

小米故障应急响应经验分享

2024年6月



<https://sre-elite.com>



王哲
小米云平台运维部

工作经历:

09-12, 新浪网, 平台架构部
12-21, 奇虎360, 云平台部
21-now, 小米, 云平台部

个人简介:

- 微博的第一版代码我上的线
- 从零构建360内部私有云平台Hulk
- 运维开发-产品经理-技术管理

目录

CONTENT

01 | 故障应急面临的挑战

02 | 构建哪些能力应对

03 | 一些实践经验

04 | 总结

01

小米故障应急面临的挑战

小米集团战略

「人车家全生态」

Human X Car X Home
All Your Needs in one smart ecosystem



人、设备、智能服务之间 相互协作，共同进化

让智能跟随人的需求而流动
同时，人把智能带给生态，生态不断学习和进化，反过来服务人

小米澎湃OS

6.99亿

可连接智能设备*

200+

品类

95%+

覆盖用户生活场景

了解详情 [👉](#)

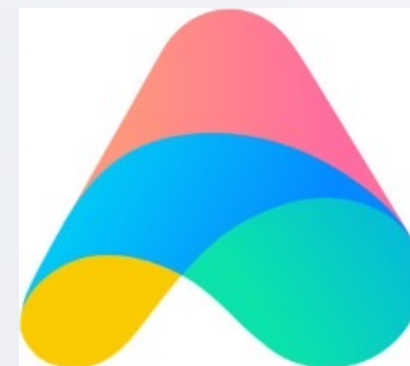
全球第三的手机制造厂商&全球领先的消费级AIoT平台



米家760万MAU



7亿联网设备



小爱1.14亿MAU

20+
业务线

12个
核心机房

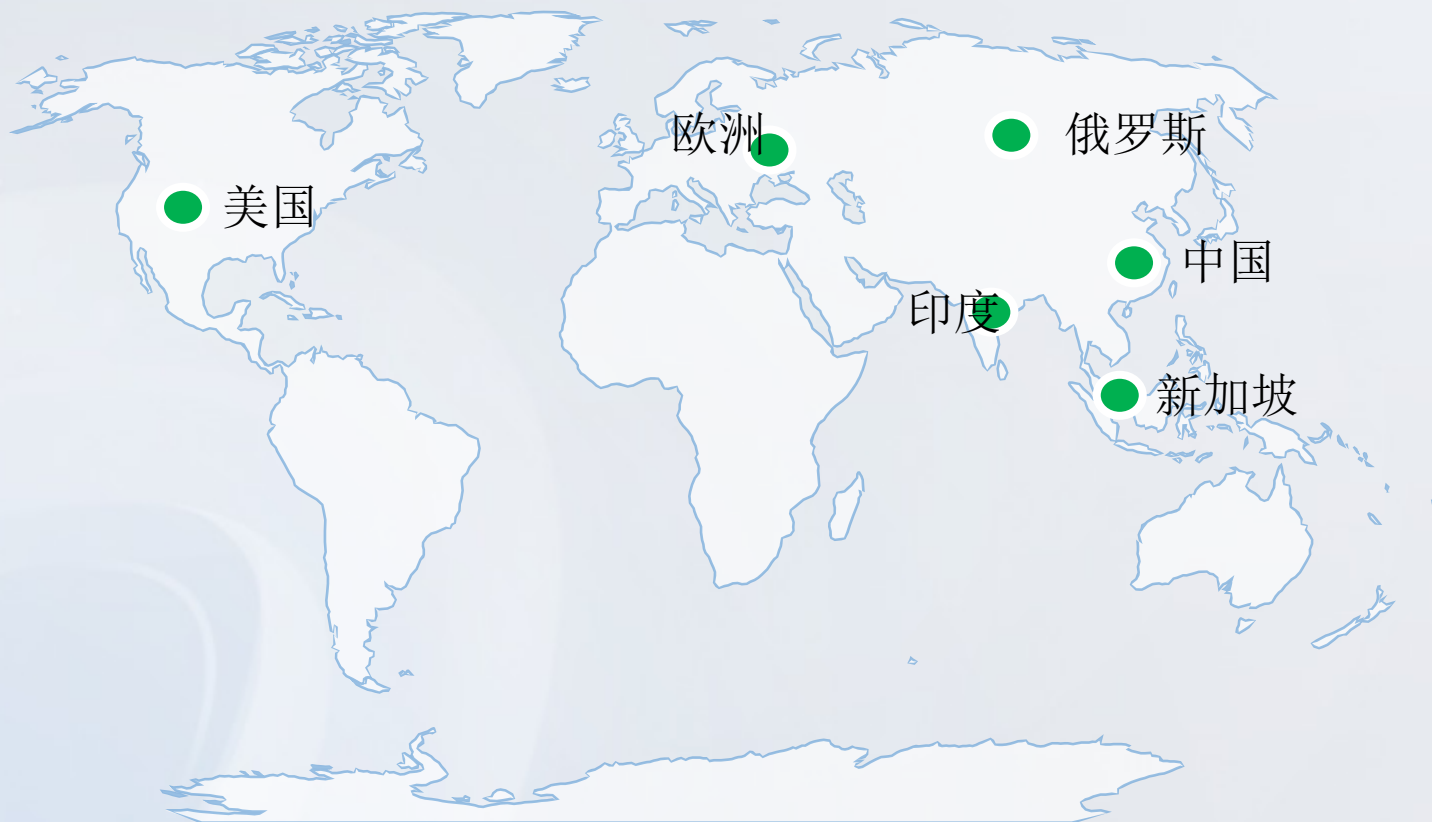
X万
服务器

100亿
日请求量

XPB
存储

业务全球布局, 数据本地存储

自有IDC, 阿里云, 金山云, 腾讯云, 火山云, AWS, Azure, GCP



02

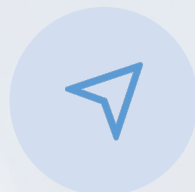
构建哪些能力应对故障响应

故障发现



一分钟内发现问题

故障诊断



五分钟内定位问题

故障恢复



十分钟内控制损失

故障复盘



02-1

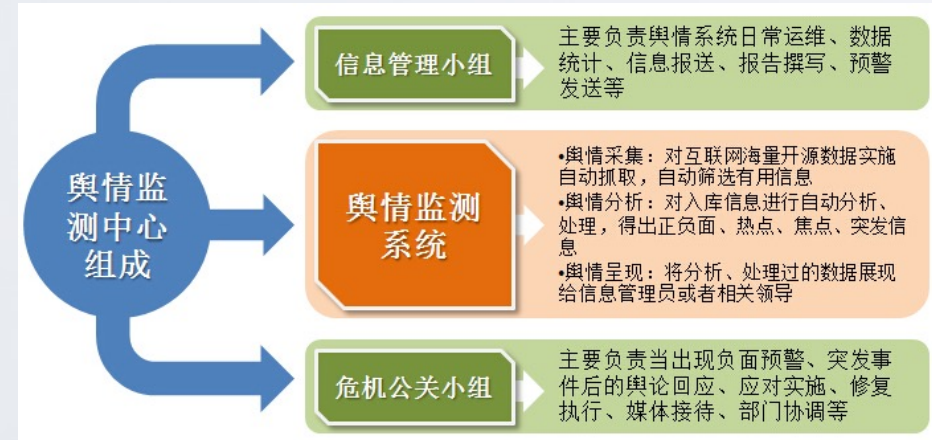
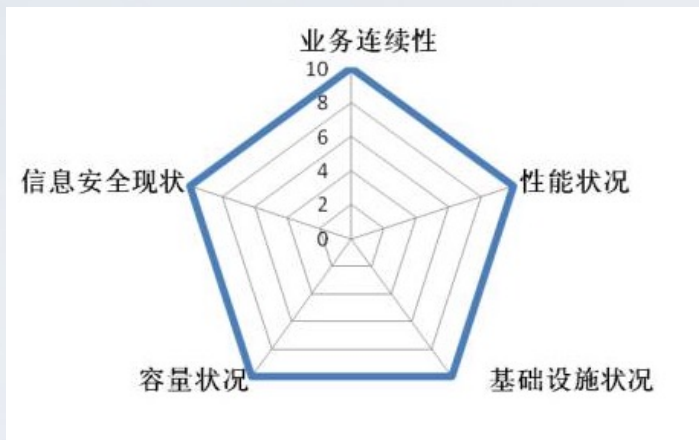
故障发现

