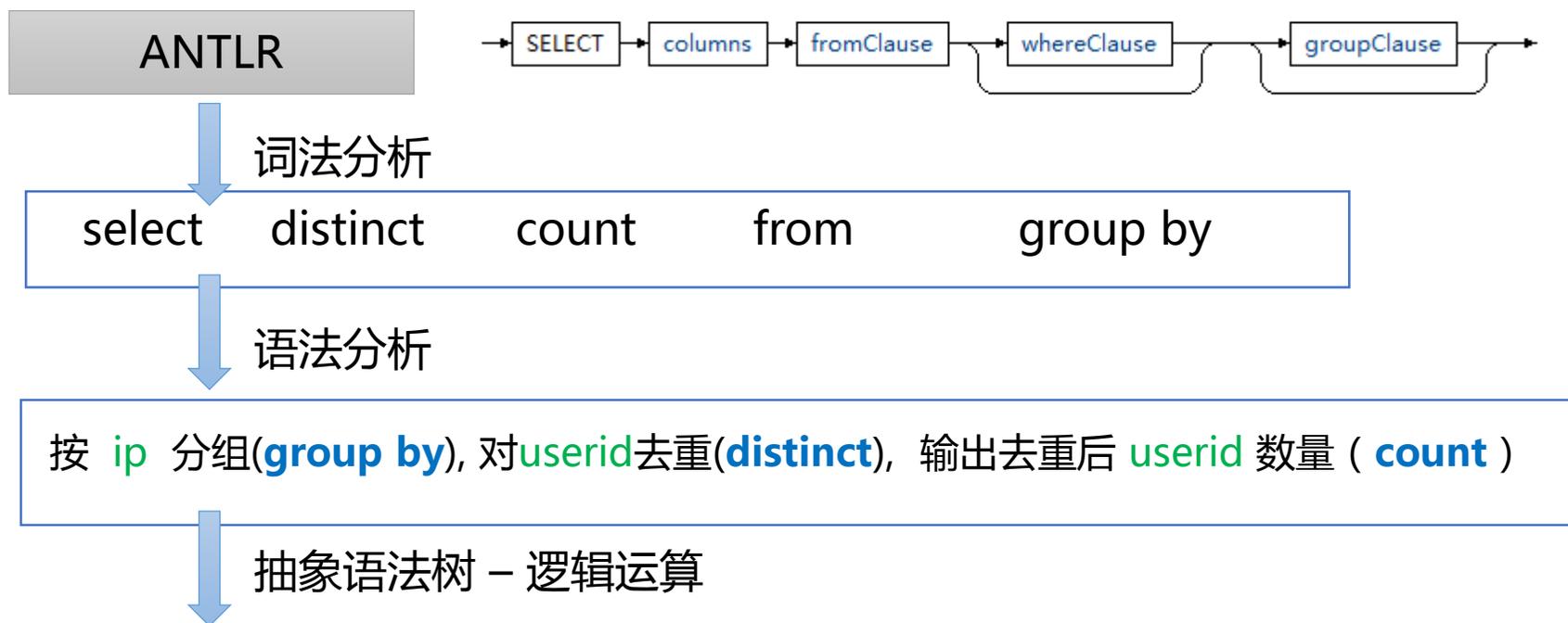
- 功能

- 支持的消息:                   格 式 化: JSON、XML  
                                  非格式化: 特殊字符间隔
- 支持的输入数据源:       Kafka、MQ
- 支持的输出数据源:       Kafka、MQ、Cassandra、Hbase、Redis、DB
- 支持的输出方式:         实时输出、定时输出。

