



GOPS 2016
Shenzhen



全球运维大会

2016

深圳站

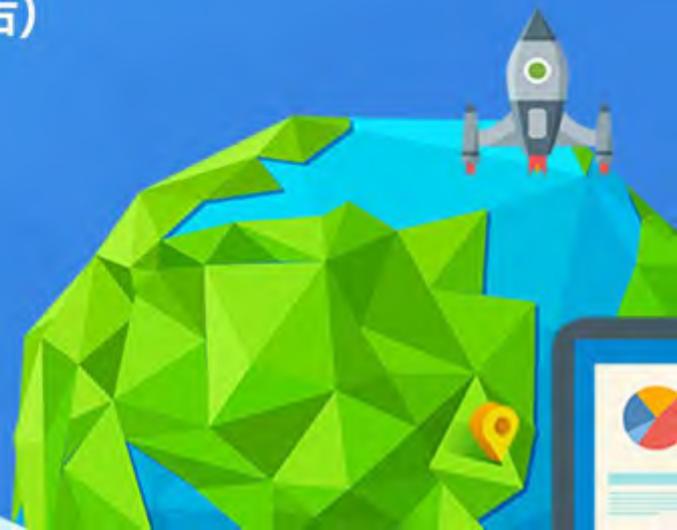
会议时间：3月25日-3月26日

会议地点：深圳·南山区 圣淘沙酒店(翡翠店)

主办单位： 开放运维联盟
OOPSA Open OPS Alliance  高效运维社区
GreatOPS Community

指导单位： 数据中心联盟
Data Center Alliance

协办单位：中国新一代IT产业推进联盟





GOPS 2016
Shenzhen



全球运维大会

2016

深圳站

由点及面，腾讯智能监控实践与思考

梁定安，腾讯



个人简介

- 梁定安（大梁）
- 10年互联网运维
- 腾讯社交平台运维负责人



社 交 网 络 运 营 部



目录

1 运维监控 in 腾讯社交

2 做好监控必须具备的要素

3 智能监控的实践分享

4 监控建设到质量管理体系建设



监控的意义和目标



- 可靠性
- 可用性
- 用户体验



监控的手段



监控的本质

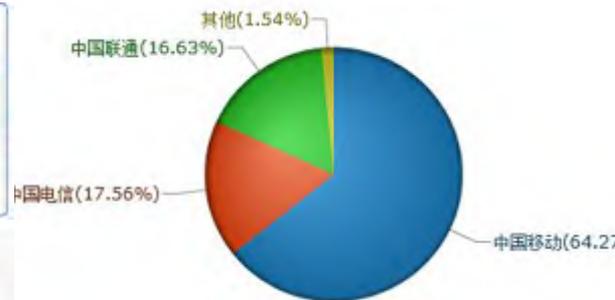
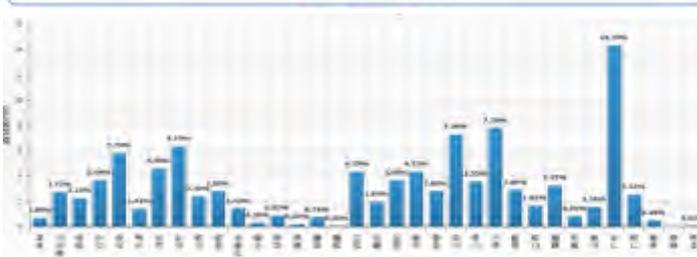
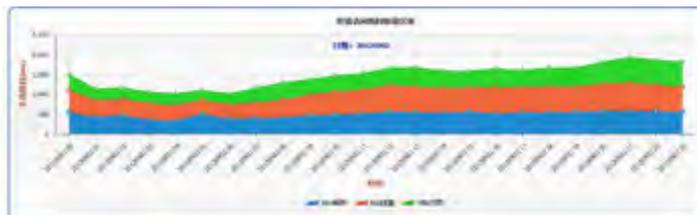
请求量
成功率
耗时