发现故障的几种手段

监控

巡检

上报



基础设施监控

- 混合云环境下的IT基础架构信息实时监控（性能/可用性）

应用程序性能监控

- 应用程序的功能/运行状况监控
- APM, 跨微服务, 主机, 容器和Serverless

日志管理

- 为应用, 系统和云服务提供日志管理, 可创建索引, 查询可视化, 报警
- 支持机器学习, 实现预测和阈值自定义等能力

用户体验监控

- 模拟客户监控产品的体验和可用性
- 真实用户监控

网络性能监控

- 云或者混合环境的网络流量进行分析和可视化处理



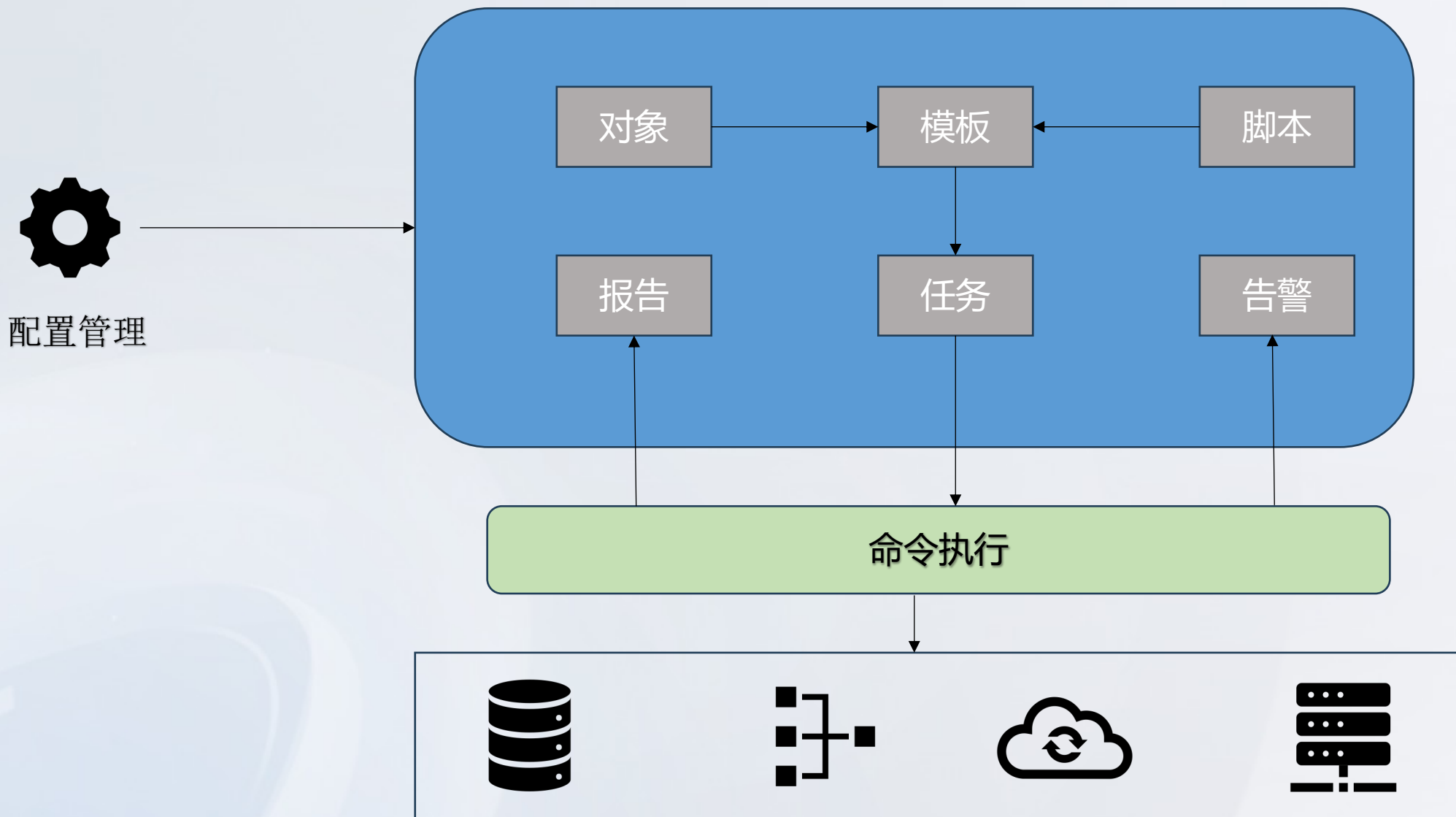
协同



易用



整合



巡检指标



关键词：自动化的监控不能绝对解决的问题，依赖终端用户的体验上报是很好的兜底手段

舆情监听

- 媒体信息收集
- 论坛讨论热点

客服反馈

- 用户投诉上报
- 客服统计数据
- VIP重点服务关注

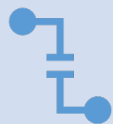
运营反馈

- 新功能上线检查
- 日常操作反馈

02-2

故障诊断

关键词：故障恢复是第一位！



故障范围：

快速确认故障影响模块，逻辑功能依赖和上下游服务依赖



根因确认：

找到关键告警和分析入口，定位系统和责任人



影响评估：

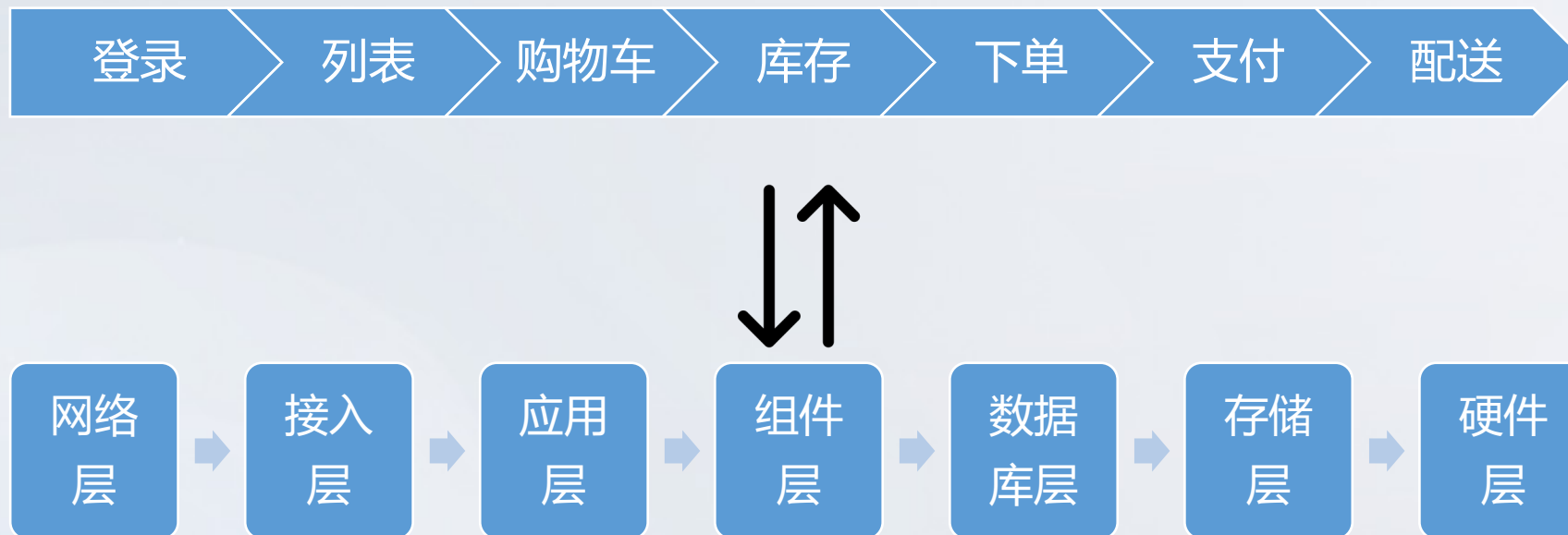
基于故障范围，确认影响情况，决策信息同步周期和上报级别



组织能力：

值守，应急响应，升级流程，资源保障，信息同步...

