

# 超大规模弹性计算节点 自动化运维稳定性实践

唐磊（小唐）

阿里云 技术专家

# InfoQ 企业会员

## 企业数字化传播一站式服务

InfoQ 企业会员是为满足企业在中国开发者群体中的品牌曝光需求而推出的一款矩阵化资源包。可为企业提供包括“企业号服务”、“企业动态宣发”、“品牌展示通道”在内的多项专属权益与服务，助力企业高效触达开发者群体，提升数字化时代影响力。



### 企业号服务

深度触达 300 万中高端开发者



### 企业动态宣发

新媒体矩阵覆盖百万粉丝



### 品牌展示通道

线上平台 10 万+ 流量曝光



# 个人简介

姓名：唐磊，花名：小唐，网名：石头

- 2008 ~ 2015：学生@中南/清华，（CG&CAD）
- 2015 ~ 2017：工程师@宜信大数据创新中心，（互金-信贷）
- 2017 ~ 2019：工程师、TL@大疆创新，（社区 SkyPixel）
- 2019 ~ 至今：工程师@阿里云神龙计算平台，（弹性计算-异常调度）

## 01 概述&背景

- 客户诉求
- 业务难点

## 02 业界方案

- 业界方案
- 发展趋势

## 03 我们的方案

- 基于专家规则的自动化运维策略
- 运维评价
- 发布熔断



# 01 概述&背景

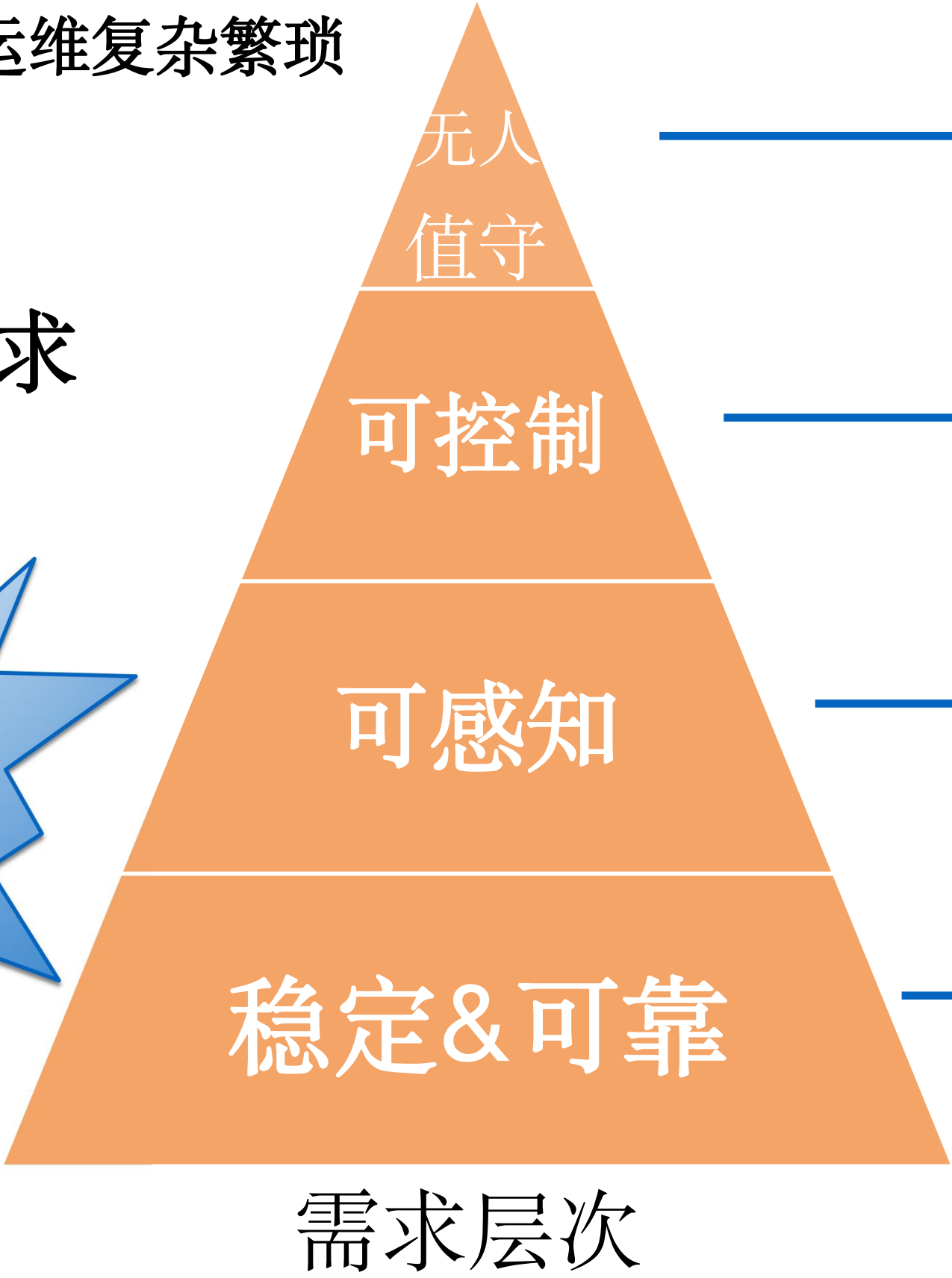
# 概述&背景 - 客户诉求

## ❑ 云下环境特点

- 维护成本高
- 资源利用率低
- IT 资源管理和运维复杂繁琐
- ...



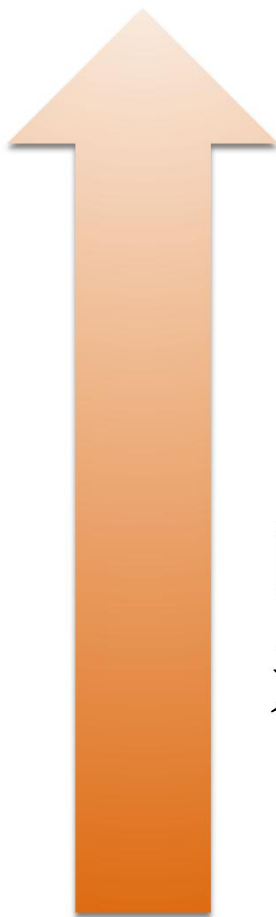
## ❑ 客户用云诉求



- 智能自愈:
- 智能编排
  - 预测&自愈
- 可控制:
- 原子操作
  - 自动化
- 异常感知:
- 监控&告警
  - 根因诊断
- 稳定&可靠:
- 稳定性 SLA
  - 性能 SLA