# 流式计算平台-Libra

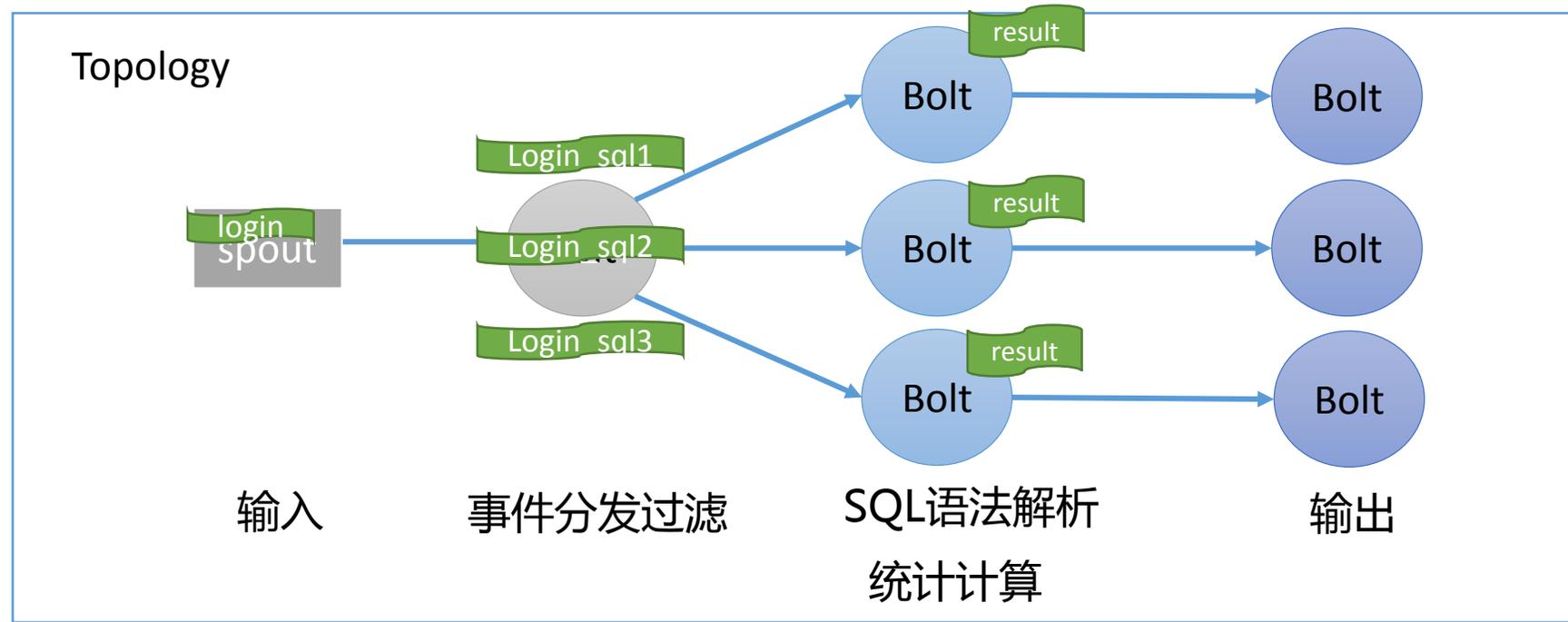
- 技术点

Select count(distinct userid) from login group by ip



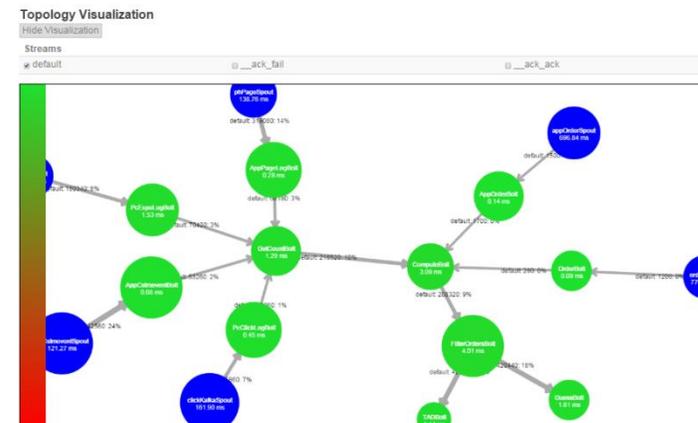
# 流式计算平台-Libra

- 将不同的算子映射到四种spout/bolt上去



# 流式计算平台-问题及解决

- Worker挂了，内存数据丢失？
  - 使用Cassandra缓存中间结果
  - 增量写入
  - 延迟加载
- Big Topology的页面刷新慢
  - 原因：可视化需要加载zk数据，耗时严重
  - 解决：只在点击显示可视化时，才触发取数据的操作。



