

B站统一监控系统的设计,演进 与实践分享

梁晓聪 devops
@lxcong

About Me

- 梁晓聪
- 2015年加入B站
- devops
- 热爱新技术,热爱开源
- 小宅男



故事的开始

B站炸了. 舆情监控(括弧笑脸)



大家正在搜：池昌旭 林允儿

首页 视频 发现 游戏 注册 | 登录

超过1000万人正在使用



今天B站炸了吗

知识就是力量，法国就是培根，B站就是爆炸。

+ 关注

私信



她的主页

她的相册

Lv9

海外 日本

丧偶

2009年6月26日

简介：知识就是力量，法国就是培根，B站就是爆炸。

个性域名：yamanasion



今天B站炸了吗

6月6日 17:48 来自 微博 weibo.com

炸了

我的内心毫无波动

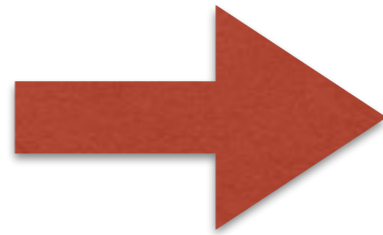


私信聊天

我们的挑战

当前情况:

- 技术栈多
- 产品模块复杂
- 业务爆发式增长
- 运维要求高



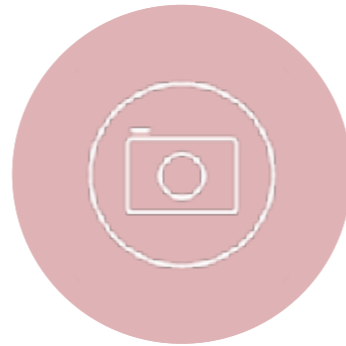
监控问题爆发:

- 覆盖率低
- 误报, 漏报多
- 告警风暴



重新定义的监控系统

◆ 完整的监控体系



◆ 科学的告警策略



◆ 统一的告警中心



完整的监控体系

用户端监控

客户端质量

- 用户端网络质量
- 劫持情况
- 崩溃&卡顿
- 返回码
- 响应时间
- 错误率

播放质量

- 点播/直播
- 播放卡顿
- 平均首帧
- 播放失败率
- 弹幕加载
- cdn质量

服务端监控

业务层

- qps/tps
- 耗时分布
- 饱和度
- 吞吐量
- 依赖响应
- 缓存命中率
- 调用链
- SLA
- 日志

应用层

- cache资源
- db资源
- mq资源
- lb资源
- es资源
- 分布式文件
- 进程监控

基础层

- 虚拟机
- 物理设备
- 容器
- 专线质量
- 机房出口质量
- 交换设备
- http
- tcp
- ping

如何推进？

分析监控场景对应监控手段

场景

类型

手段

服务端监控

metric类型

时间序列数据

日志类型

日志处理流

自定义类型

自研

.....
用户端监控

客户端

apm

播放器

自研

如何推进？

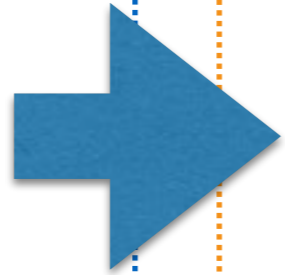
分析监控场景对应监控手段

场景	类型	手段
服务端监控	metric类型	时间序列数据
	日志类型	日志处理流
	自定义类型	自研
用户端监控	客户端	apm
	播放器	自研

metric方案选型

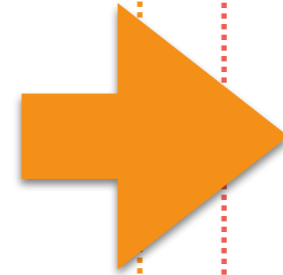
metric数据特征

- 能覆盖大部分监控场景
- 固定几种数据类型
 - ◆ Counter
 - ◆ Gauge
 - ◆ 等..
- 时序数据
 - ◆ 具有统计特性
 - ◆ 具有规律性



选型原则

- 基于开源方案, 二次开发
- 具备现代时间序列数据库的特性
- 活跃项目, 具有成熟的生态环境



结论

- **prometheus**
- 支持任意维度label
- cncf基金会

metric

现状:

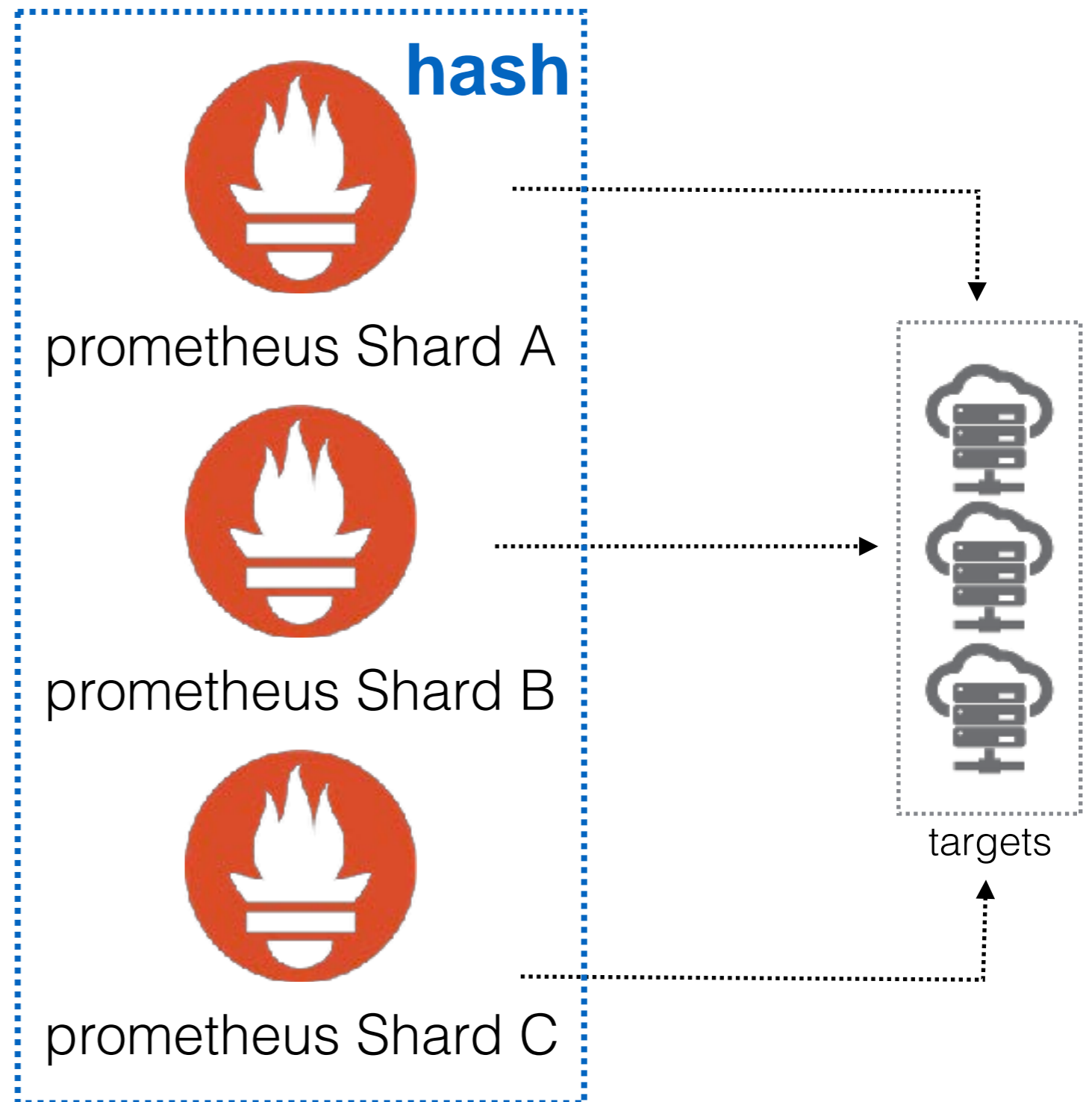
- 40w+/s的指标采集
- 10k+ 监控目标
- 10+ prometheus节点