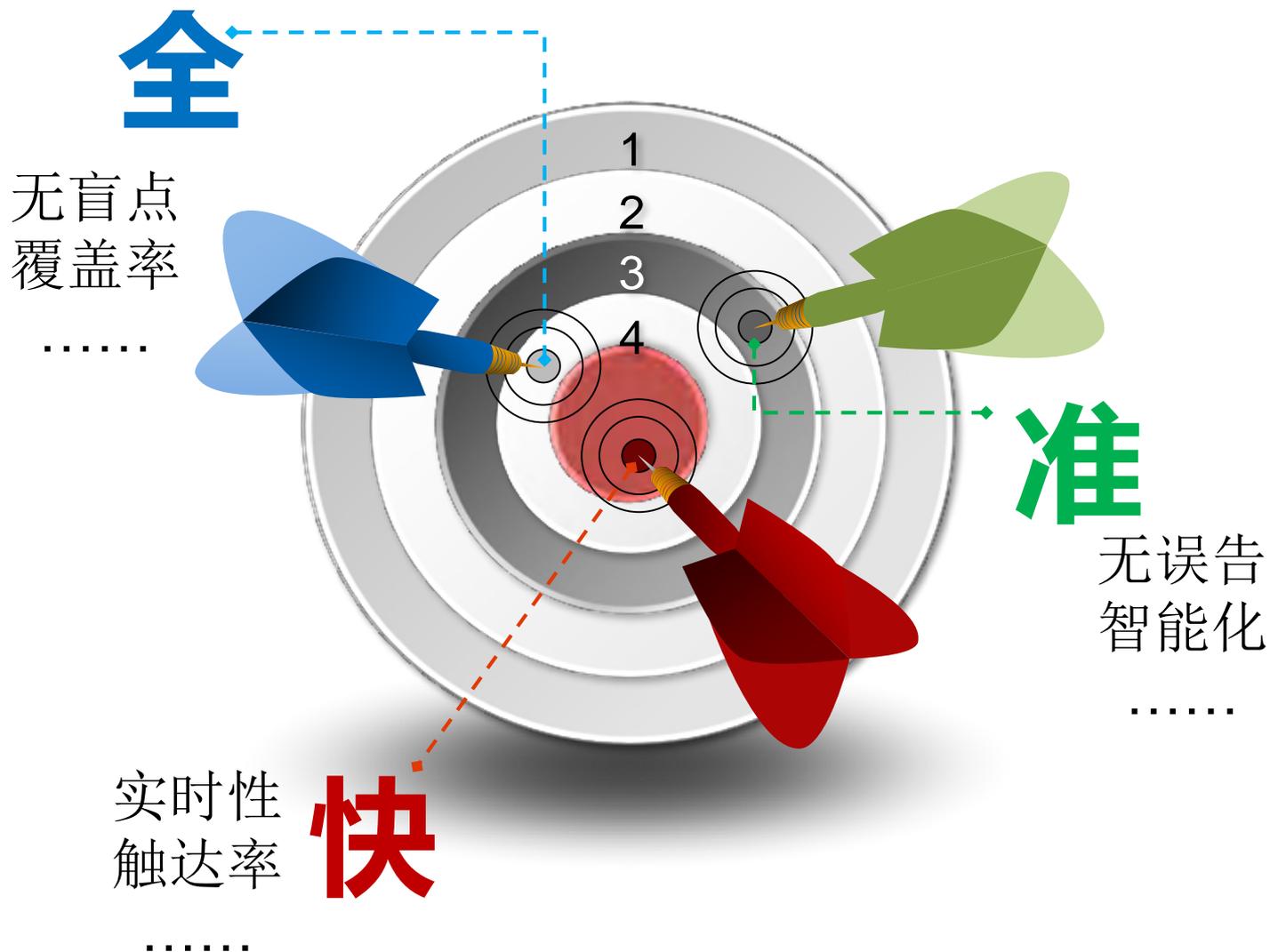
趋势
对比
波动
阈值
分布
聚类
区间



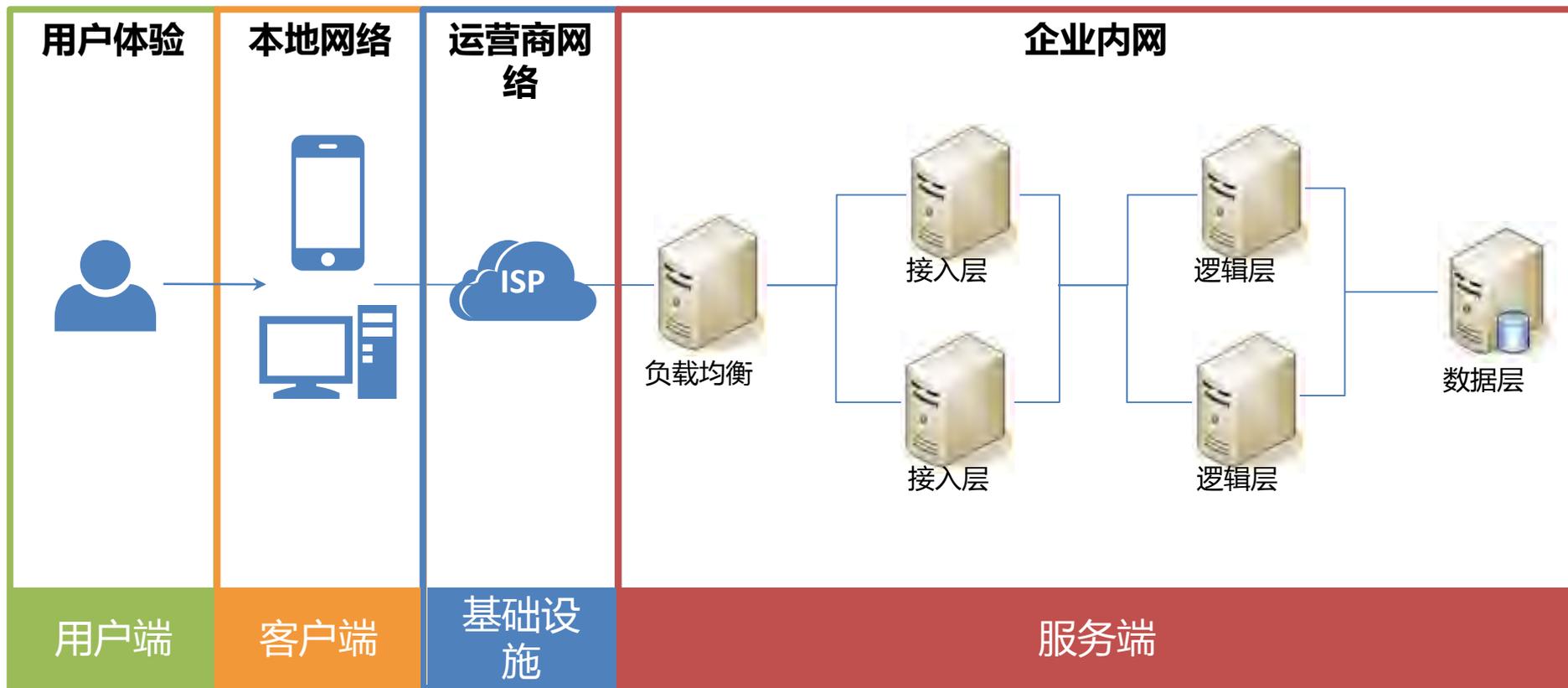
图
表
告
警



监控的目标



全链路监控



SNG监控全景图

TEG服务监控：

N: 网络质量监控
C: CDN监控
D: 数据层监控

SNG服务监控：

Y: 业务染色监控
R: 返回码监控
S: 测速系统
A: 自动化测试
M: 模块间调用
C: 组件监控

基础监控：

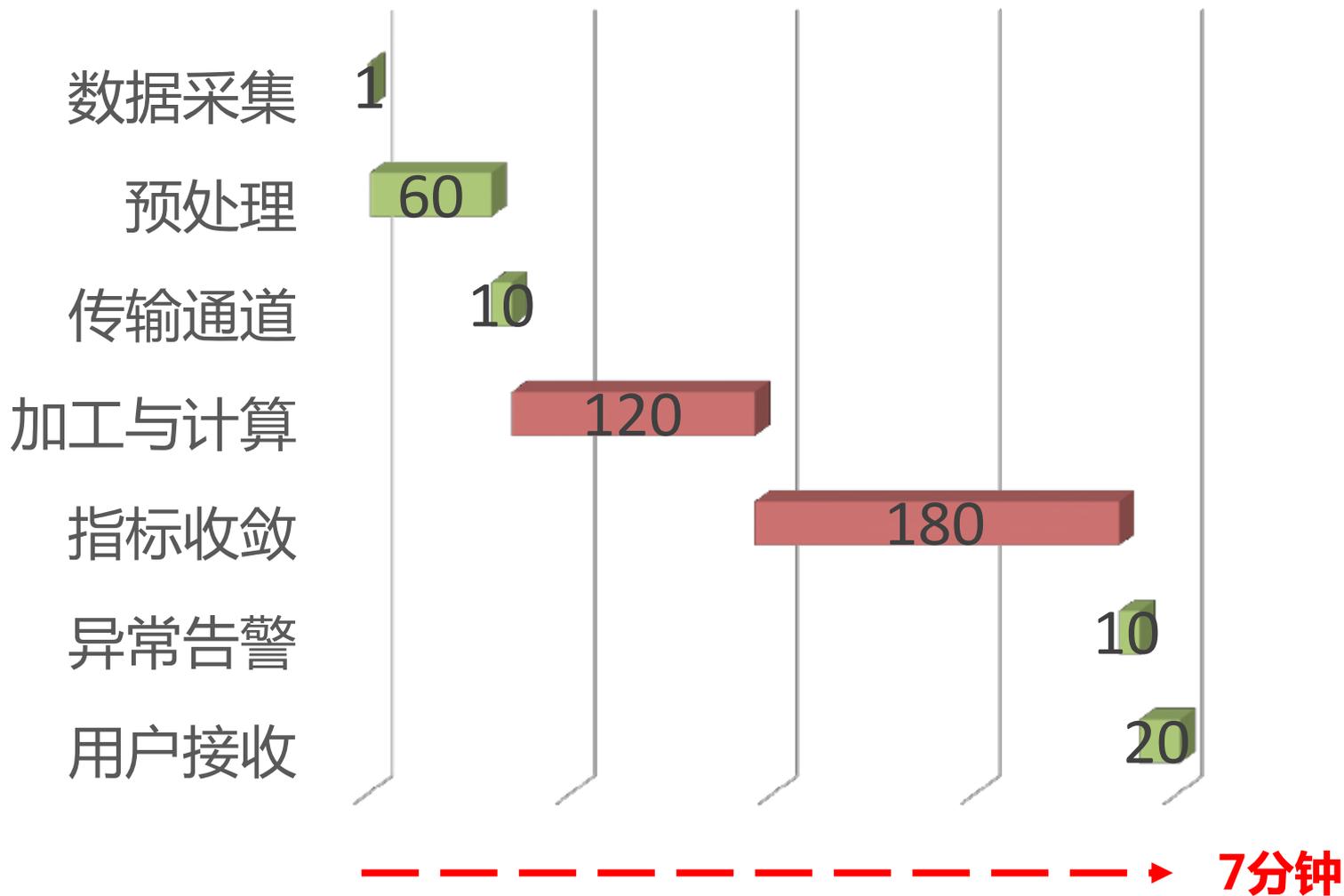