# 流式计算平台-问题及解决

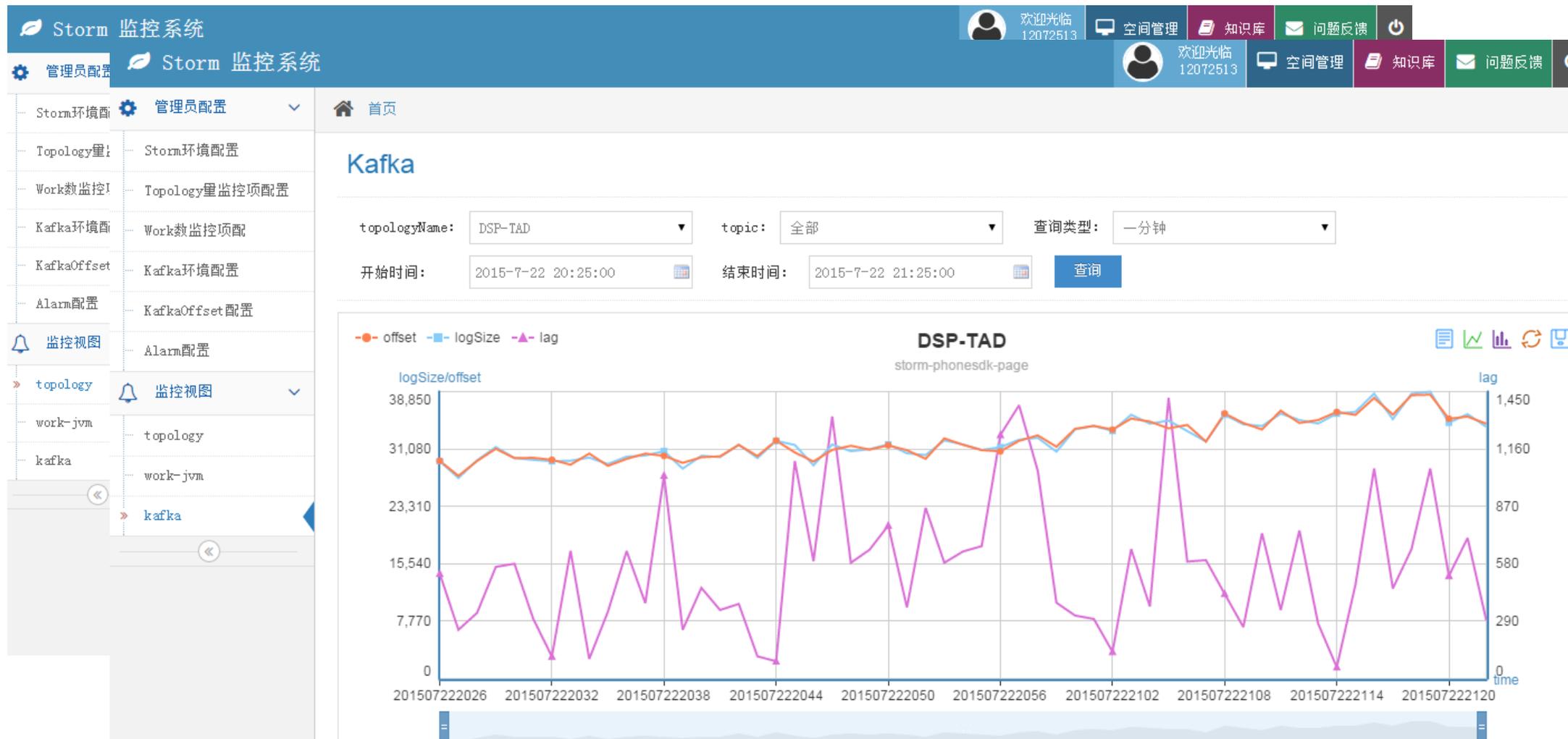
- Jar包太大导致Topology启动慢？
  - 原因：jar包太大，导致supervisor下载jar包时对nimbus带宽压力大，下包时间比较长。
  - 解决：代码包及依赖包分离，依赖包提前分发。
- Supervisor找不到jar包，导致挂掉？
  - 原因：
    1. supervisor两个线程不同步。
    2. supervisor与nimbus状态不同步。
  - 解决：捕获异常并处理。