问题:

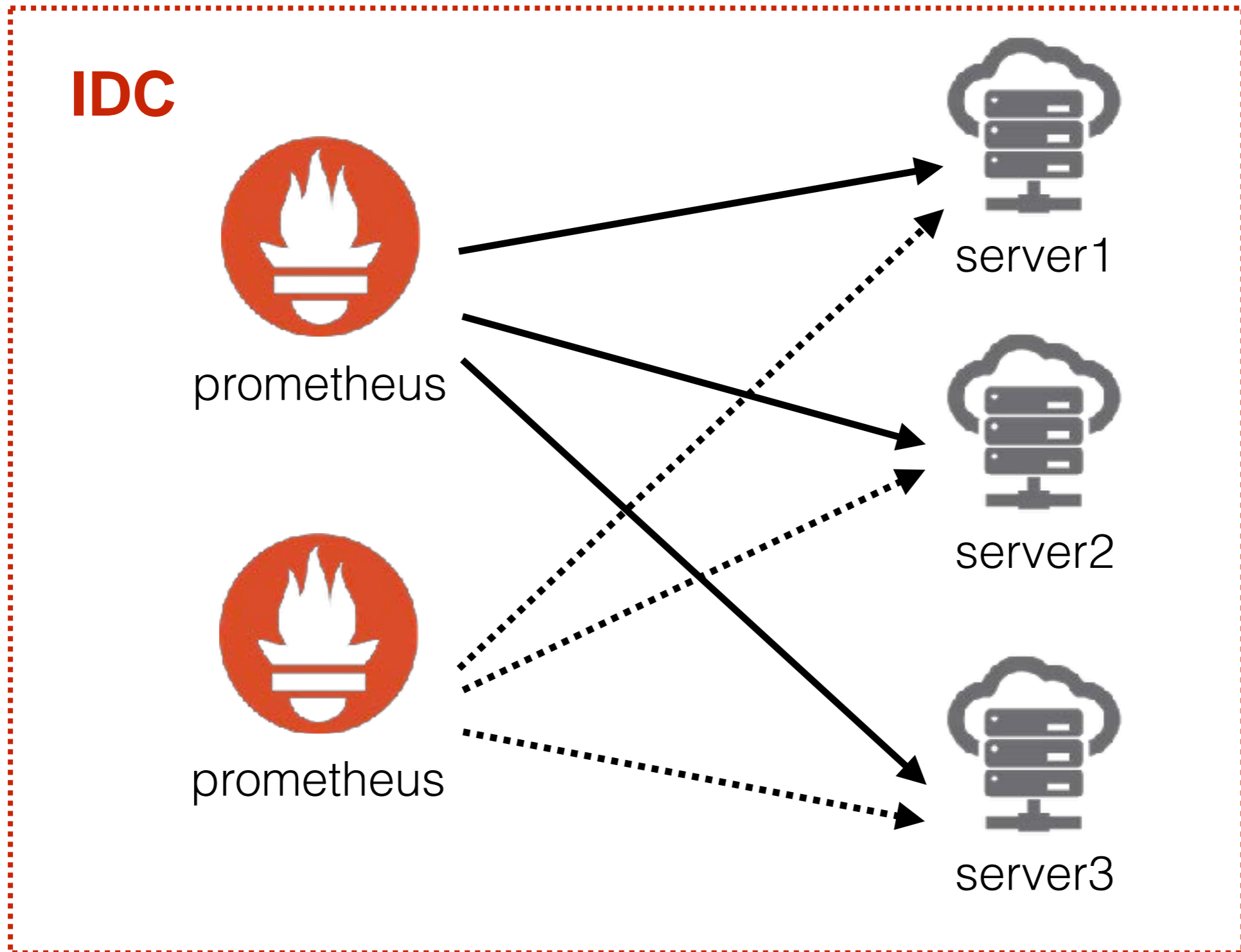
- 性能
- 高可用 ?
- 分布式
- 使用成本

性能问题

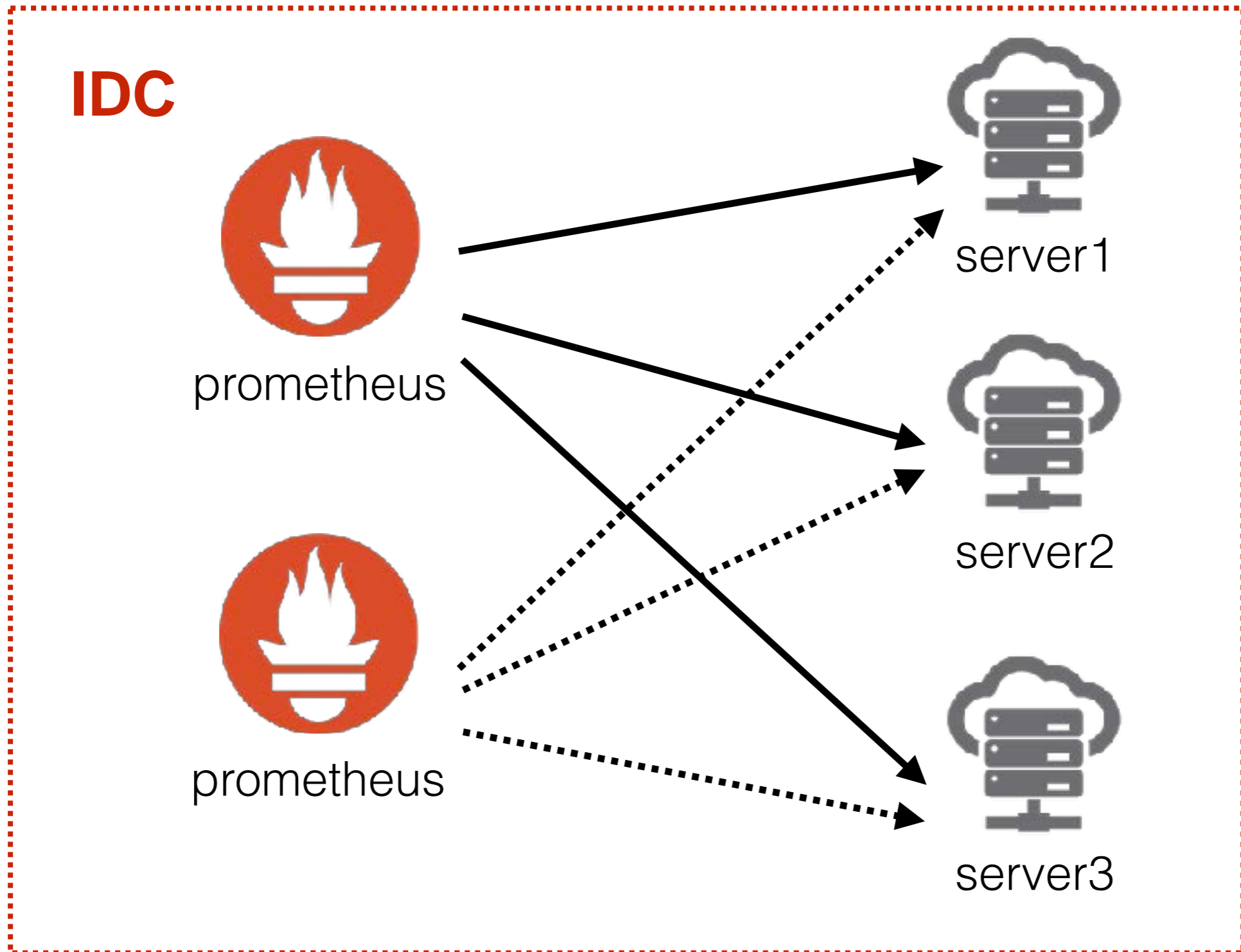
- 本地ssd
- horizontal sharding
(实验性质使用)
- prometheus 2.0 (tsdb)



