

苏宁大数据平台运维实践

王志强

苏宁云商.大数据平台技术总监

QCon

全球软件开发大会

10月17-19日 上海·宝华万豪酒店



扫码锁定席位

九折即将结束

团购还享更多优惠，折扣有效期至9月17日

扫描右方二维码即可查看大会信息及购票



如果在使用过程中遇到任何问题，可联系大会主办方，欢迎咨询！

微信：qcon-0410

电话：010-84782011

ArchSummit

全球架构师峰会 2017



扫码锁定席位

12月8-9日 北京·国际会议中心

七折即将截止立省2040元

使用限时优惠码AS200，

以目前最优惠价格报名ArchSummit

仅限前20名用户，优惠码有效期至9月19日，

扫描右方二维码即可使用



如果在使用过程中遇到任何问题，可联系大会主办方，欢迎咨询！

微信：aschina666

电话：15201647919

极客搜索

全站干货，一键触达，只为技术

s.geekbang.org



扫描二维码立即体验

有没有一种搜索方式，能整合 InfoQ 中文站、极客邦科技旗下12大微信公众号矩阵的全部资源？

极客搜索，这款针对极客邦科技全站内容资源的轻量级搜索引擎，做到了！

扫描上方二维码，极客搜索！

这里只有 技术领导者

EGO会员第二季招募季正式开启



E小欧

报名时间：9月1日-9月15日

扫描添加E小欧，
邀您进入EGO会员预报名群

立即报名



你认为以下哪些关键字对运维最重要？

智能

易用

高效

性能

稳定

MTTR

流程

自动化

安全

规范

审计

SLA

故障预测

TABLE OF CONTENTS

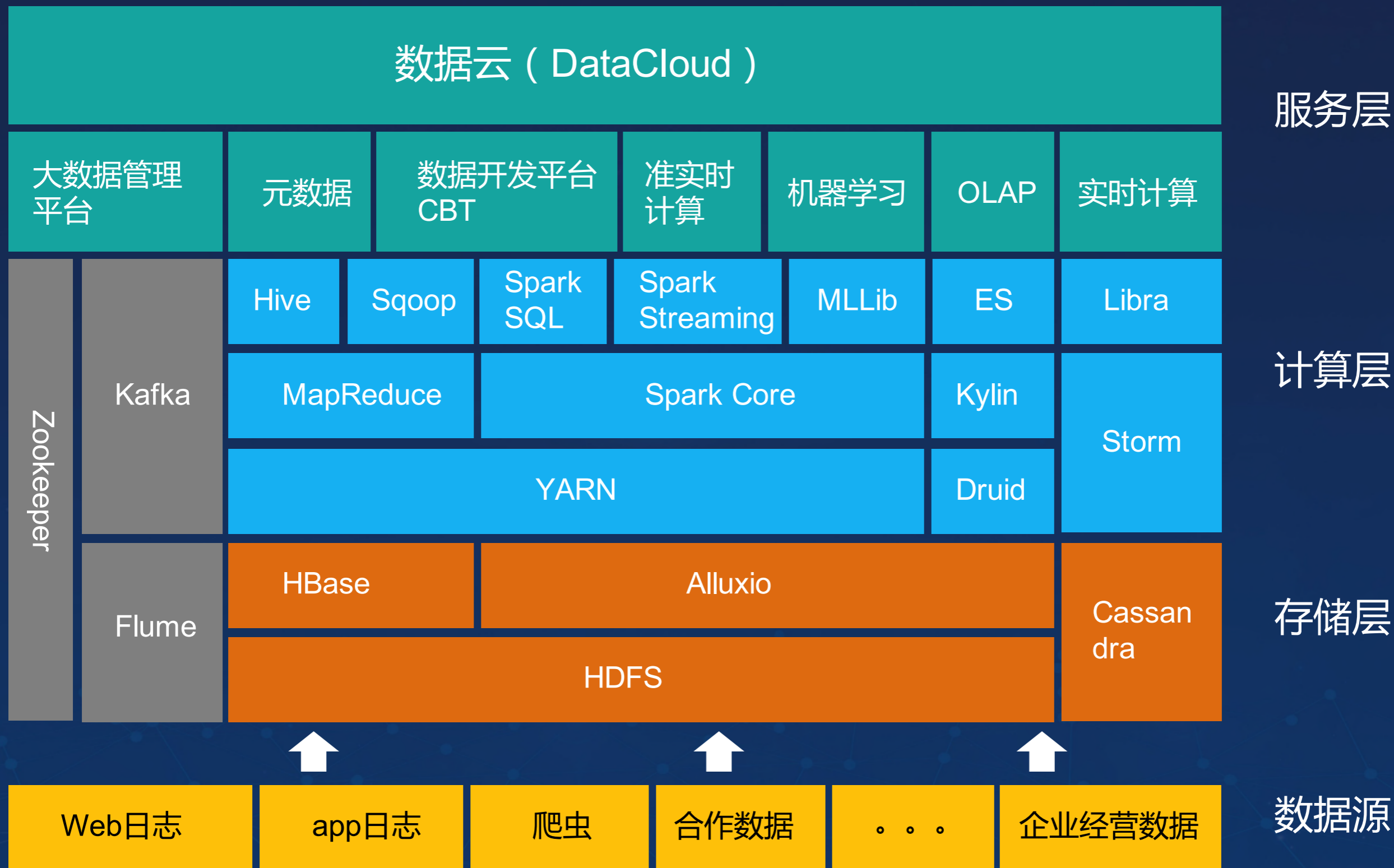
苏宁大数据平台基本介绍

大数据平台运维的痛点及解决方案

平台优化及增强

DOING & TO DO

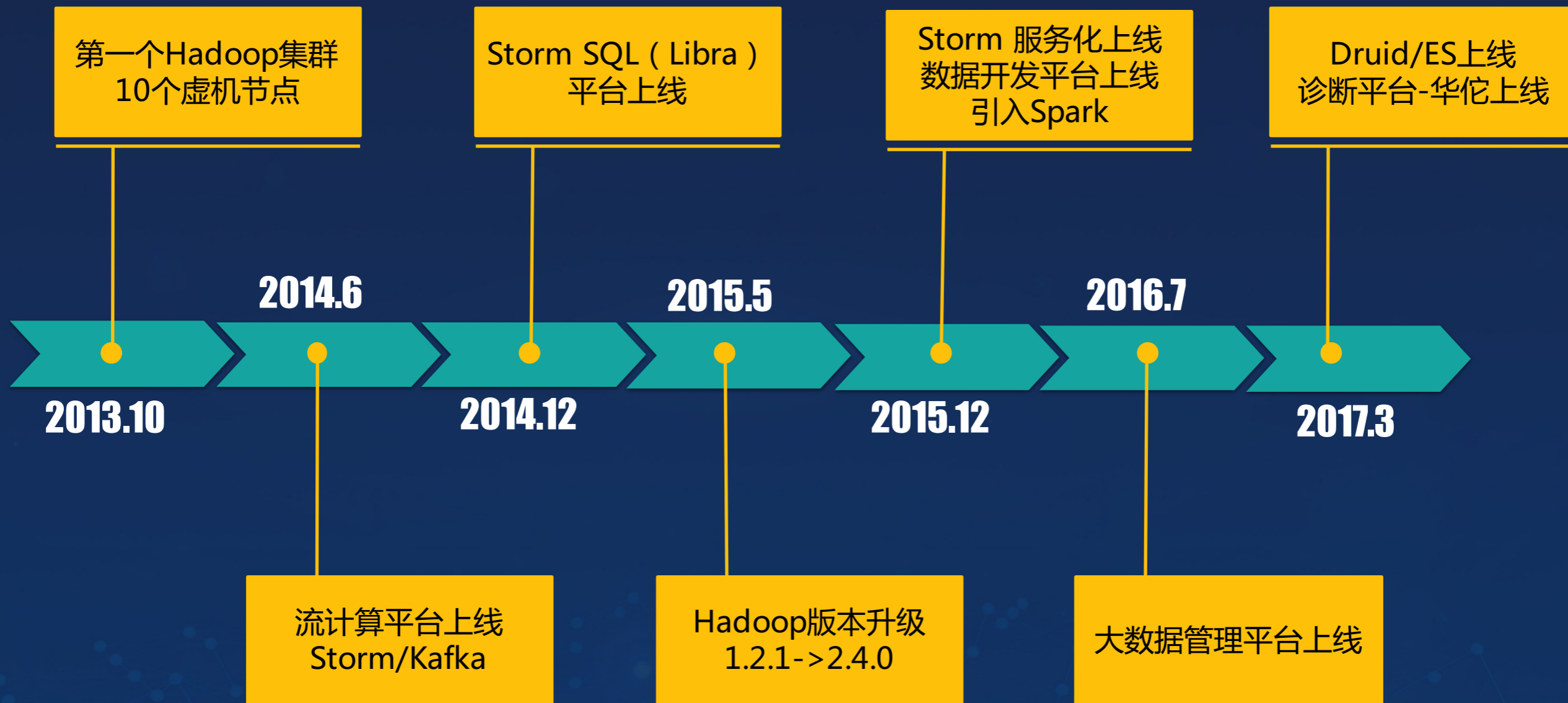
苏宁大数据平台软件栈介绍



典型数据流向



平台发展历程



平台规模

离线计算

1500节点
5万任务流, 30万任务/天
处理500TB/天
新增30TB/天
150账户, 500+开发人员
Hive:MR:Spark=90:3:7

KV存储(HBase)