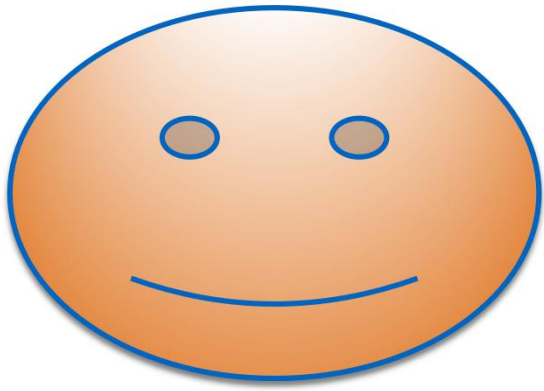


幸福感



自动化运维  
提升幸福感

已发生的不可用: 及时止损  
未发生的不可用: 提前规避



安全感

## 业务难点 – 基础设施规模大



- 云计算基础设施规模决定了其运维复杂度
- 没有现成的产品可以借鉴，需要探索出自己的道路



# 业务难点 - 产品形态多、业务领域广

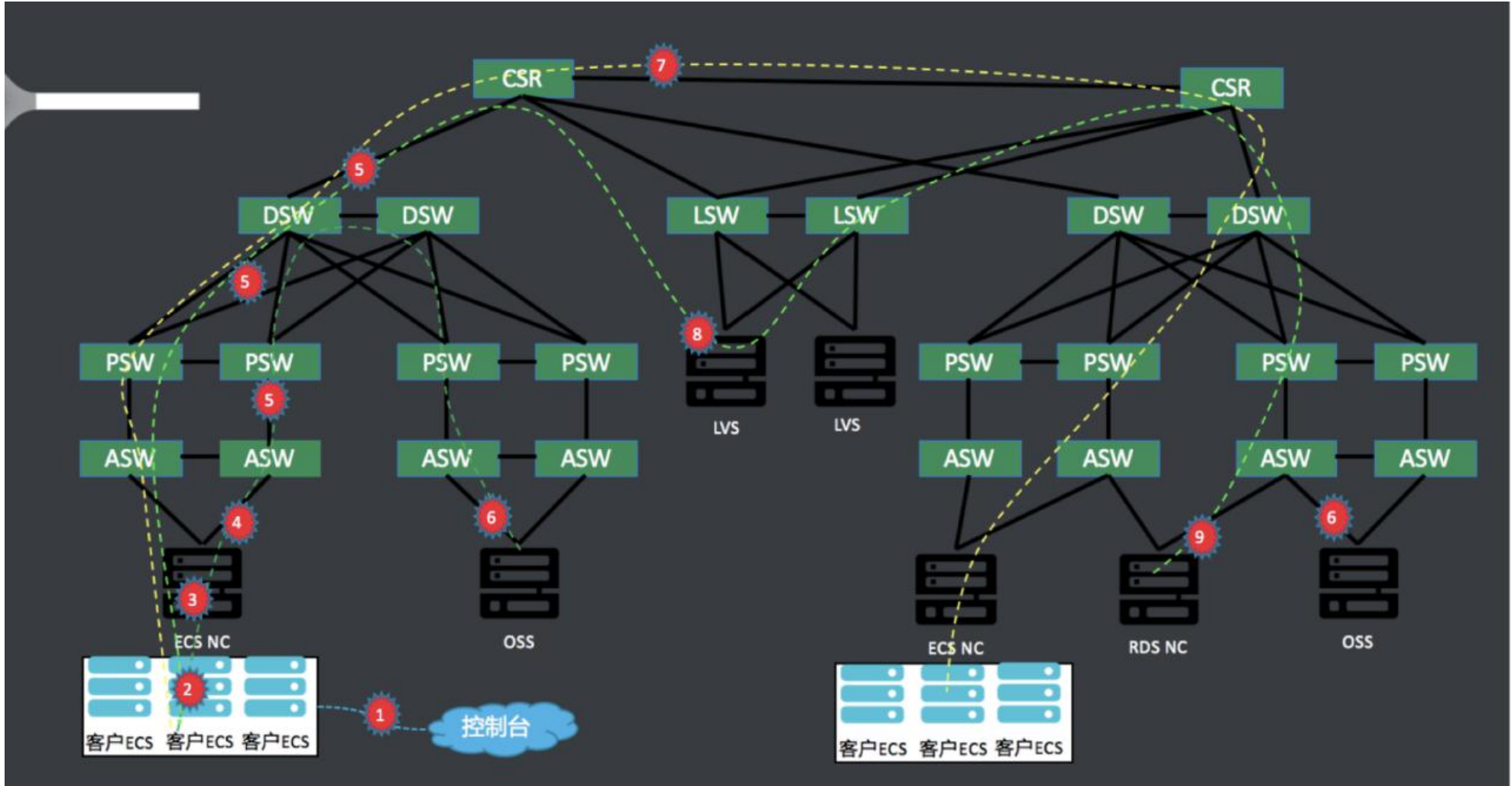
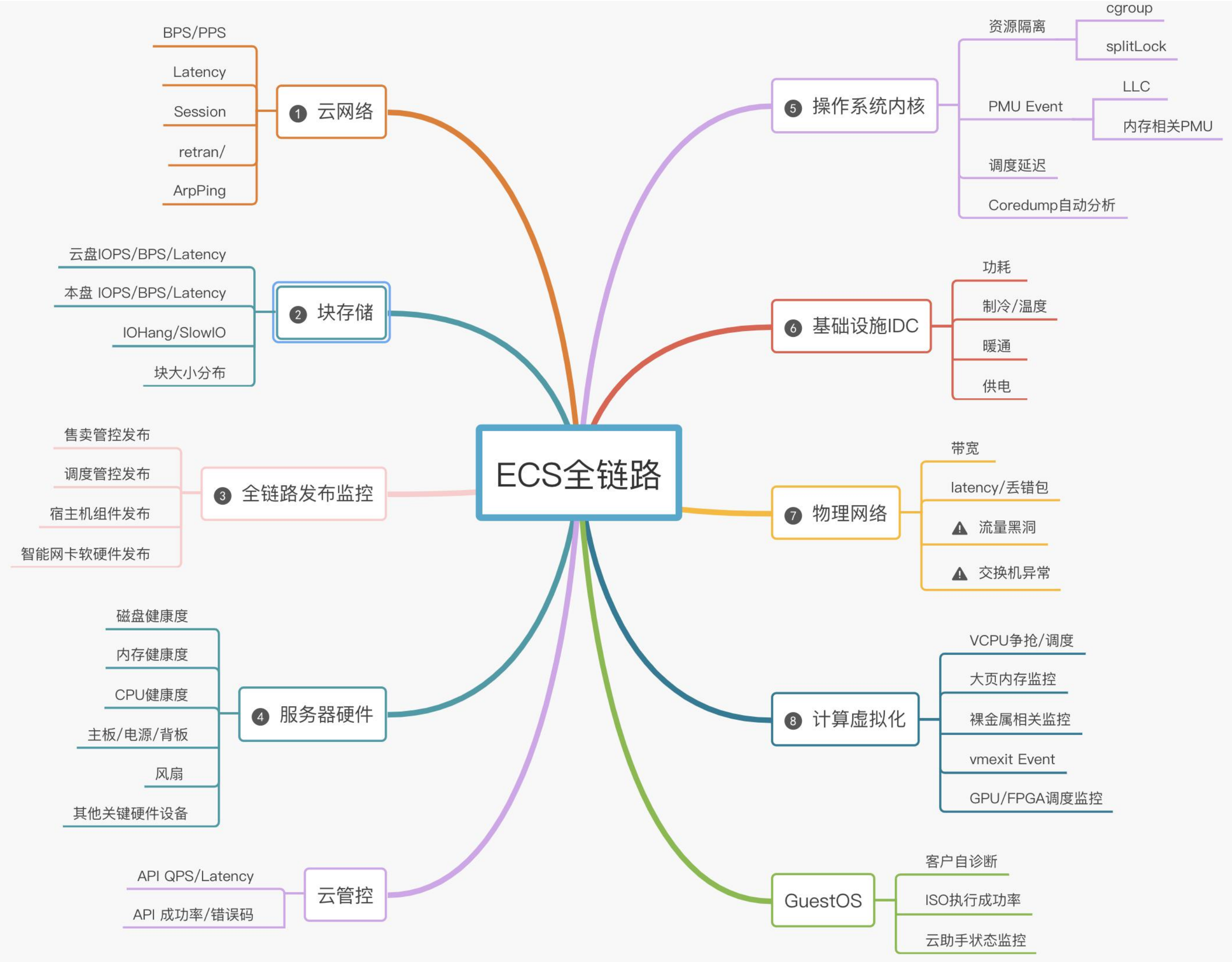