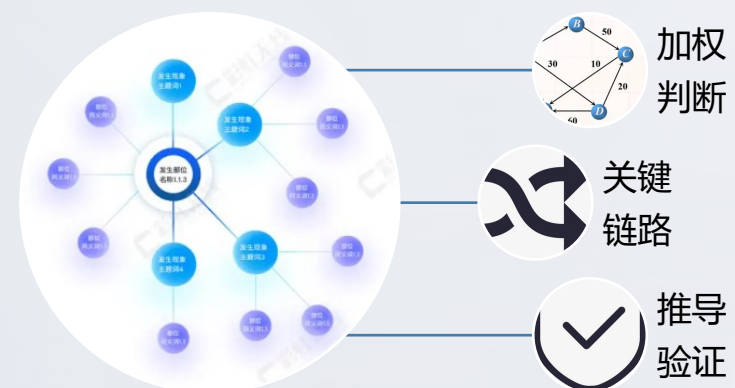
关键词：SRE向上要理解业务功能模块，向下要精通云基础技术组件



关键词：AIOps最卷的就在这里

归纳分析（点：聚类）

演绎推理（线和面：关联）

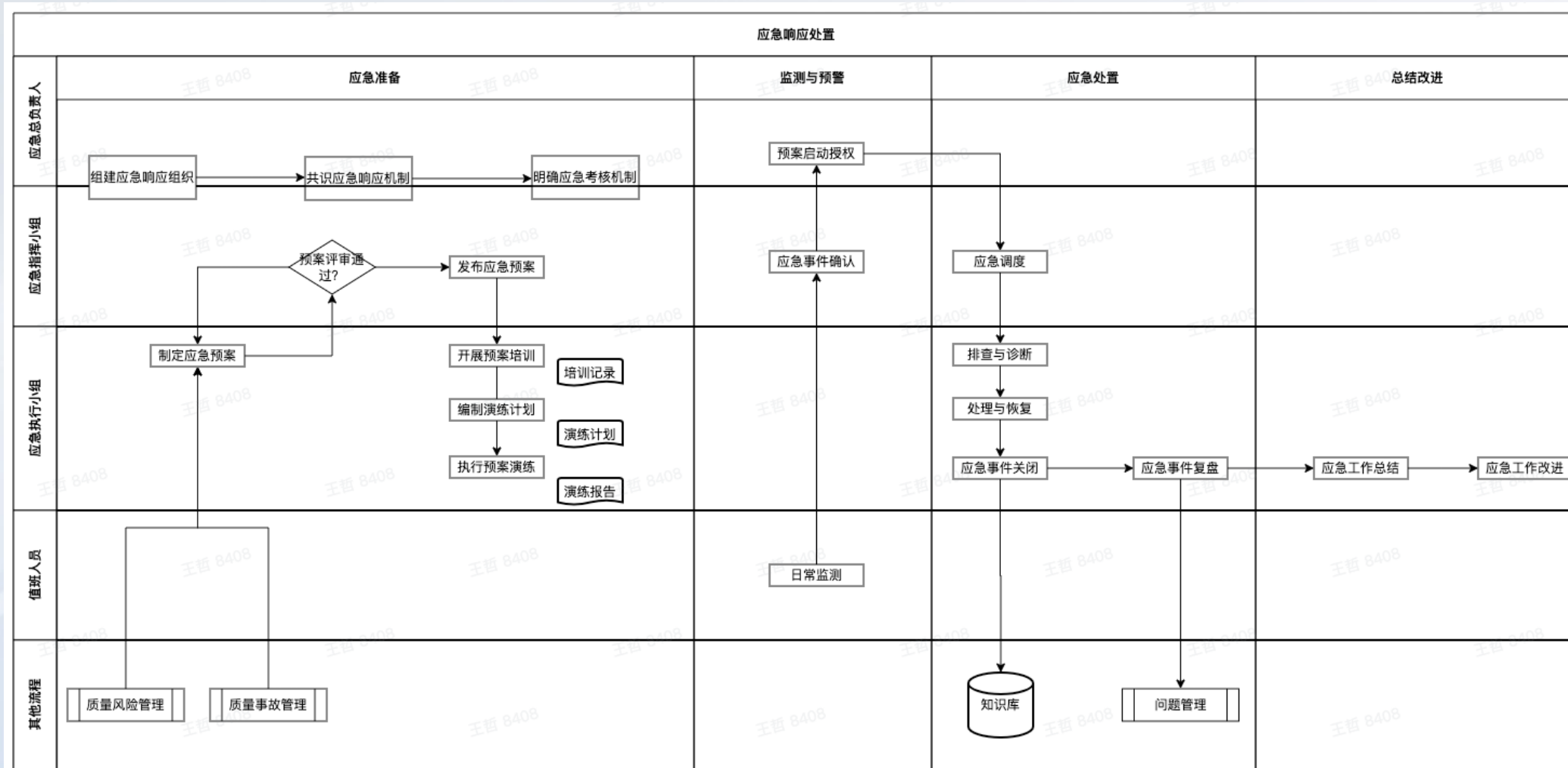


知识图谱

关键词：及时升级！

用户请求损失	用户无法使用功能造成的损失
影响用户数量	实际影响的用户数量
影响用户时长	精确的故障时长
客服来访数量	用户通过客服的报障数
用户反馈数量	用户通过反馈渠道的投诉反馈
直接经济损失	因为故障造成的经济损失
政府监管	抽检, 召回, 立案
市场舆情	维权事件, 媒体报道, 社会话题
售后维修	维修工单/销量
集团经营影响	产能, GMV金额, 核心业务时长