HA



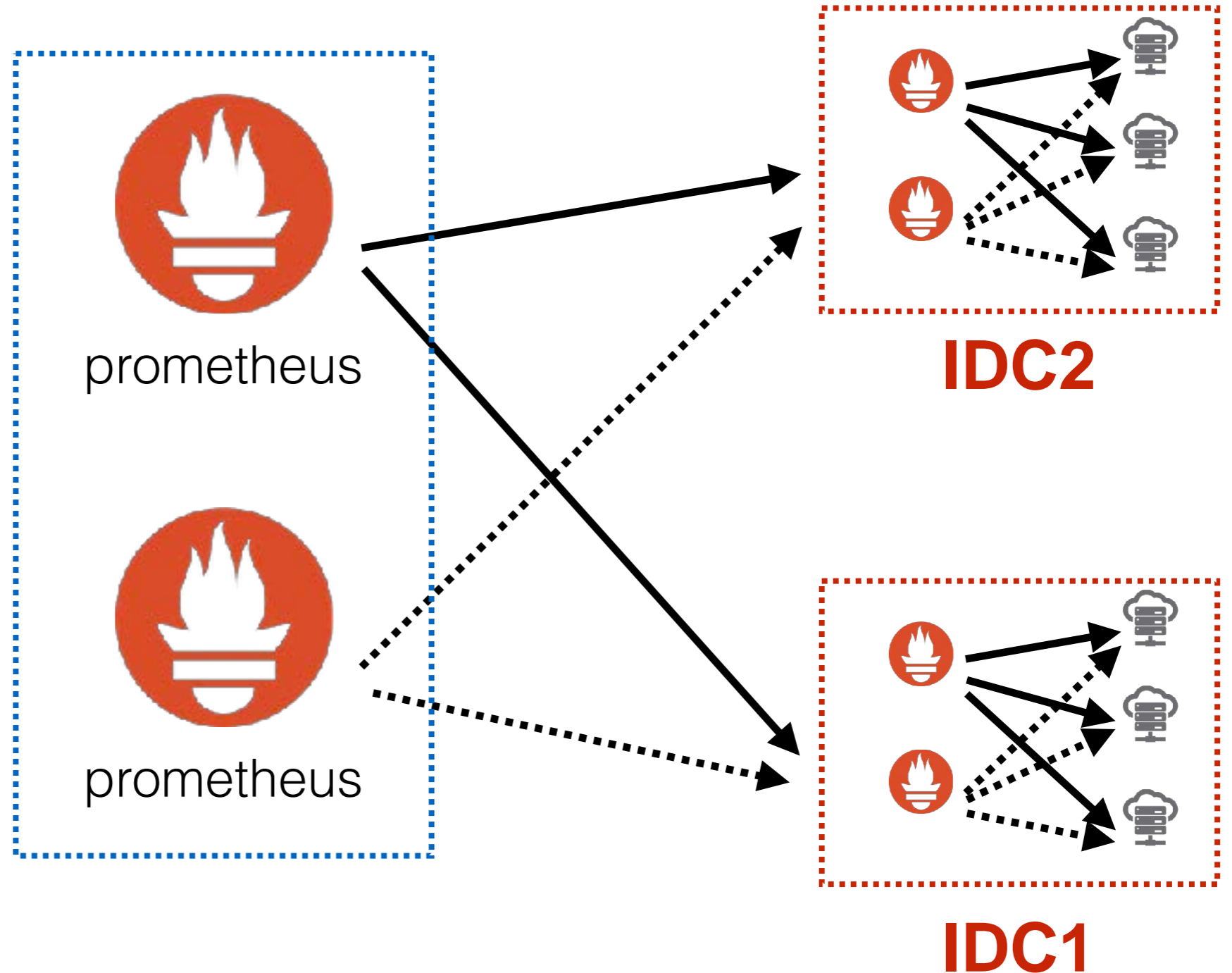
HA



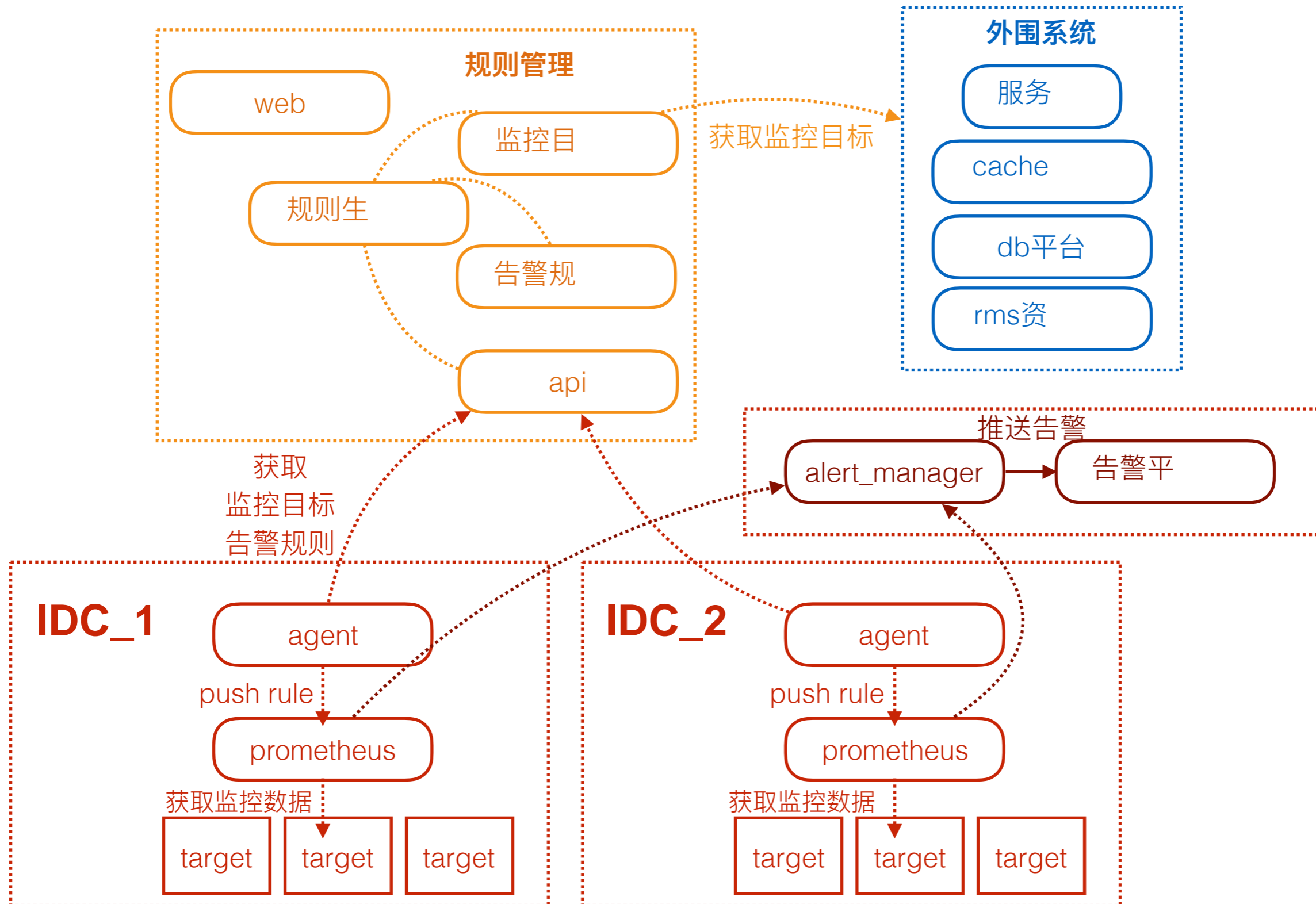
