

# 腾讯海量监控包袱与创新

聂鑫

腾讯社交网络运营部运维负责人

# SPEAKER INTRODUCE

---

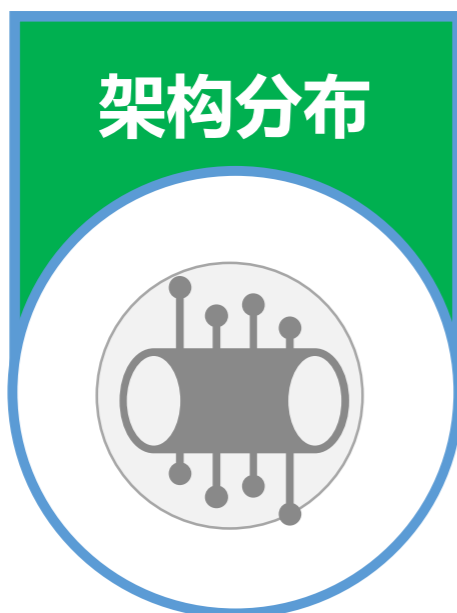
## 聂鑫

- 从开发到运维，伴随腾讯社交网络运营部成长的十年，负责过腾讯社交产品所有业务运维工作，目前主要负责QQ、空间等产品运维团队管理工作。
- 经历多个业务产品的诞生到蓬勃，伴随着运维团队的成长和成熟，见证着腾讯一代代运营技术的创新和发展。作为运维界老兵有好多故事想和大家讲，也特别愿意听听各位经历的酸甜苦辣。

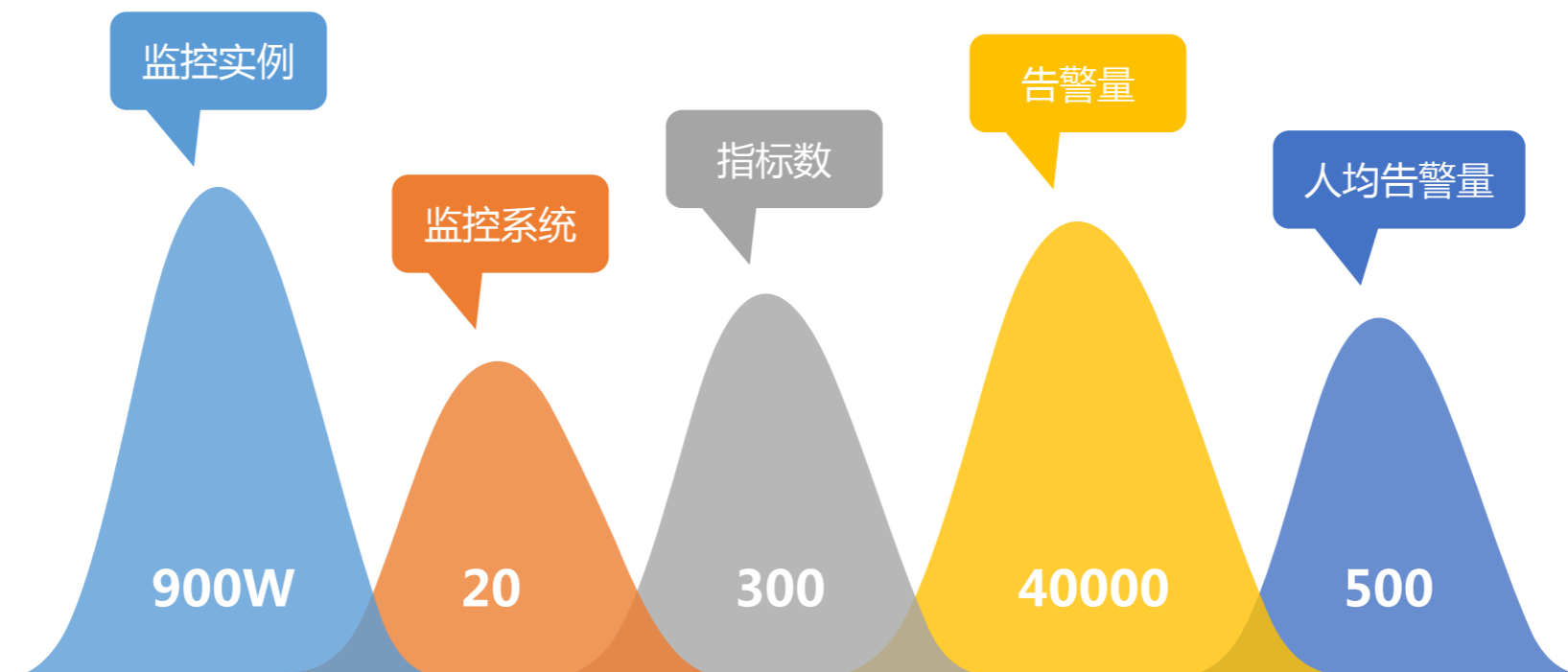


SPEAKER  
ArchSummit 2017 ShenZhen

# 腾讯 技术运营重点



短信告警 **5** 万条/天  
单人最高**1500**条/天



- 腾讯正在做哪些监控
- 有哪些不一样的地方
- 有哪些值得关注点

# 正在做哪些监控

在监控领域有三个主题 { 快、准、全 }  
他们永远是矛盾的，调和矛盾成了运维  
技术或艺术



# 监控体系演进

2006

N:网络质量监控  
B:网管基础监控  
A:自动化测试  
M:模块间调用  
S:测速系统

2007

W:站点分析系统  
m:模块监控  
L:容量管理

2009

I:L5组件监控  
Y:一致性  
Monitor特性监控

2011

S:QZ组件监控  
F:设备特性监控  
R:返回码监控

2013

C:CDN监控  
P:Ptlogin监控  
D:存储质量  
H:客户端环境  
R:ROOT根源分析

2014

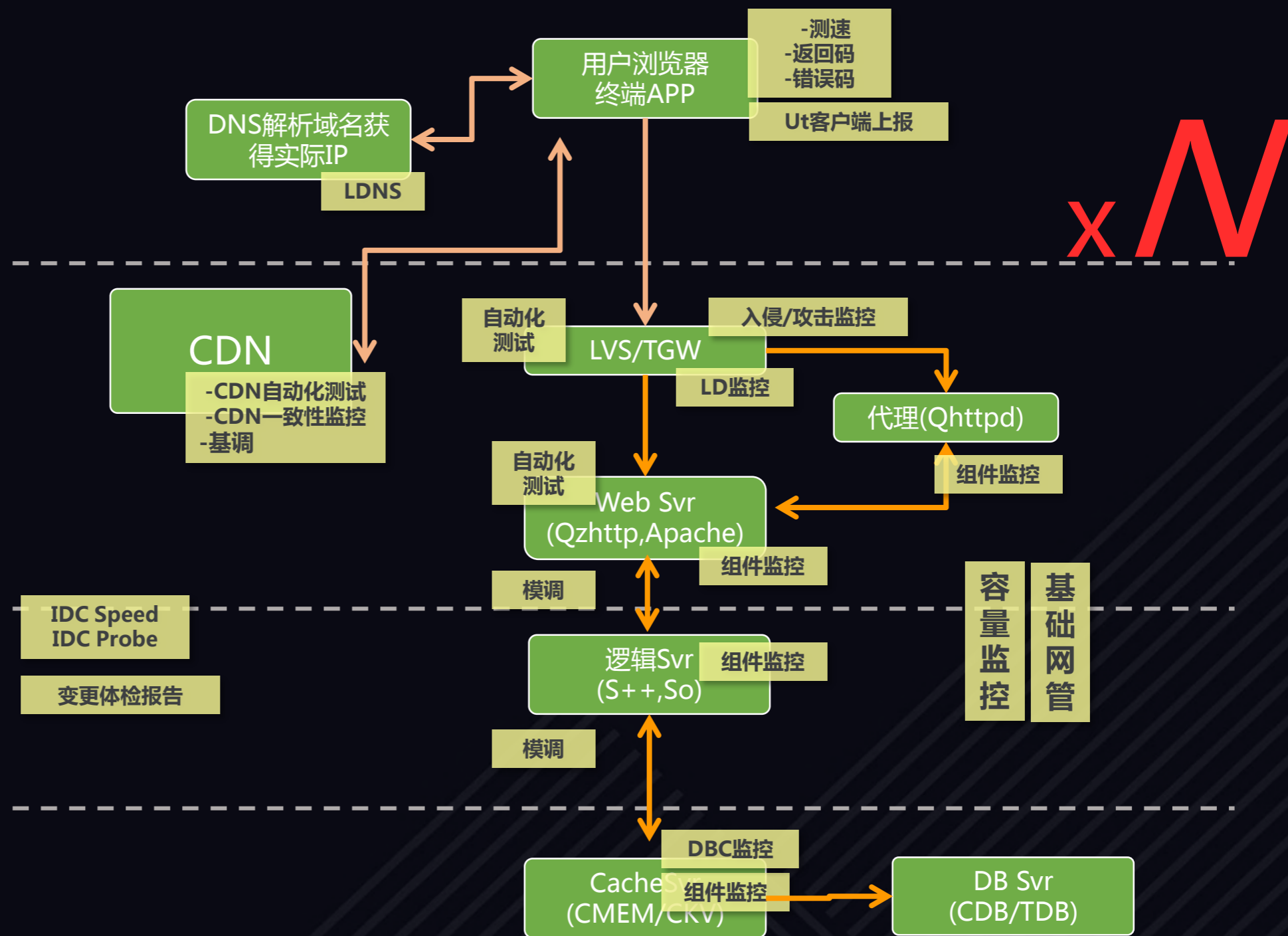
I:华佗移动端分析  
Q:舆情监控  
G:日志大数据分析  
U:UIN染色分析

2015

DLP核心指标  
Habo多维分析

2016

Q调拨测监控  
移动端卡慢  
全链路日志



# 覆盖完整

## 用户端监控

- 测速
- 返回码
- 自动化测试
- 基调
- 移动分析(mua)
- html5

## 业务侧监控

- 核心产品指标
- 各纬度业务指标
- 攻击防御
- 舆情监控

## 服务内监控

- 模块间调用
- L5失败率
- 组件监控
- 强制一致性
- Gslb、lvs

## 基础资源

- 丢包断线
- 死机重启
- 硬件故障
- 容量监控

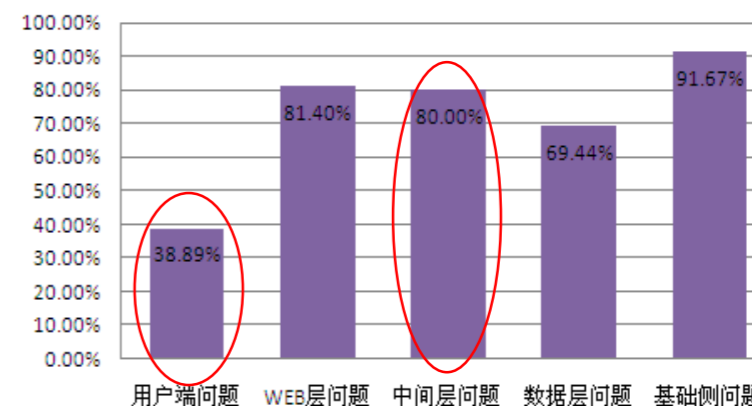


# 多、大、杂

## 业务增长 VS 监控系统发展

	2009年	2010年	2014年	2017年
主要监控系统数	9	11	20+	18
主要监控指标数	132	178	300+	400
监控实例数	-	45w	900w	2000w
基础告警数/天	192	300	3000+	5000+
业务告警数/天	362	390	3.9w	4.97w
个人告警量/天 (包含运维开发)	-	最大:177 平均:13	500+	1500 184

### 基于架构层的监控覆盖率



- 客户端、数据层监控覆盖率弱;
- 系统建设离散, 监控数据分散;
- 综合分析能力弱, 定位时间长;
- 告警数量过多;