160节点, 6个集群
50个NameSpace
600+张表
高峰期200万QPS

流式计算

Storm :
1400虚拟机, 41个集群
单集群最大174节点
1500流计算任务

Spark Streaming :
93节点, 80个任务

OLAP引擎

Druid :
39节点
30个DataSource

Elastic Search :
15节点
4个Index

TABLE OF CONTENTS

苏宁大数据平台基本介绍

大数据平台运维的痛点及解决方案

平台优化及增强

DOING & TO DO

大数据平台运维的痛点

- 痛点1. 部署及运维复杂
- 痛点2. 无资源使用视图
- 痛点3. 任务相互影响，资源隔离性差
- 痛点4. 排查问题耗时长，应用优化门槛高

▶ 痛点1. 部署及运维复杂

平台管理员

集群组件辣么多，部署环境辣么多，累成狗。。。。

堡垒机->跳板机->运维机，就为了开个权限。。。。

谁懂Flume运维的苦？

业务开发

我开发用到一个用户行为的数据，在哪个表里？

配置crontab的方式串联任务，伤不起。。。。

不同账户间的任务相互依赖严重，跨账户依赖难实现

解决1. 平台化、自动化

大数据管理平台：主机管理，集群管理自动化

大数据管理平台

帮助 王志强

集群管理

主机

客户端

账户

权限

资源

监控

资源展板

任务展板

主机管理

输入ip或者hostname



+ 添加

编号	IP	Host	主机类型	操作
1	10.27.1.138	namenode1-sit.cnsuning.com	cluster	详情 销毁
2	10.27.1.141	namenode2-sit.cnsuning.com	cluster	详情 销毁
3	10.27.1.142	slave01-sit.cnsuning.com	cluster	详情 销毁
4	10.27.1.143	slave02-sit.cnsuning.com	cluster	详情 销毁
5	10.27.0.241	slave03-sit.cnsuning.com	cluster	详情 销毁
6	10.27.0.242	slave04-sit.cnsuning.com	cluster	详情 销毁

主机列表

解决1. 平台化、自动化

大数据管理平台：主机管理，集群管理自动化

The screenshot displays a web-based management interface for a Big Data platform. The top navigation bar includes the title '大数据管理平台' and a user profile for '王志强'. A left sidebar lists various management functions such as '集群管理', '主机', '客户端', '账户', '权限', '资源', '监控', '资源展板', '任务展板', 'OLAP', 'HBase HA管理', and 'HBase数据探查'. The main content area is titled '主机管理 > 详情' and is divided into three sections: '主机基本信息', '主机资源信息', and '主机部署组件信息'. The '主机基本信息' section shows details for 'namenode1-sit.cnsuning.com' with IP '10.27.1.138'. The '主机资源信息' section lists hardware specifications like '24 Intel(R) Xeon(R) CPU E5-2620 v2 @ 2.10GHz' for CPU and '264484056' for memory. The '主机部署组件信息' section provides a table of installed services and their status.

主机部署组件信息		
HMaster	common-hbase-sit	启动 停止
NameNode	common-hdfs-sit	启动 停止
JournalNode	common-hdfs-sit	启动 停止

Annotations in the image include a green box labeled '主机详情' with arrows pointing to the '主机基本信息' and '主机资源信息' sections, and another green box labeled '主机运行组件' with an arrow pointing to the '主机部署组件信息' table. A '刷新' button is located to the right of the resource information, and '全部启动' and '全部停止' buttons are located to the right of the component information table.

解决1. 平台化、自动化

大数据管理平台：主机管理，集群管理自动化

The screenshot displays the '大数据管理平台' (Big Data Management Platform) interface. The top navigation bar includes the platform name, a '帮助' (Help) link, and a user profile for '王志强'. The left sidebar contains a menu with '集群管理' (Cluster Management) highlighted in red, along with other options like '集群看板', '长任务进度', '软件部署', '配置管理', '主机', '客户端', '账户', '权限', '资源', '监控', and '资源看板'. The main content area is titled '集群看板' (Cluster Dashboard) and features a '+ 创建集群' (Create Cluster) button in the top right. Below this, there are tabs for 'HDFS集群', 'YARN集群', 'HBase集群', and 'ZooKeeper集群'. A green arrow points from a teal box labeled '集群列表' (Cluster List) to the 'HBase集群' tab. The 'HDFS集群' tab is active, showing a cluster named 'common-hdfs-sit' with the description 'sit环境HDFS集群' and '详情 节点'. A blue status bar at the bottom of the cluster card indicates '运行正常' (Running normally).

解决1. 平台化、自动化

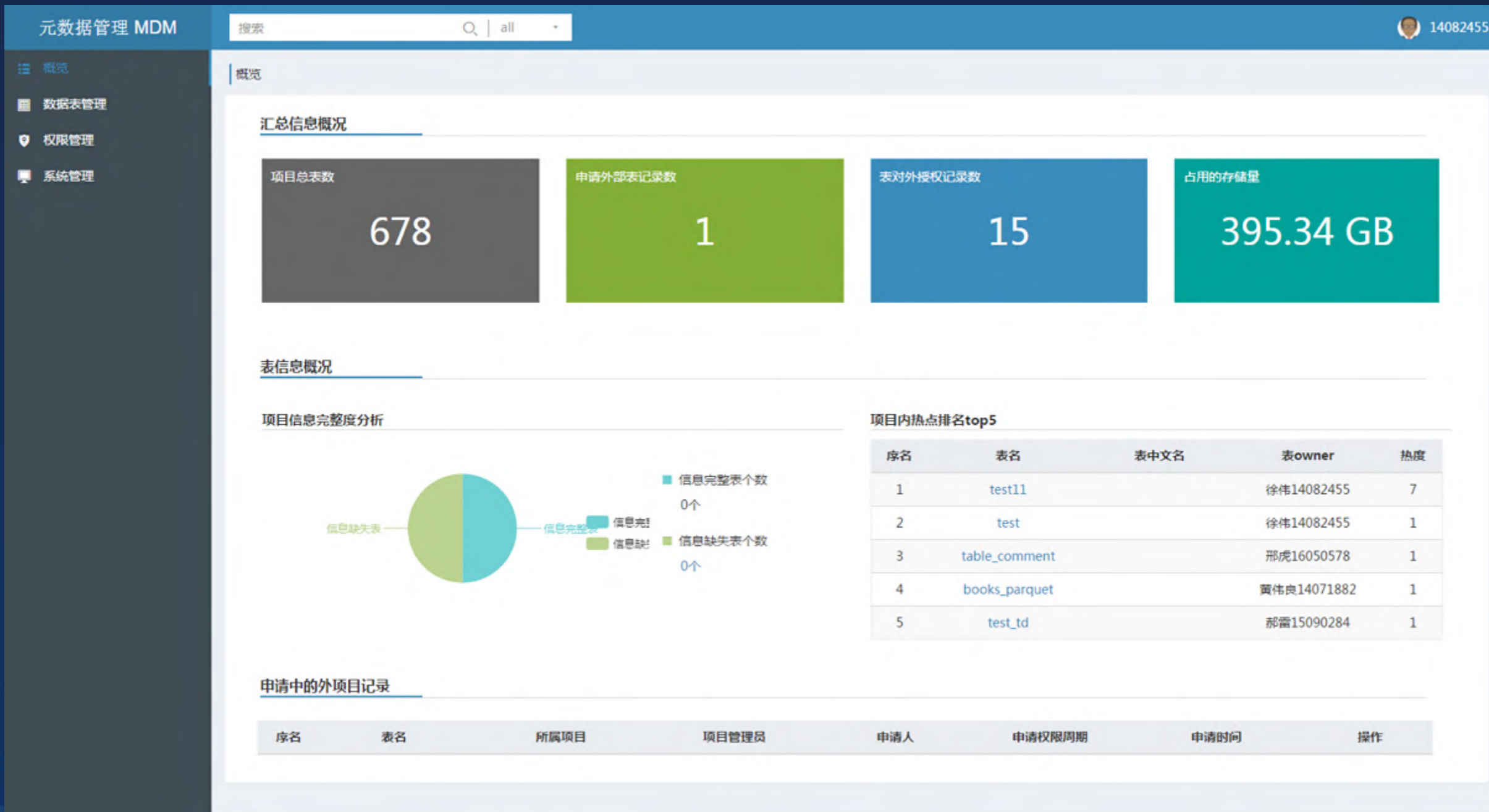
大数据管理平台：主机管理，集群管理自动化

The screenshot displays the '大数据管理平台' (Big Data Management Platform) interface. The top navigation bar includes '帮助' (Help) and the user '王志强' (Wang Zhiqiang). The left sidebar lists various management functions: 集群管理 (Cluster Management), 主机 (Hosts), 客户端 (Clients), 账户 (Accounts), 权限 (Permissions), 资源 (Resources), 监控 (Monitoring), 资源看板 (Resource Dashboard), 任务看板 (Task Dashboard), OLAP, HBase HA管理 (HBase HA Management), and HBase数据探查 (HBase Data Exploration). The main content area is titled '集群看板 > common-yarn-sit详情' (Cluster Dashboard > common-yarn-sit Details). It features a '集群详细信息' (Cluster Detailed Information) section with the following data:

common-yarn-sit的基本信息		查看页面(使用window运维机访问)>>	
sit环境YARN集群			
配置信息			
Version	hadoop-2.4.0.5	ResourceManagers	namenode2-sit.cnsuning.com,namenode1-sit.cnsuning.com
Zookeepers	namenode1-sit.cnsuning.com,namenode2-sit.cnsuning.com,slave01-sit.cnsuning.com	NodeManager节点数	10
依赖hdfs	common-hdfs-sit		
运行信息			
运行正常	40%	34	查看节点>>