L: 容量管理
P: 进程监控
F: 特性监控

移动端监控：

T: 舆情监控
K: 卡慢监控
D: 多维监控



监控的速度



统一上报协议



ID, IP, 时间, 值



按ID/IP/时间聚合统计



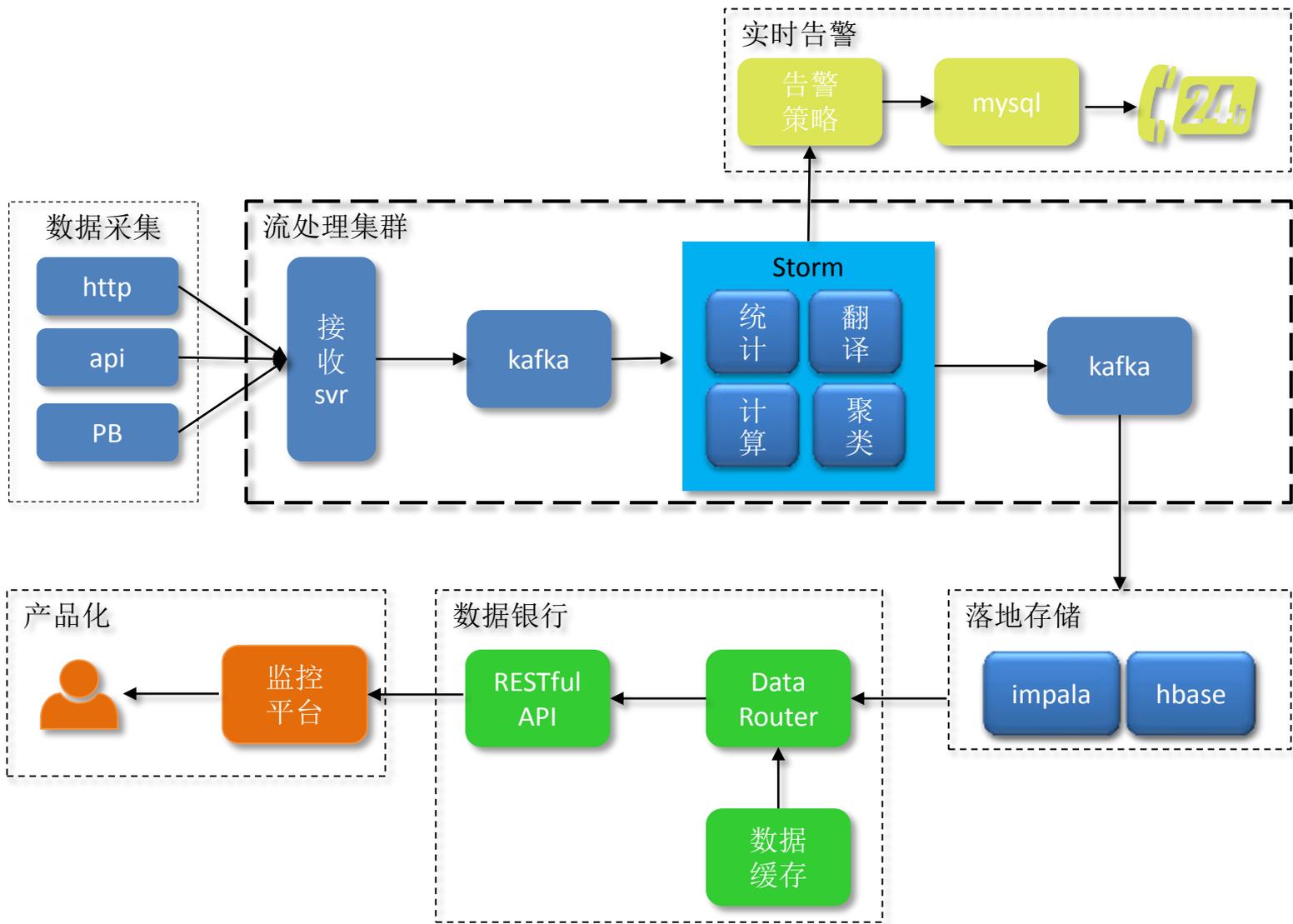
ID, 时间, 运营商, 版本号,
QQ号, 业务类型, 播放状态, 播
放页面url, 请求段播放时长, 完
整播放时长, 视频文件下载地
址, 播放id, 请求id, 用户下载
速度, 播放器版本号, 命令字.....