# 有哪些不一样的地方

腾讯织云  
CLOUD  
MANAGEMENT

放下包袱来 **创新**

不是破旧立新而是尊重历史

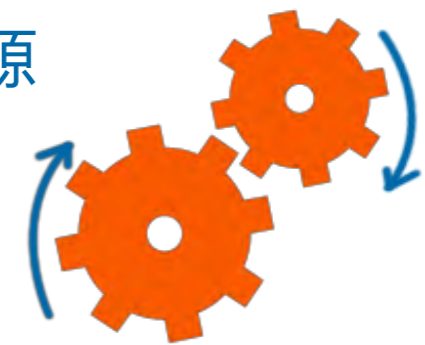
**坚决优化**历史演进中的架构落后

# ROOT

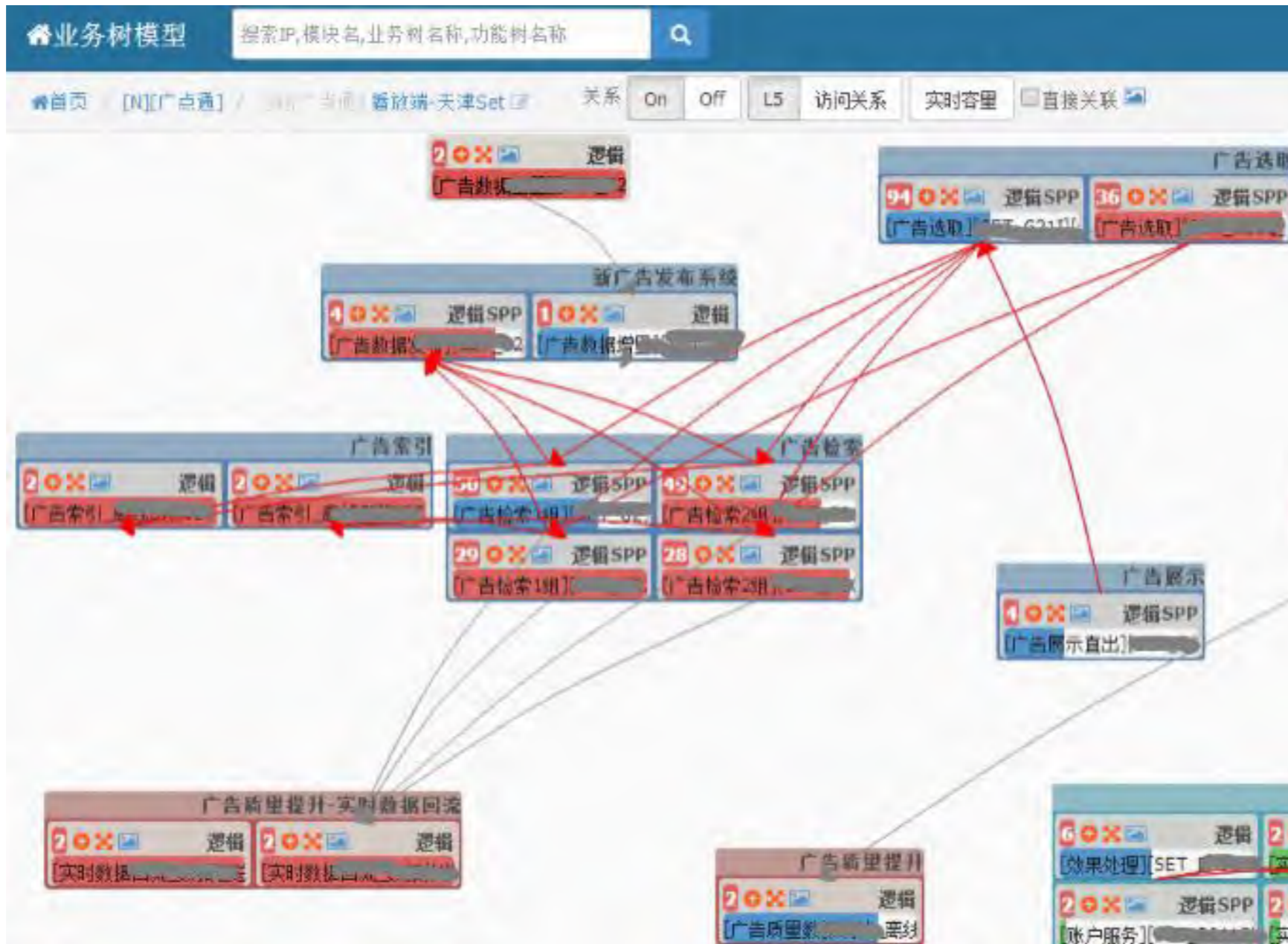
根源智能分析法



基于业务架构，结合数据流关系，通过时间相关性、面积权重等算法，将监控告警进行筛选分类，发掘有业务价值的告警，并直接分析给出告警根源

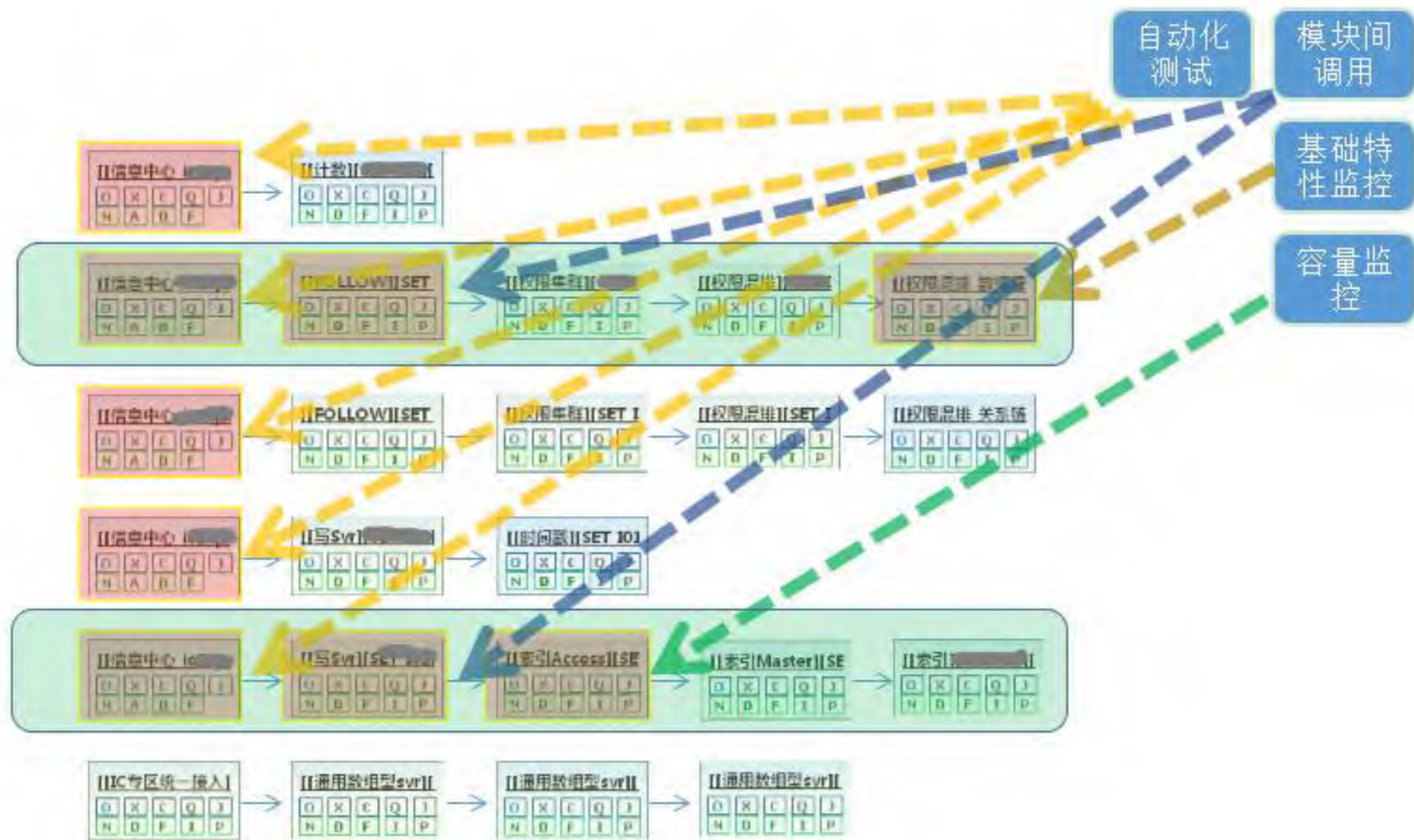


# ROOT原理：多维关系降维

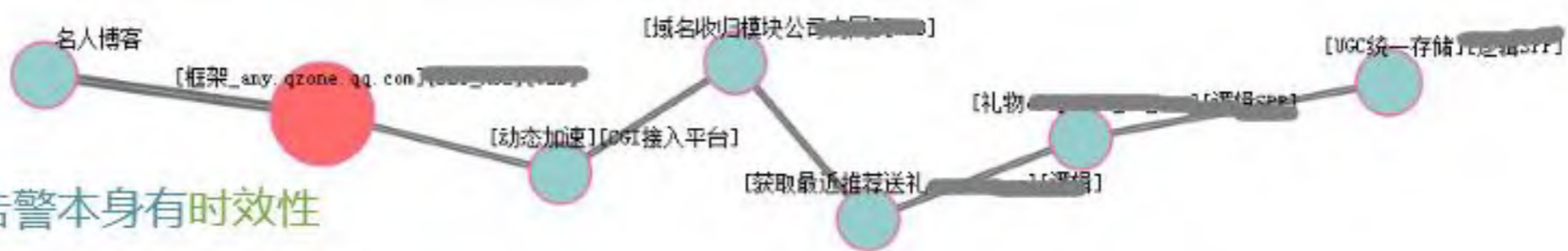




# ROOT原理：告警叠加



# 时间片与时间相关性



- 告警本身有时效性
- 时效性源于告警延时
- 连续性可能是干扰
- 链路相关性和时间相关性一起决定准确性

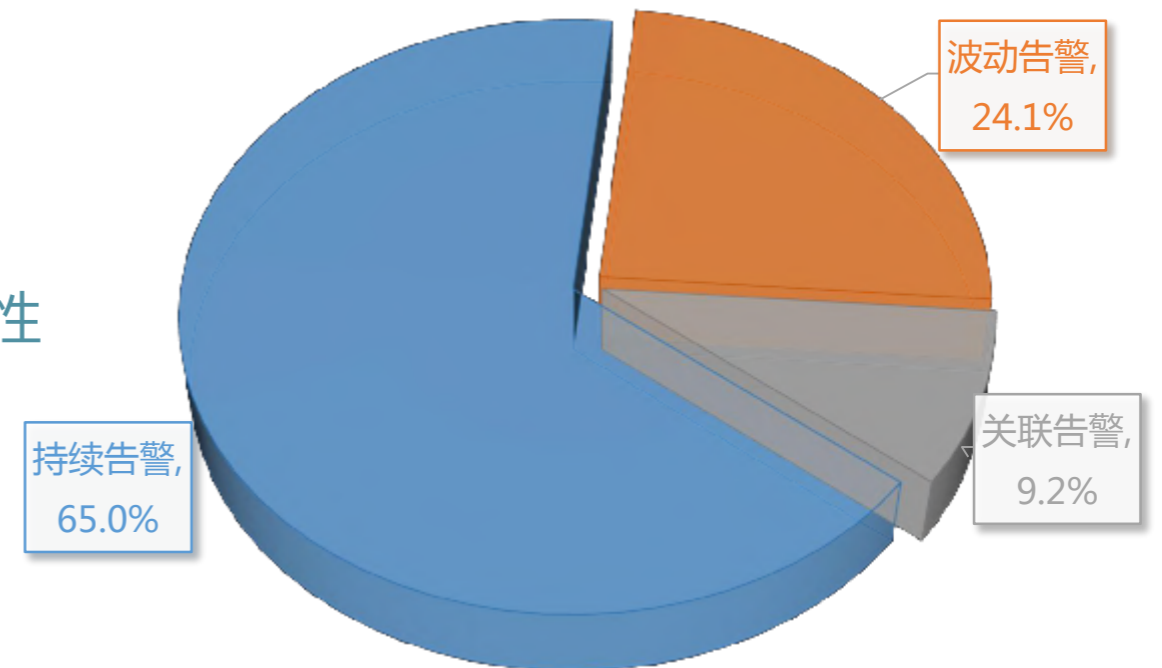
# 分类筛选和选择性处理

- **原因告警**、**现象告警**

- **原因告警** : 往往是造成故障的根源, 却往往无需处理
- **现象告警** : 故障的结果, 往往看不出根源, 需要分析

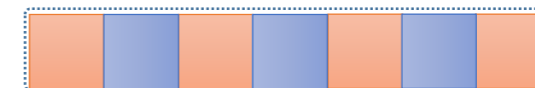
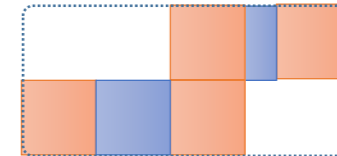
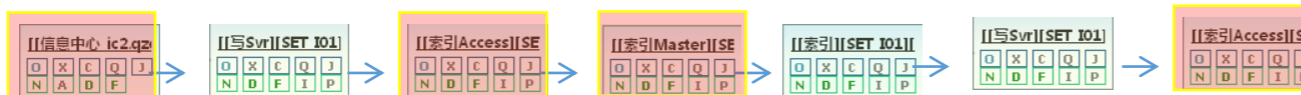
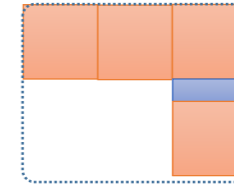
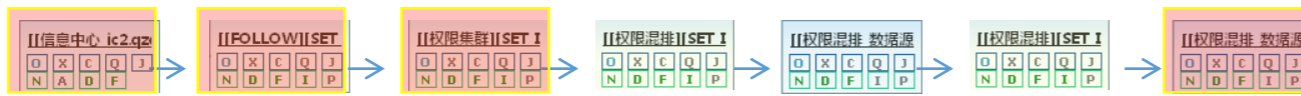
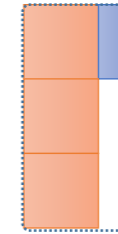
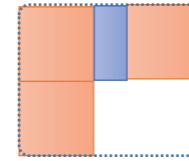
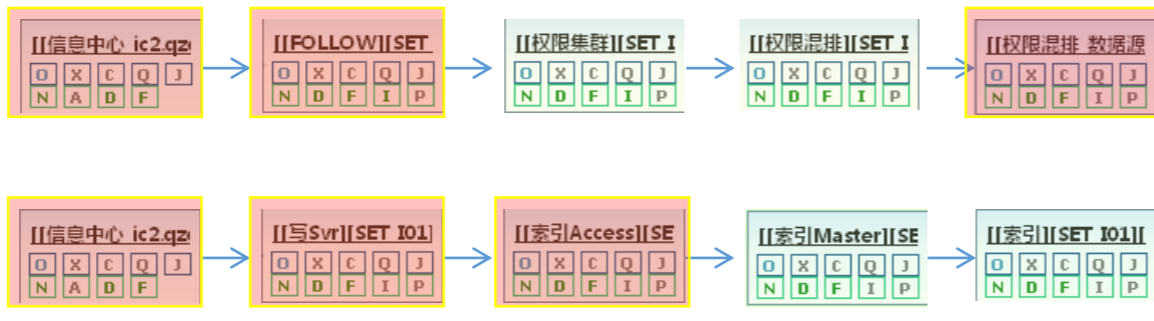
- **持续告警**、**波动告警**、**关联告警**