# 业务难点 - 知识面广&技术难度深

覆盖子系统众多

链路长



技术难度深

CPU子系统举例

**core**

- Intel: SKL/ICL/Atom
- AMD: MiLan/ROME
- YiTian710
- kunpeng

**cache**

- LLC 一致性
- LLC 容量 QoS
- LLC 争抢

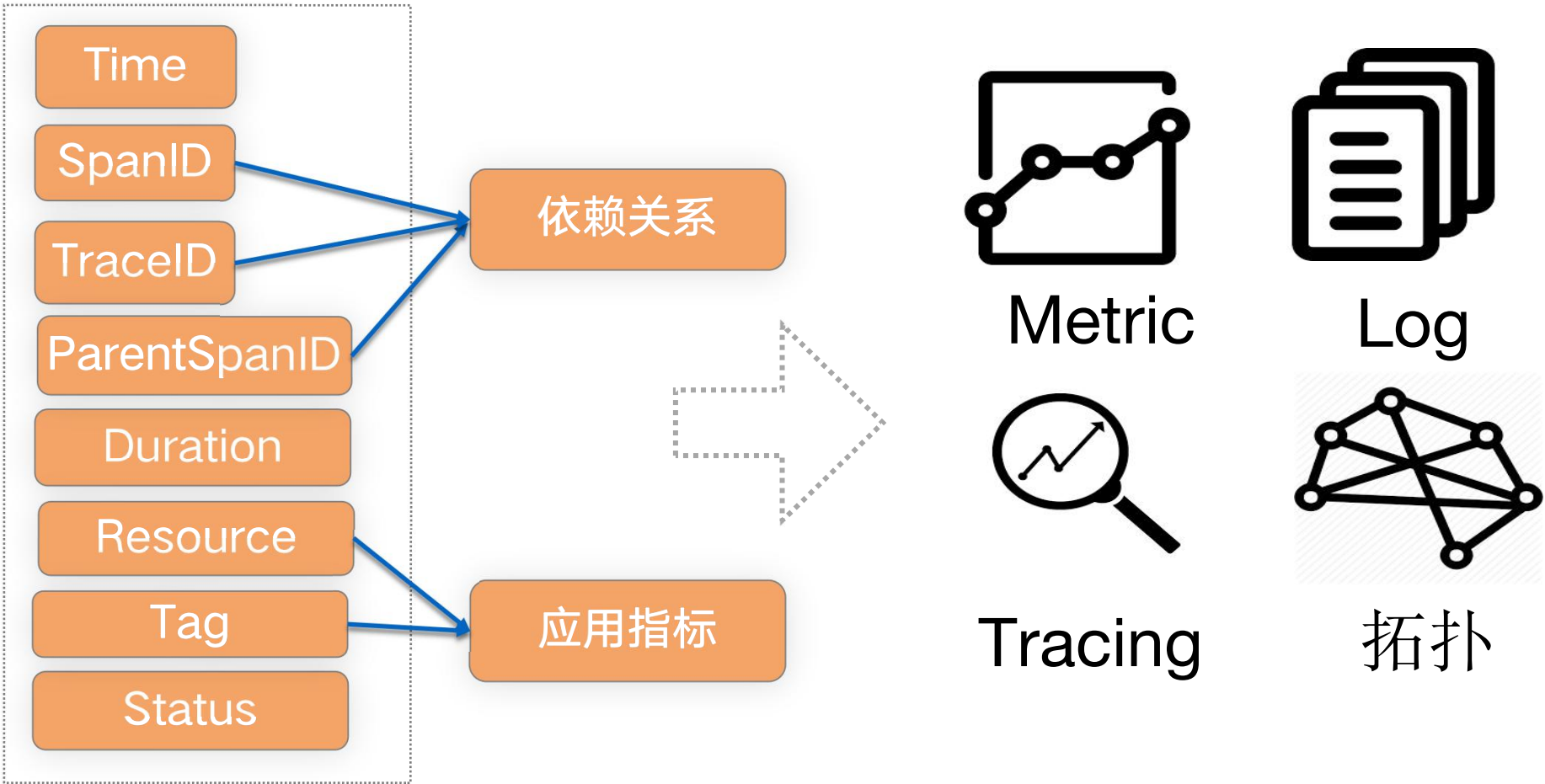
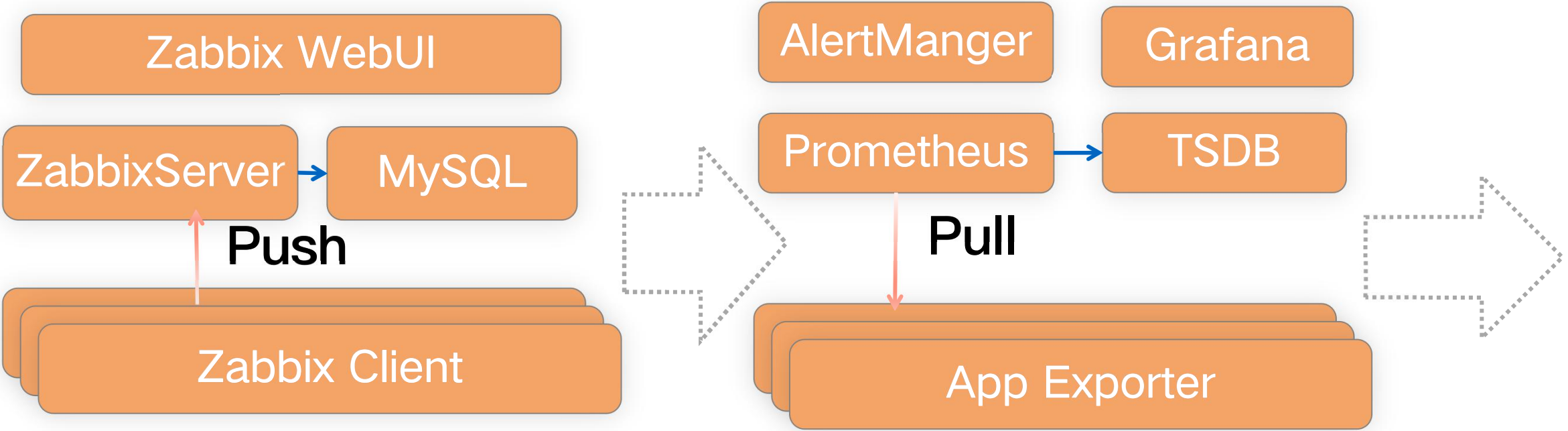
**IMC**