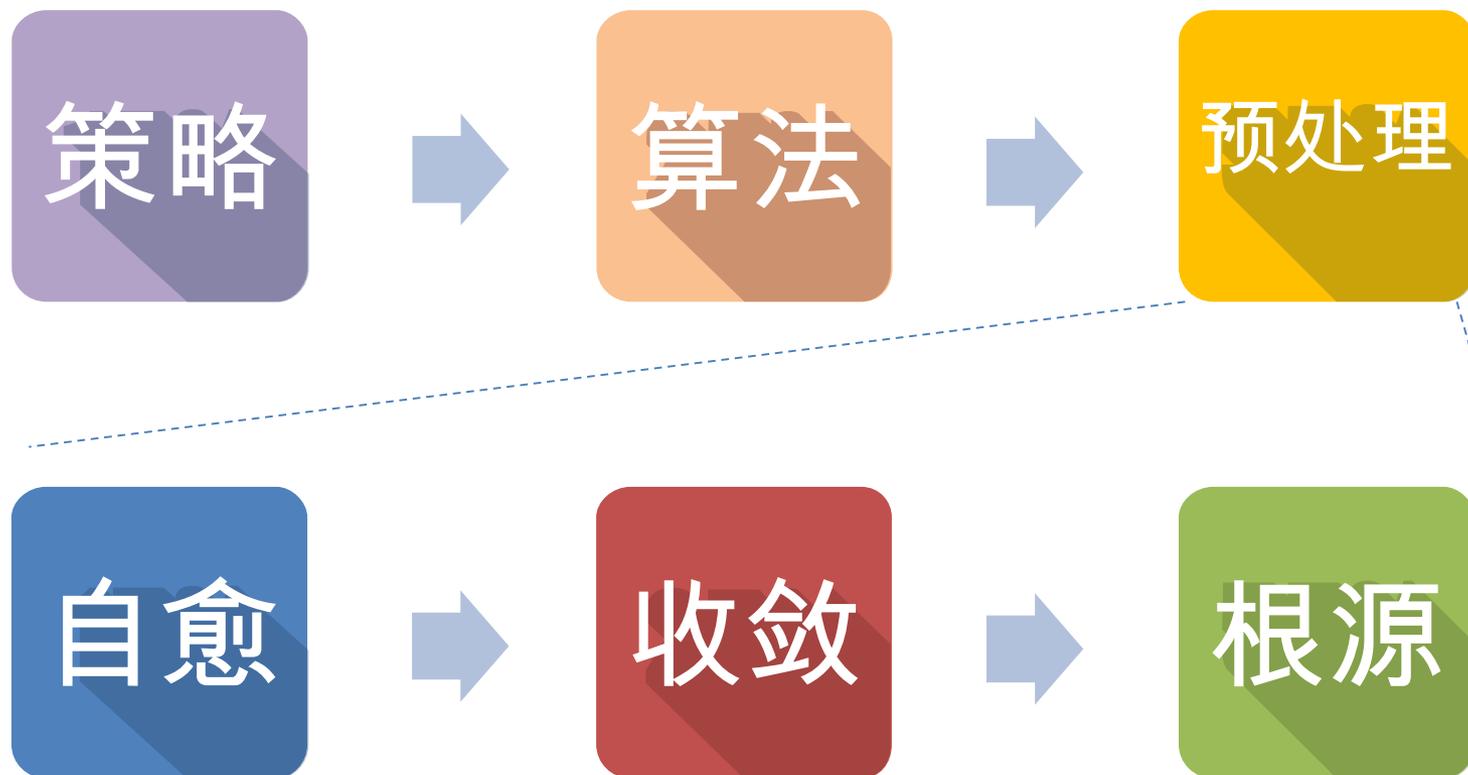
按场景分类
多维度组合



监控集群能力



准：智能监控



自愈



质量多
告警多
需求
海量服务
变更杂
设备多
故障

设备多
监控
质量
需求
变更杂
质量
变更杂
监控
海量服务
设备多

告警多
设备多
变更杂
故障
监控
需求
海量服务
变更杂
设备多
故障

设备多
变更杂
监控
海量服务
需求
故障
海量服务
需求
告警多

监控
设备多
质量
故障
需求
海量服务
变更杂
设备多
故障

故障
设备多
海量服务
需求
变更杂
故障
设备多
故障

基础监控

服务端监控



客户端监控

用户端监控



ROOT智能监控

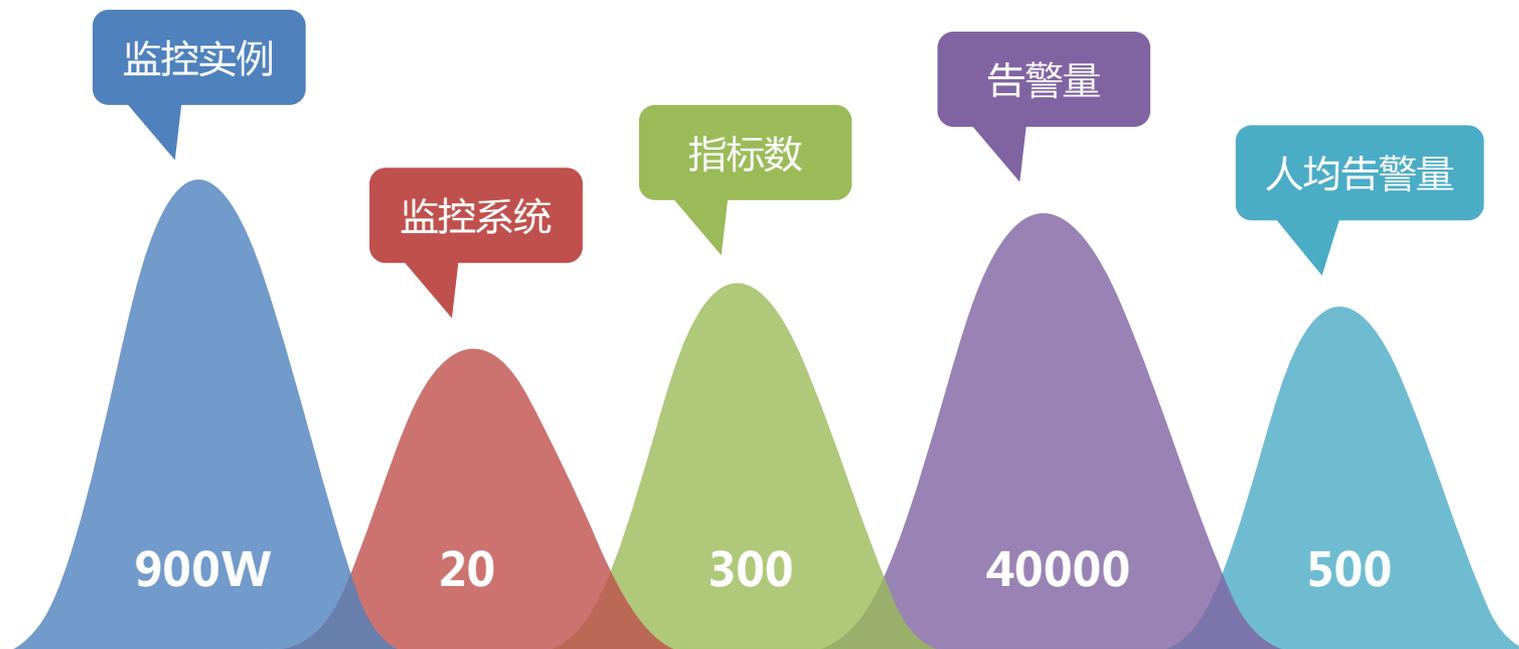
包袱

- 业务架构庞大而复杂
- 大量现象告警（点）
- 告警收敛无法最大化