- **持续告警** : 不紧急、不重要
- **波动告警** : 业务重要性决定告警重要性
- **关联告警** : 有因有果, 即时处理





# 权重与面积算法



# 算法案例

## 1、链路中告警模块数=1

长=1 (只有一个模块告警时固定为1), 宽=(1+告警模块所在链路的序号除以链路总模块数), 面积=长\*宽= $1 * (1 + (iarr+1) / lnkcout) * 100$

- a、1-0-0-0, 权重面积= $1 * (1 + (0+1) / 4) * 100 = 125$  ;
- b、0-1-0-0, 权重面积= $1 * (1 + (1+1) / 4) * 100 = 150$  ;
- c、0-0-0-1, 权重面积= $1 * (1 + (3+1) / 4) * 100 = 200$  ;

备注：链路中只有一个模块告警，并且结合业务链路生成的特性，告警模块越靠后，权重面积越大；

## 2、链路中告警模块数>1

长=链路中连着告警模块的最大个数(iarrmax), 宽=连着或不连着告警模块宽都为 $1 + 1 / (连着不告警的模块个数)$ , 面积=长\*宽= $iarrmax * (1 + 1 / N + \dots) * 100$

- a、1-0-0-0-1, 权重面积= $1 * (1 + 1/3 + 1) * 100 = 233$  ;
- b、1-0-0-1-0, 权重面积= $1 * (1 + 1/2 + 1) * 100 = 250$  ;
- c、1-1-0-0-1, 权重面积= $2 * (1 + 1/2 + 1) * 100 = 500$  ;
- d、1-1-0-1-0, 权重面积= $2 * (1 + 1/1 + 1) * 100 = 600$  ;
- e、1-1-0-0-1-1-0-0-1, 权重面积= $2 * (1 + 1/2 + 1 + 1/2 + 1) * 100 = 800$  ;
- f、1-1-0-0-1-1-0-1-1, 权重面积= $2 * (1 + 1/2 + 1 + 1/1 + 1) * 100 = 900$  ;
- g、1-1-1-0-1-0-0-1-1, 权重面积= $3 * (1 + 1/1 + 1 + 1/2 + 1) * 100 = 1350$  ;

## 3、特殊情况：

1、链路中，前面模块都没有告警，但最后模块连着告警（相当于链路中全模块告警），权重面积\*10；

2、链路中，模块全告警，权重面积\*10；

- a、0-0-0-1-1, 权重面积= $(2 * 1 * 100) * 10 = 2000$  ;
- b、1-1-1-1-1, 权重面积= $(5 * 1 * 100) * 10 = 5000$  ;

# DLP

## 业务生死指标

衡量业务死、活的指标

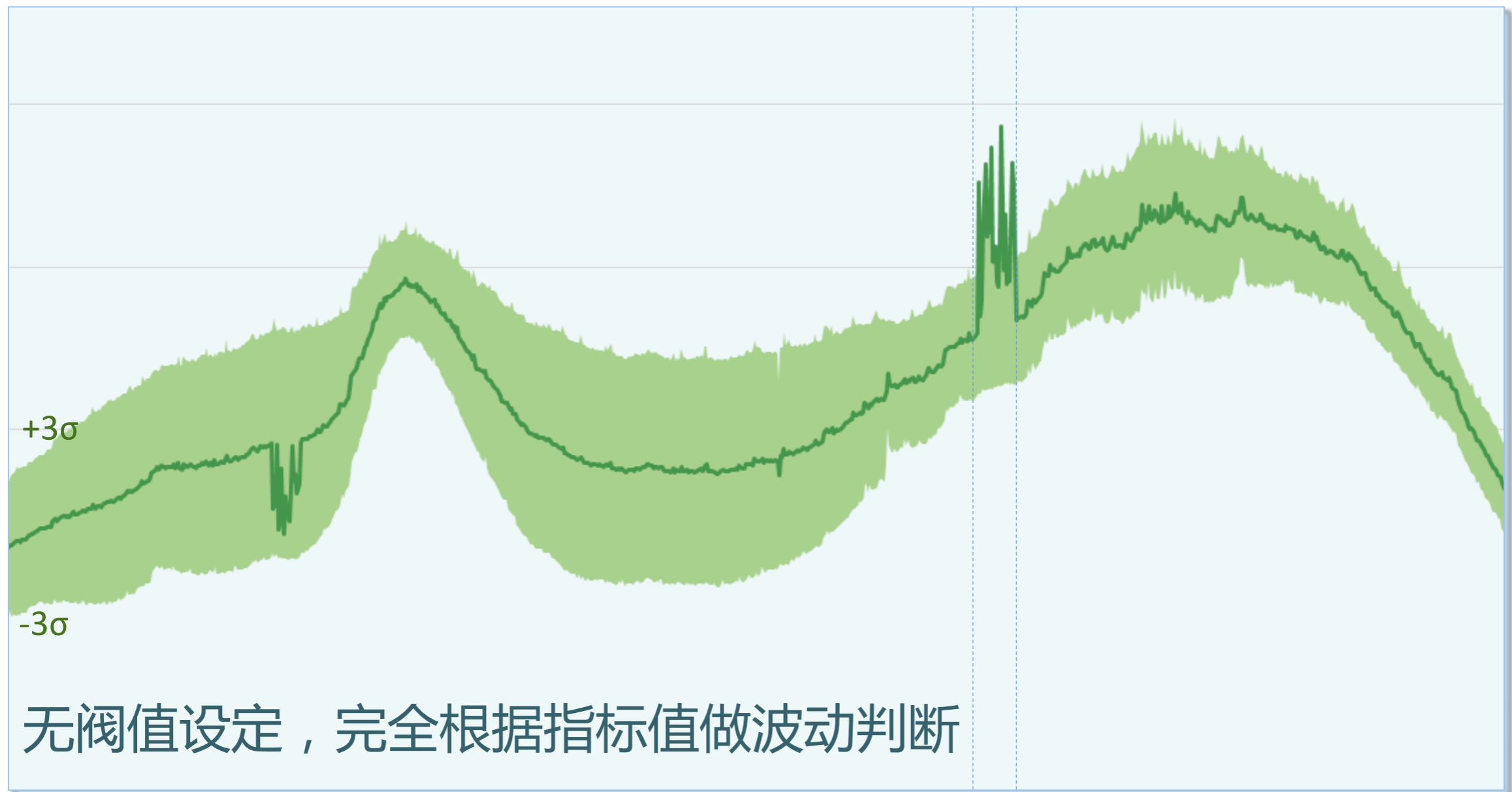


无阈值设定，完全根据指标值做波动判断

一个服务只能有一个生死指标

不建议用业务指标做生死指标

# 自定义-告警泛滥的罪魁祸首



# 一个服务只能有一个生死指标

日期: 2016-12-01	当天图	可疑属性 1小时
132628	PushData收到SyncData量 (无线侧量)	
111086	延时统计:最高延时 111086 延时统计:本轮计数 165616 延时10s以上的SyncData量_warning	
201739	已分配文件句柄数 201739 系统最大文件句柄数	
283451	日志客户端从Channel取数据成功 283549 发送日志包成功 283550 发送日志包时连接未建立	
283817	发送成功日志量 424627 日志发送流量	
222659	albs寻址成功 231520 发送流水日志 283448 日志写入Channel成功	
205427	US_GetUser失败-指定app查询 205428 US_GetUser失败-拉定制实例	
205429	US_GetUser失败-拉所有实例 205478 US_GetUser失败-拉主显状态 205664 非法优先级数据	
205688	不识别的状态值 205939 从US中读到不识别的App 210896 US_GetUser失败-推状态	
205515	设置msf push/幸push改变时间 205518 msf离线转app在线 205517 msf离线转push在线	
205518	msf app在线转push在线 205519 msf push在线转app在线	
187462	接入配置API: 拉取配置超时 196613 接入配置api: 拉配置发回成功	
196616	接入配置api: 拉配置回包成功 196625 接入配置api: ip寻址成功 222130 主循环次数	
222133	需要拉取配置次数 222134 拉取配置失败 222135 拉取配置成功 222136 TopSendRecv超时	
222140	OC返回配置数据有效 222141 Agent更新配置 222145 ClientAPI上报版本信息成功	
222147	ClientAgent上报版本信息成功	
251103	agent上报状态失败 251115 agent上报状态正确 251116 agent上报状态错误, iCoId=0	
251117	agent上报状态错误, iCoId=1 251134 上报状态成功 251135 更新状态上报	
251136	收到下发APP配置	
330388	recv收包超时丢弃数	
447804	插件msr同步成功 447805 插件msr同步失败	
666602	入包总数 666605 入包送达上层协议 666606 出包数 666608 分片重组超时	
666609	入包需重组 666610 分片重组成功 666611 分片重组失败 666612 分片成功	
666614	创建分片数	

DLP

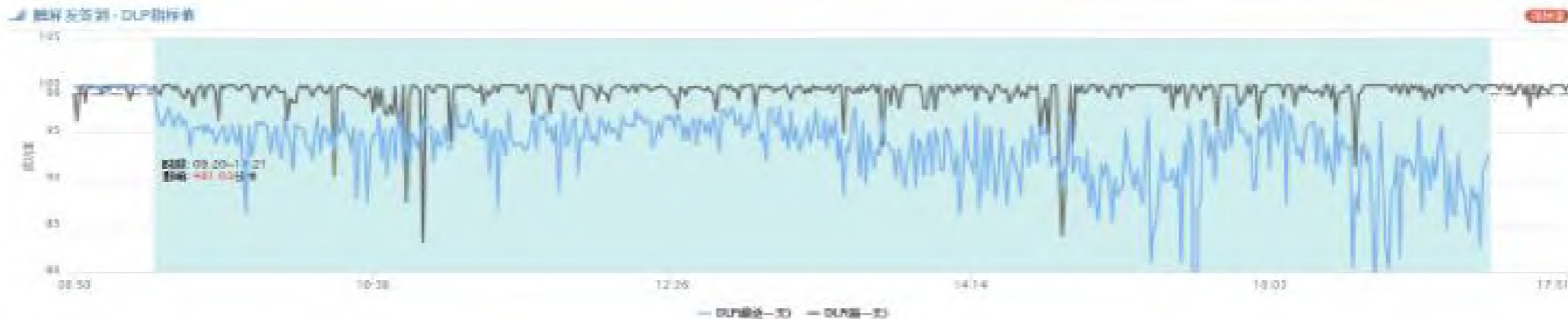


# 案例跟踪

告警时间	DLP名称	告警级别	告警详情	处理状态
2015-09-16 09:20:00	数据互连	数据汇聚-4001:19126,0.921,数据互连1		已处理

viewall | selected | warning

变更	基础告警	公告	主网IP	公网IP	返回码
			IP, 失败占比, 失败次数	IP, 失败占比, 失败次数	返回码, 失败占比, 失败次数
			10.213 3.01%, 625	10.213 2.31%, 19180	-4001, 92.1%, 19126
			10. 2.96%, 614	10. 56%, 116	-10305, 5.6, 1156
			10.2 2.2, 2.92%, 606	10.21 54%, 112	-11174, 1.3, 222
			10.21 4, 2.89%, 600	10. 5.51%, 106	-19997, 0.3, 98
			10.213 2.87%, 597	10.2 40%, 84	-10085, 0.4, 92