Federation

建议

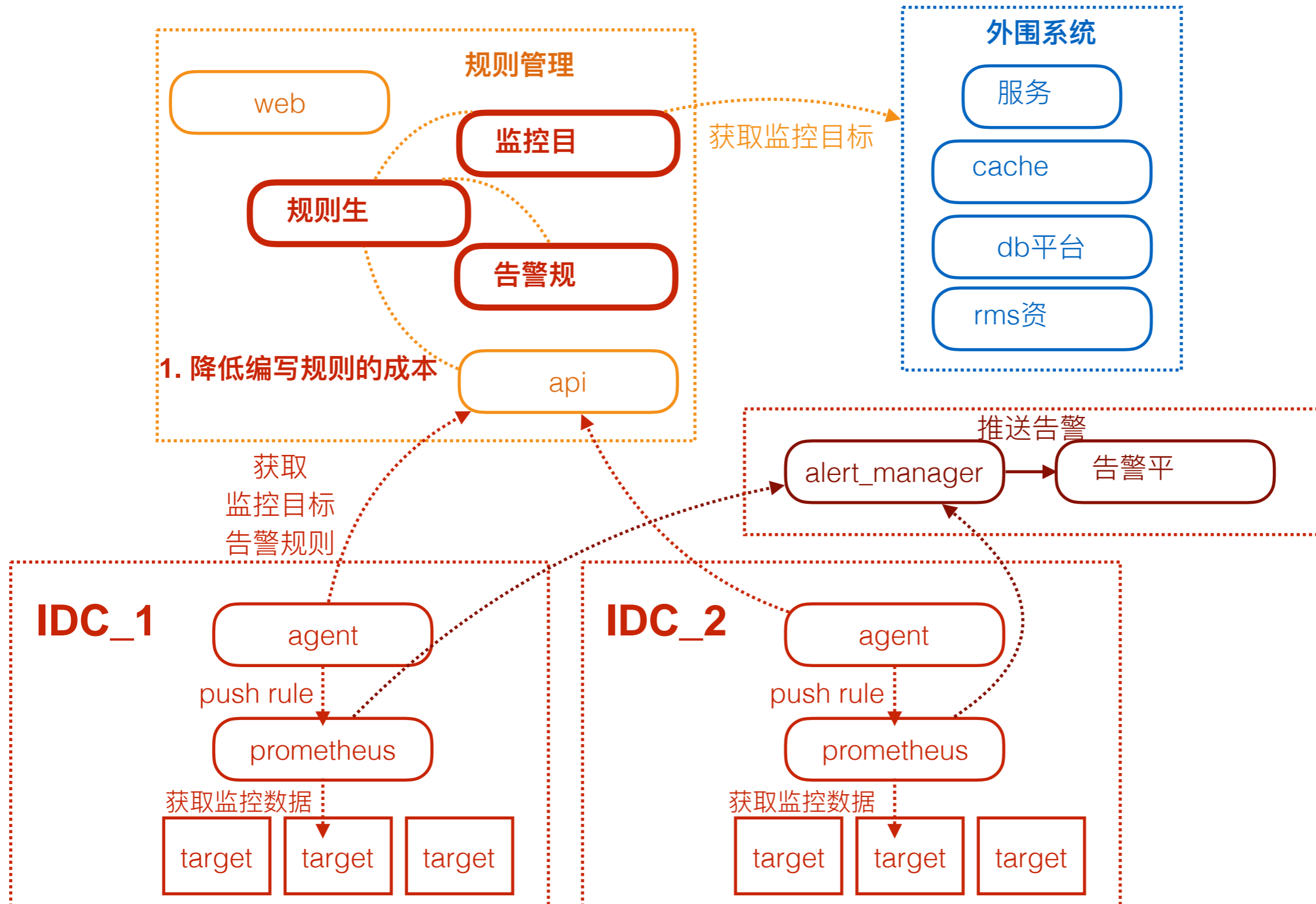
filter数据
精度降低