- IMC freq
- IMC channel
- memory buffer

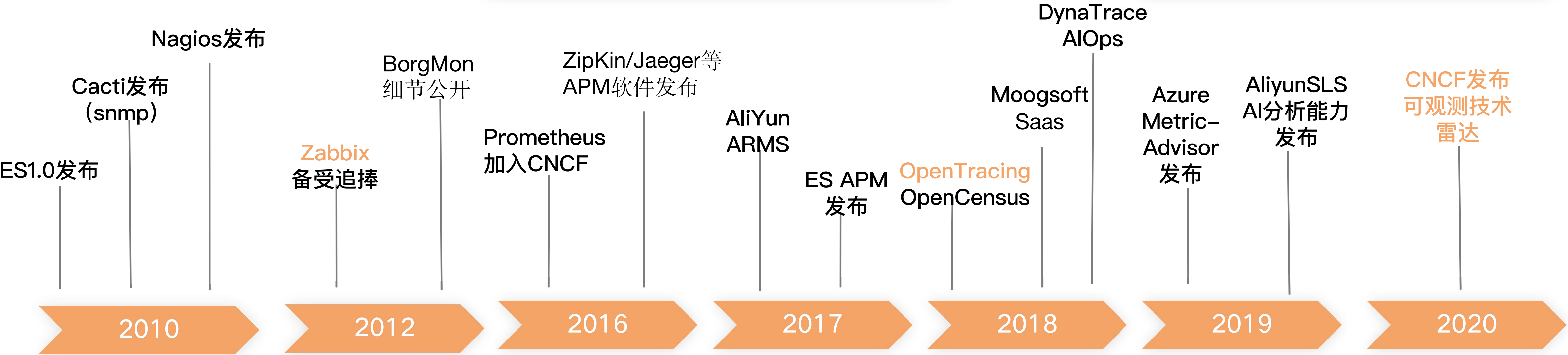
## 02 业界方案



# 业界方案及发展



- 拓扑自动分析
- OneAgent通用采集
- 动态阈值告警
- 影响面自动分析
- APITrace分析
- 故障根因分析



## 03 我们的方案

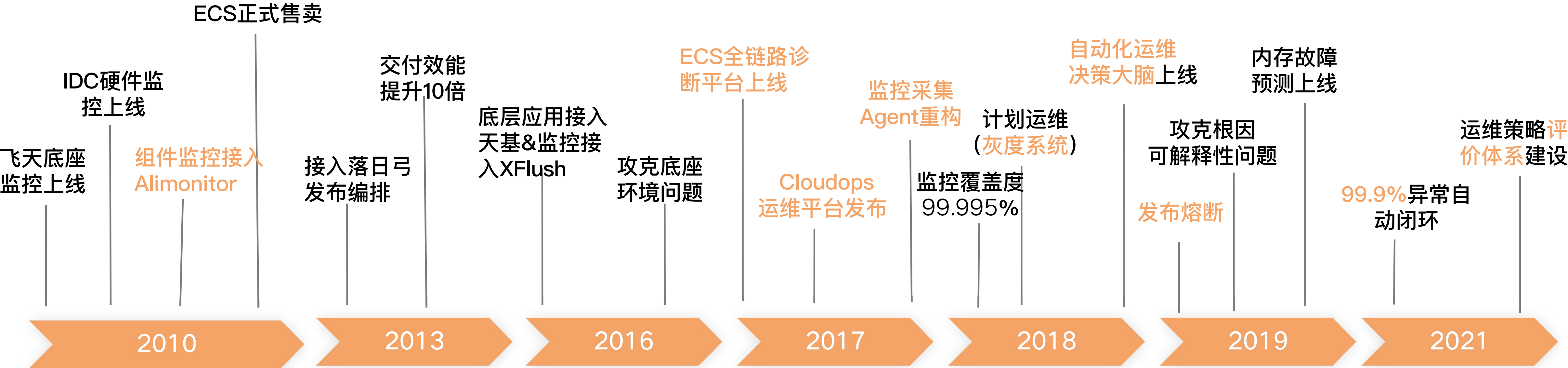
# 我们的方案 - ECS 监控运维体系 发展历程



工具+人工时代

平台+半自动编排

数据化智能化





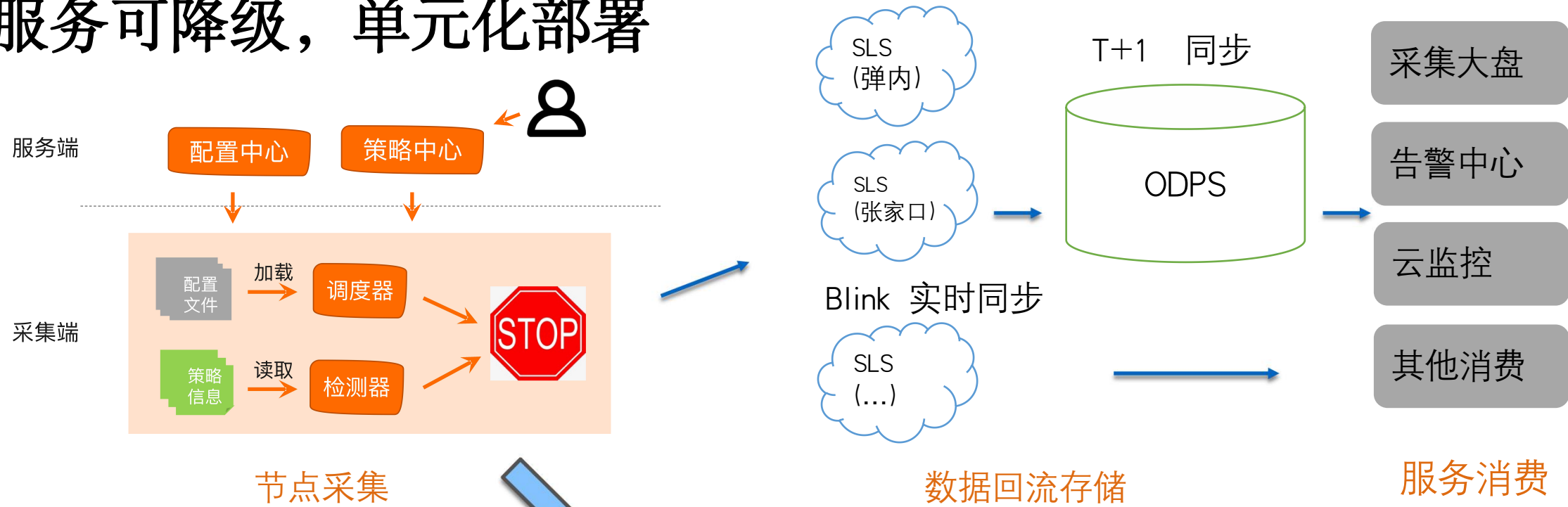


# 我们的方案 – 监控数据采集

## 规模效应

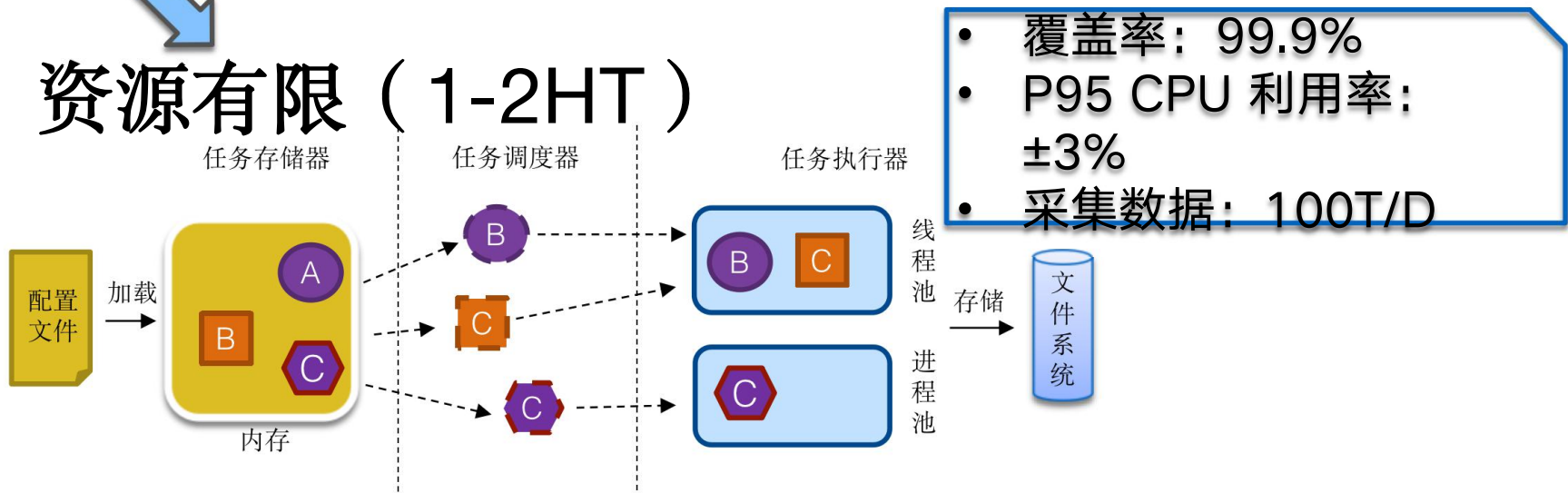
对比项	AlarmAgent	Prometheus	Nightingale
来源	自研	K8s开源	滴滴开源
触发方式	采集侧 服务侧	服务侧	采集侧
采集策略	周期采集 按需采集	周期采集	周期采集
支持规模	百万级节点	千级别	万级别
结果完整度	高 采集回流分离	中 结果实时回流	中 结果实时回流
采集自监控	完善 报表和告警	有限 节点丢失提醒	有限 节点丢失提醒
灰度发布	支持	不支持	不支持