## 访问链路



告警

汇聚

关联

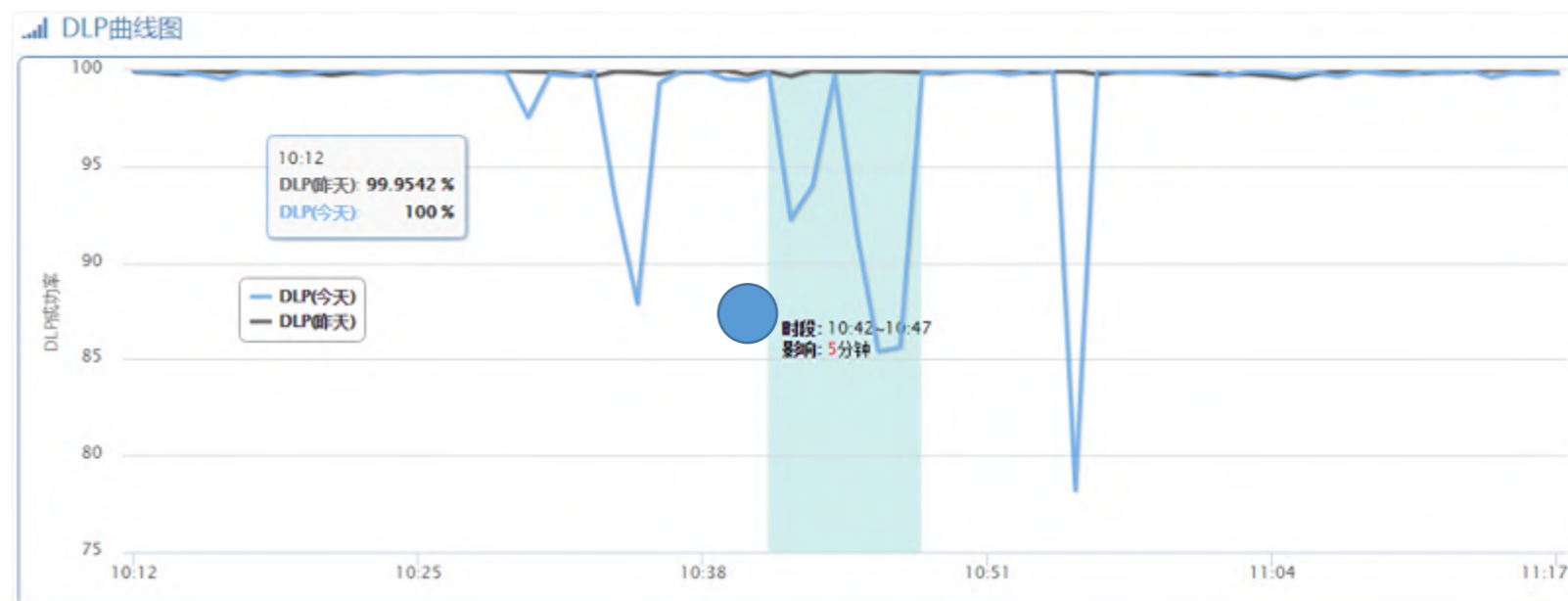
对比

架构

# 关联计算

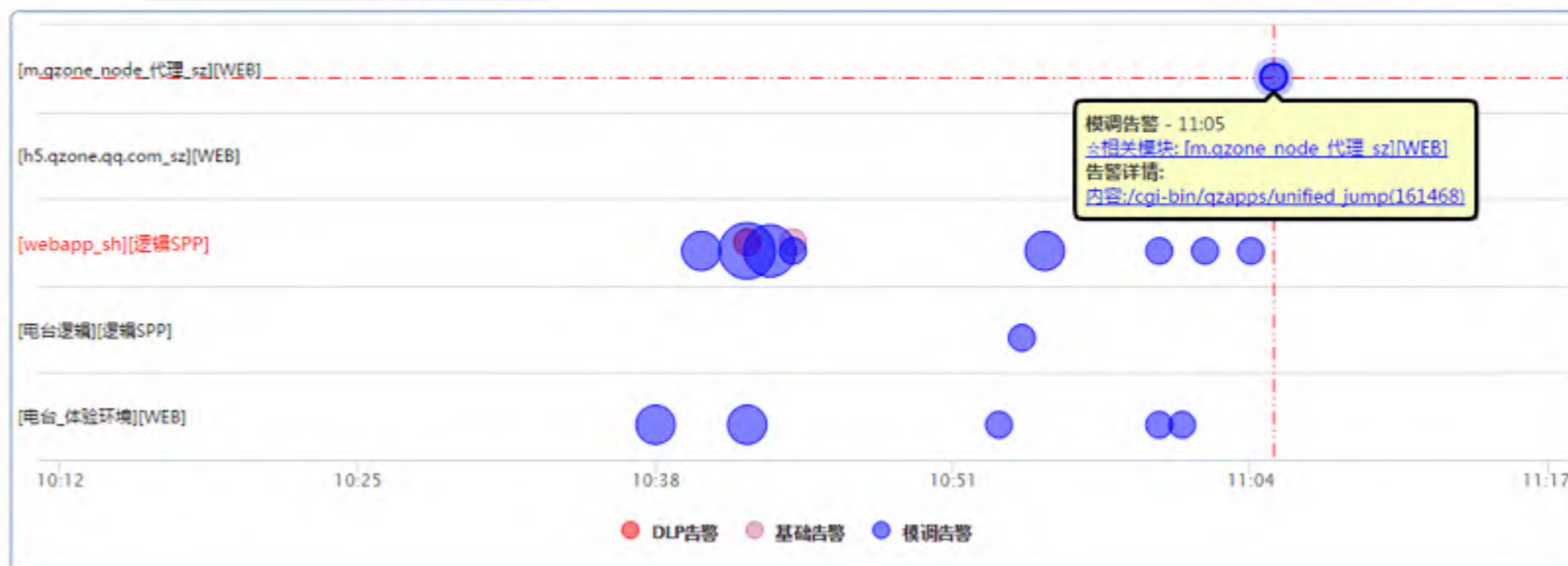
自身故障/变更  
网络故障公告  
DBC故障公告

关联链路告警



## 访问链路分析

链路选择:





# 全链路监控

帮业务组织数据

无接入门槛的数据组织方式

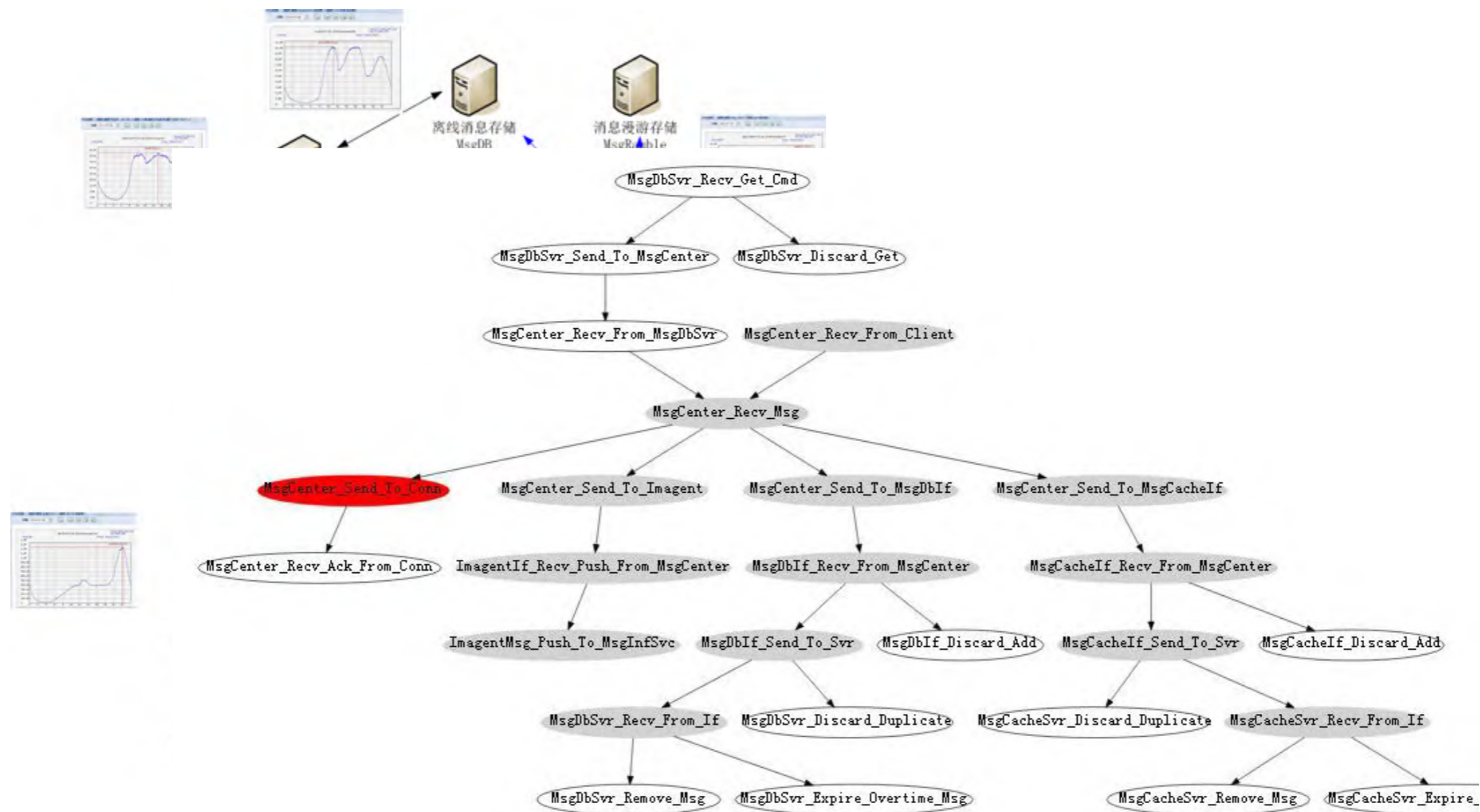


不把数据多当包袱

界定数据的生产者与消费者

帮助生产者消费数据

# 早期染色监控



一条消息在系统内的51个状态

# 兼容各种数据源

织云基础监控  
织云特性监控  
其他已有数据源



织云多维监控



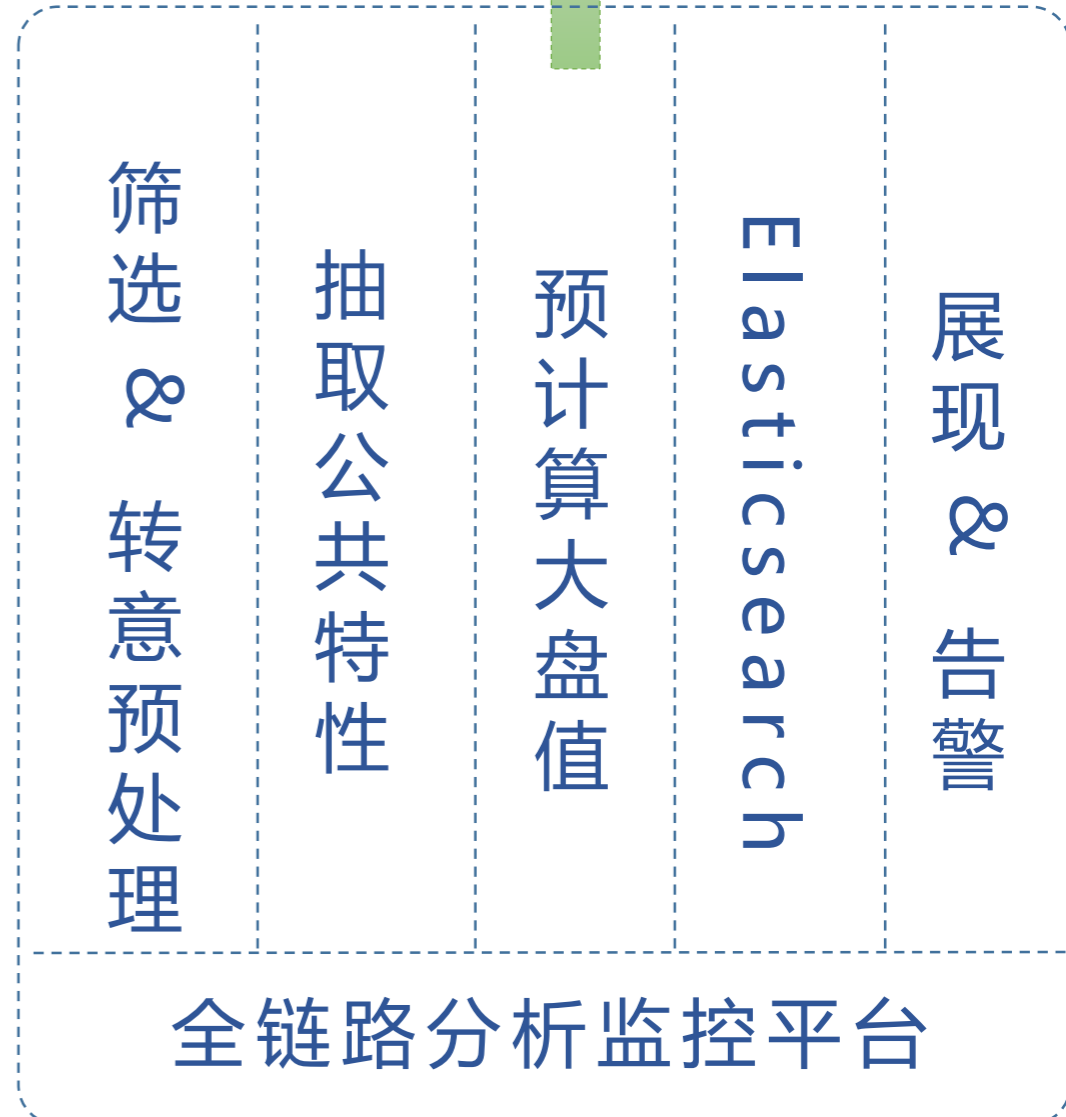
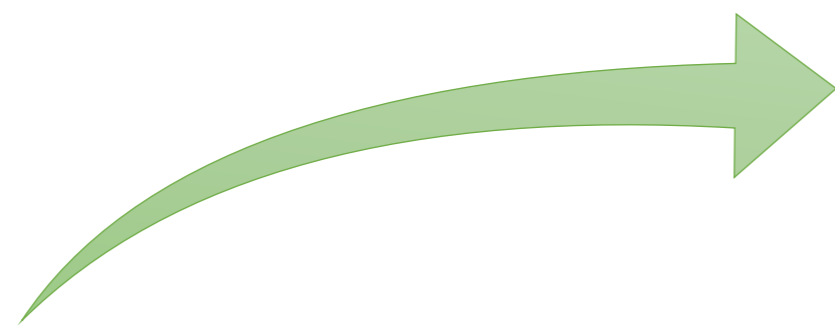
现网各类日志