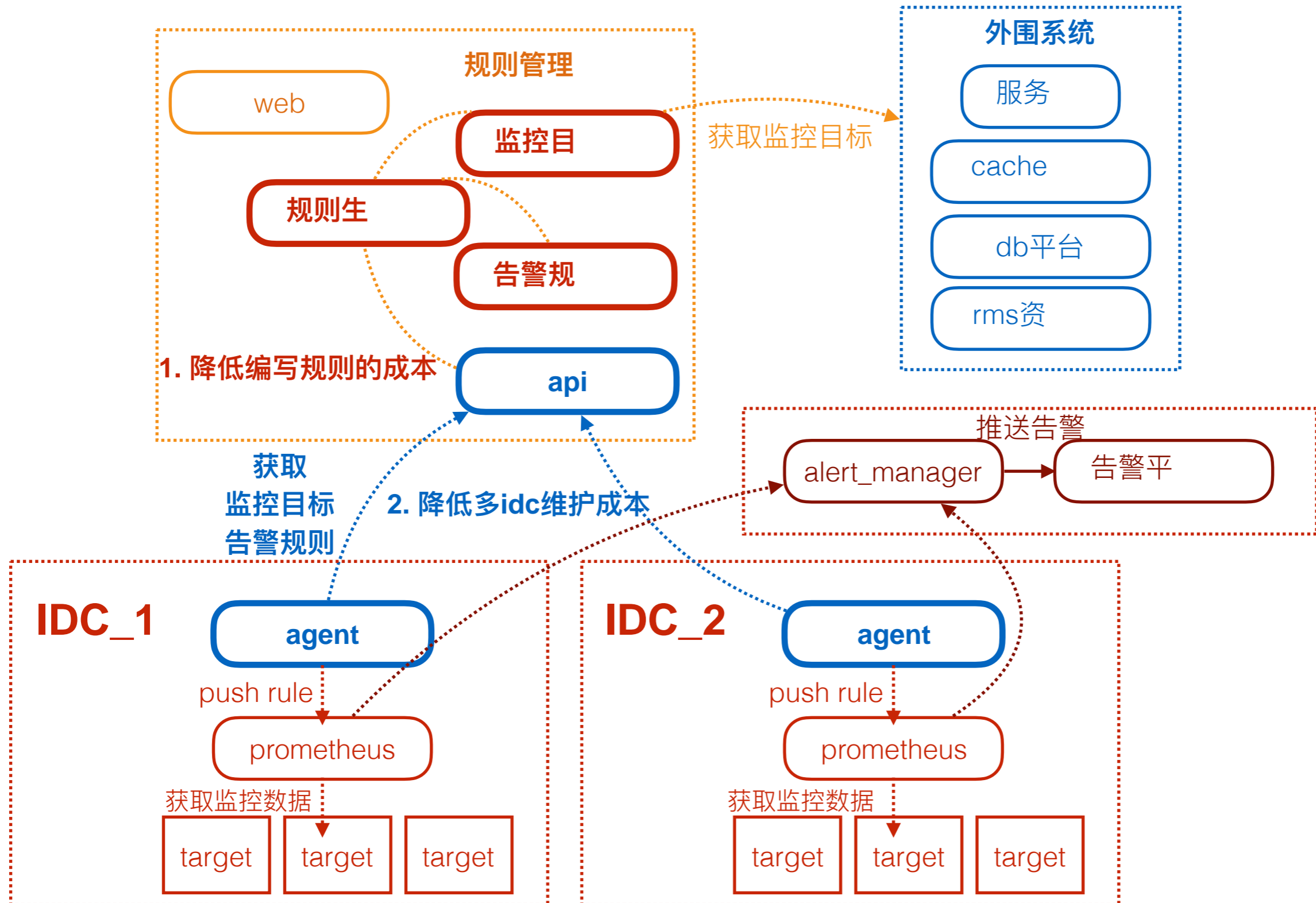
降低使用成本



降低使用成本



降低使用成本



规则管理页面

1 关联资源

* 产品:

* 资源类型:

应用分组:

2 设置报警规则

* 规则描述: % [收起](#) [删除](#)

* 规则描述: % [拓展](#) [删除](#)

+ 添加报警规则

持续时间: (分钟)

生效时间: 至

3 通知方式

通知对象: 联系人通知组

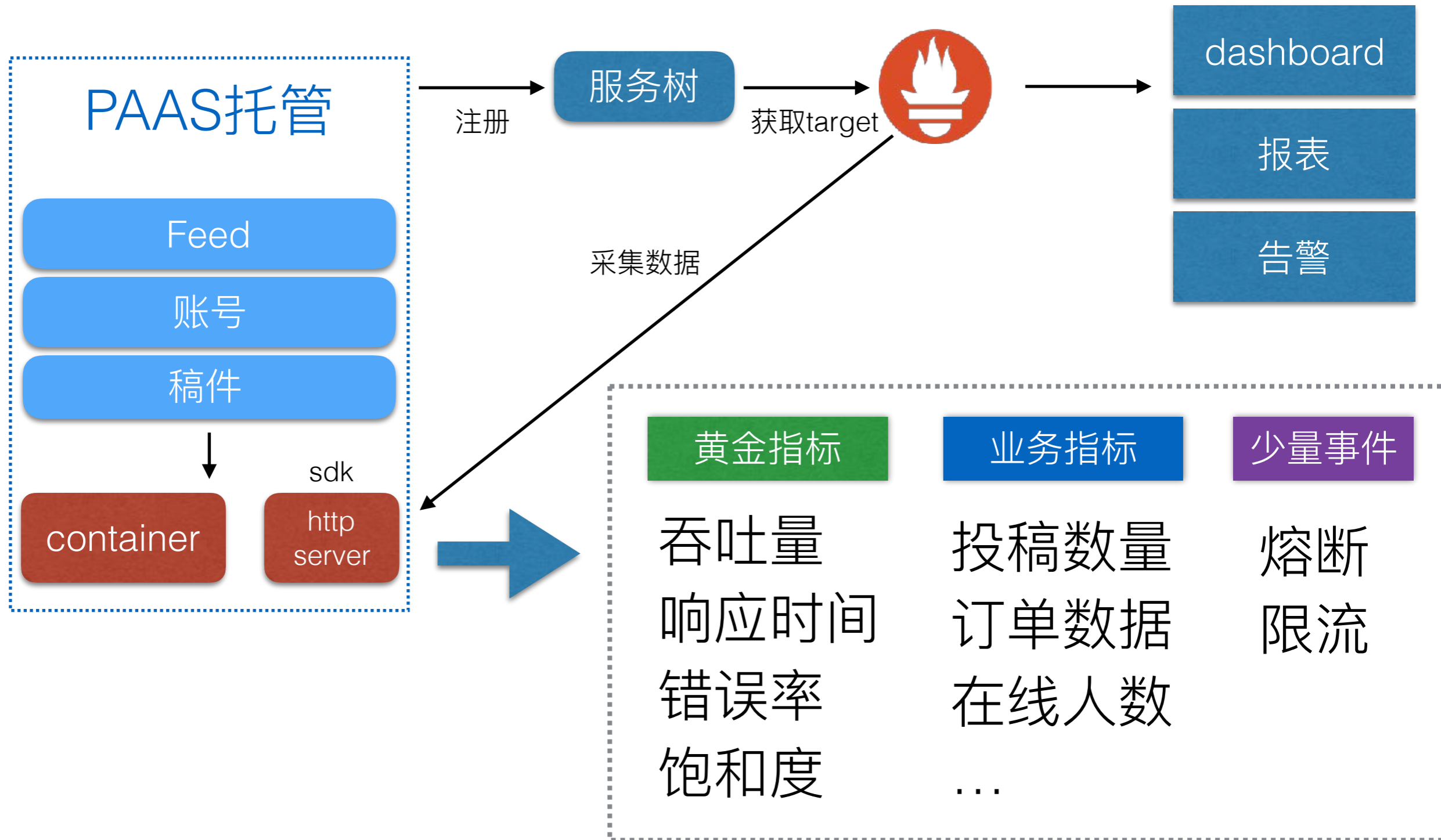
- dev_alarm_group
- /game/fate/yc/运维人员
- /game/fate/its/运维人员
- /game/fate/test/运维人员
- /game/fate/bc/运维人员
- /game/ags/bx/运维人员
- /game/ags/ch2-all/运维人员