# 流式计算平台-监控

- 实时计算监控报警尤其重要！！
- 平台监控：
  - 进程监控，worker监控，zk服务监控，worker内存监控。
- 业务监控：
  - 失败，堆积，数据质量监控等。

Storm Monitor

# 流式计算平台-监控



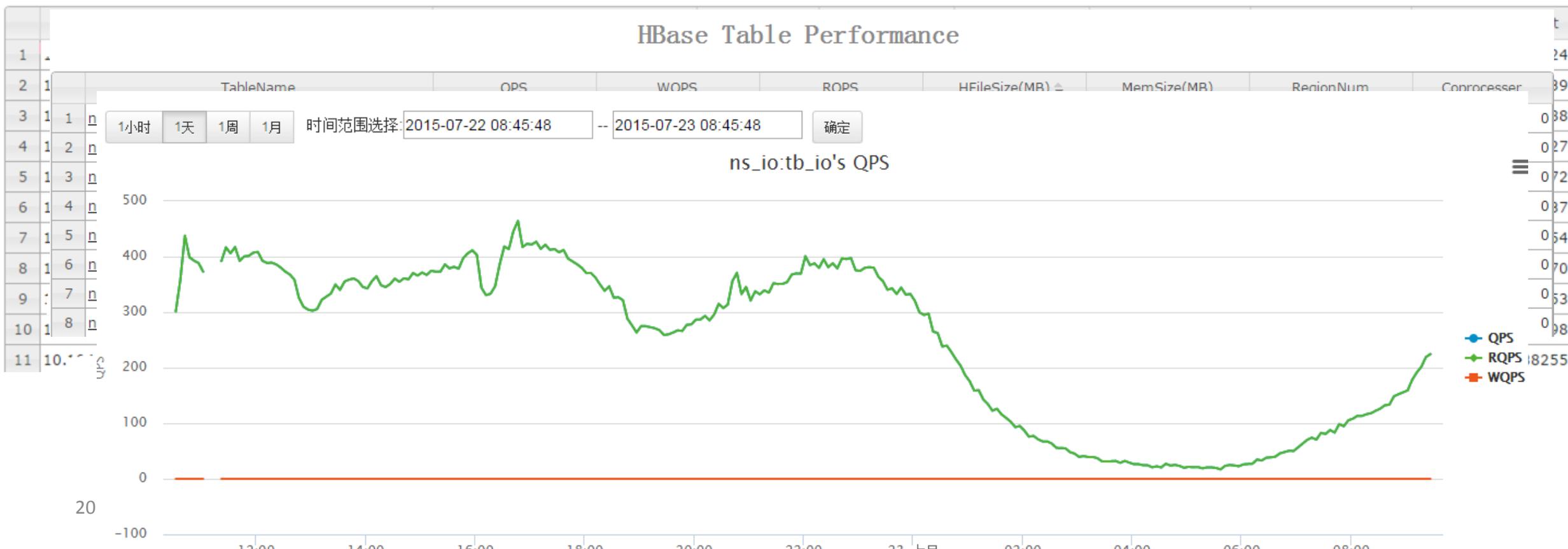
# KV存储平台

- 包括HBase, Cassandra
- HBase:
  - 海量KV存储
  - 应用场景: 搜索query分析、商品分析、广告联盟、云台纵览等。
  - 目标: 在线KV平台
- Cassandra:
  - 流式计算缓存中间结果数据和存储结果数据。