- 原因告警（端到端）被淹没

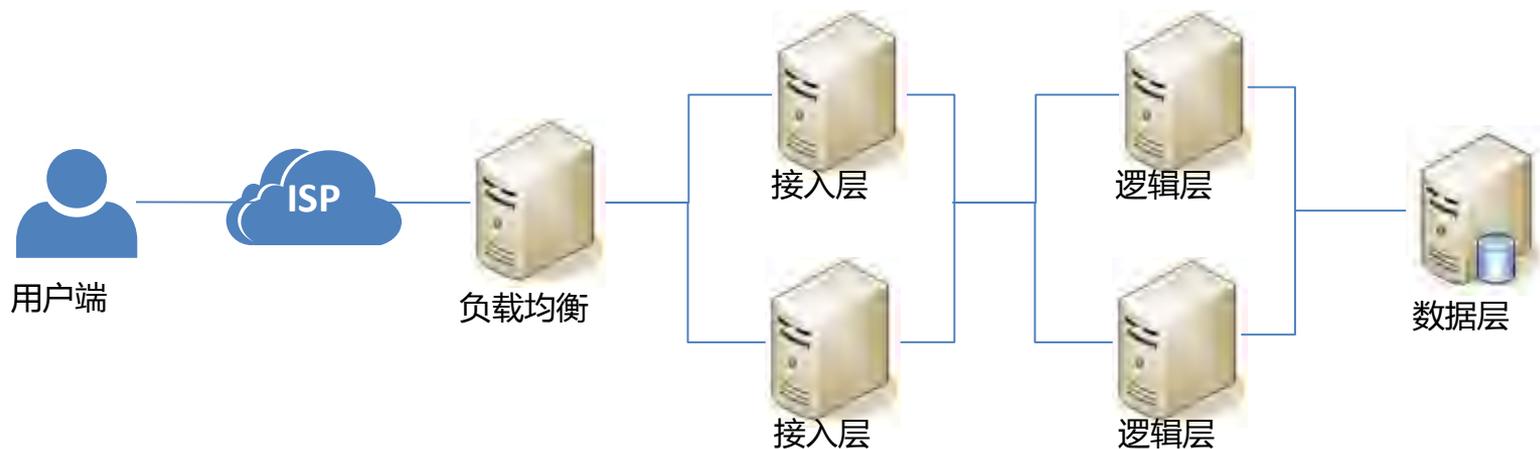


ROOT

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示意图



假如：DB宕机。

现实：用户端、接入层、逻辑层、数据层的监控点均有 **N** 个告警产生。

理想：智能定位到数据层监控，只发出 **1** 个告警。

现象告警

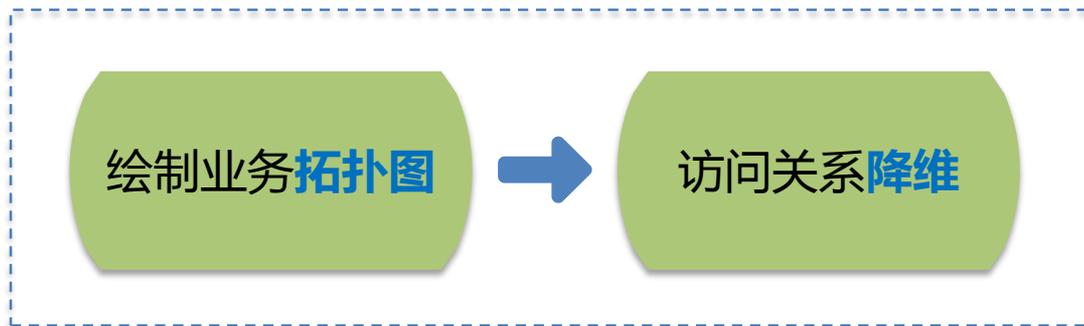


原因告警

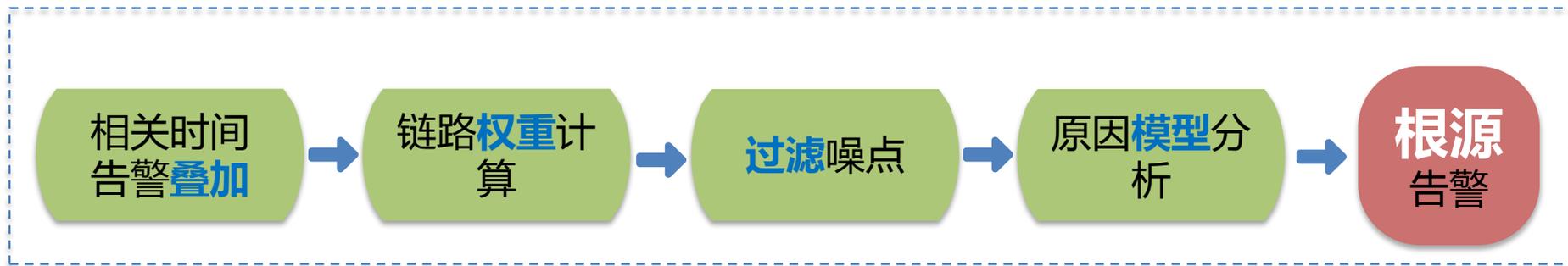


ROOT分析原理

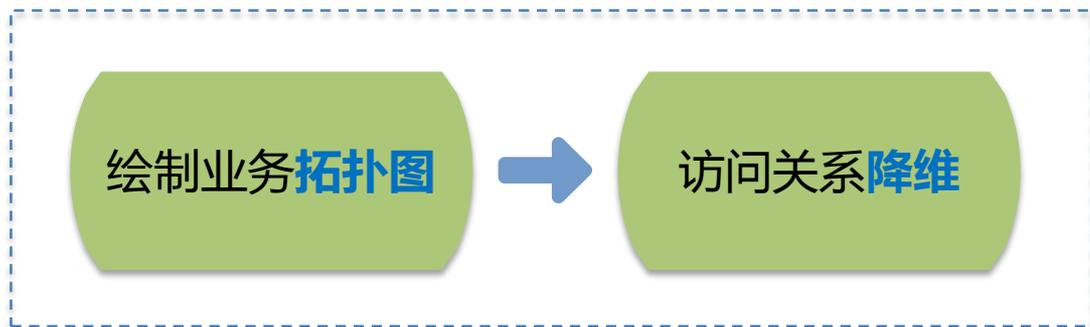
基础数据



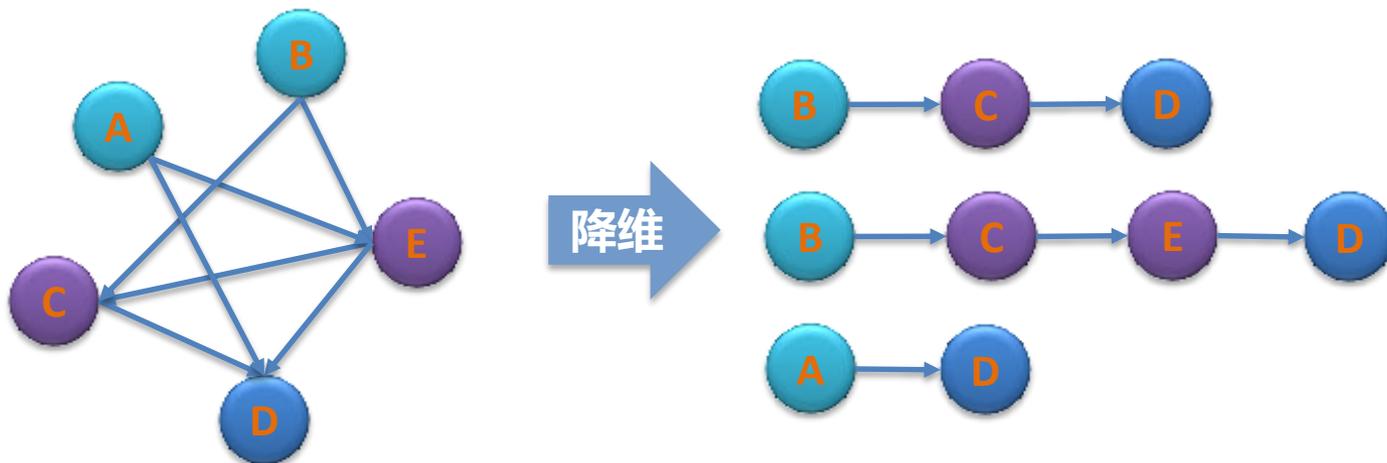