不同的层次都由不同的团队来负责运维管理
同层次不同的硬件/系统/应用都由不同的小组来负责运维管理



02-3

故障恢复

关键词：每一个动作前都应该知道预期的结果，每一步真正的操作后都需要充分检查

1

重启：

单个或多个机器上的服务出现响应问题，先重启就能先恢复，能恢复就能止损

2

回滚：

确认跟发布相关，回滚是最好的选择，最起码减少变量，更快定位

3

扩容：

确认硬件跑满，压力增大，立即扩容也是最佳选择之一

关键词：每一个动作前都应该知道预期的结果，每一步真正的操作后都需要充分检查

1

限流：

控制入口流量，减少超出后端承载量的请求，自动丢弃新来的请求

2

降级：

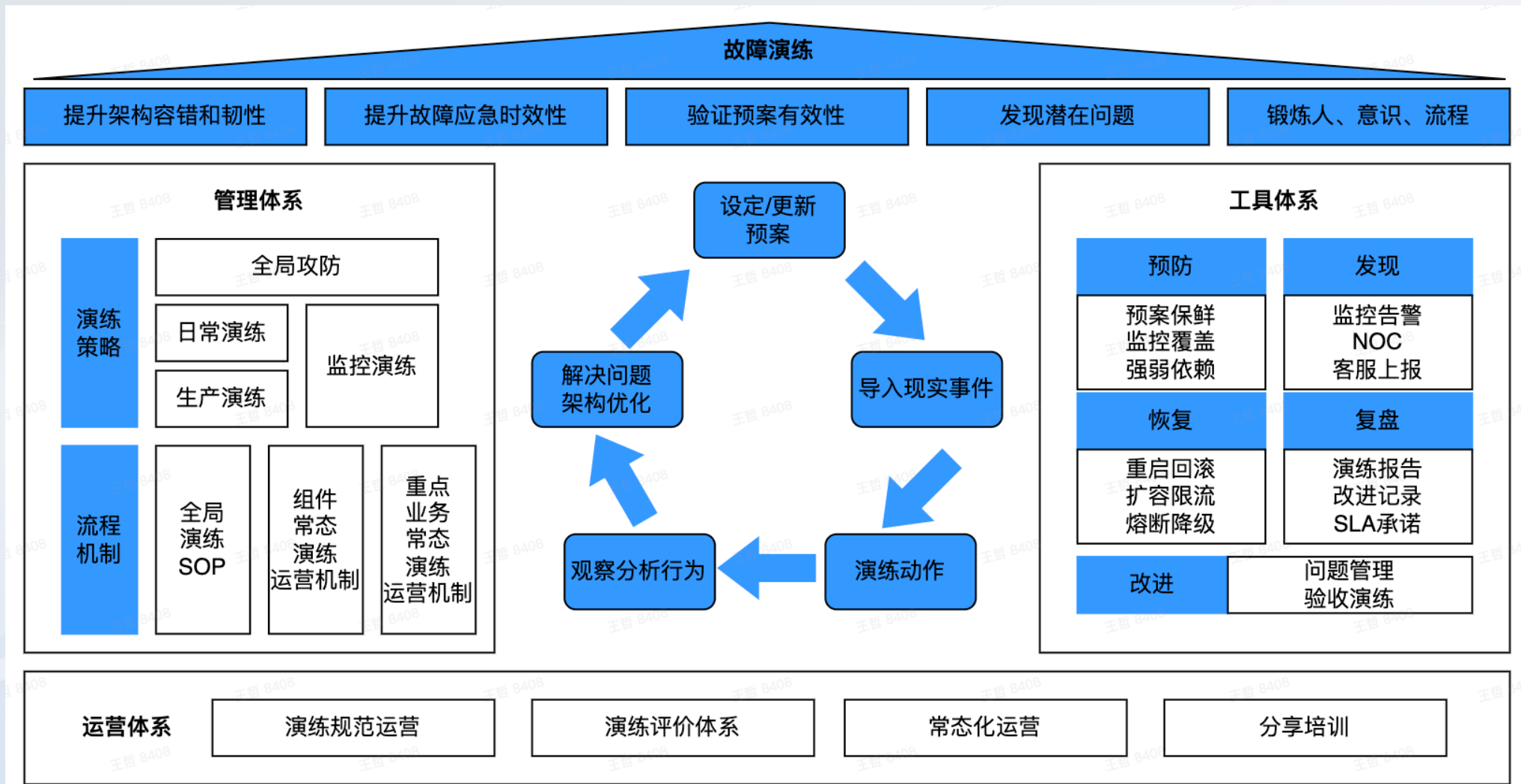
有意识的降低系统的部分功能和服务质量，确保系统的核心功能和关键服务继续运行

3

熔断：

依赖第三方服务异常，导致自身响应时长超长情况下快速返回失败而不是等待，阻止级联效应

关键词：预案和爱情一样都有保鲜期



02-4

故障复盘

关键词：避免类似的问题重复发生

还原事实

- 过程中做了哪些动作
- 依据的流程机制标准是啥

问题反思

- 存在什么问题
- 做对了什么做错了什么

原因分析

- 主观原因
- 客观原因

重来一遍

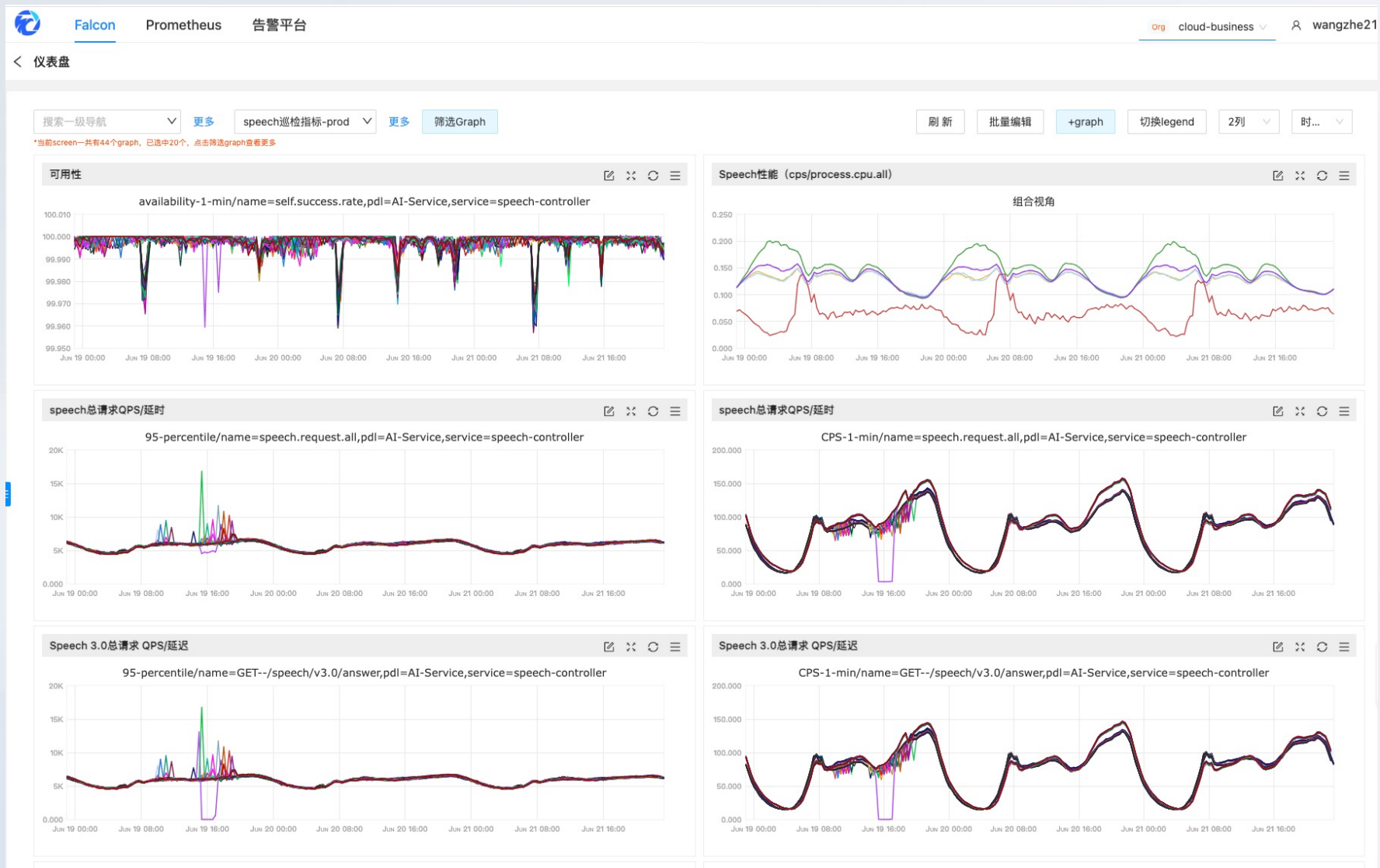
- 重来一遍怎么做
- 明确重点改进计划

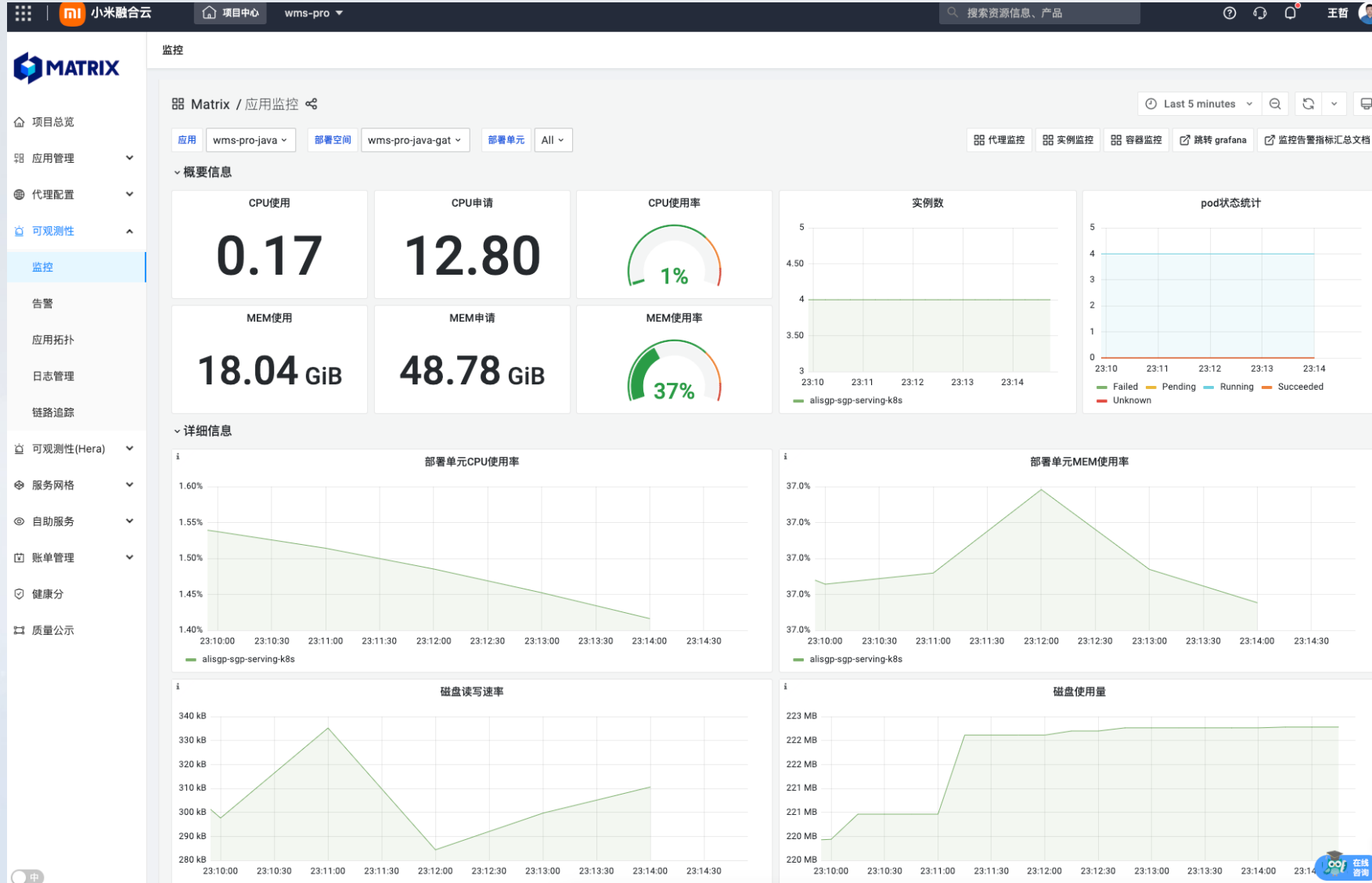
关键词：避免类似的问题重复发生



03

小米的一些实践经验





【PROBLEM】【P1】可用性 mibi.api.xiaomi.com 给商户提供的服务端接口服务

域名: mibi.api.xiaomi.com 米币.米币手机端APP接口 给商户提供的服务端接口服务

集群: cnbj4-talos

报警条件: 可用性 all(#3) 93.617 < 99.95

报警时间: 2022-09-22T23:50:00

报错数量: count: 6, rate: 6.383

http状态码: ('500', 6)

nginx分布(top 10):

('c3-miui-mibi-api04.bj', 4)

('c3-miui-mibi-api03.bj', 2)

路径分布(top 5):

('/_j_security_check', 6)

客户端IP(top 5):

('121.62.63.178#湖北', 6)

后端服务(top 10):

('c3-miui-mibi-fe03.bj:9025#mibi-web', 5)

('c3-miui-mibi-fe04.bj:9025#mibi-web', 1)

[看图](#) [关联分析](#) [编辑路径](#)

处理人: @叶金鑫

认领

屏蔽2小时

屏蔽1天

封禁IP

容量巡检

CPU机器 MICE容器 Matrix容器 GPU机器 主机信息(达标/过保/故障)