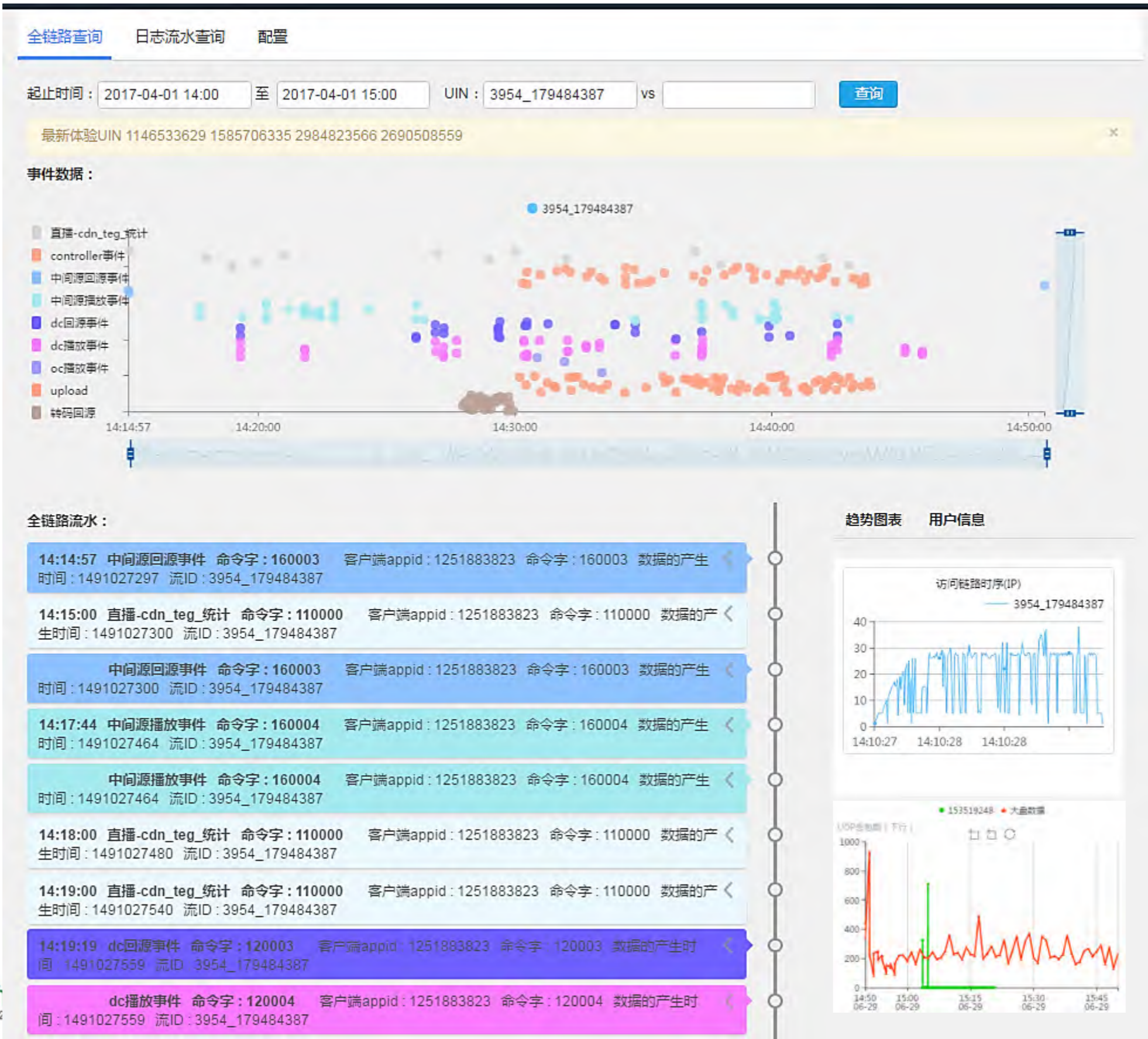
织云日志中心



业务格式数据  
织云舆情监控



# 各类数据的再利用





# 数据中挖掘各种纬度

数据源: 空间点播APP(dc00895)

查询

## 成功率TOP10

事件	成功率	事件量	耗时
video_play_ret	99.96700	1409087485	0.000
video_downloa...	99.98400	1116061822	0.000
video_hotlink_...	99.98900	1382539093	0.388
video_downloa...	100.00000	220615050	0.000
video_link_redi...	100.00000	861936568	0.070
video_link_redi...	100.00000	867437443	0.000
video_downloa...	100.00000	1082542483	0.046

## 事件量TOP10

事件	成功率	事件量	耗时
video_play_ret	99.96700	1409087485	0.000
video_hotlink_...	99.98900	1382539093	0.388
video_downloa...	99.98400	1116061822	0.000
video_downloa...	100.00000	1082542483	0.046
video_link_redi...	100.00000	867437443	0.000
video_link_redi...	100.00000	861936568	0.070
video_downloa...	100.00000	220615050	0.000

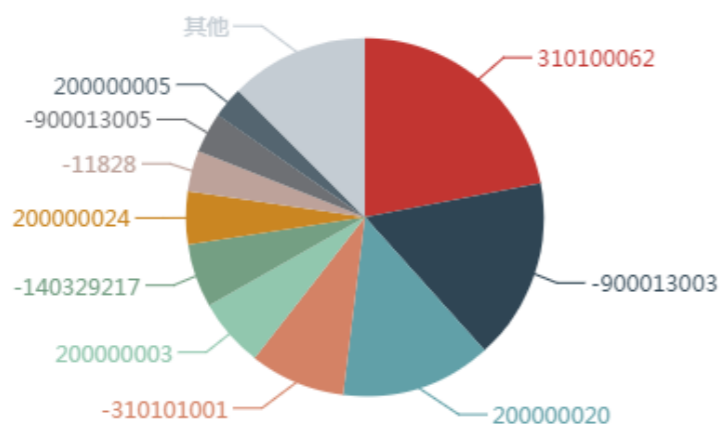
## 耗时TOP10

事件	成功率	事件量	耗时
video_hotlink_...	99.98900	1382539093	0.388
video_link_redi...	100.00000	861936568	0.070
video_downloa...	100.00000	1082542483	0.046
video_downloa...	99.98400	1116061822	0.000
video_downloa...	100.00000	220615050	0.000
video_link_redi...	100.00000	867437443	0.000
video_play_ret	99.96700	1409087485	0.000

## 事件video\_play\_ret异常返回码占比分析



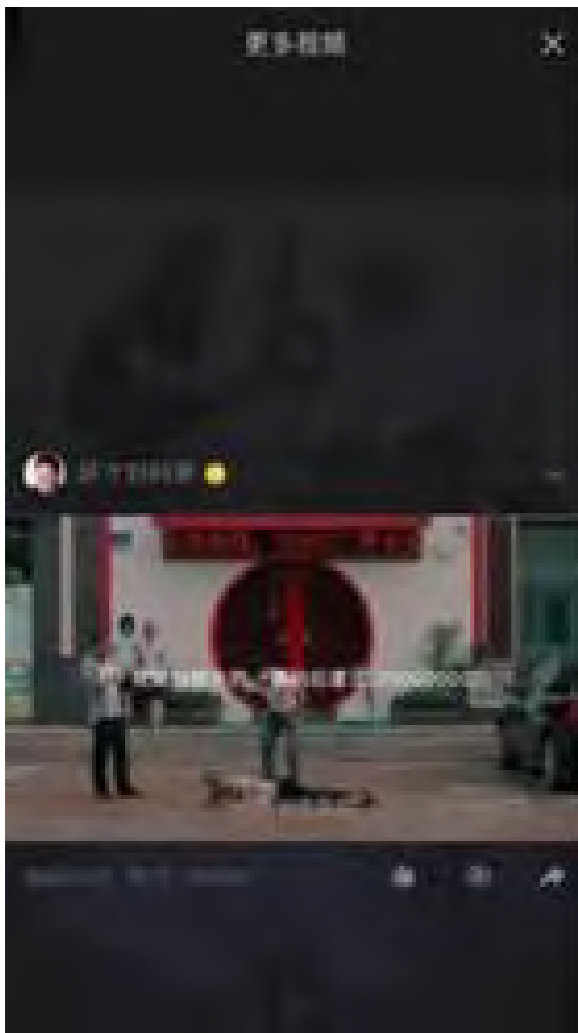
- 310100062
- 900013003
- 200000020
- 310101001
- 200000003
- 140329217
- 200000024
- 11828
- 900013005
- 200000005
- 其他



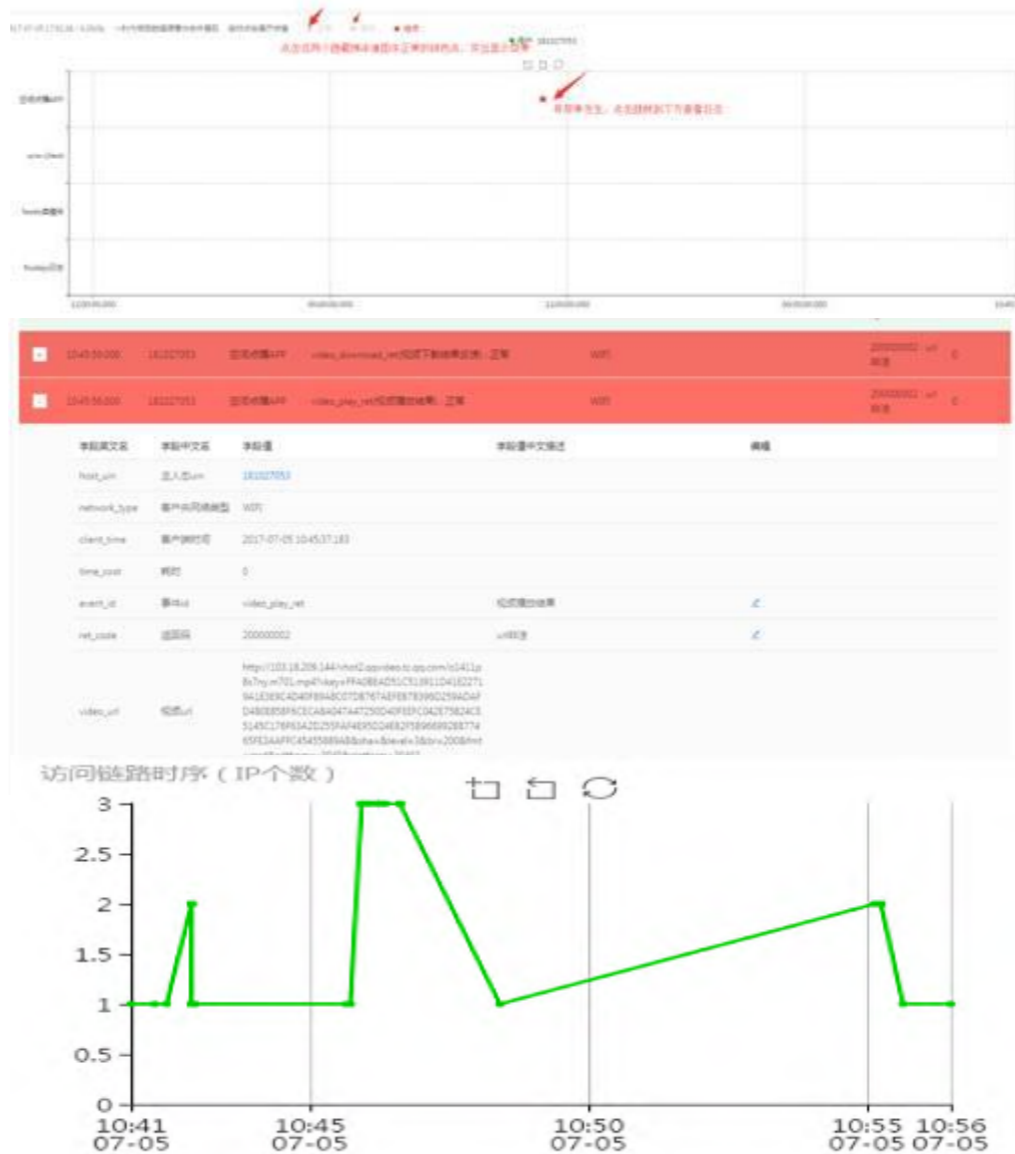
返回码	大盘占比	返回码描述
310100062	0.2195475	视频状态不合法
-900013003	0.1638776	播放器加载超时
200000020	0.1360407	socket错误
-310101001	0.0863267	http连接超时
200000003	0.0618855	DNS解析失败
-140329217	0.0578005	vkey校验失败
200000024	0.0477247	发情请求错误