解决1. 平台化、自动化

元数据管理：数据字典，权限申请审批实施自动化



解决1. 平台化、自动化

元数据管理：数据字典，权限申请审批实施自动化

元数据管理 MDM

搜索: dpa_ | all

概览
数据表管理
权限管理
系统管理

Hive(1000) Kafka(3)

项目: ETL系统(ETL) 类目: 清空

dpa_ 查询 重置

表名	项目	所属数据库	负责人	最后更新时间
dpa_br_outer_click_d	ETL系统(ETL)	bi_dpa	彭虎(11080108)	2014-10-08 15:29:39
dpa_ord_bill_tf	ETL系统(ETL)	bi_dpa	彭虎(11080108)	2017-07-18 11:38:39
dpa_ord_bill_tf_d	ETL系统(ETL)	bi_dpa	彭虎(11080108)	2017-05-09 20:40:51
dpa_ord_bill_tf_spark	ETL系统(ETL)	bi_tmp	彭虎(11080108)	2015-12-22 14:48:40
dpa_ord_bill_tf_tmp	ETL系统(ETL)	bi_dpa	彭虎(11080108)	2015-12-25 09:09:19

数据字典，支持搜索

解决1. 平台化、自动化

元数据管理：数据字典，权限申请审批实施自动化

The screenshot displays the '元数据管理 MDM' (Metadata Management MDM) interface. The main content area shows details for the table 'dpa_ord_bill_tf_d' in the 'Hive' database. The '项目管理员' (Project Administrator) field is highlighted with a red box, showing '彭虎 11080108'. A modal dialog box titled '权限申请' (Request Permissions) is open, showing the table name 'dpa_ord_bill_tf_d', a dropdown for '权限有效期' (Permission Validity) set to '永久' (Permanent), and a text input for '申请理由' (Request Reason) with the placeholder '请输入理由,不超过200字'. The dialog has '确定' (Confirm) and '取消' (Cancel) buttons.

表基本信息	字段信息	分区信息	变更历史	授权信息
表名: dpa_ord_bill_tf_d				
中文名:				
所属数据库: bi_dpa				
所属项目: ETL系统(ETL)				
项目管理员: 彭虎 11080108				
表owner: 彭虎 11080108				
描述:				
权限状态: 无				
是否敏感信息表: 否				
表来源: Hive仓库				

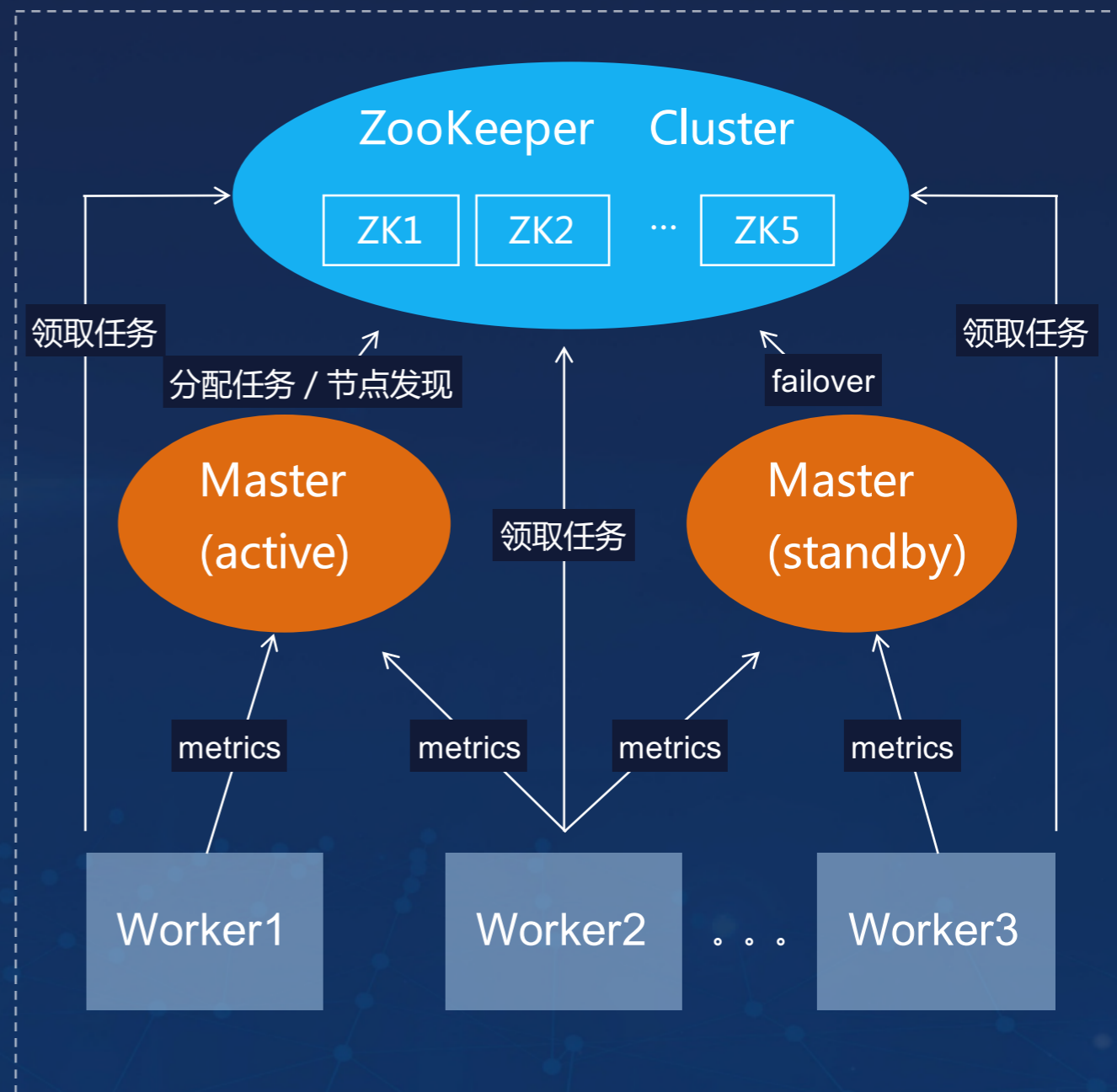
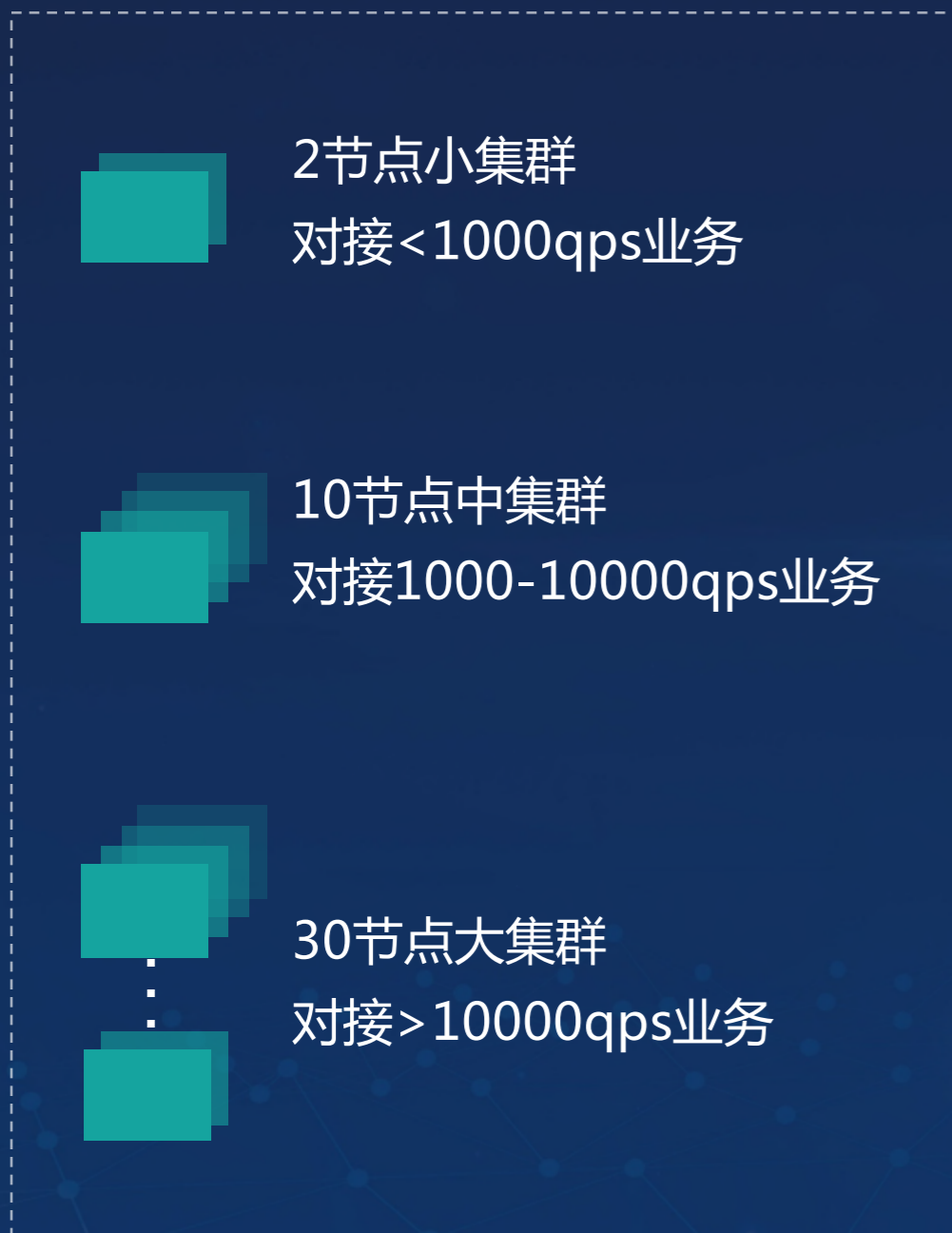
授权人所属系统	授权人	授权时间	批准人
---------	-----	------	-----

Copyright © 2016 苏宁云商. All rights reserved.