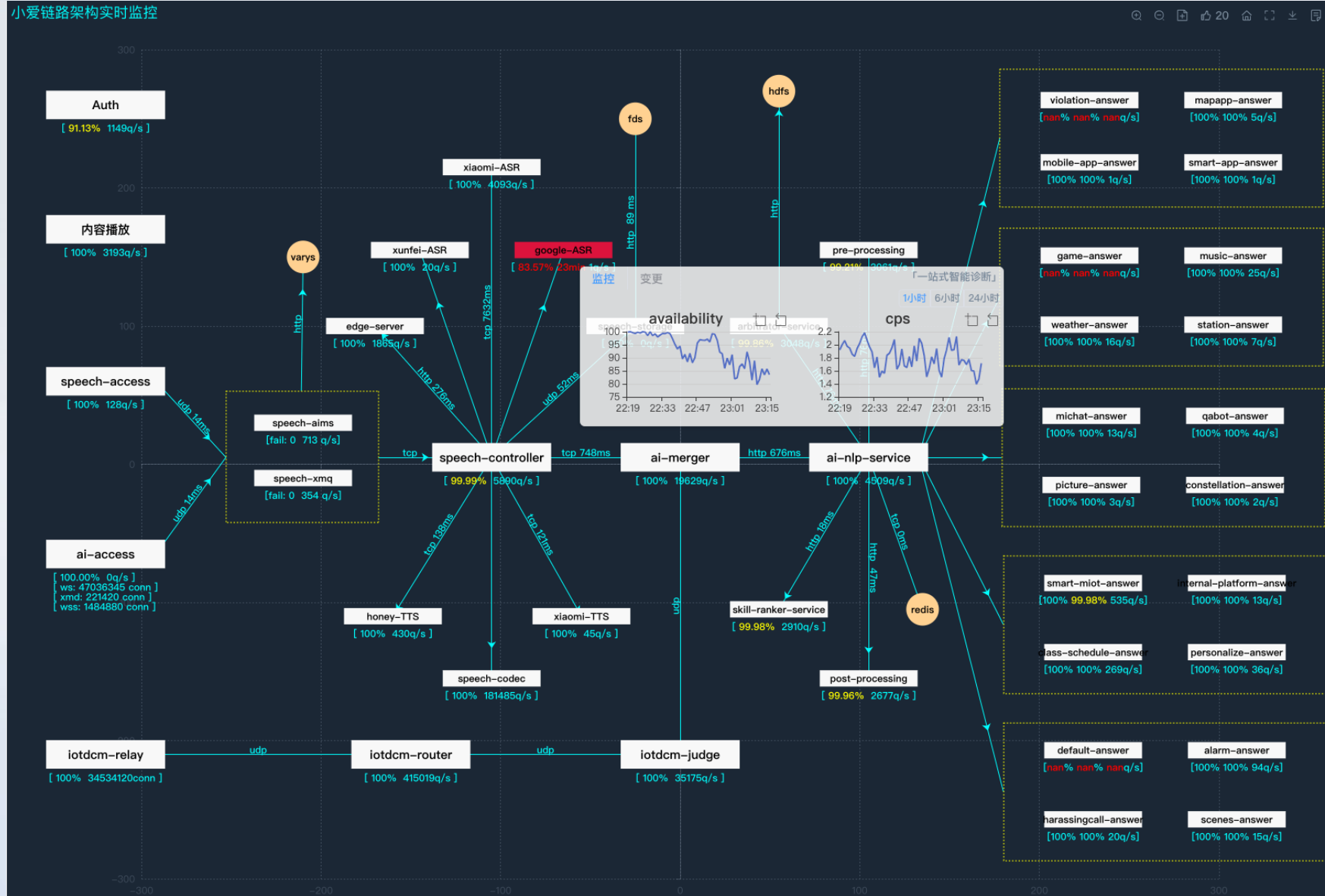
分组: 111 机器: 1178

cpu.busy 近一周 分位选择 阈值: > 70 % 均值: 29.8% 机房选择

组名	用量	数量	c3	c4	c5	qrc6	tj1	sh	others
ai-archtest	cpu: 99.7% mem: 85.4%	4	c3-ai-archtest00.cj: 24.4%	c4-ai-archtest00.cj: 2.8%		qrc6-ai-archtest00.cj: 8.0%			sgp1-ai-archtest00.cj: 0.0%
ai-ai-service-ksc...	cpu: 97.9% mem: 38.8%	1					tj1-ai-ai-service-ksc-va2-2-10341-9540vhd0.kscn: 97.9%		
ai-prod-mon-it	cpu: 96.8% mem: 57.2%	2	c3-ai-prod-mon-it01.cj: 96.8% c3-ai-prod-mon-it02.cj: 87.9%						
ai-prod-smartmiot	cpu: 84.1% mem: 88.0%	31	c3-ai-prod-smartmiot01.cj: 52.3% c3-ai-prod-smartmiot02.cj: 52.2% c3-ai-prod-smartmiot03.cj: 48.2% c3-ai-prod-smartmiot04.cj: 47.7% c3-ai-prod-smartmiot05.cj: 47.1% c3-ai-prod-smartmiot06.cj: 46.9% c3-ai-prod-smartmiot07.cj: 46.6% c3-ai-prod-smartmiot08.cj: 46.0% c3-ai-prod-smartmiot09.cj: 45.4% c3-ai-prod-smartmiot10.cj: 42.8%	c4-ai-prod-smartmiot01.cj: 84.1% c4-ai-prod-smartmiot02.cj: 55.5% c4-ai-prod-smartmiot03.cj: 51.5% c4-ai-prod-smartmiot04.cj: 51.0% c4-ai-prod-smartmiot05.cj: 50.7% c4-ai-prod-smartmiot06.cj: 48.7% c4-ai-prod-smartmiot07.cj: 48.2% c4-ai-prod-smartmiot08.cj: 47.9% c4-ai-prod-smartmiot09.cj: 46.8% c4-ai-prod-smartmiot10.cj: 45.9%					