计算逻辑



降维策略



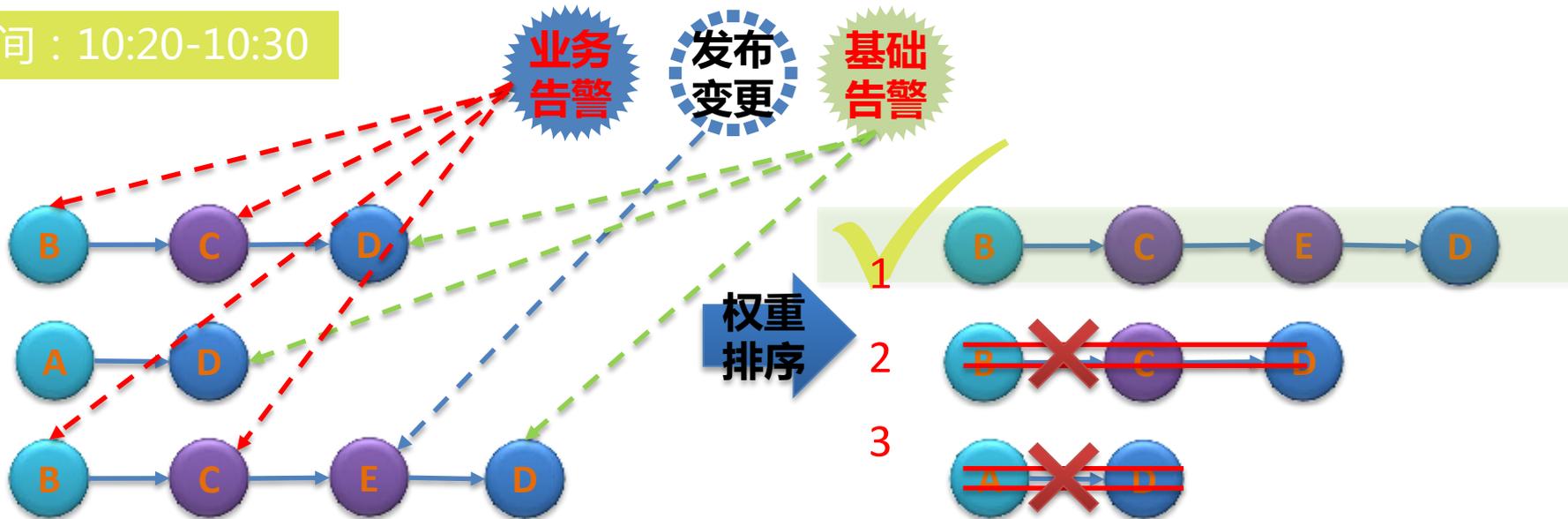
- L5访问关系
- 模调关系
- IP间抓包



关联分析

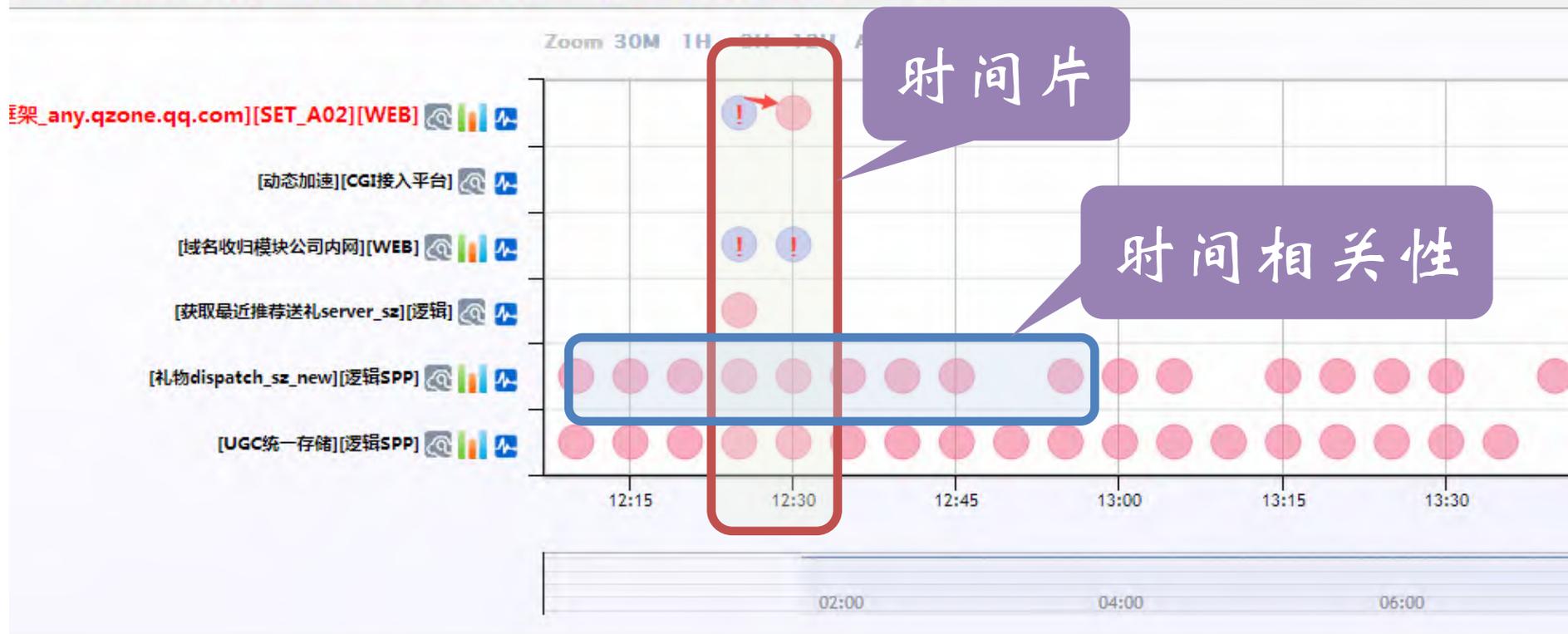


时间：10:20-10:30



时间相关性分析

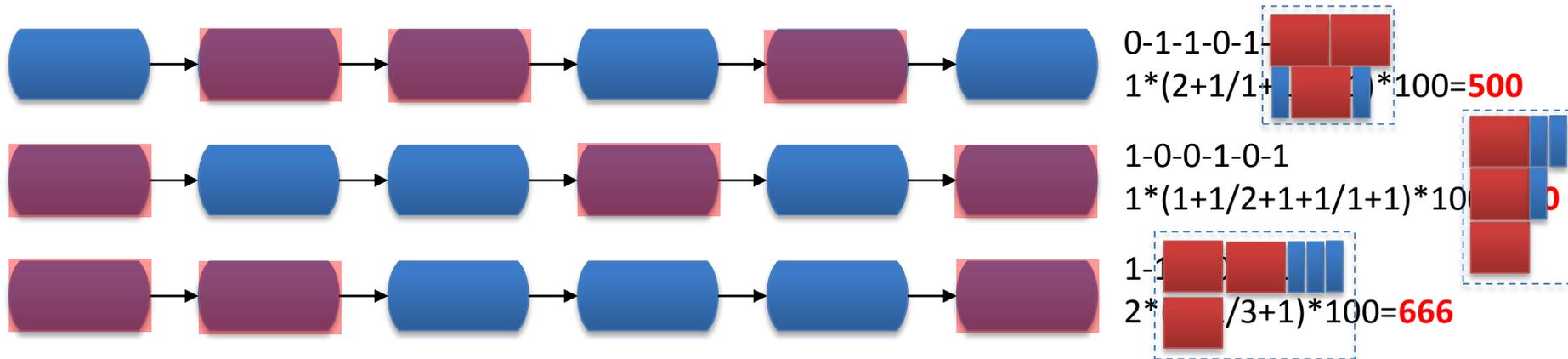
搜索模块/[N][Qzone]-[Qzone平台_基础平台][框架]-[框架_any.qzone.qq.com][SET_A02][WEB]



关联告警准确性：时间有效性，持续性，告警时延，链路相关性。



权重面积分析



✓ 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；

✓ 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= $iarrmax*(1+1/N+...)*100$

- 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、1-1-1-0-1-0-0-1-1, 权重面积= $3*(1+1/1+1+1/2+1)*100=1350$;

✓ 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= $(5*1*100)*10=5000$;



ROOT架构





举个栗子!!



6个时间片内



- A : 用户打开超时
- B : 服务调用延时高
- C : 组件失败率突增
- D : 有版本变更发布

关系链路中叠加告警信息计数



相关时间片告警分类

C : 组件失败率突增

波动

是否恢复?