## 服务可降级，单元化部署

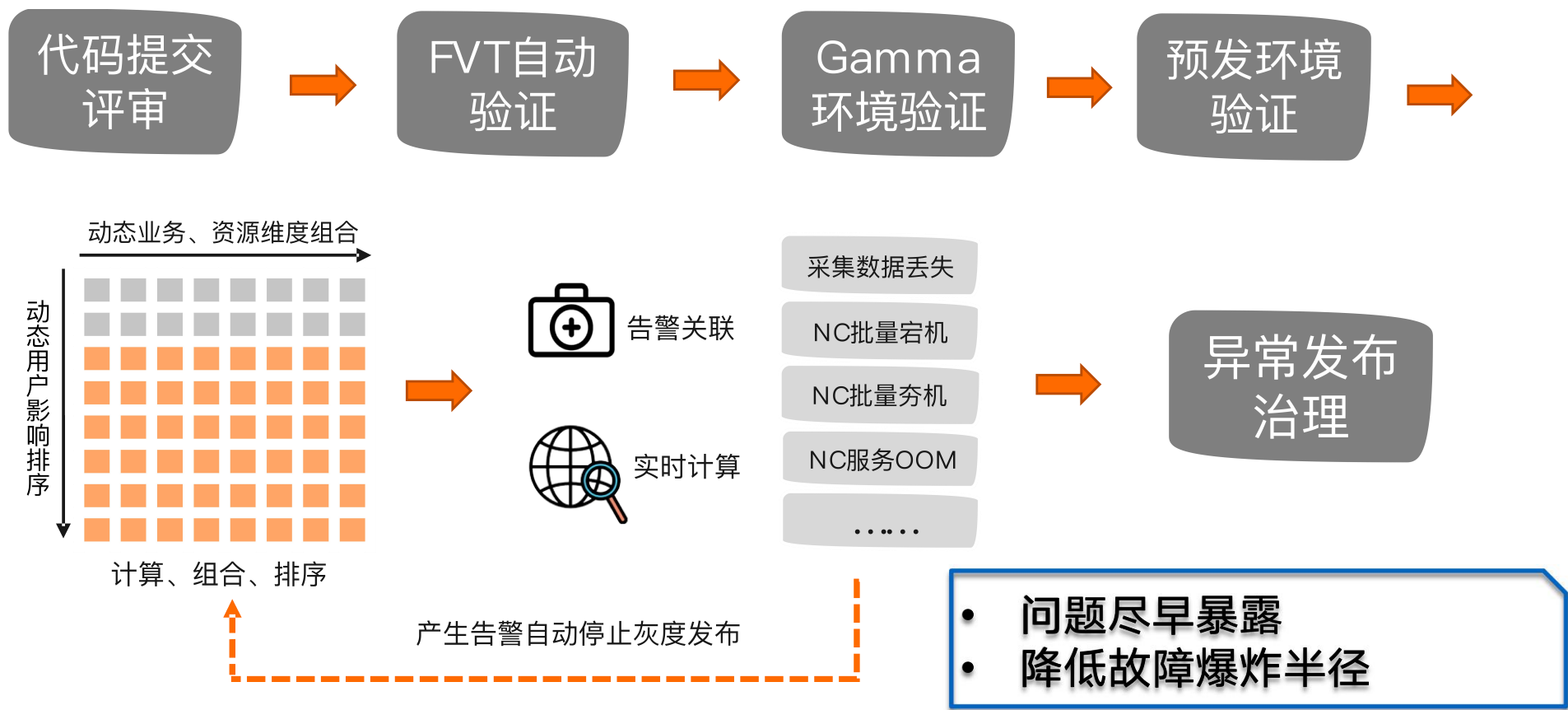


## 资源最大化利用：资源有限（1-2HT）

- 调度器轻量化
- 采集之间数据共享



## 代码发布灰度可控

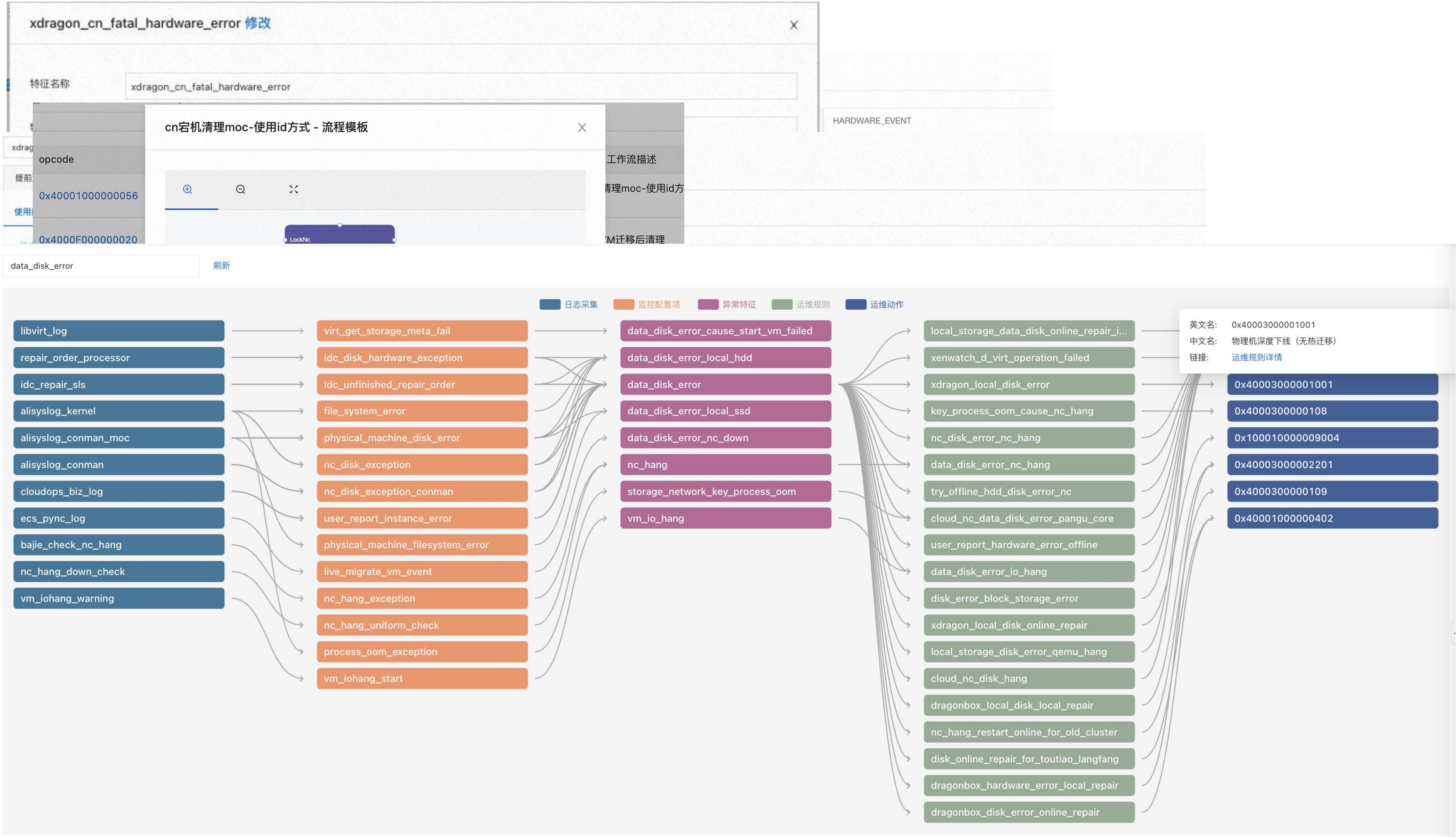




# 我们的方案 – 运维策略相关定义

## 相关定义：

- 监控异常
- 特征定义
- 运维规则
- 运维动作





# 我们的方案 – 运维评价（理论）

## 产生背景：

- 区分运维动作的“好坏优劣”？锁机器、热迁移、下线
- 是否存在过度运维的问题？
- 对客户真实体感是什么？

## 评价度量\*：

- 性能度量
- 不可用度量
- 控制面度量

$$KeyMetric_{perf} = \frac{\sum_{l=1}^l F_l \sum_{i=1}^n F_e * Vm_i Perf_T}{\sum_{i=1}^n Vm_i TotalLifeTime_T}$$

$$KeyMetric_{down} = \frac{\sum_{i=1}^n F_e * Vm_i Down_T}{\sum_{i=1}^n Vm_i TotalLifeTime_T}$$

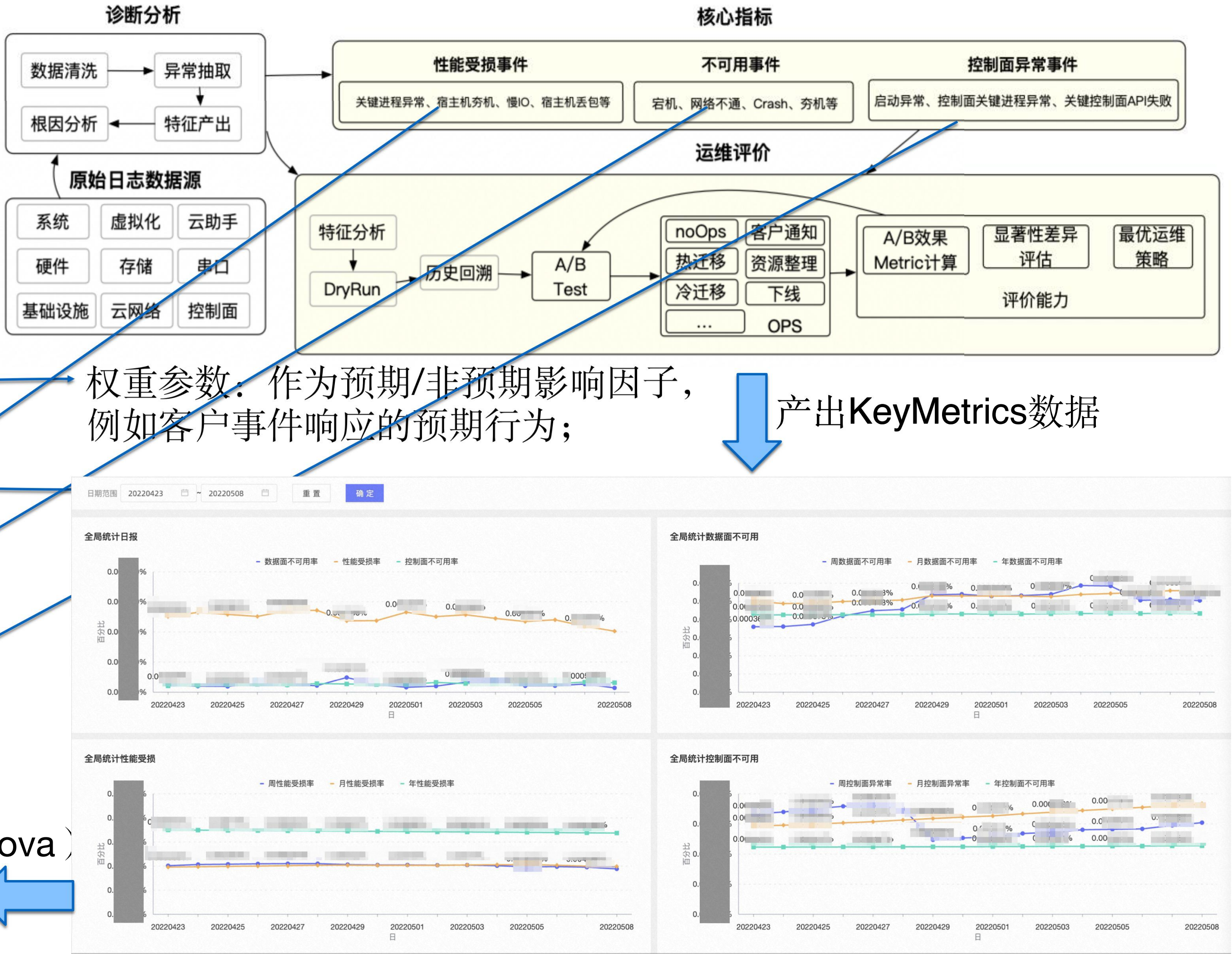
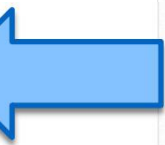
$$KeyMetric_{ctrl} = \frac{\sum_{l=1}^l F_l \sum_{i=1}^n F_e * Vm_i Ctrl_T}{\sum_{i=1}^n Vm_i TotalLifeTime_T}$$