# 举个栗子吧

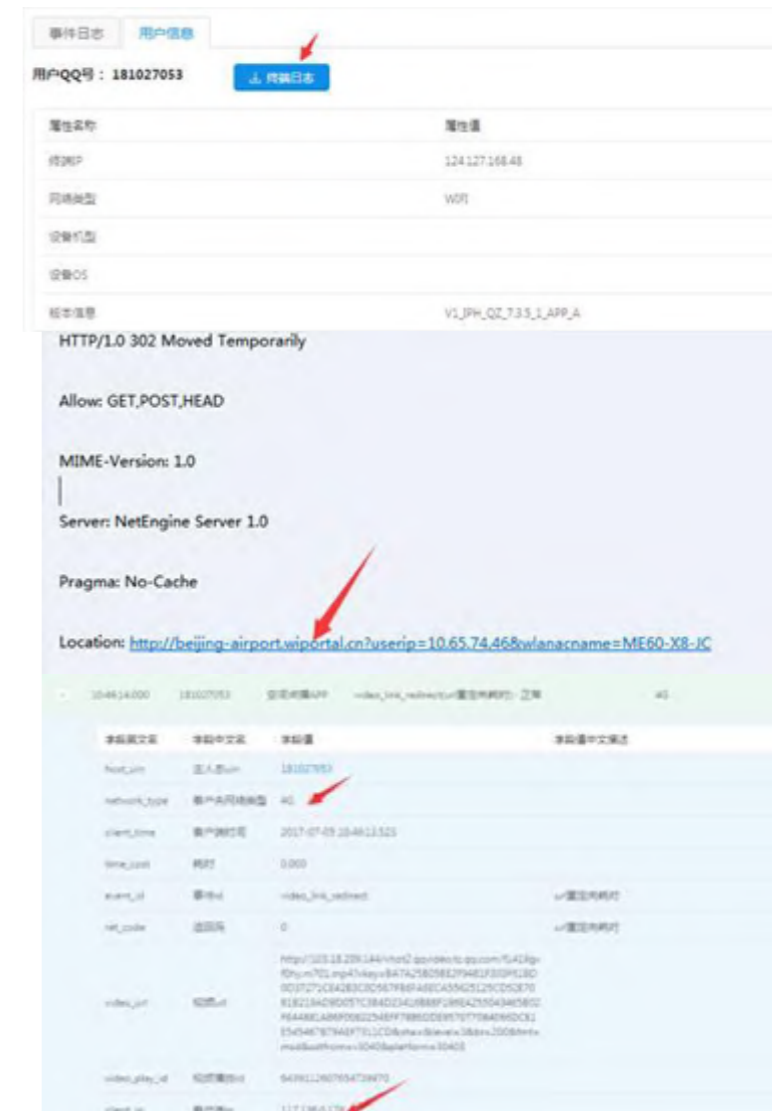
## 收到用户反馈



## 织云查该日志



## 验证分析结果



至此：根因是用户进了机场后，因wifi开关开启，自动试连，并跳转登陆界面

# 跟进时代，践行机器学习

CLOUD  
MANAGEMENT



海量业务的监控**优势**

机器该学习**什么**

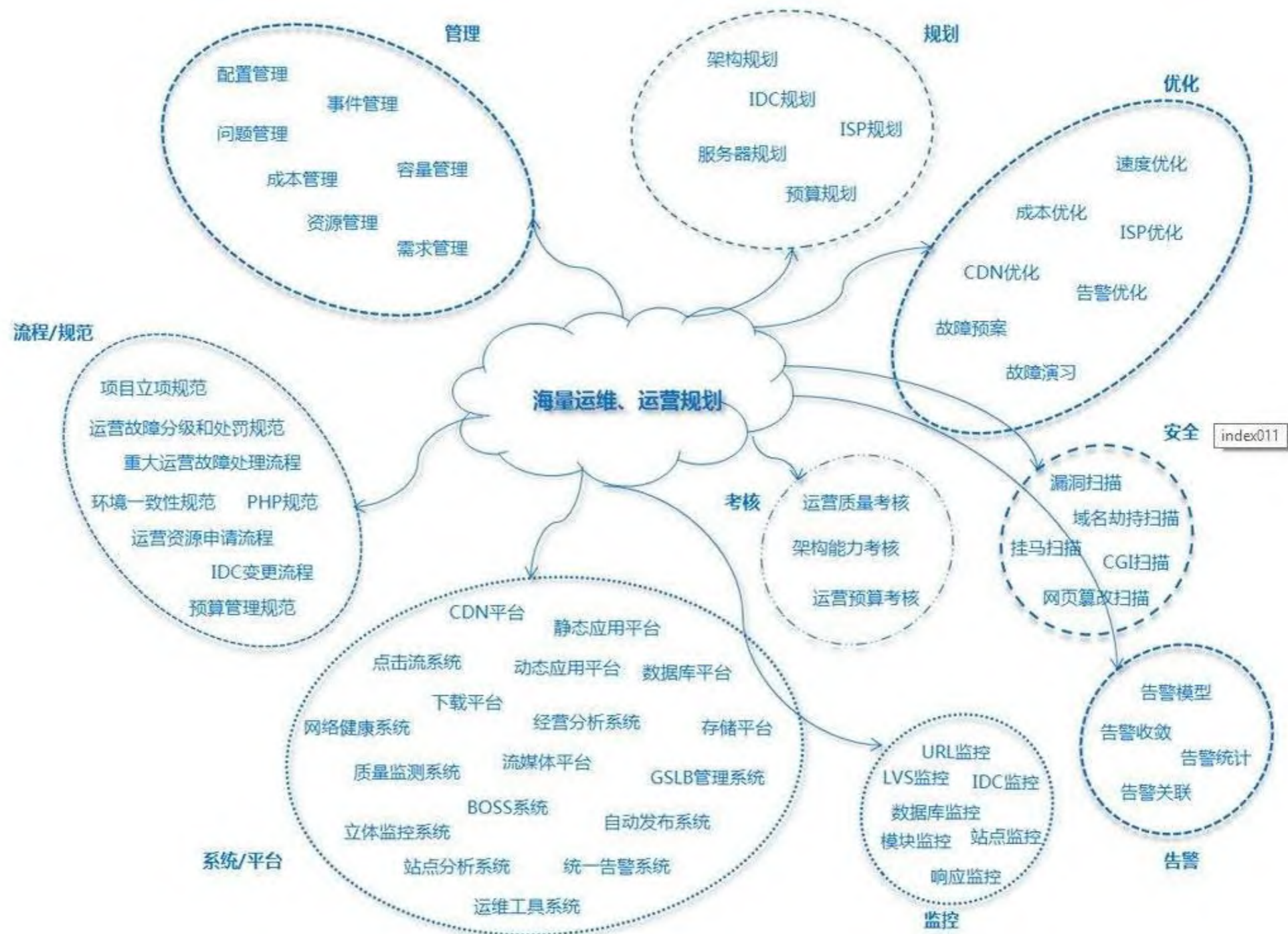
**教**机器正确学习



# 咖啡运维

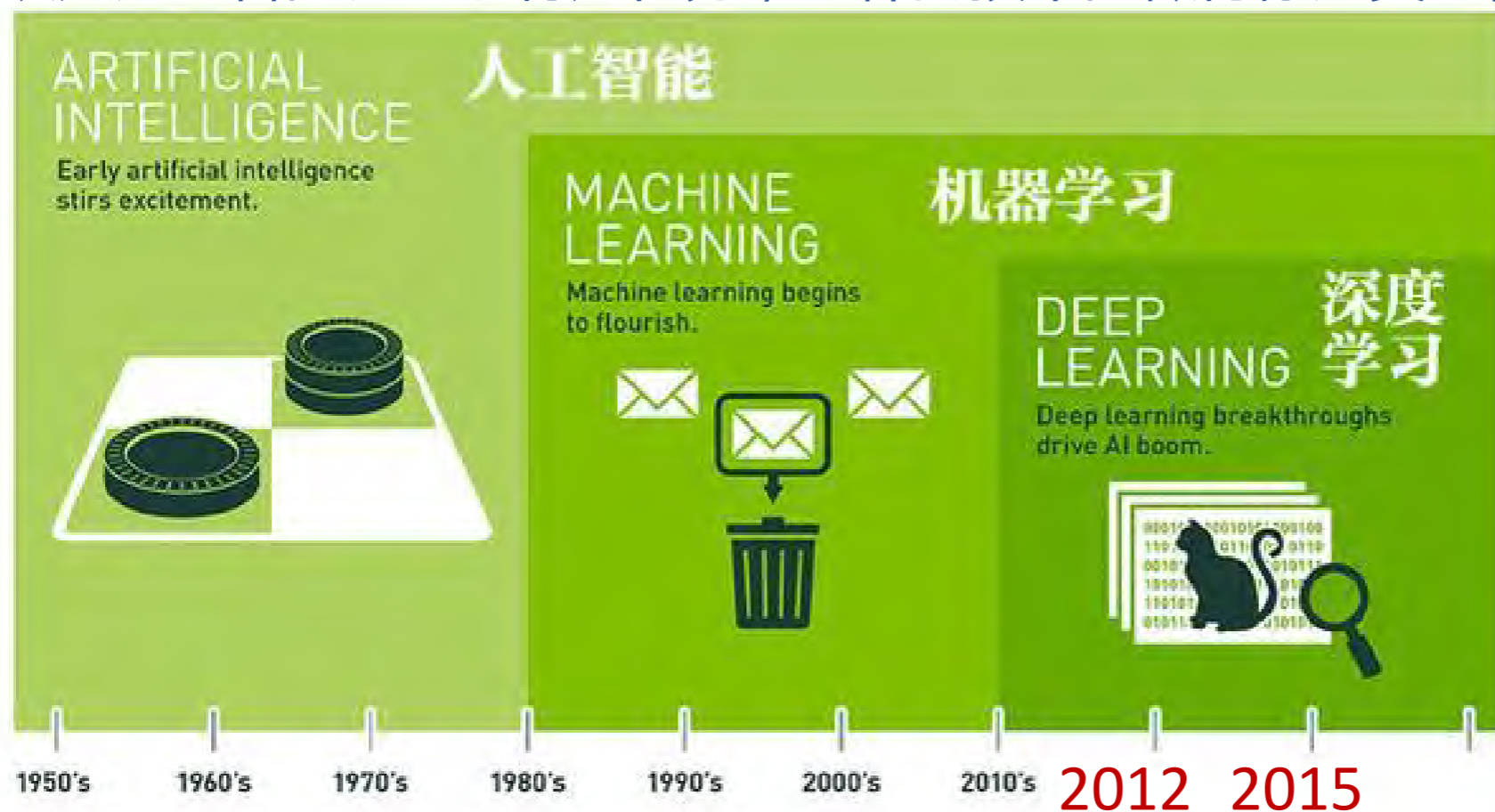


# 运维需要这么多技能吗



# AI 走向咖啡运维之路

狭义人工智能。对于特定任务，这样的技术能做得像人类一样好，甚至更好

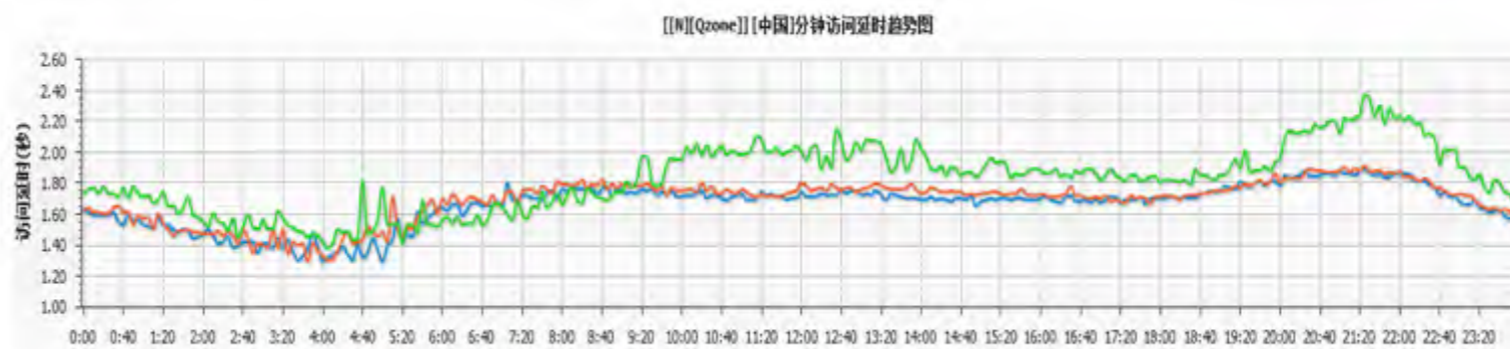
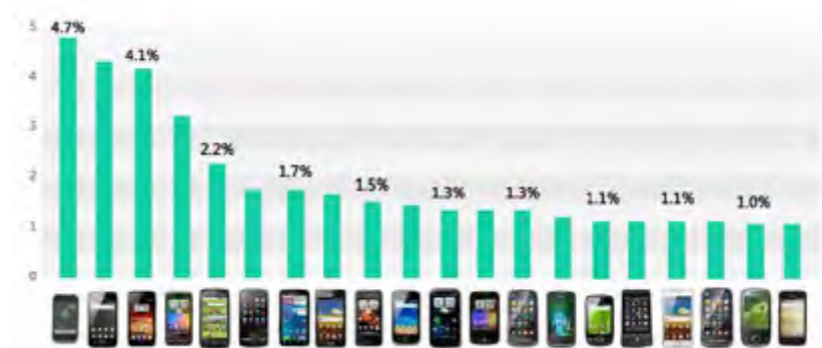
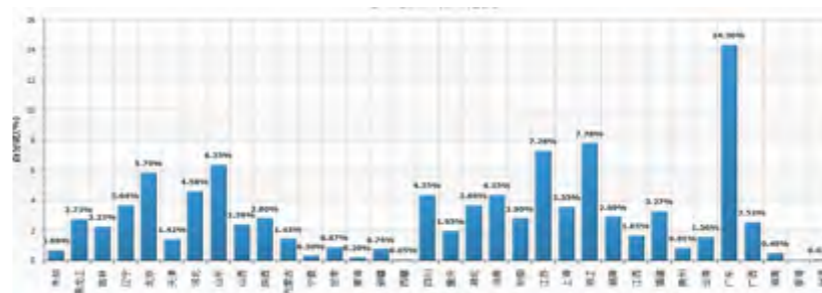
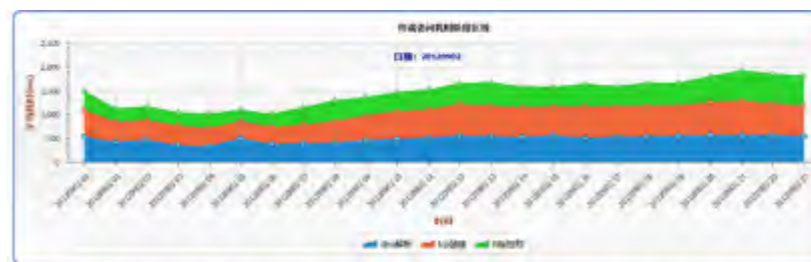
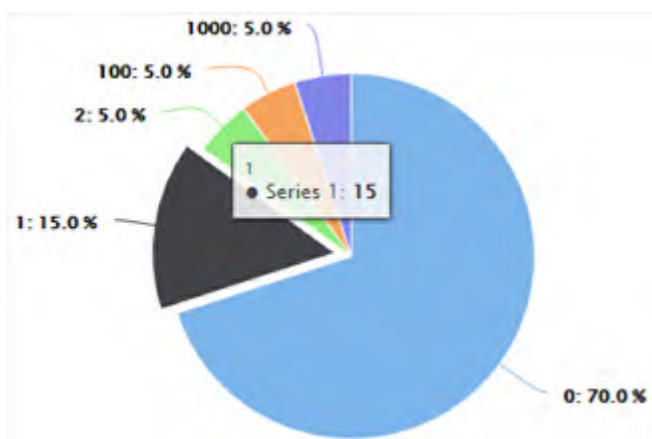


Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

视频  
图像  
文本  
推荐



# 常见的分析模型



✓ 趋势、对比、波动、阈值、分布、聚类

# 重新检视ROOT、DLP、全链路

ROOT

基于架构  
基于经验  
基于概率



收敛告警事件

DLP

基于规范  
基于分工



产生告警事件

全链路

基于数据  
基于模型



提高事件处理能力

# 第一个阶段 机器学习之文本

数据积累

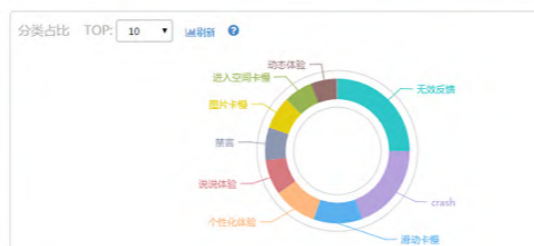
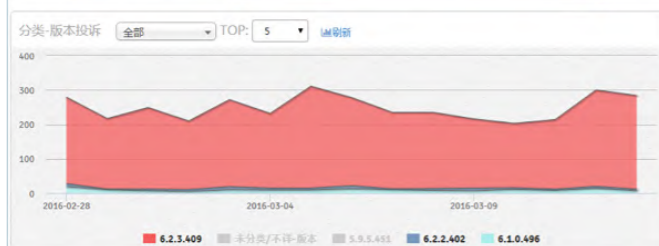
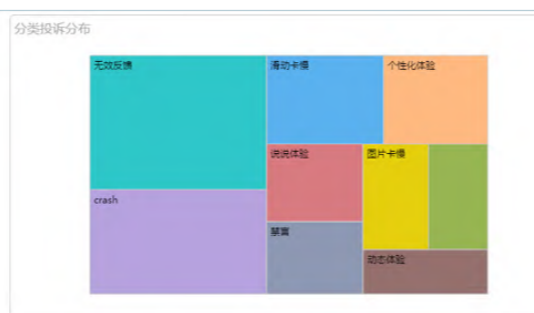
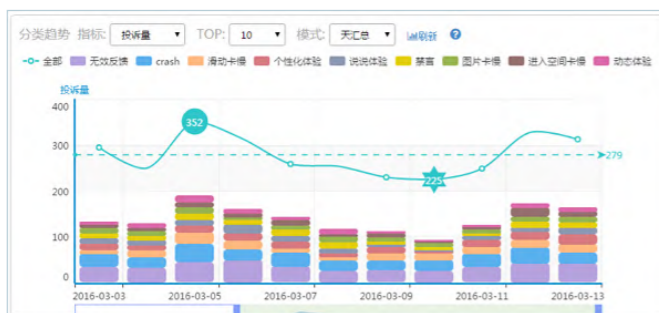


问题发现



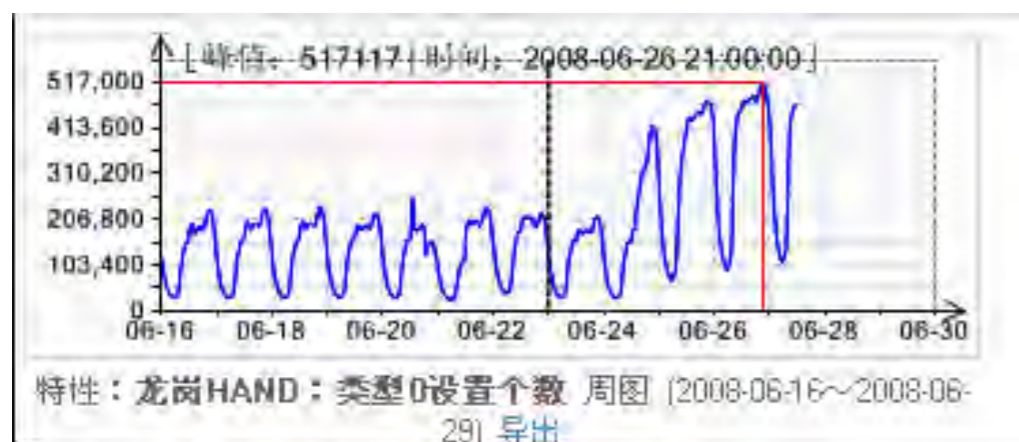
问题处理

## 织云舆情监控 + AI客服

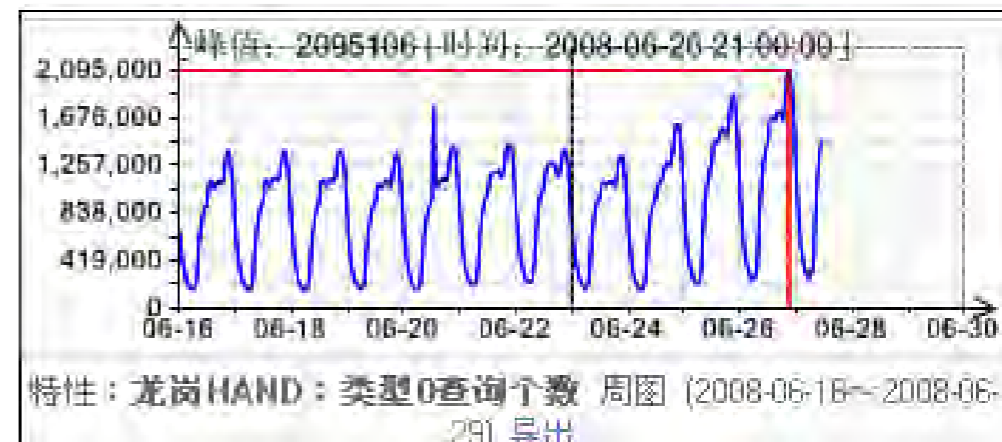


## 第二个阶段

## 机器学习 之 图像

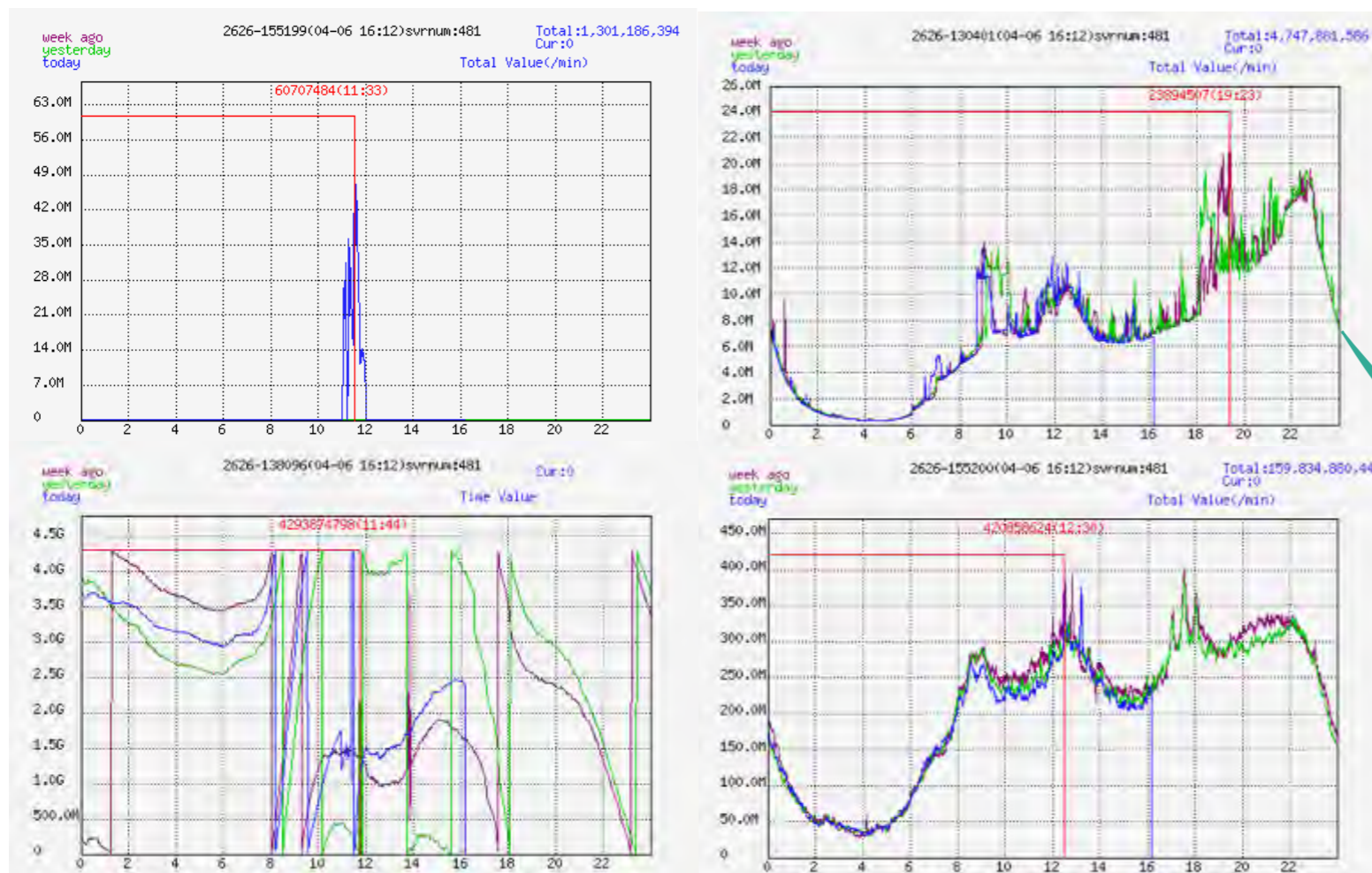


图像相似性





# 第三个阶段 如何告诉AI规则是什么



有监督学习



告警

如果用了自动找出来的历史告警，模型就会学成历史告警的策略 = nothing

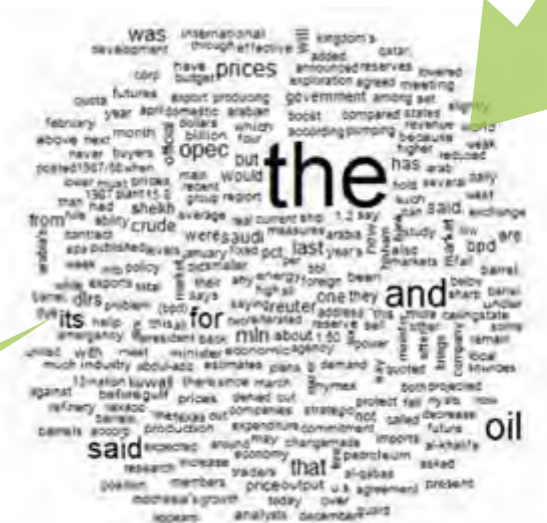
# 第四个阶段 告诉AI数据的意义



```

61.105.249.217 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
Mozilla/5.0 (compatible; Googlebot/2.1; http://www.google.com/bot.html)
66.249.66.193 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
(compatible; Googlebot/2.1; http://www.google.com/bot.html)