miot-camera-ec2...	cpu: 82.5% mem: 17.5%	27					tj1-miot-camera-ec2-master07.kscn: 82.5% tj1-miot-camera-ec2-master03.kscn: 74.9% tj1-miot-camera-ec2-master06.kscn: 74.8% tj1-miot-camera-ec2-master02.kscn: 73.7% tj1-miot-camera-ec2-master05.kscn: 69.9% tj1-miot-camera-ec2-master04.kscn: 68.0% tj1-miot-camera-ec2-master08.kscn: 66.7% tj1-miot-camera-ec2-master01.kscn: 34.6% tj1-miot-camera-ec2-master00.kscn: 0.7%	sh1-miot-camera-ec2-master03.kscn: 77.1% sh1-miot-camera-ec2-master02.kscn: 69.0% sh1-miot-camera-ec2-master01.kscn: 64.7% sh1-miot-camera-ec2-master06.kscn: 63.3% sh1-miot-camera-ec2-master05.kscn: 59.5% sh1-miot-camera-ec2-master04.kscn: 59.2% sh1-miot-camera-ec2-master07.kscn: 57.4% sh1-miot-camera-ec2-master00.kscn: 53.3% sh2-miot-camera-ec2-master00.kscn: 0.4%	gz1-miot-camera-ec2-master05.kscn: 65.1% gz1-miot-camera-ec2-master07.kscn: 63.9% gz1-miot-camera-ec2-master02.kscn: 61.6% gz1-miot-camera-ec2-master01.kscn: 60.6% gz1-miot-camera-ec2-master03.kscn: 59.4% gz1-miot-camera-ec2-master06.kscn: 57.6% gz1-miot-camera-ec2-master08.kscn: 55.9% gz1-miot-camera-ec2-master04.kscn: 55.5% gz1-miot-camera-ec2-master00.kscn: 1.2%

选择机房, 最少选择1个机房

- 全选
- c3 c4 c5 qrc6
- tj1 sh mos mb
- or sgp fr



查看全部 cop.xiaomi_owt.AI_pdl.AI-Service



磁盘总量
5.6 PB



磁盘剩余
3.4 PB / 60.7%



覆盖机器
2381



清理次数
531,543

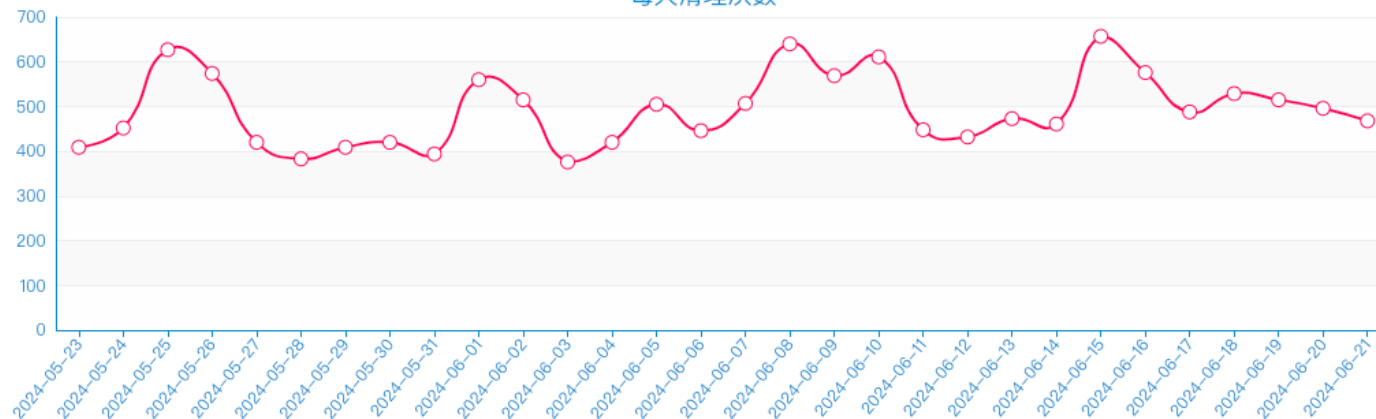
清理空间
118.6 PB

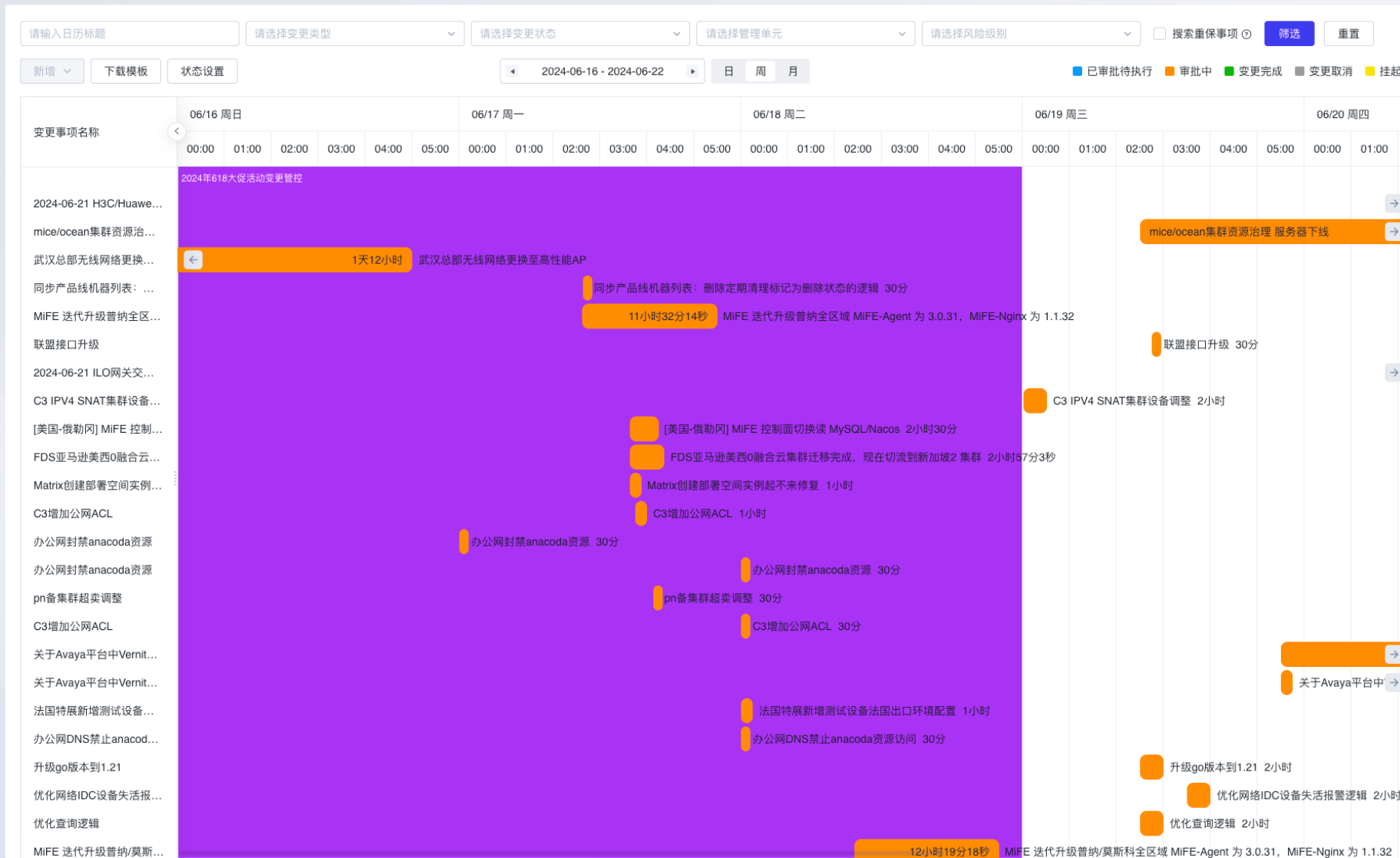
清理次数 Top 10

Hostname	Count
[redacted]	8611
[redacted]	8561
[redacted]	7557
[redacted]	7467
[redacted]	7423
[redacted]	7266
[redacted]	6924
[redacted]	6613
[redacted]	6140
[redacted]	5983