已选组

- /game/fate/dev/运维人员

例子 - 业务监控

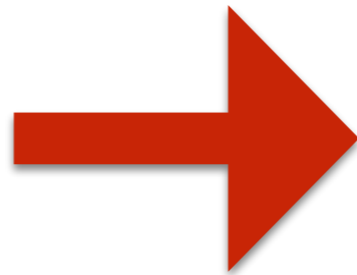


统一的告警中心

解决什么问题？

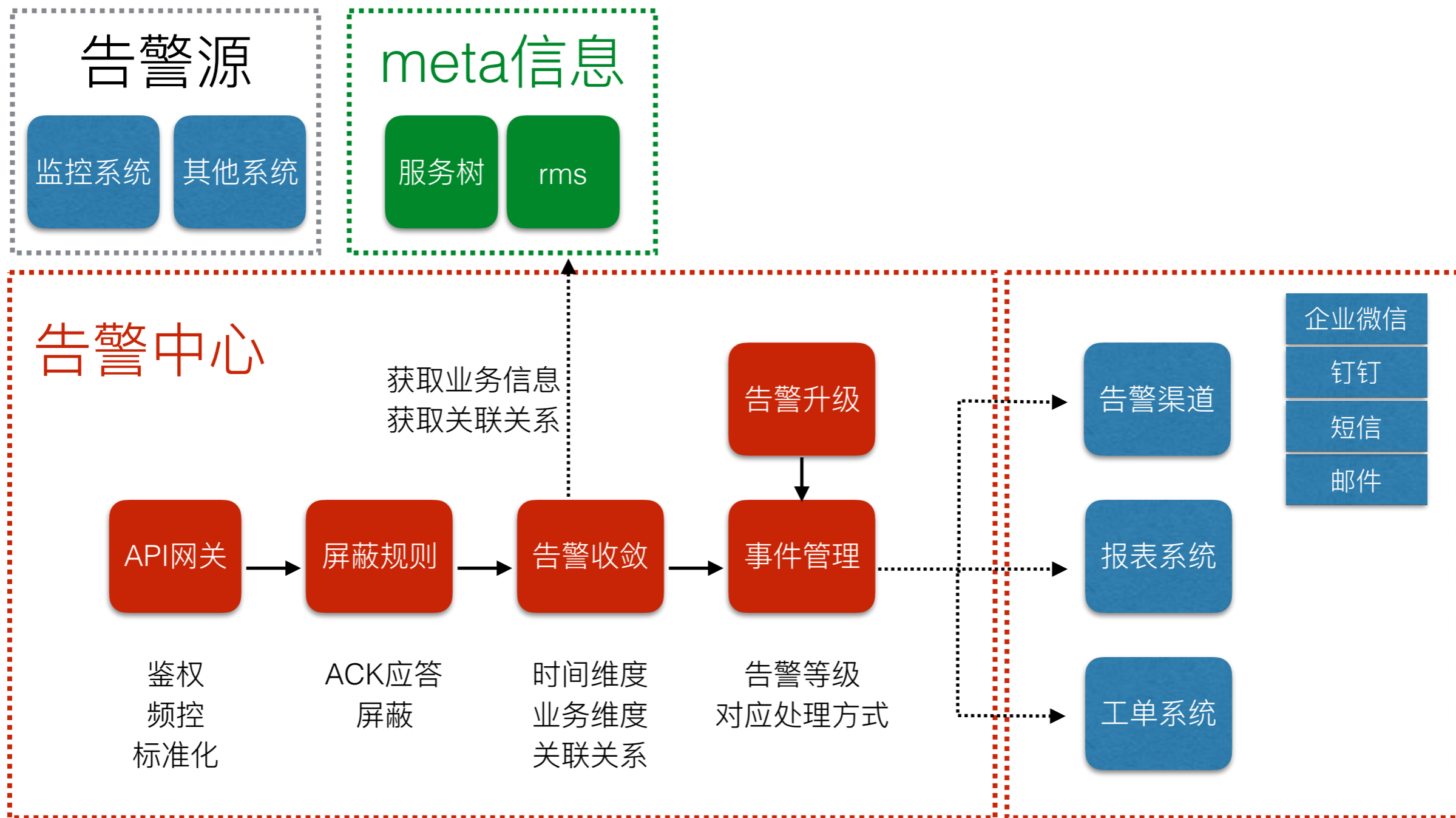
问题

- 告警源头多
- 告警风暴, 大量重复告警
- 发送告警渠道多
- 重要告警没有及时到达
- 优化告警没有数据依据



核心功能

- 告警标准化
- 告警收敛
- 告警渠道管理
- 告警升级
- 告警报表



有意思的尝试

新版[告警]

waf shd-wm-nginx-03 ip_conntrack表使用率 99.66%

[http://www.10000.com/.../.../...](#)

下午2:00

新版[恢复]

livedm shd-shanghai-dmz-01 节点宕机(无法Ping通)

livedm shd-shanghai-dmz-01 节点宕机(无法Ping通)

restart shd-wm-nginx-03

确认重启

Server Code: 80000000

Server Name: shd-wm-nginx-03

请回复 你懂的暗号

8000

提交任务成功,返回结果:Chassis Power Control: Reset

科学的告警策略

科学？

~~machine learning?~~

~~deep learning?~~

不要盲目的使用机器学习

先让告警有意义

可读的

- 时间
- 源头
- 规则
- 影响
- 状态

正确的

- 正确反映现实

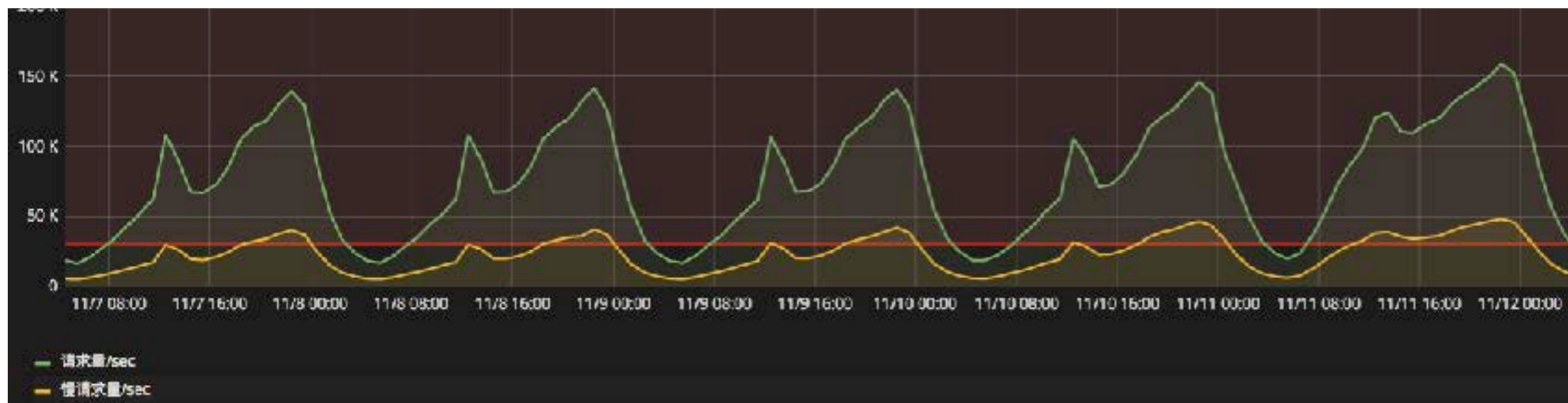
有价值的

- 发现问题

案例1

告警规则: 业务A 慢请求量 > 10k/s

wrong



固定阈值

告警阈值需要随着流量变化而调整

建议:

告警规则: 业务A 慢请求比例 > 80%

案例2

告警规则: 磁盘容量可用率 < 10%

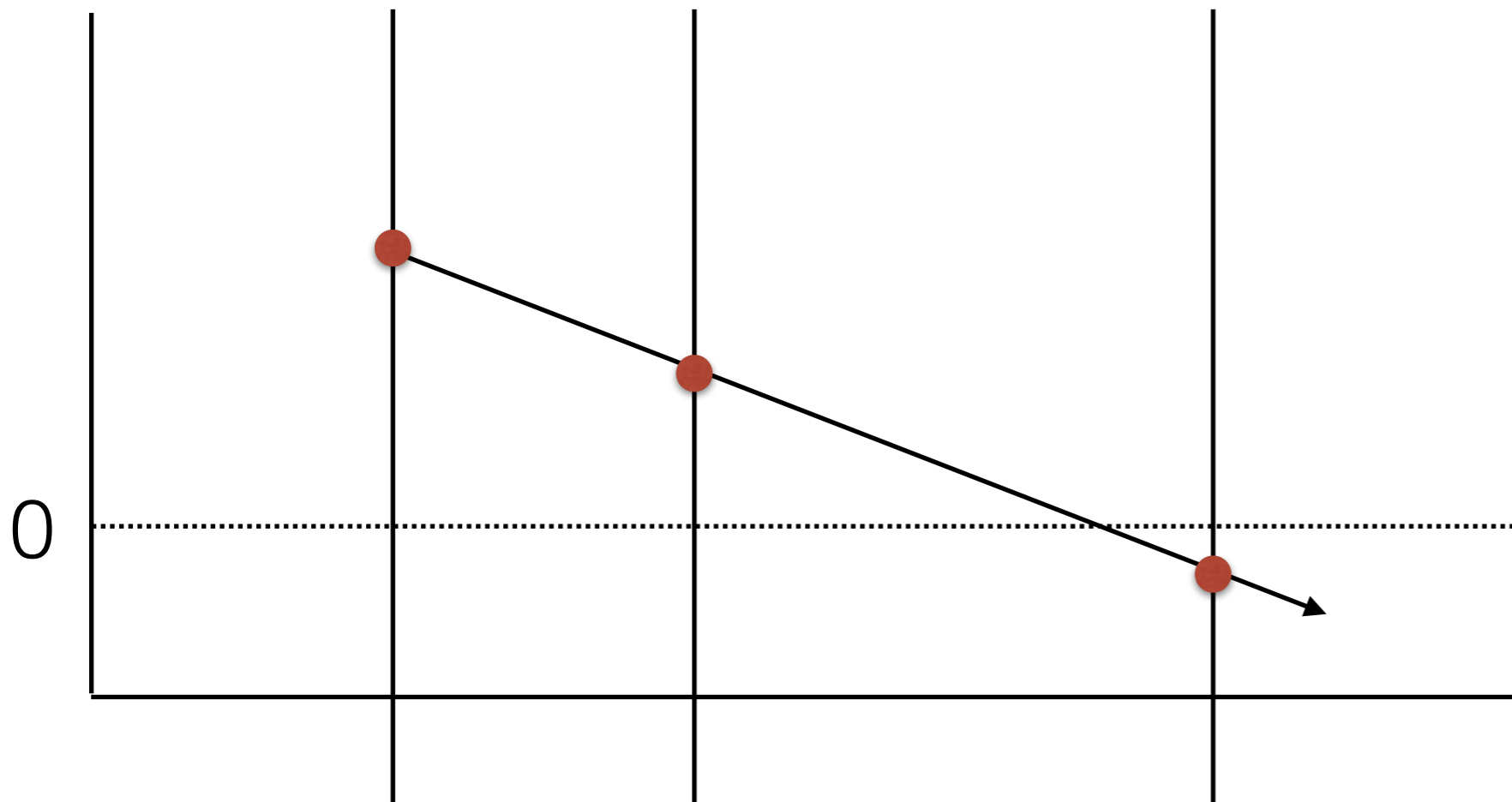
告警规则: 磁盘容量预计将于3小时后饱和

`predict_linear(node_filesystem_free{}[1h], 3 * 3600) < 0`

-1h

now

+3h

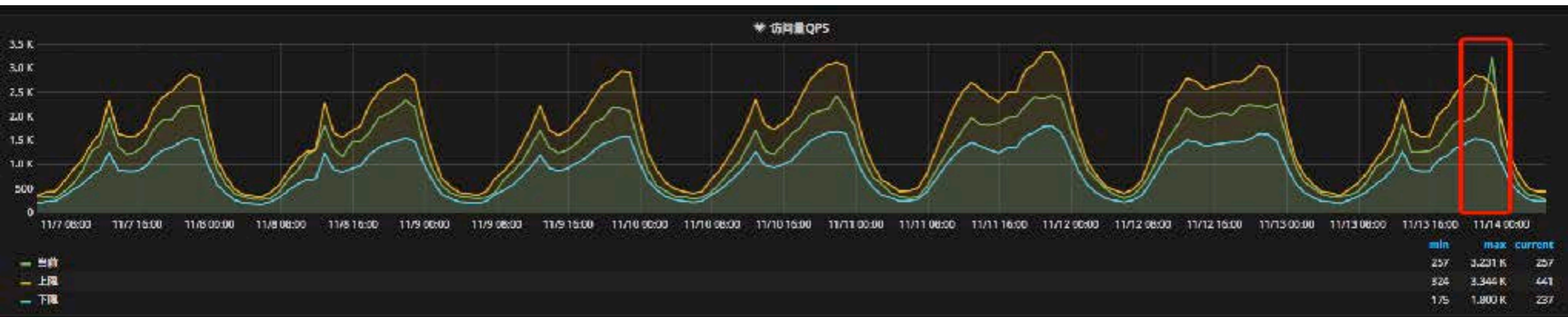


异常检测

异常流量

告警规则: 预测业务A请求量异常

$\text{abs}(\text{requests} - \text{requests:holt_winters_rate1h offset 7d}) > 0.3 * \text{requests:holt_winters_rate1h offset 7d}$



异常响应

```
(
  instance:latency_seconds:mean5m
> on (job) group_left()
  (
    avg by (job)(instance:latency_seconds:mean5m)
    + on (job)
      2 * stddev by (job)(instance:latency_seconds:mean5m)
  )
)
> on (job) group_left()
  1.2 * avg by (job)(instance:latency_seconds:mean5m)
and on (job)
  avg by (job)(instance:latency_seconds_count:rate5m) > 1
```


todo

针对历史事件

- 异常事件关联关系挖掘
- 全链路模块调用分析
- 瓶颈分析

针对当前事件

- 异常检查(动态阈值)
- 异常定位(根因分析)
- 快速止损

针对未来事件

- 故障预测
- 容量预测
- 趋势预测

Thank You!

哔哩哔哩 - (° - °)つ□ 乾杯~ - bilibili