权限申请

解决1. 平台化、自动化

数据流管理平台：集成Flume，智能扩缩容，插件式



解决1. 平台化、自动化

数据开发平台：支持10种不同的任务类型，支持任务流/任务管理，解决复杂依赖问题，可扩展

The screenshot displays the '数据开发IDE' (Data Development IDE) interface. The top navigation bar includes '数据开发IDE', '首页', '任务设计', '资源中心', '运维中心', and '支持中心'. A left sidebar contains a search bar and a list of task components: '任务流', '组件', '草稿', '数据交换', 'Hive任务', 'MapReduce任务', 'DATASTAGE', 'RPC任务', 'Sqoop任务', 'Java任务', 'Python任务', 'Spark任务', and 'SparkSql任务'. The main workspace shows a '新建任务流' (New Task Flow) button. A '新建草稿' (New Draft) dialog box is open, featuring a '任务类型选择' (Task Type Selection) section with radio buttons for: '数据交换任务' (selected), 'Hive任务', 'MapReduce任务', 'DATASTAGE', 'RPC任务', 'Java任务', 'Sqoop任务', 'Python任务', 'Spark任务', 'SparkSql任务', and '机器学习'. The '基本信息' (Basic Information) section includes input fields for '任务名称*' (Task Name) and '任务描述*' (Task Description), both with placeholder text '请输入任务名称' and '请输入描述文本' respectively. At the bottom of the dialog are '确定' (Confirm) and '取消' (Cancel) buttons.

解决1. 平台化、自动化

数据开发平台：支持10种不同的任务类型，支持任务流/任务管理，解决复杂依赖问题，可扩展

The screenshot displays a web-based IDE for data development. The top navigation bar includes '数据开发IDE', '首页', '任务设计', '资源中心', and '运维中心'. The user is logged in as '王志强'. The main workspace shows a task configuration form for a task named 'sangqiang_spark_config'. The form includes fields for 'Jar名称' (JavaOK.jar), 'Main Class' (com.suning.java.Test), 'Args' (--Xms200M), and 'spark版本' (spark-1.5.2). Below the form is a 'spark opts' section with a table for parameter configuration.

参数名	参数值	
	3333	X
	222	X
	1111	X

解决1. 平台化、自动化

数据开发平台：支持10种不同的任务类型，支持任务流/任务管理，解决复杂依赖问题，可扩展

The screenshot displays a data development IDE interface. At the top, there is a navigation bar with tabs for '数据开发IDE', '首页', '任务设计', '资源中心', and '运维中心'. On the right side of the navigation bar, there are links for '支持中心' and '王志强'. Below the navigation bar, there is a search bar with the text '输入搜索内容'. The main workspace is divided into two panels. The left panel shows a list of task flows, including '9_onlyonce_6...', 'sync_bi_dpa.t...', 'ISE_HBASE_O...', 'HP_MSA_H', 'test_liuyh', '智能客服离线...', '111', 'HP_MSA_D', 'HP_CCBS_INIT', and '小时任务测试'. The right panel displays a complex task flow diagram with multiple nodes and dependencies. The nodes are represented by boxes with labels like 'HP_MSA_TANA_SMPACT' and 'HP_MSA_TANA_SMPACT'. The diagram shows a hierarchical structure with arrows indicating dependencies between tasks. At the bottom of the interface, there is a '日志查看' (Log View) section and a '任务流面板' (Task Flow Panel).

▶ 痛点2. 无资源使用视图



为什么我的任务跑不动了，是资源不够了么？



给我多分点资源吧？



另外，我昨天跑了多少任务，占了多少资源？

已读



已读



解决2. 资源可视化、人民币化

- ✓ 存储/计算资源计量计费
- ✓ 资源池使用可视化
- ✓ 任务展板

解决2. 资源可视化、人民币化

存储/计算资源计量计费



解决2. 资源可视化、人民币化

资源使用可视化



解决2. 资源可视化、人民币化

任务展板

大数据管理平台

帮助

王志强

- 集群管理 >
- 主机
- 客户端
- 账户
- 权限
- 资源 >
- 监控 >
- 资源展板 >
- 任务展板 >
- 排行榜
- 任务概览
- 任务详情

任务概览

账户 All

集群 common-yarn-sit

组件 All

2017-08-09

总任务数:

10753

成功任务数:

10747

失败任务数:

1

主动杀死任务数:

5

总耗费

222元

处理数据总量

1.69TB

平均处理数据量

165MB

平均运行时间

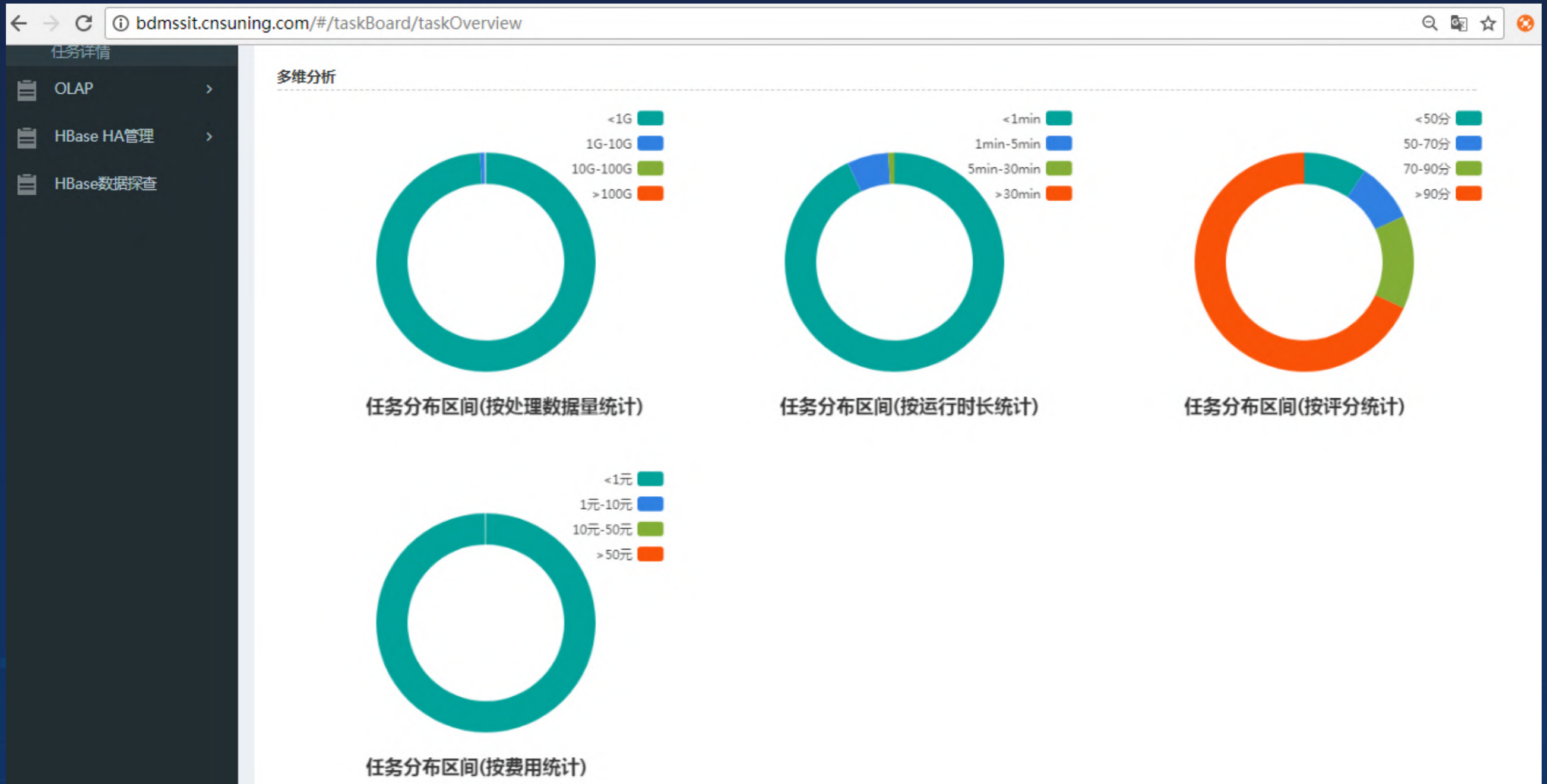
31s

MR/SPARK任务占比

10753 / 0


解决2. 资源可视化、人民币化

任务展板



解决2. 资源可视化、人民币化

任务展板

大数据管理平台 帮助  王志强

集群管理 > 主机 客户端 账户 权限 资源 > 监控 > 资源展板 > **任务展板 >** 排行榜 任务概览 任务详情 OLAP > HBase HA管理 > HBase数据探查

排行榜

组件: All 📅 2017-08-09

任务总量Top10

排名	账户	任务数
1	csi	3347
2	bigdata	2232
3	erp	1274
4	lbi	1086
5	bi	779
6	srdss	361
7	scpdm	240
8	smmb	205
9	sopdm	191
10	pos	155