7天 30天

每天清理次数





异常管理

异常上报

复盘任务

一般任务

系统管理

管控期

业务管理

空间管理

异常上报

新建异常上报

请选择业务

发生时间 2024-05-22 → 2024-06-21

创建时间 2024-05-22 → 2024-06-21

参与的业务

ID	主题名称	创建时间	异常级别	待办任务	业务名称	流程状态	责任部门	创建人	操作
<input type="checkbox"/>	12457 【中国区-汽车交付中心网...	2024-06-21 16:17:51	质量异常 - 轻微异常	-	中心防火墙	进行中	-		🔍 📄 🗑️ ...
<input type="checkbox"/>	12455 【平台研发-容器-Matrix...	2024-06-21 14:12:25	待定	-		进行中	-		🔍 📄 🗑️ ...
<input type="checkbox"/>	12426 【IT基础服务-基础服务-网盘...	2024-06-20 16:26:46	质量异常 - 轻微异常	-		进行中	-		🔍 📄 🗑️ ...
<input type="checkbox"/>	12456 【IoT业务运维-物联网-米...	2024-06-21 14:43:15	待定	-		进行中	-		🔍 📄 🗑️ ...
<input type="checkbox"/>	12342 【电商业务运维-无-keyce...	2024-06-17 18:23:03	质量异常 - 一般异常	-		完成	集团技术委员会/基础技术平台部/信息安全部/安全研发组		🔍 📄 🗑️ ...
<input type="checkbox"/>	12335 【网络-网络设备-IDC-核心...	2024-06-17 09:58:45	质量异常 - 轻微异常	-		完成	小米集团外		🔍 📄 🗑️ ...
<input type="checkbox"/>	12332 【平台研发-公有云-资源管理] ...	2024-06-15 15:50:42	质量异常 - 一般异常	-		进行中	-		🔍 📄 🗑️ ...
<input checked="" type="checkbox"/>	12334 【网络-网络设备-IDC-核心...	2024-06-17 09:51:03	质量异常 - 一般异常	-		完成	小米集团外		🔍 📄 🗑️ ...
<input type="checkbox"/>	12327 【国际运营中心-国际业务SRE...	2024-06-14 18:54:49	质量异常 - 一般异常	-	联网业务	进行中	-		🔍 📄 🗑️ ...
<input type="checkbox"/>	12322 【IT基础服务-基础服务-负载...	2024-06-14 17:12:25	待定	-		进行中	-		🔍 📄 🗑️ ...
<input type="checkbox"/>	12315 【Miks-Miks研发-Mi...	2024-06-14 13:04:21	质量异常 - 轻微异常	-		进行中	-		🔍 📄 🗑️ ...
<input type="checkbox"/>	12275 【基础平台部-基础组件运维-K...	2024-06-13 16:03:27	质量异常 - 轻微异常	-	er	进行中	-		🔍 📄 🗑️ ...
<input type="checkbox"/>	12326 【国际运营中心-国际业务SRE...	2024-06-14 18:44:38	质量异常 - 一般异常	-	联网业务	进行中	-		🔍 📄 🗑️ ...
<input type="checkbox"/>	12313 【国际运营中心-国际业务SRE...	2024-06-14 10:59:13	质量异常 - 轻微异常	-	家&可穿戴	完成	AWS		🔍 📄 🗑️ ...
<input type="checkbox"/>	12317 因AVAYA G450宕机导致...	2024-06-14 15:24:12	质量异常 - 轻微异常	-	A	进行中	集团技术委员会/基础技术平台部/云平台部/基础架构部 北京华胜天成科技股份有限公司		🔍 📄 🗑️ ...

统一故障上报入口，建立大部及公司级故障应急响应中心，提高故障处理效率。



故障 小米服务故障通报群
659 | 1

朱建平-云平台质量管理

@所有人 【Gitlab升级通告】
04/09 23时00分 Gitlab将进行停服升级，本次升级开启gitlab FDS集成功能，预计操作时间5min，操作间与代码仓库有关服务(例如部署系统、上线审批等)会出现拉取不到分支情况。如遇意外情况即刻回滚有问题请及时联系
联系人员：
郭斐 17362980125
陆文龙 18600668903

云平台质量值班-雷振 4月9日 19:19

【紧急通知】@所有人
接莫斯科专线01 & 02 提供商通知，境外ISP计划北京时间04/10 04:00 ~12:00对莫斯科专线进行维护，官方通告中断时长8小时，割接期间北京到莫斯科专线01 & 02存在同时中断的可能性，导致莫斯科机房过专线访问其他IDC中断，以上请大家知晓，并做好相应预案；
线下尝试反复与莫斯科专线01 & 02 提供商沟通协调割接改期，或避开同时割接，目前协调沟通未果，们仍在进一步升级尝试沟通协调割接改期，或避开同时割接，如有新进展再与大家同步；