67.194.113.210 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
(compatible; Googlebot/2.1; http://www.google.com/bot.html)
118.190.29.199 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
Mozilla/5.0 (compatible; Googlebot/2.1; http://www.google.com/bot.html)
119.187.28.184 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
Mozilla/5.0 (compatible; Googlebot/2.1; http://www.google.com/bot.html)
66.249.66.193 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
(compatible; Googlebot/2.1; http://www.google.com/bot.html)
67.194.113.210 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
(compatible; Googlebot/2.1; http://www.google.com/bot.html)
240.140.179.150 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
(compatible; Googlebot/2.1; http://www.google.com/bot.html)
81.158.248.217 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
Mozilla/5.0 (compatible; Googlebot/2.1; http://www.google.com/bot.html)
81.158.248.217 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
Mozilla/5.0 (compatible; Googlebot/2.1; http://www.google.com/bot.html)
106.113.193.4 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
(Windows-NT-5.1.1)
106.113.193.4 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
(Windows-NT-5.1.1)
118.190.29.199 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
Mozilla/5.0 (compatible; Googlebot/2.1; http://www.google.com/bot.html)
118.190.29.199 [2017/04/01 15:00:00] GET /... HTTP/1.1 200 512 0.000
Mozilla/5.0 (compatible; Googlebot/2.1; http://www.google.com/bot.html)
    
```