处理数据量Top10

排名	账户	处理数据量
1	spider	706.78GB
2	srdss	271.88GB
3	ztbd	210.14GB
4	bigdata	165.41GB
5	pcms	148.15GB
6	mobdss	118.66GB
7	sousuo	29.57GB
8	pos	27.91GB
9	scpdm	15.15GB
10	lbi	12.13GB

任务评分Top10 (由高到低)

排名	账户	评分
1	bigdata	100

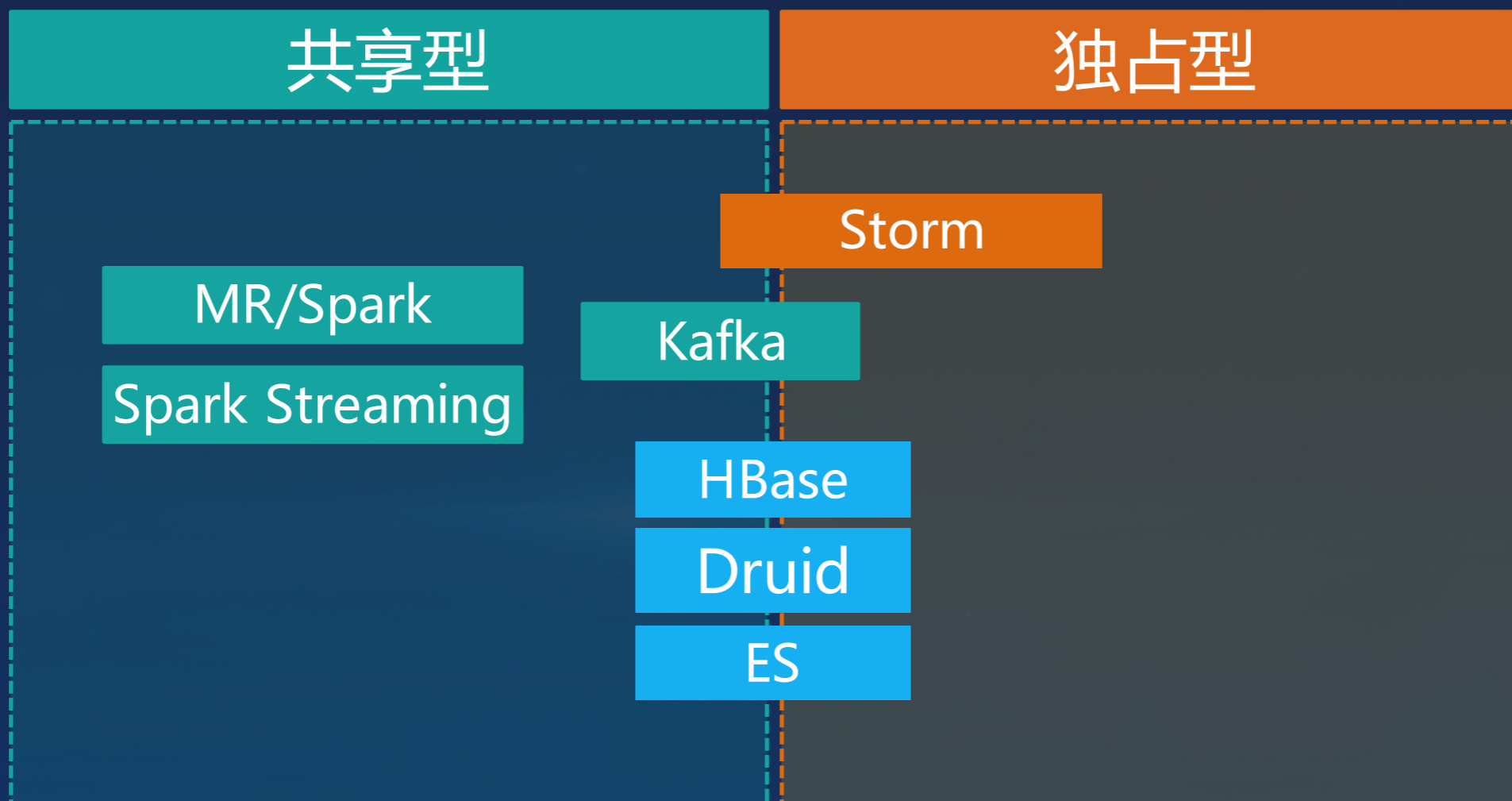
任务评分Top10 (由低到高)

排名	账户	评分
1	safe	48

▶ 痛点3. 任务互相影响，资源隔离性差

我的任务为什么突然变慢了？
是不是别人影响我了？

解决3. 差异化服务、物理隔离



资源浪费：离线、流本来就是错峰的

解决3. 差异化服务、物理隔离

共享型-YARN

FairScheduler：每个账户一个Pool，配置最大最小资源；最小资源是能保证的资源，最大资源是上限；

优先级: backport from 2.8；改进为7个优先级，每提升1个优先级，获得分配权重增加1倍；

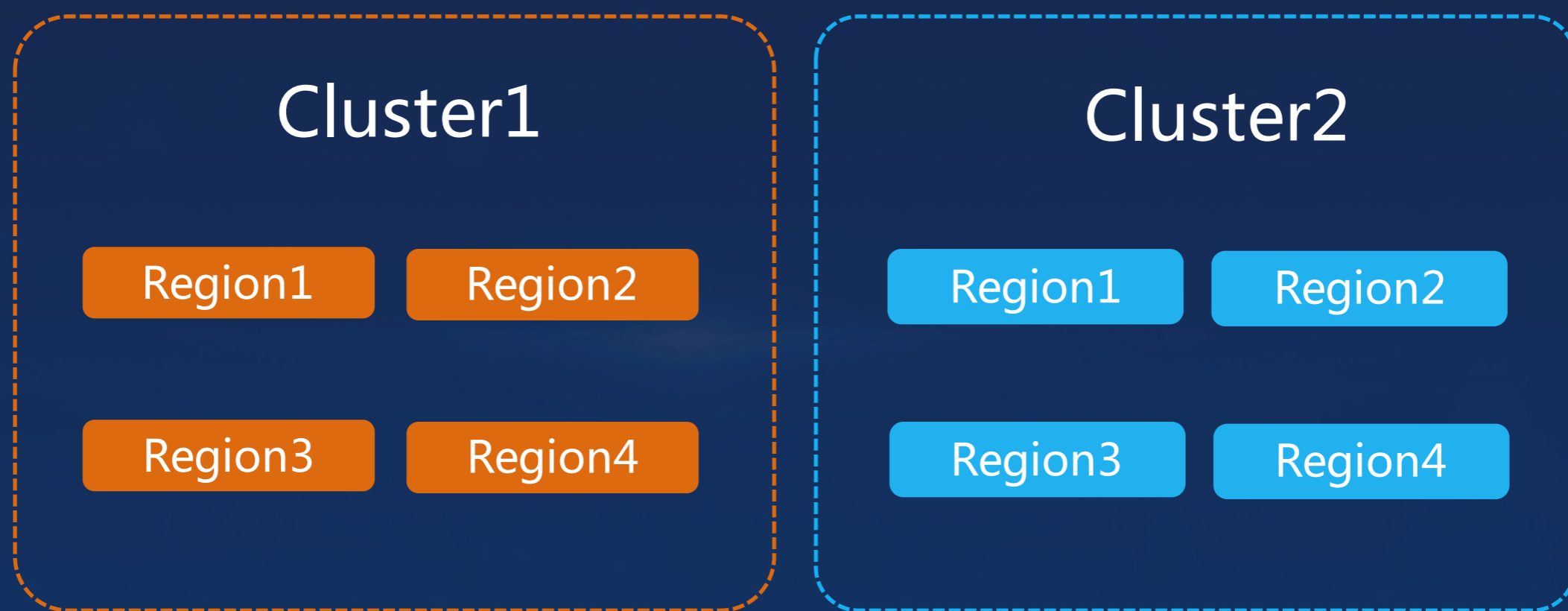
降低长GC的影响：尽快失败，合理设置 GCTimeLimit=90 GCHeapFreeLimit=10 ParallelGCThreads=8；