权重参数：作为预期/非预期影响因子，例如客户事件响应的预期行为；

产出KeyMetrics数据

## 差异化分析：

- 显著性差异检验 - 单因素方差分析F检验（Welch's anova）
- 精准控制切流比例 - 功效分析（Cohen's f）



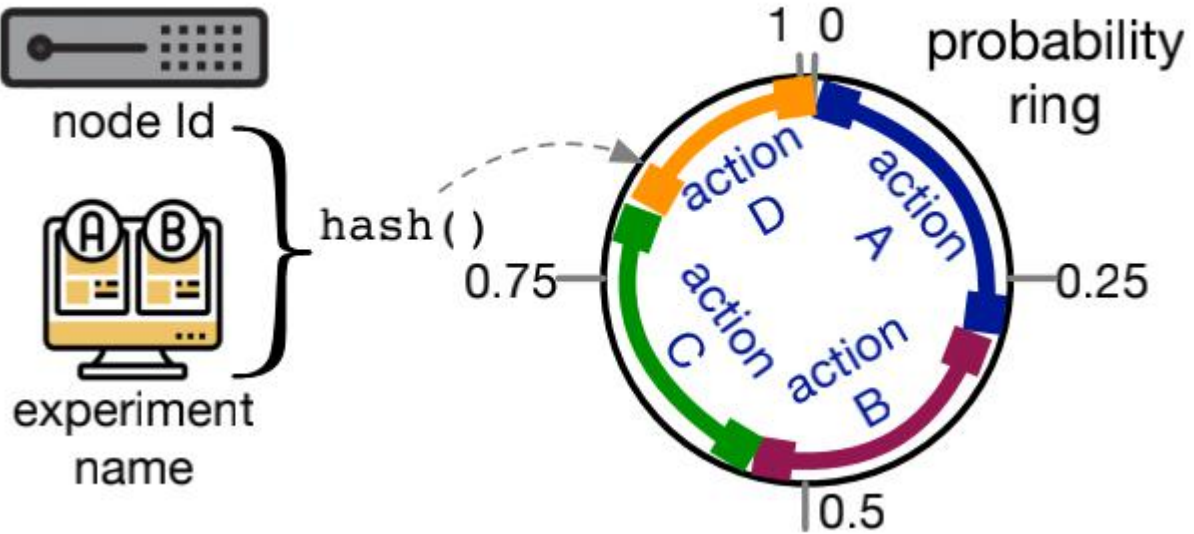
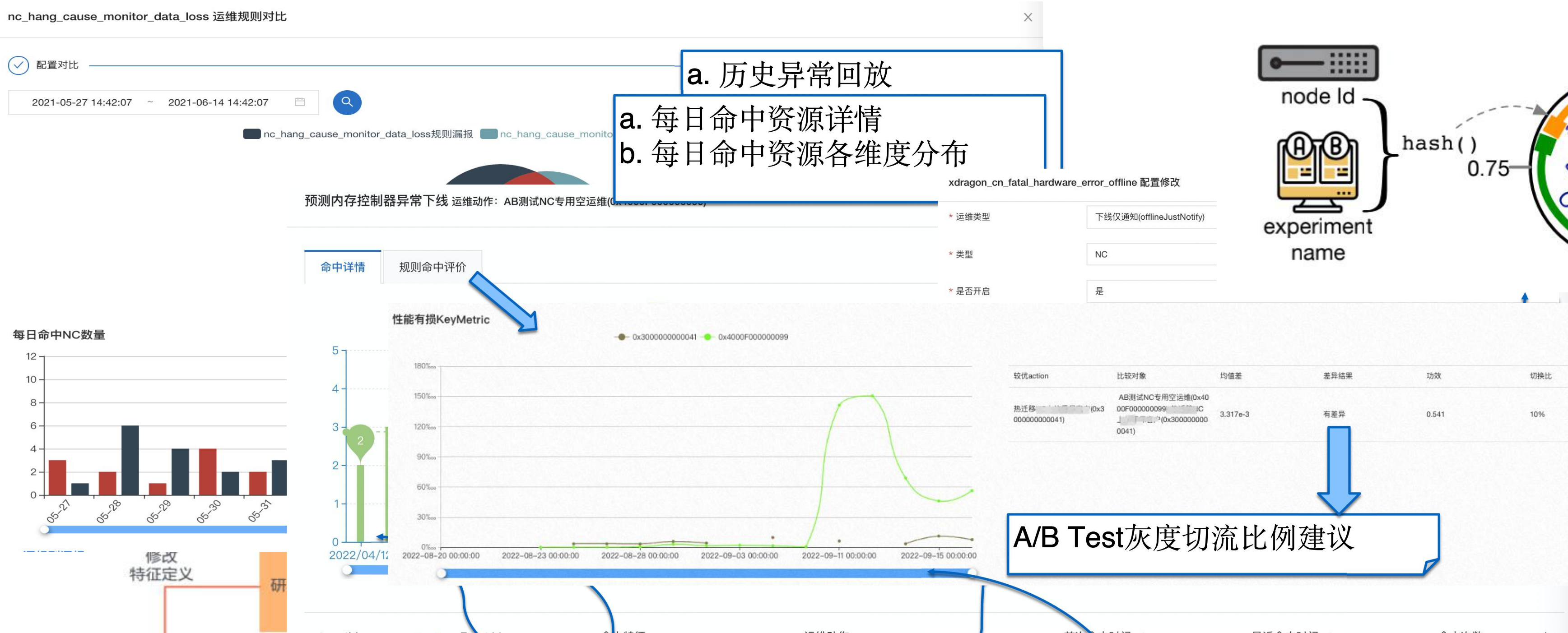
[\*] Levy S, et al. Predictive and Adaptive Failure Mitigation to Avert Production Cloud VM Interruptions.[C]// Operating Systems Design and Implementation. 2020.