机器学习



# 有哪些有值得关注点

腾讯织云  
CLOUD  
MANAGEMENT

监控是 平台  
也是 产品  
更重要是 运营



如果监控是产品

快  
即时性  
告警快

准  
告警准  
误告少

全  
无遗漏  
覆盖广

如果监控是平台

稳

强

易

如何运营监控

指标

闭环

生态

# 指标

## SMART

- 很具体
- 可衡量
- 可达到
- 可观察
- 有时间

DLP告警 = 服务异常

服务异常 x DLP告警时间 = 服务不可靠性

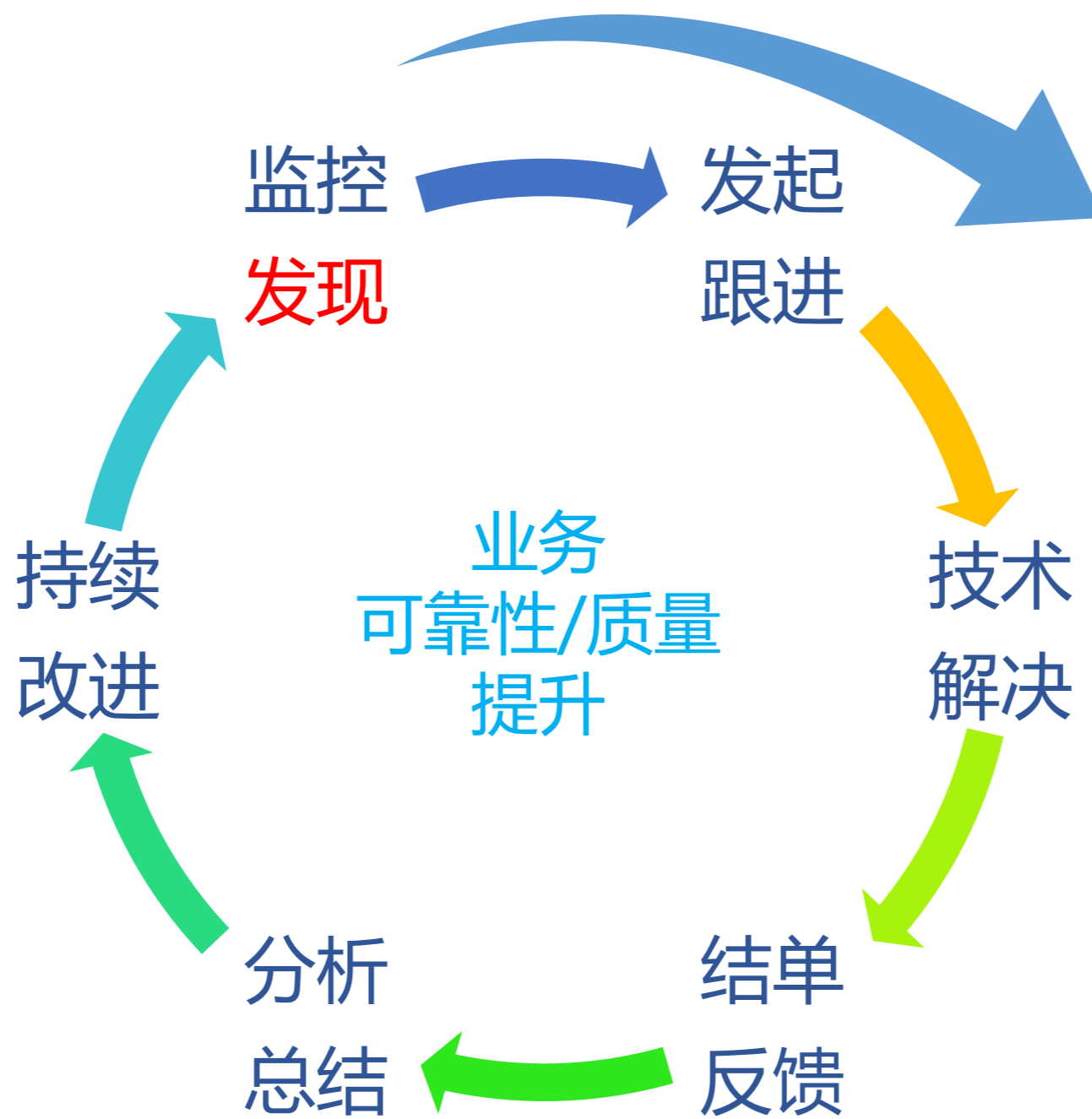
SUM ( 服务可靠性,.....) = 业务不可靠性

业务不可靠性/一段时间(如月,年) = 业务可靠性**指标**(年,月)

- 可横向对比
- 有趋势变化
- 可以目标明确



# 闭环



# 生态



客户端卡慢



舆情监控



H5/测速

自动化测试

摸调



摸调

自动化测试

摸调



摸调

自动化测试

摸调



摸调

## Crash /卡慢

- 移动端监控

## 速度

- APP H5测速
- web测速

## 体验

- 多媒体图片
- 海外速度
- 运营商

## 成功率

- ATT
- 摸调
- monitor
- 业务特性告警
- 组件监控

✓ 体验有提升

✓ 成功率改进

✓ 可用性衡量

横向业务**指标**对比

分析手段

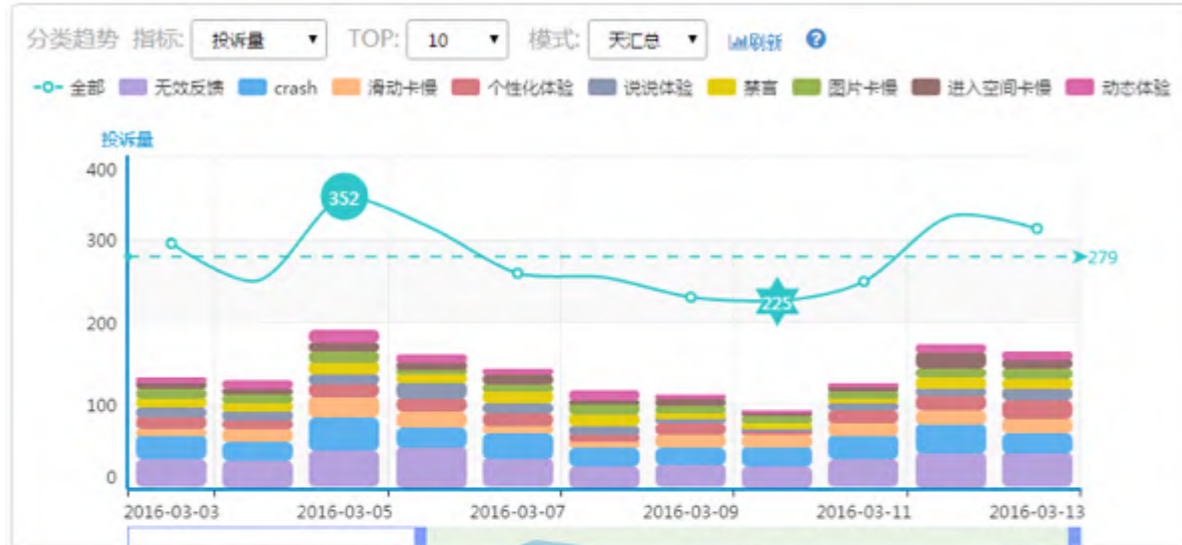
目标

# TIPS

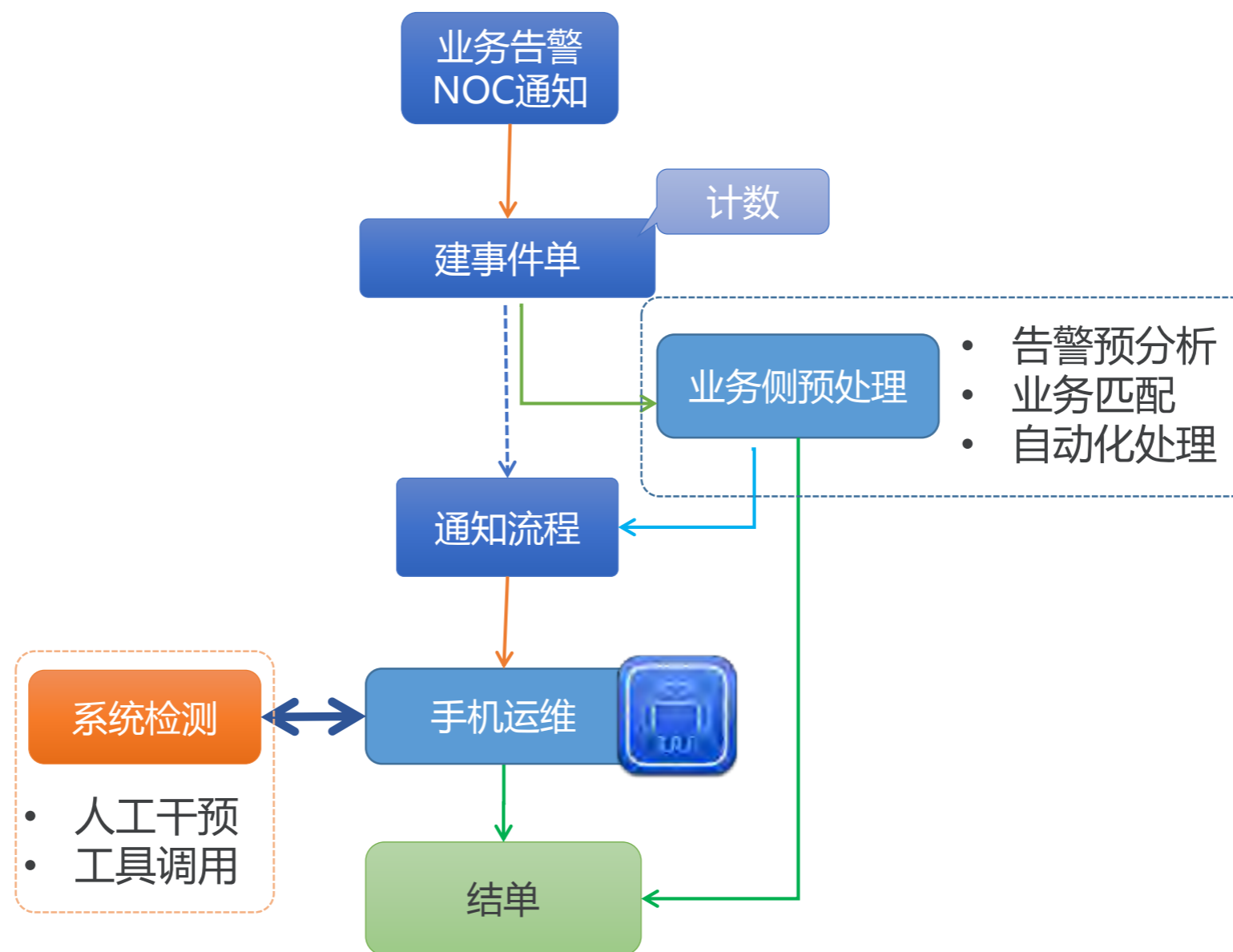
腾讯织云  
CLOUD  
MANAGEMENT

自动处理、  
移动办公、  
舆情监控、  
告警分级

# 织云舆情监控



# 告警自动告警处理





# 织云移动运维办公



# 监控如何分级

告警级别

时间分级

范围分级

推送规范



请交流，请携手