Cgroup

开启抢占

解决3. 差异化服务、物理隔离

共享型-Kafka



避免大集群：对ZK的压力，对controller的压力

逻辑分区：按照不同的业务场景分配，避免不同保障级别业务的相互影响

解决3. 差异化服务、物理隔离

独占型-Storm



独立的物理区域

物理CPU与虚拟机CPU 1:1 亲和性绑定

CPU ST从40%+降低至1%

痛点4. 排查问题耗时长，应用优化门槛高

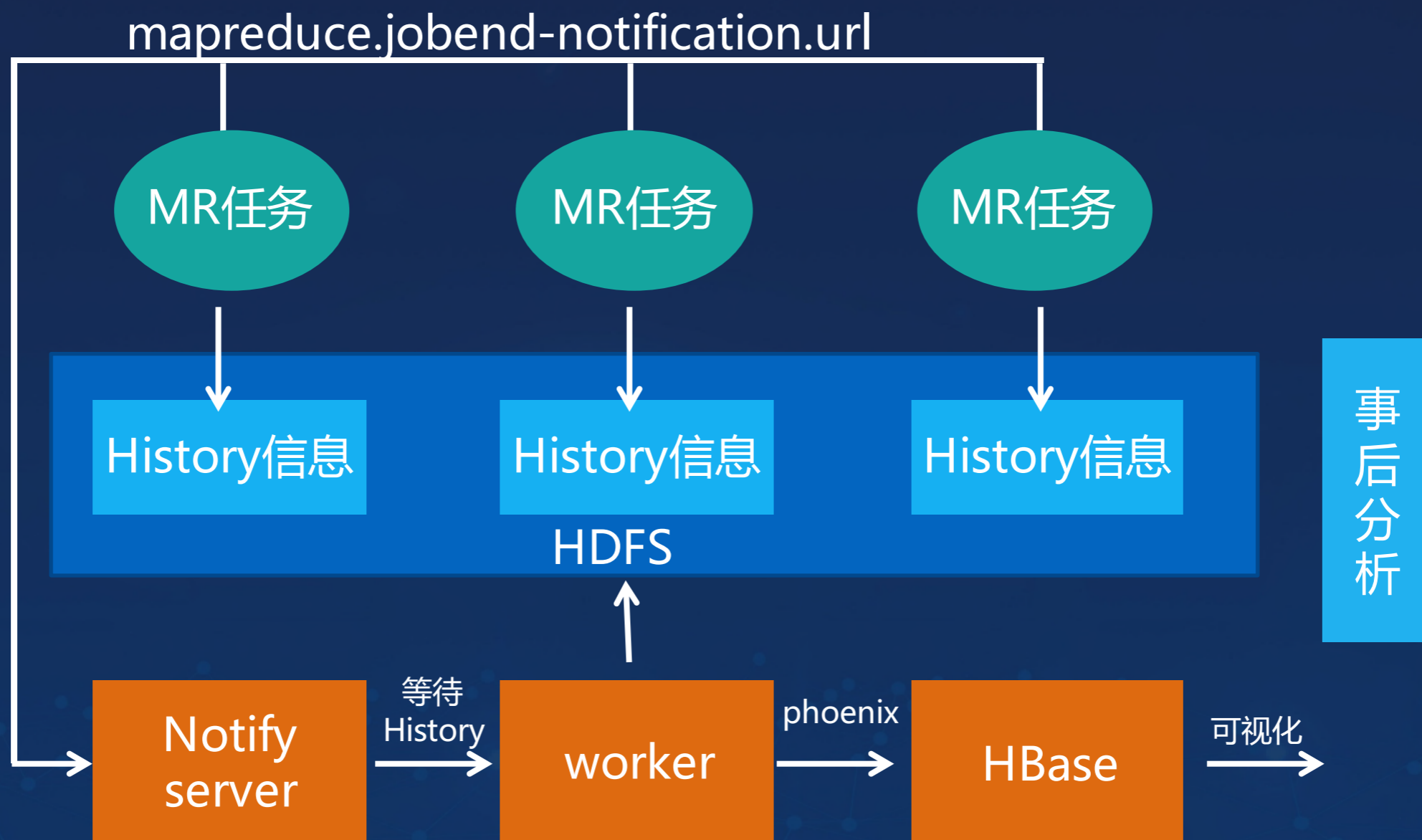


用户不清楚任务为什么失败

用户不清楚该如何优化任务

解决4. 智能诊断、优化建议

离线计算平台-MR任务分析平台



解决4. 智能诊断、优化建议

离线计算平台-MR任务分析平台

任务排行榜

评分最低Top20				处理数据量最多Top20			
排名	账户	jobID	评分	排名	账户	jobID	处理数据量
1	bi	job_1482378778761_19995963	44	1	sousuo	job_1482378778761_20032979	2.57TB
2	aps	job_1482378778761_19986972	44	2	sousuo	job_1482378778761_20039489	2.57TB
3	cloudytrace	job_1482378778761_19882779	44	3	srds	job_1482378778761_20002279	2.57TB
4	aps	job_1482378778761_20003396	44	4	spider	job_1482378778761_19879587	2.44TB
5	smmb	job_1482378778761_20040120	44	5	spider	job_1482378778761_19875138	2.39TB
6	smmb	job_1482378778761_19987516	44	6	safe	job_1482378778761_19953908	1.99TB
7	smmb	job_1482378778761_19987516	44	7	safe	job_1482378778761_19950502	1.99TB
8	lbi	job_1482378778761_19987516	44	8	safe	job_1482378778761_19961139	1.99TB
9	lbi	job_1482378778761_19888904	44	9	srds	job_1482378778761_19996728	1.46TB
10	birec	job_1482378778761_20032814	44	10	osmos	job_1482378778761_19879911	1.41TB

不同维度的排行榜

总CPU时间最长Top20				Task数排名Top20			
排名	账户	jobID	耗时	排名	账户	jobID	Task数
1	bi	job_1482378778761_19896644	3442184	1	sopdm	job_1482378778761_19987327	22908
2	bi	job_1482378778761_19903213	2259819	2	aps	job_1482378778761_19979424	20925
3	erp	job_1482378778761_20004500	2040751	3	aps	job_1482378778761_19977850	20391
4	erp	job_1482378778761_20010608	1533820	4	ztbd	job_1482378778761_19916949	17344

解决4. 智能诊断、优化建议

离线计算平台-MR任务分析平台



解决4. 智能诊断、优化建议

离线计算平台-MR任务分析平台

任务详情 > job_1499408566046_8795319 输入任务ID

基本信息

账户名称	bjredbaby
集群名称	common_yarn_prd3
任务名称	create table tmp_zq_zs_order_06 stored ...)(Stage-1)
任务状态	FAILED
提交时间	2017-08-10 17:14:18
结束时间	2017-08-10 17:15:37
总CPU时间	1.25min
处理数据量	0MB
map/reduce	1 / 4



失败原因

失败原因 数字格式化异常
建议解决方案 ["检查数字格式化"]

失败诊断信息

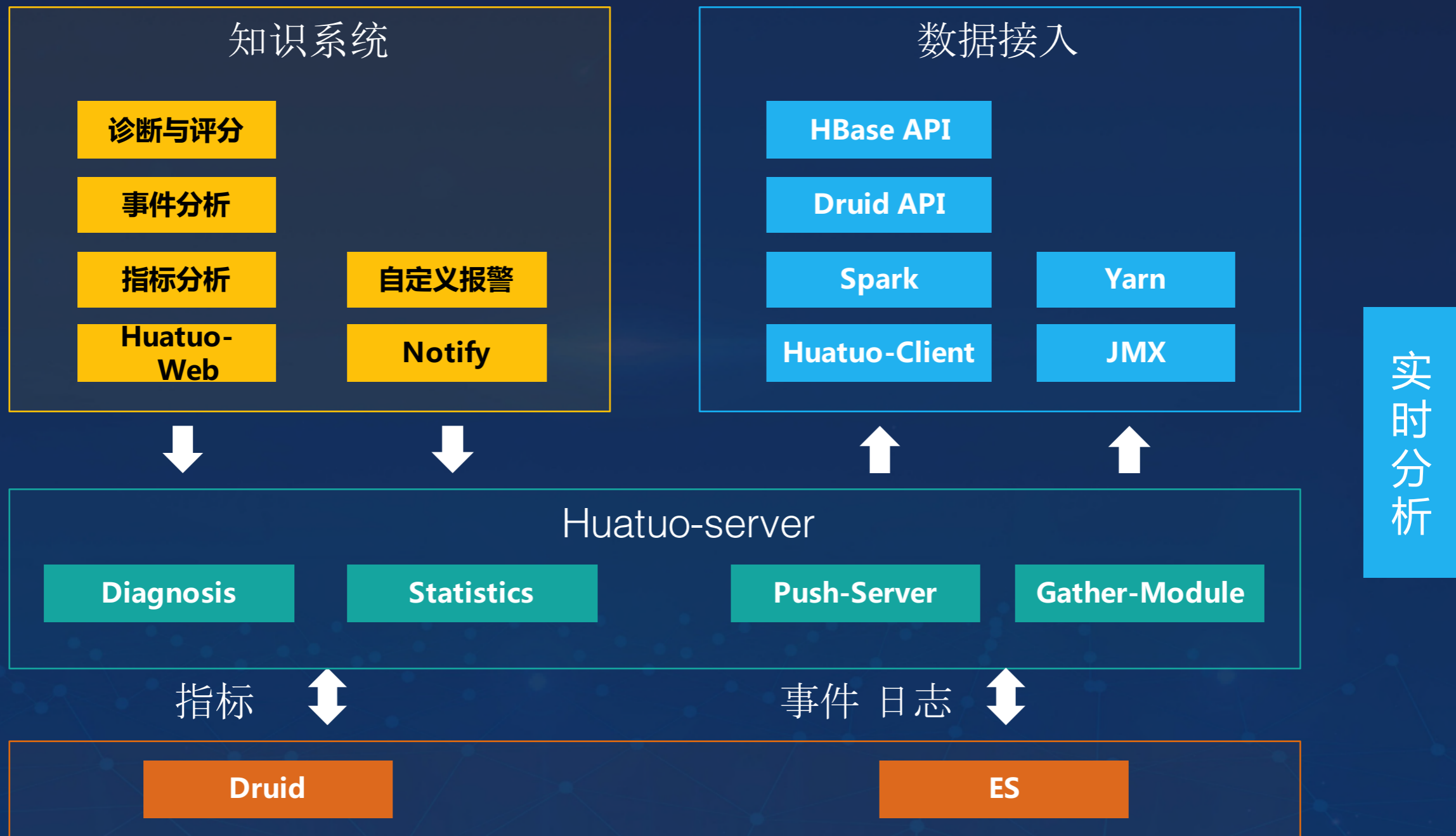


异常日志信息

Log Type: stderr Log Length: 243 log4j:WARN No appenders could be found for logger (org.apache.hadoop.metrics2.impl.MetricsSystemImpl). log4j:WARN