# 我们的方案 – 运维评价（工程落地）

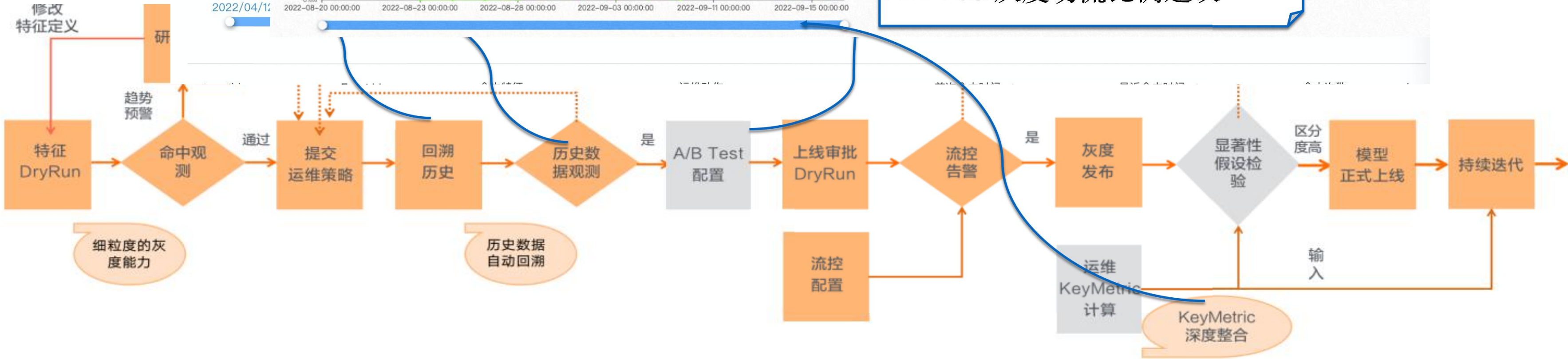
## 落地方案

- 评价模块
- 数据加速层
- KeyMetric
- 异常感知模块



## 面临的问题

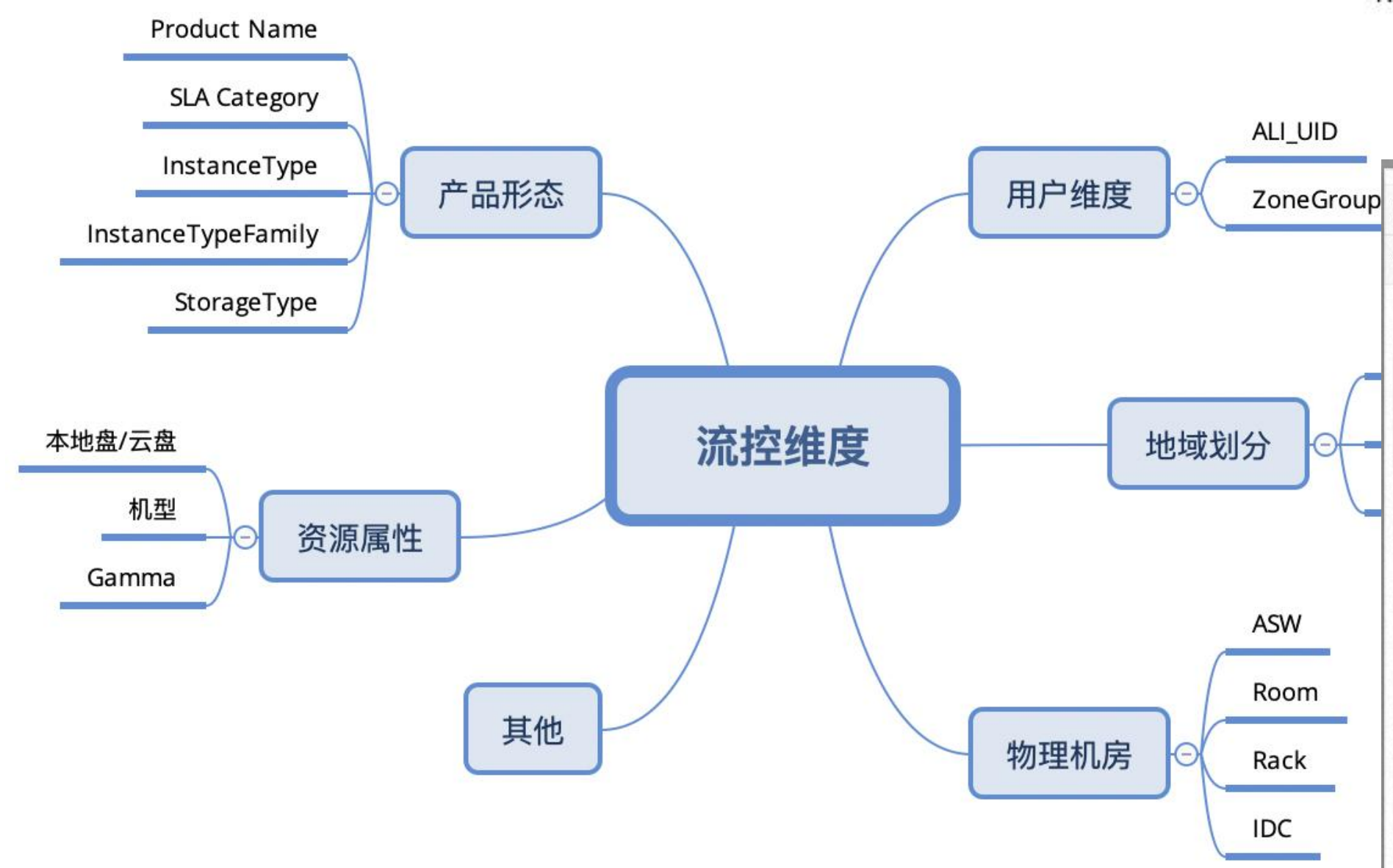
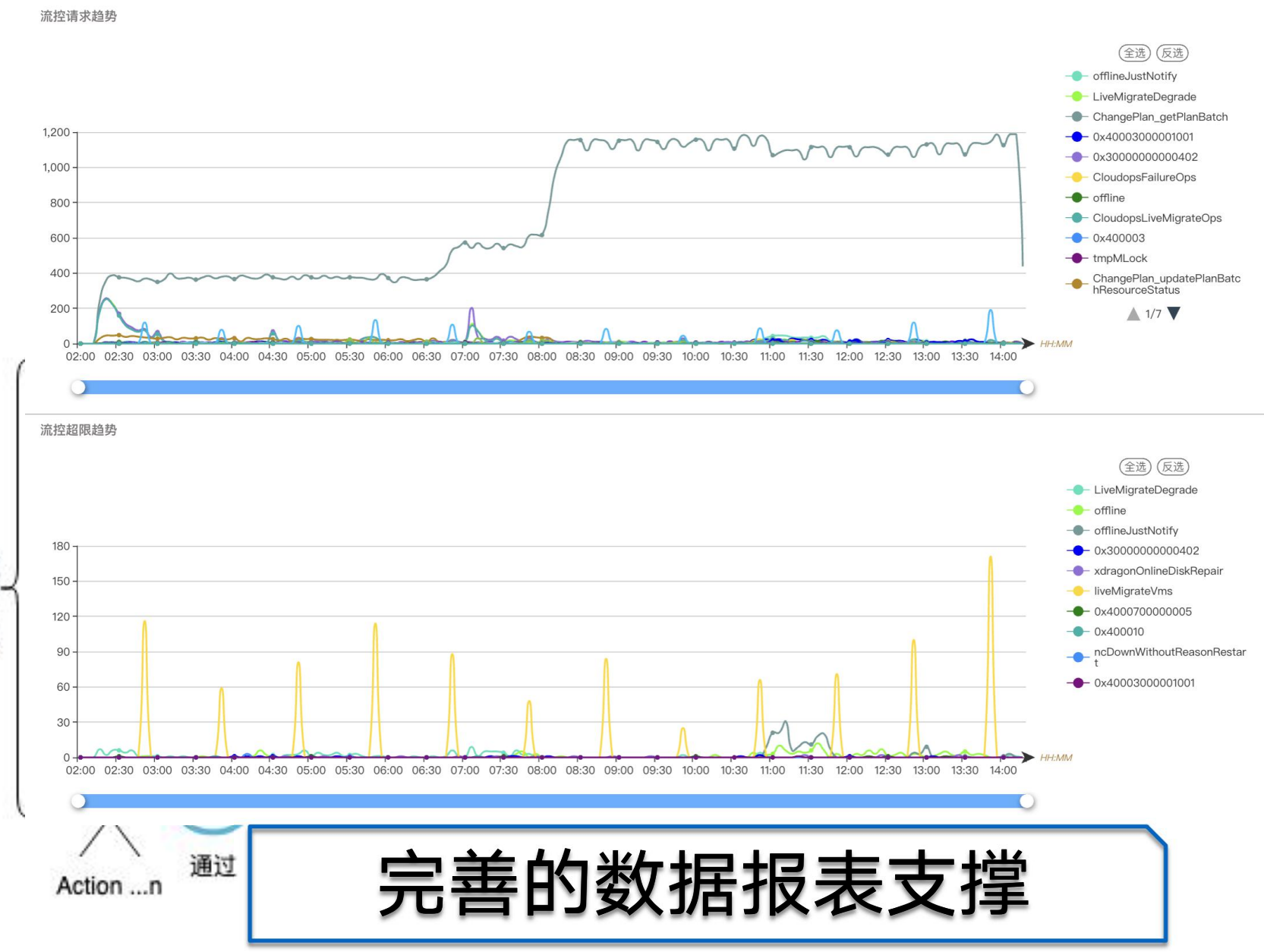
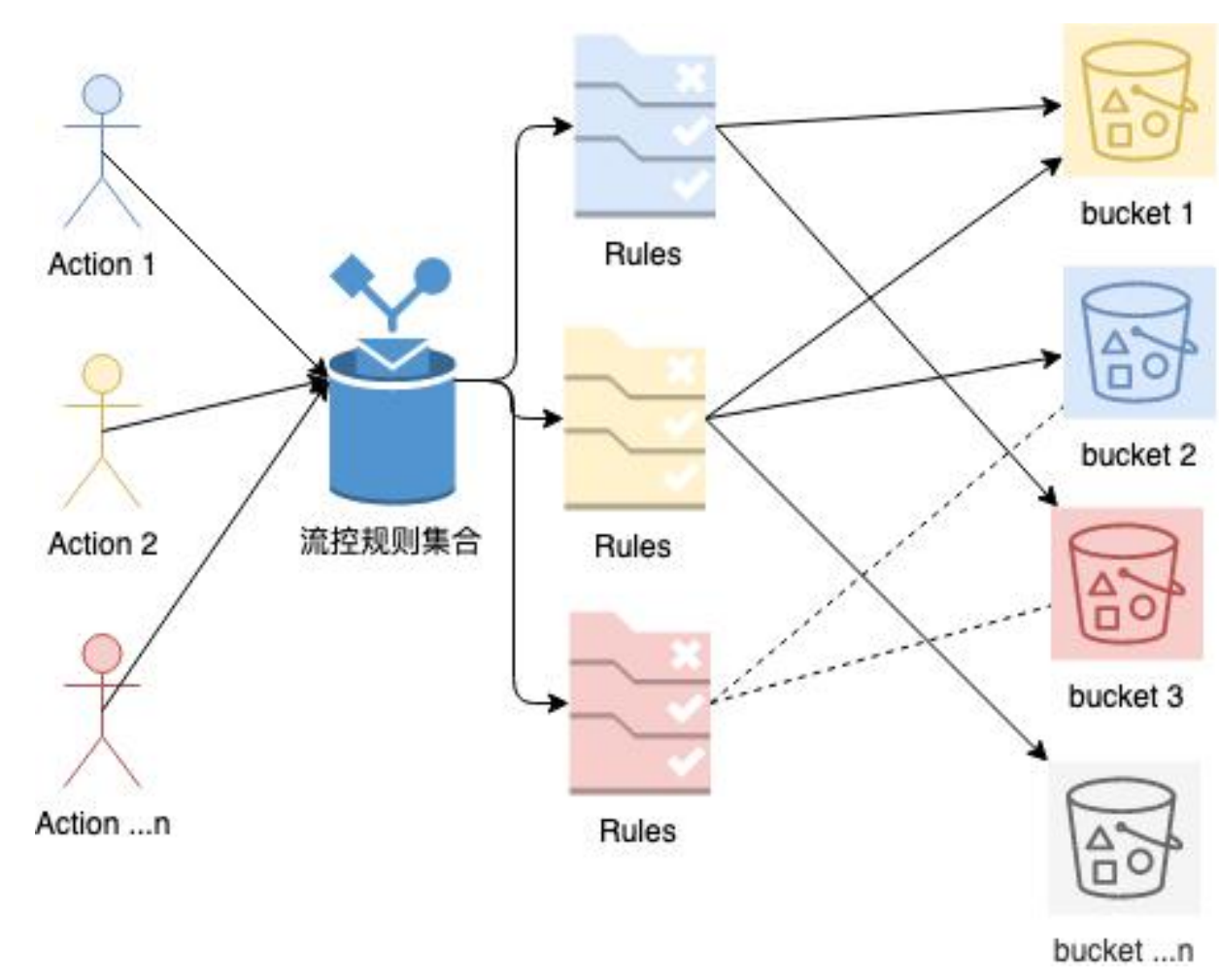
- 如何与现有的运维体系整合？
- 如何安全高效的灰度发布上线？