寻找关联告警

A : 用户打开超时
关联到
 D : 有版本变更发布

模型+根源+历史数据

结论

发现 : [A]用户打开超时
 是由于[D]有版本变更发布造成
 上次同类问题发生在 上周二
 责任人 王小宝

B : 服务调用延时高

长期告警

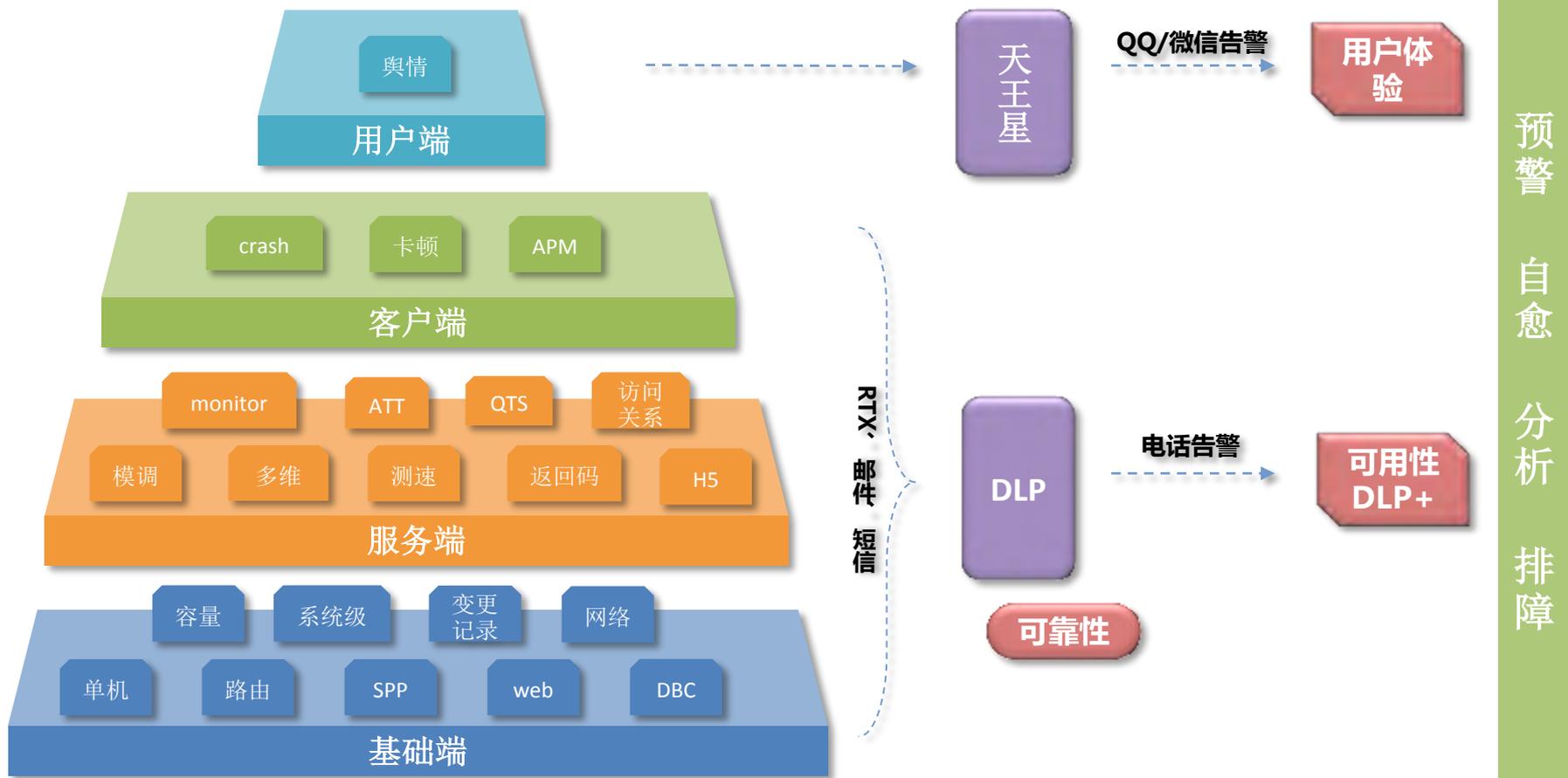
大数据分析



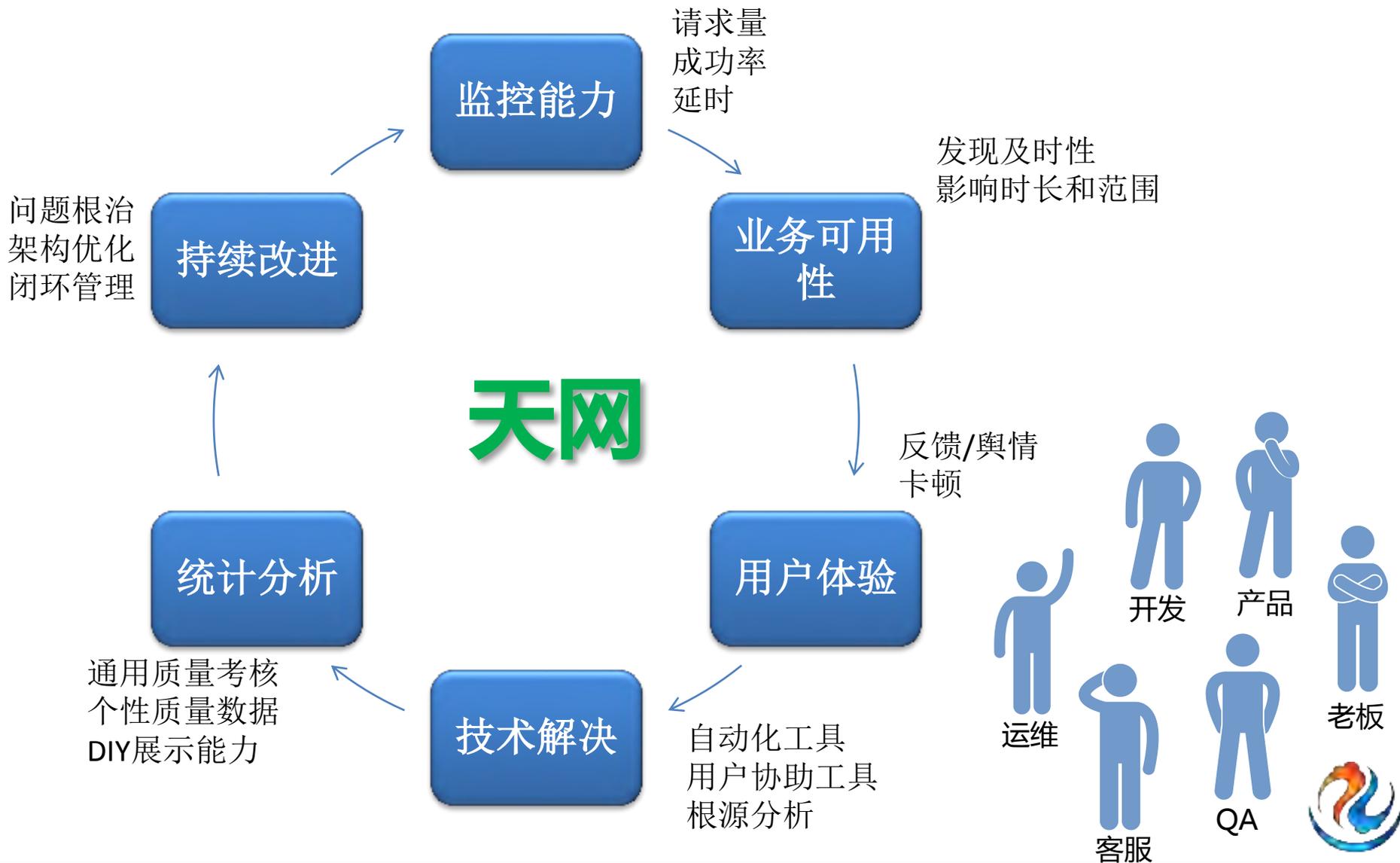
质量体系：生态构建



天网体系介绍



天网：质量体系



攀登探索