解决4. 智能诊断、优化建议

离线计算平台-华佗智能诊断平台(spark)



解决4. 智能诊断、优化建议

离线计算平台-华佗智能诊断平台(spark)

资源角度

- 宿主机状态分析
- HDFS资源使用分析

- Driver/Executor进程状态分析
- 资源利用率分析

- Cache 利用率分析
- Shuffle内存利用率分析

性能角度

- Task耗时链分析
- 长尾Task分析

- 任务调度Overhead分析
- Reduce并发度分析
- JDBC并发度分析Kafka读并发度分析

故障角度

- Shuffle数据倾斜
- HDFS Commit阻塞
- Executor丢失
- Spill事件分析
- 高维Parquet写性能诊断
- RDD Size Estimator耗时诊断
- 任务事件流

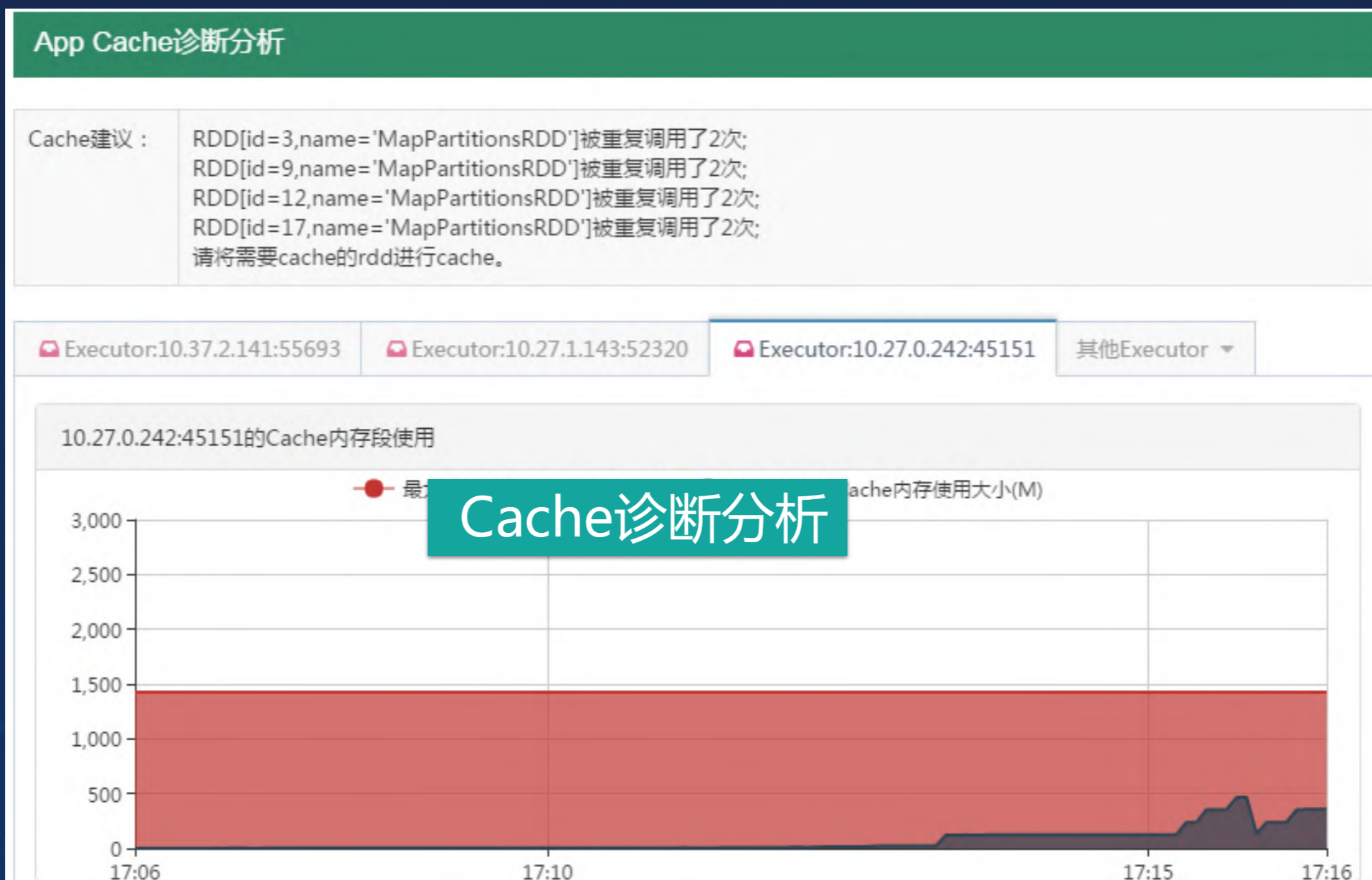
解决4. 智能诊断、优化建议

离线计算平台-华佗智能诊断平台(spark)



解决4. 智能诊断、优化建议

离线计算平台-华佗智能诊断平台(spark)



解决4. 智能诊断、优化建议

离线计算平台-华佗智能诊断平台(spark)



解决4. 智能诊断、优化建议

离线计算平台-华佗智能诊断平台(spark)