# KV存储平台-HBase

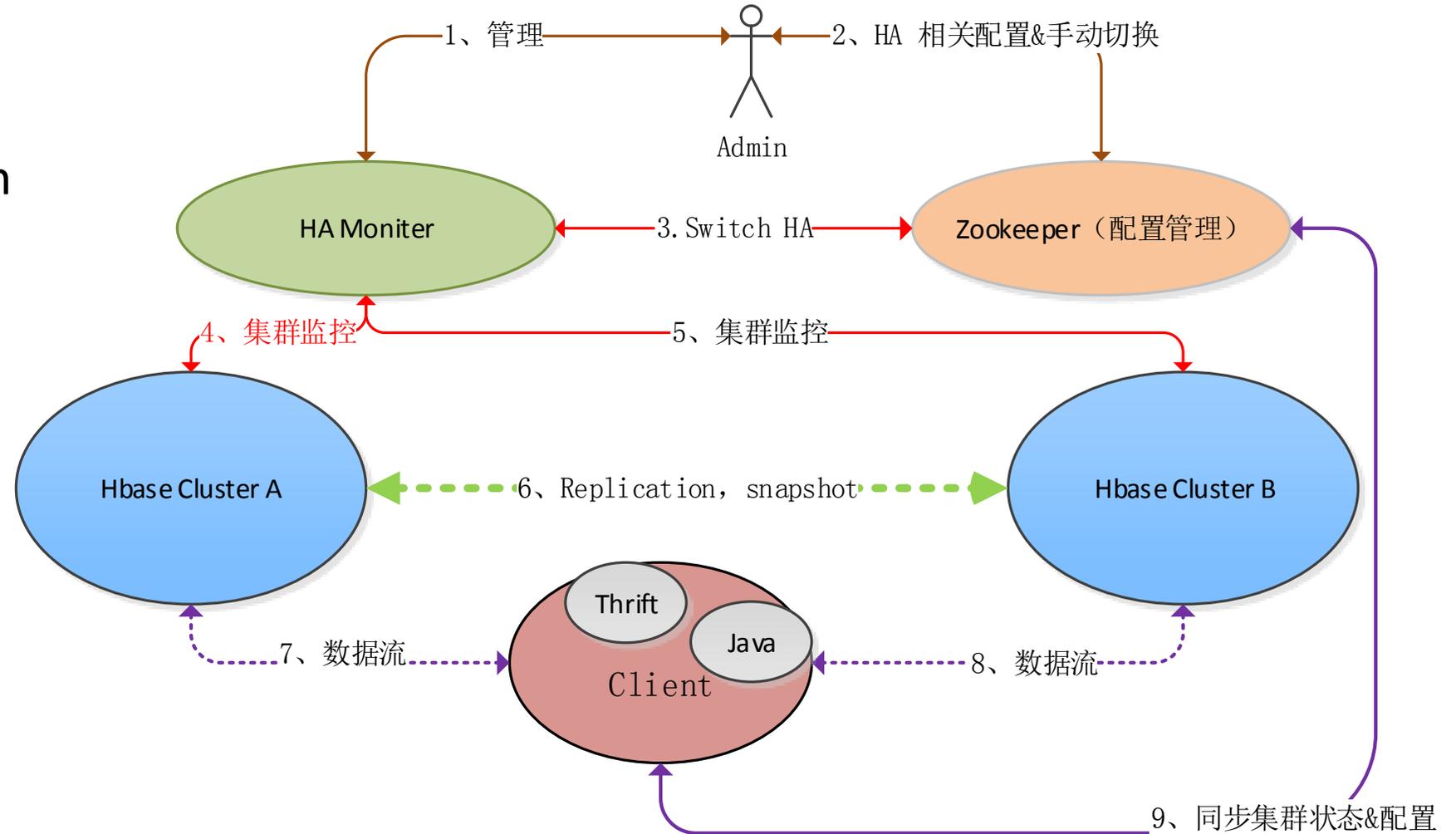
- 监控

## HBase Basic Information



# KV存储平台-HBase

- Cluster HA
  - Replication
  - Monitor
  - Failover



# 数据开发平台(CBT)

CBT 苏宁大数据开发平台

你好, 欢迎光临

空间管理 知识库 问题反馈 退出

管理员首页

数据源管理

任务流管理

任务管理

告警管理

任务运行状态

事件管理

任务组管理

公共资源管理

用户管理

用户映射

问题反馈

任务流总数 8

任务总数 288

正在执行 1

等待执行 0

执行成功 287

执行失败 0

### 今日的任务流列表

任务流名	描述	上线状态	优先级	执行频率	开始时间
DIM_SYNC	维度表同步任务流	上线	低	每天	2015-07-01 05:00:00
SSH_SIT_REPORT	苏宁云居测试补数专用	上线	低	每天	2015-07-21 15:53:41
huwangliang_testinfo	PC流量站内外分析	上线	低	仅此一次	2015-07-01 00:00:00
FTP_TO_OTHER_DAY	同步BI数据给其他系统	上线	低	每天	2015-06-01 06:30:00
PC_TASK_TEST	PC任务流测试	上线	一般	仅此一次	2015-06-05 09:08:39
PC_DAY_TASK_RUNTIME	关于PC相关的日任务	上线	一般	每天	2015-05-01 02:00:00
TRMNL_REPORT_MON	移动研发中心移动端报表(月报)	上线	一般	每月	2015-05-01 02:00:00
TRMNL_QR_2	移动研发中心销售报表	上线	一般	每天	2015-05-01 02:00:00