2条回复



故障 小米服务故障通报群
659 | 1

云平台质量值班-童盼盼

【通知】监控07:35 ~当前 香港 到 新加坡专线02质量下降，正常情况下不影响业务，目前已向ISP报障，ISP仍在排查中，后续有进展再向大家同步；

Haiyan1 Shao 邵海彦 通过 杨文强 分享的名片进入此群，新成员入群可查看所有历史消息

郝丰澧 专注、细节、在关键点处深入到底；多想
飞书 在线文档有问题，工程师正在处理

2条回复

朱建平-云平台质量管理 12:14

| 回复 郝丰澧：飞书 在线文档有问题，工程师正在处理
@郝丰澧 @所有人

朱建平 更新了群头像

郝丰澧 专注、细节、在关键点处深入到底；多想
| 回复 郝丰澧：飞书 在线文档有问题，工程师正在处理
已经全部恢复

从发现、响应、复盘、改善、预防进行全方位质量异常闭环管理，确保服务稳定性持续提升

故障管理目标

降低故障数量

缩短故障时长

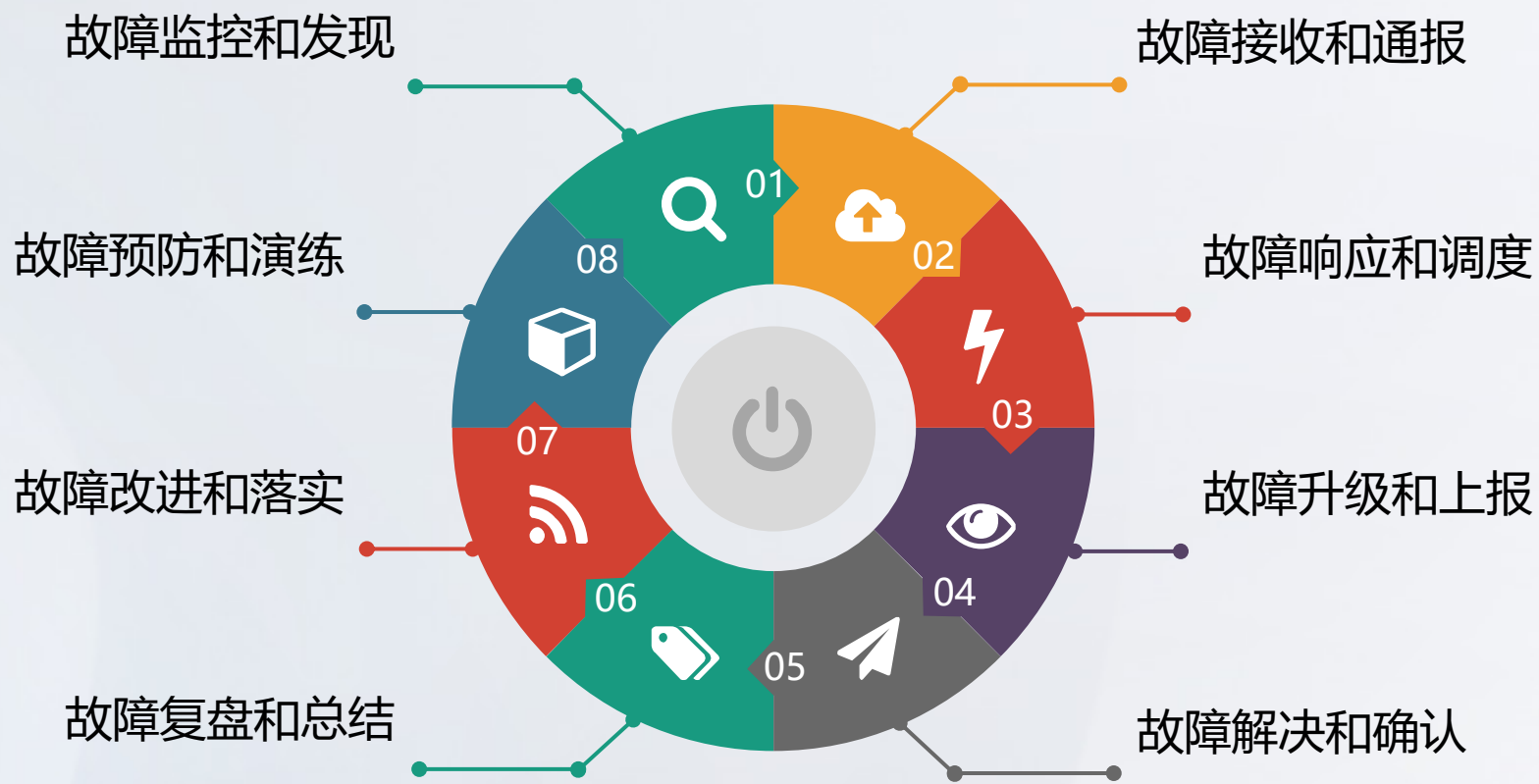
降低故障对产品影响

预防同类故障再次发生

确保各平台和产品服务质

保障良好的用户体验

维护公司品牌形象



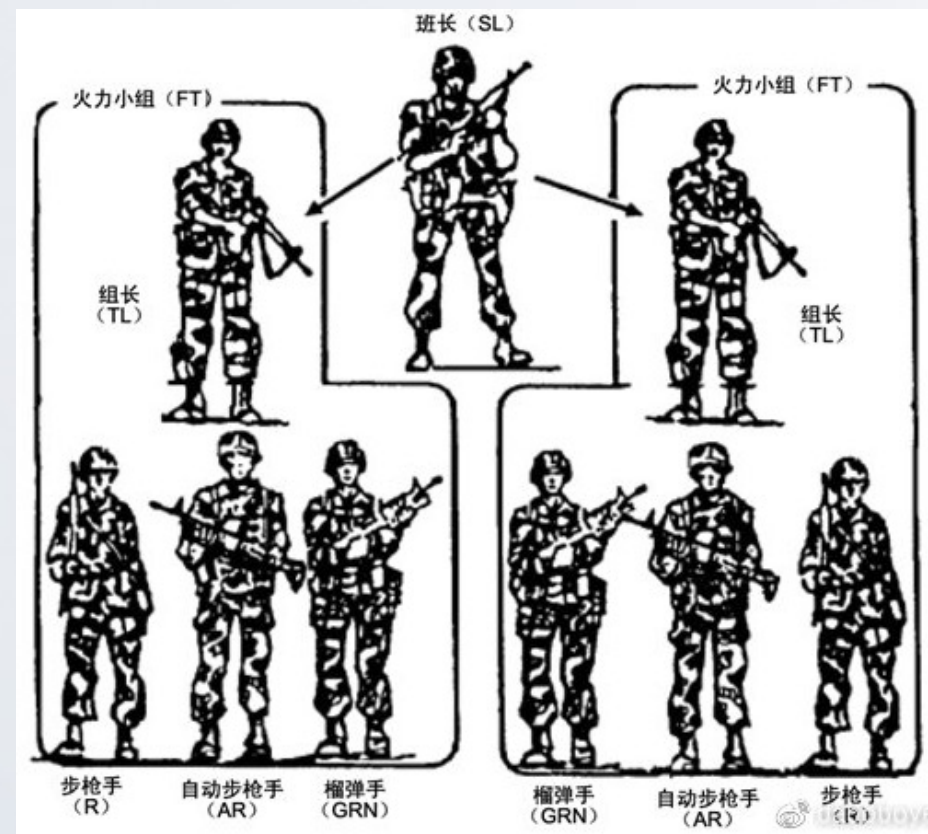
04

总结

故障应急 考较的是SRE体系构建的综合能力，包括丰富的信息采集，快速的反应和决策，稳定准确的处置措施，超强的资源整合，日常以练养战的流程体系以及经验丰富能力强大的专家团队



2004年在伊拉克的一个美国海军陆战队步兵班，里面有3挺M249。照片中有12人，第13人是照相机的

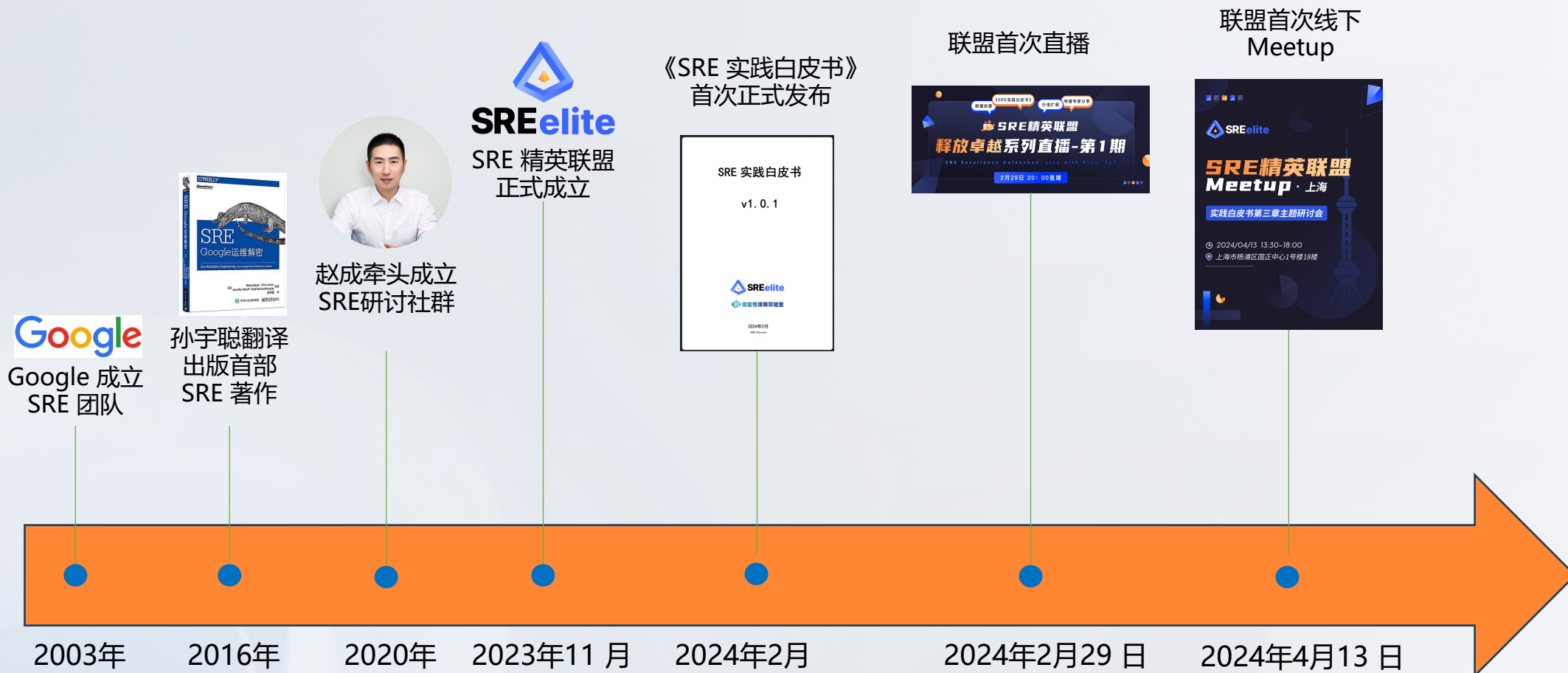


Q&A



<https://sre-elite.com>

“SRE精英联盟”概述



SRE 实践白皮书

v1.0.1



2024年2月
SRE-Elite.com



经历数年，20 多位一线专家协作编写。



扫码下载 v1.0.1。版本持续更新迭代中。



在官网 <https://sre-elite.com/notice/> 下载最新版。



公众号



视频号



B 站



YouTube