Shuffle数据倾斜诊断

事件：[INFO] ShuffleWriteStage[9.17]执行结束。任务耗时：4918MS。ShuffleWrite的Meta大小为：228字节。



解决4. 智能诊断、优化建议

流计算平台-健康状态监测系统

进程重启

代码异常

OutOfMemory

宿主机/虚拟机宕机

Kafka消息堆积

数据量暴增

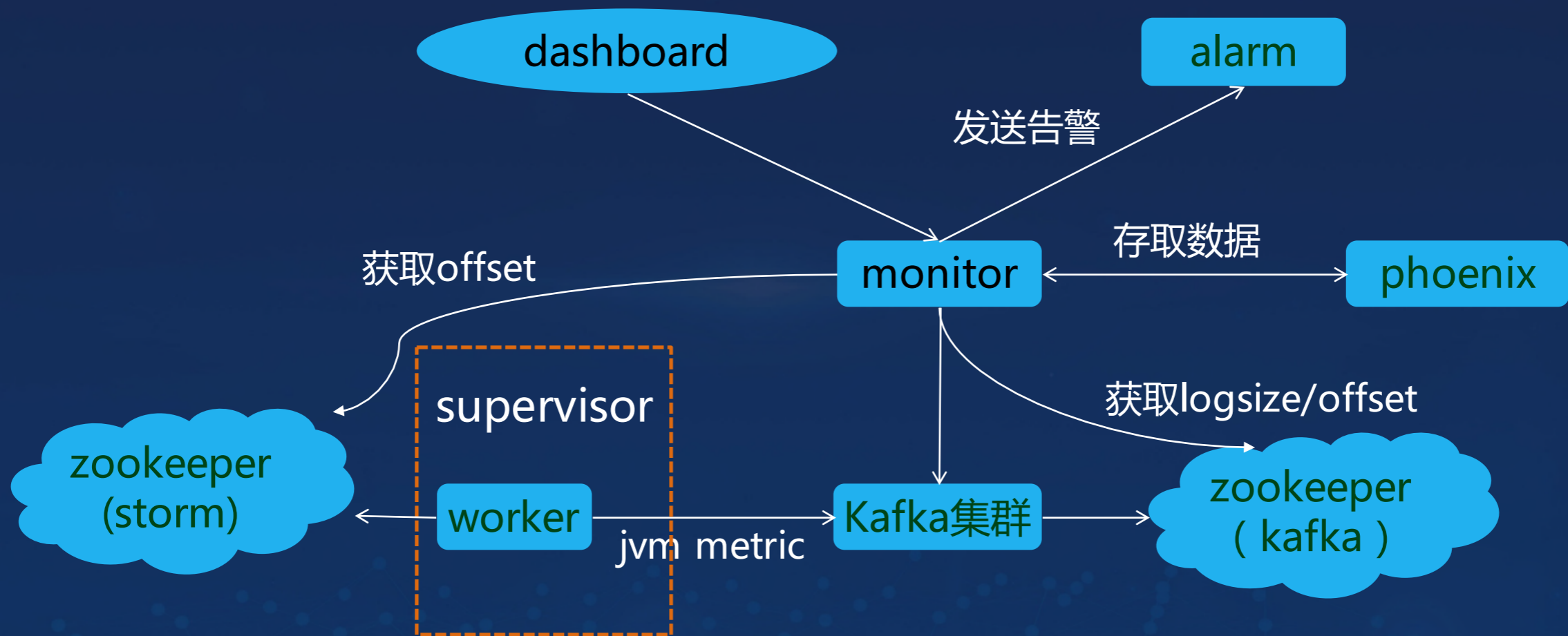
FGC/Capacity

并发不够

任务异常停止

解决4. 智能诊断、优化建议

流计算平台-健康状态监测系统



解决4. 智能诊断、优化建议

流计算平台-健康状态监测系统

ppcsHomePageB健康状态

2017年08月10日15时04分~2017年08月10日22时59分

1、Topology基本信息

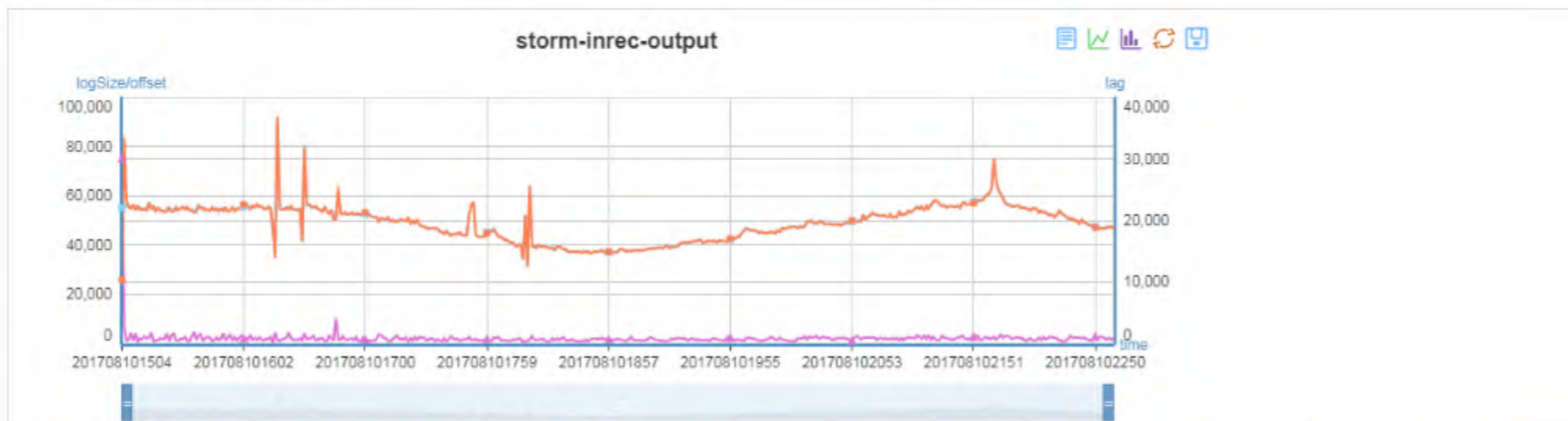
启动时长	Num workers	Num executors	Num tasks	spout.pending	message.timeout(s)	集群ID
7h54m52s	140	449	449	1000	360	ppcs-v-storm

2、Topology配置诊断

- topology的bolt并发设置合理，无需调整。此次启动也无worker发生迁移

3、数据消费视图（2017年08月10日15时04分~2017年08月10日22时59分）

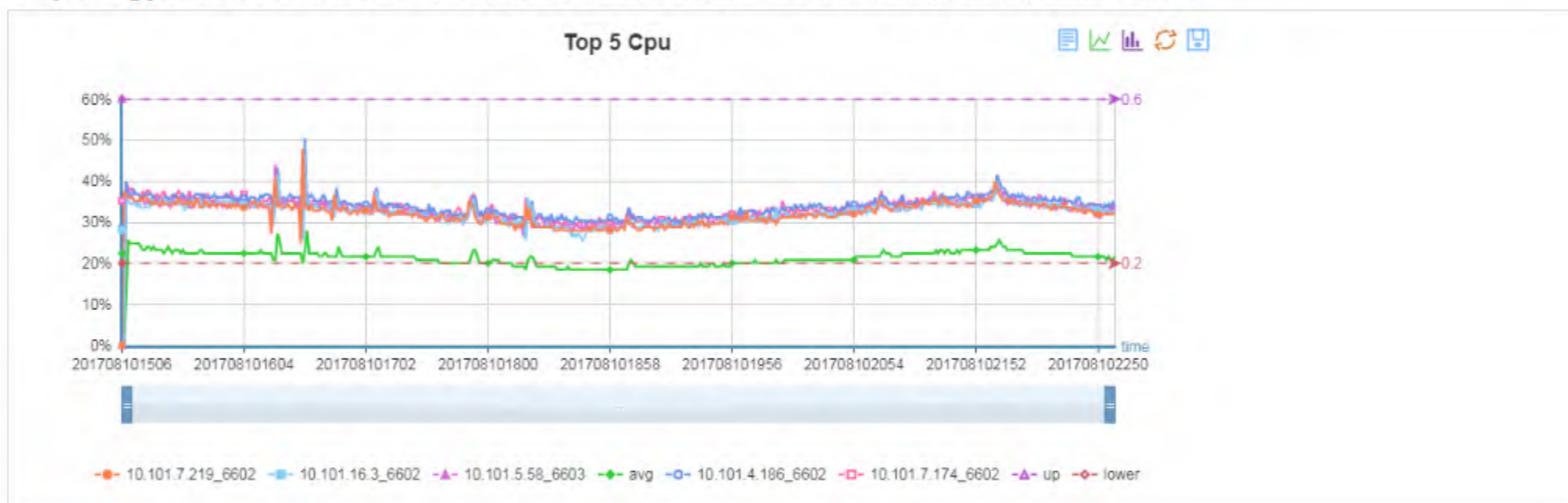
topic名: storm-inrec-output



解决4. 智能诊断、优化建议

流计算平台-健康状态监测系统

4、Topology运行CPU (2017年08月10日15时04分 ~ 2017年08月10日22时59分)



- worker:10.101.7.219_6602的cpu使用率远远高于平均使用率。
- worker:10.101.5.58_6603的cpu使用率远远高于平均使用率。
- worker:10.101.16.3_6602的cpu使用率远远高于平均使用率。
- worker:10.101.4.186_6602的cpu使用率远远高于平均使用率。
- worker:10.101.7.174_6602的cpu使用率远远高于平均使用率。

解决4. 智能诊断、优化建议

流计算平台-健康状态监测系统

5、Topology运行MEM (2017年8月10日15时4分 ~ 2017年8月10日22时59分)



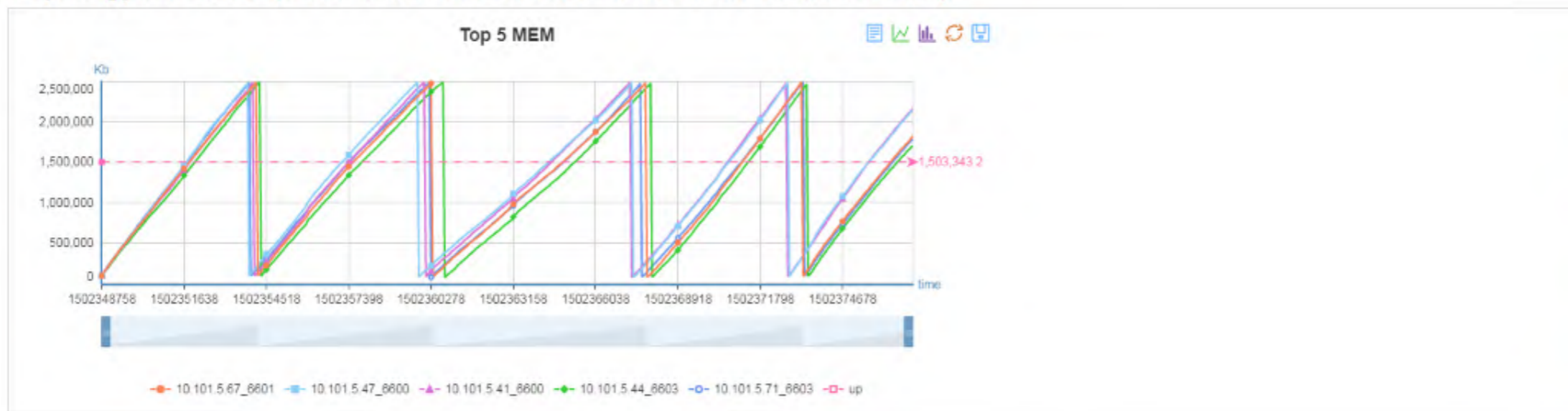
6、Topology运行Capacity (2017年08月10日15时09分 ~ 2017年08月10日22时59分)



解决4. 智能诊断、优化建议

流计算平台-健康状态监测系统

5、Topology运行MEM (2017年8月10日15时4分 ~ 2017年8月10日22时59分)



6、Topology运行Capacity (2017年08月10日15时09分 ~ 2017年08月10日22时59分)



TABLE OF CONTENTS

苏宁大数据平台基本介绍

大数据平台运维的痛点及解决方案

平台优化及增强

DOING & TO DO

平台优化及增强-稳定

- Hive metasever 连接数过高的问题

修改bonecp的配置：`maxConnectionsPerPartition=1`

- Spark Streaming & Druid System CPU过高的问题

设置`vm.zone_reclaim_mode=0`

- 透明大页导致System CPU过高的问题

`echo never >/sys/kernel/mm/transparent_hugepage/defrag`

平台优化及增强-安全

- 账户/权限体系：每个系统一个账户，不允许跨账户写。

- Hive metasever 密码加密

```
<property>
  <name>javax.jdo.option.ConnectionPassword</name>
  <value>OL2vkw+gDtQ=</value>
</property>
```

- 基于User/IP的访问控制策略：RPC层面控制，白名单

- skipTrash禁用：防止误删数据

平台优化及增强-扩展性

结合HDFS的瓶颈问题逐步优化

- 程序优化，扫全表: Hive慎用unix_timestamp方法
- 小文件合并
- YARN日志降低副本至1
- YARN日志单独放在另一个集群
- Federation + Alluxio 实现统一命名空间

TABLE OF CONTENTS

苏宁大数据平台基本介绍

大数据平台运维的痛点及解决方案

平台优化及增强

DOING & TO DO

DOING & TO DO

- ✓ Flink推广
- ✓ OLAP平台建设
- ✓ 流计算消息回溯
- ✓ 多活&灾备
- ✓ 资源统一管理（明年）

你认为以下哪些关键字对运维最重要？

智能

易用

高效

性能

稳定

流程

MTTR

自动化

安全

规范

审计

SLA

故障预测

个人号



公众号



THANKS!

智能时代的新运维

CNUTCon 2017