« 1 2 »

# 数据开发平台(CBT)

定位：为各个计算平台提供统一的、易用的开发平台。

功能模块：

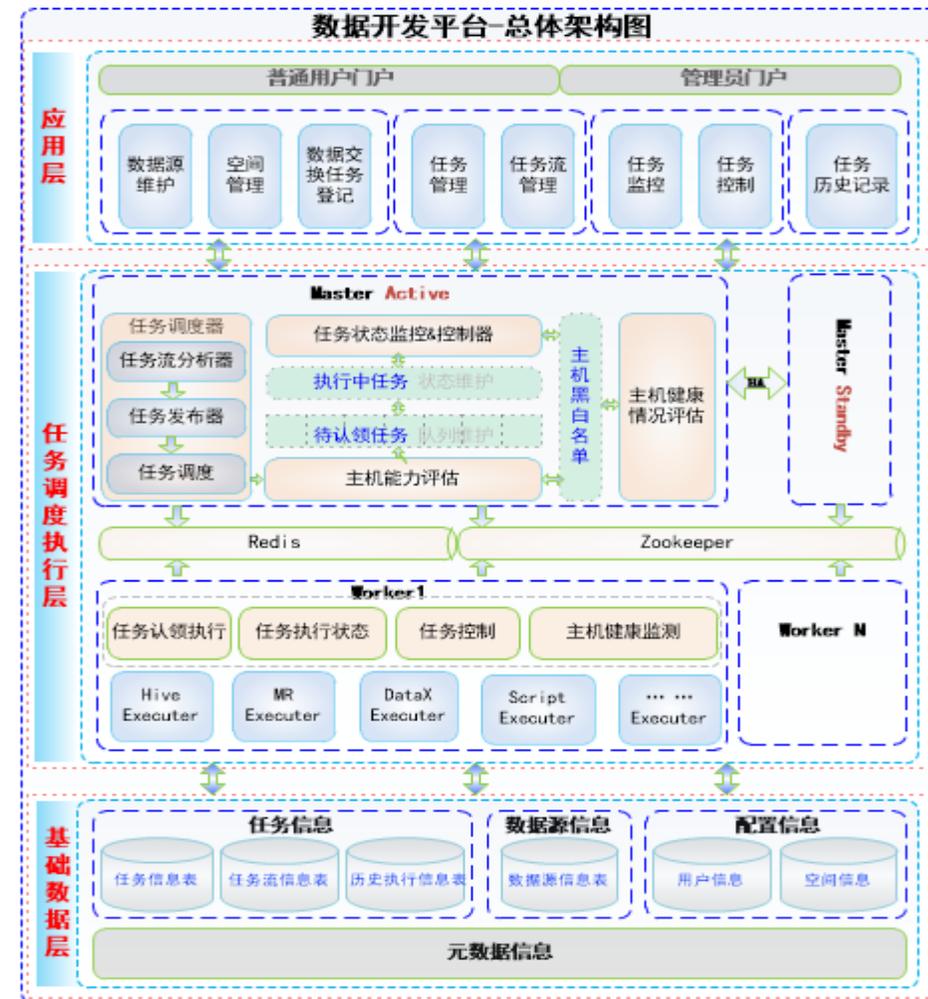


# 数据开发平台(CBT)

## • 特点

- 集群式结构
- 主节点HA
- 高可靠
- 易扩展
- 多平台支持

## • 总体架构



# TO DO

- 流式计算服务化（Storm As A Service）
- 公有云开放服务
- 完善数据开发平台，优化用户体验
- Spark推广应用
- 参与社区

The End  
Thanks

如果你对苏宁大  
数据技术有兴趣，  
欢迎来聊聊~~~~  
你懂的 $O(n_n)O\sim$