# 我们的方案 - 业务流控

- 业务流控：
  - 自动运维维持正常水位
  - 有效阻止故障发生



丰富的流控维度

The screenshot shows the '流控 action' (Flow Control Action) configuration interface. It includes a table with columns for 'name', '维度' (Dimension), '时间窗口 (秒)' (Time Window (s)), '阈值类型' (Threshold Type), '流控机器数/比例' (Flow Control Machine Count/Proportion), '最小阈值' (Minimum Threshold), '后处理' (Post-processing), '是否启用' (Whether to Enable), '是否去重' (Whether to Deduplicate), '触发条件' (Trigger Condition), '排除条件' (Exclusion Condition), '描述' (Description), and '操作' (Operation). The table lists several rules for different dimensions like '全网' (All Network), 'aswid', 'azone', 'azone-regionAlias', 'clusterAlias', and 'clusterAlias-productName'. A blue box at the bottom right contains the text '灵活的流控规则' (Flexible flow control rules).



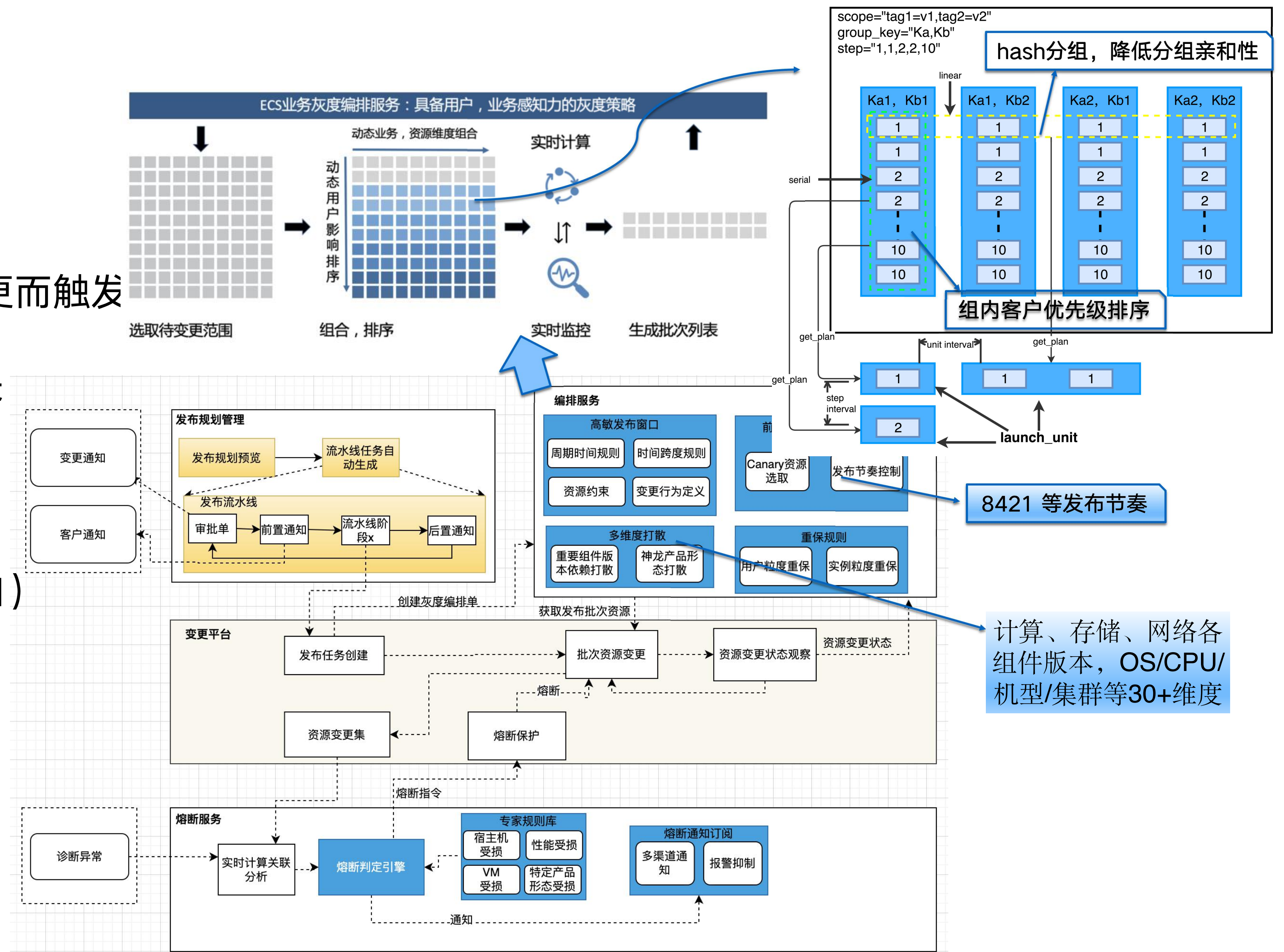
## 我们的方案 - 灰度发布

## □ 背景:

- 业界大概 XX % 的生产事故由变更而触发
- 集团全部故障中 xx %+和变更相关

## □ 业务：

- 支撑百万级资源的发布（千万级 VM）
- 发布业务方数百
- 累计变更次数 n 亿
- 发布次数 n 万





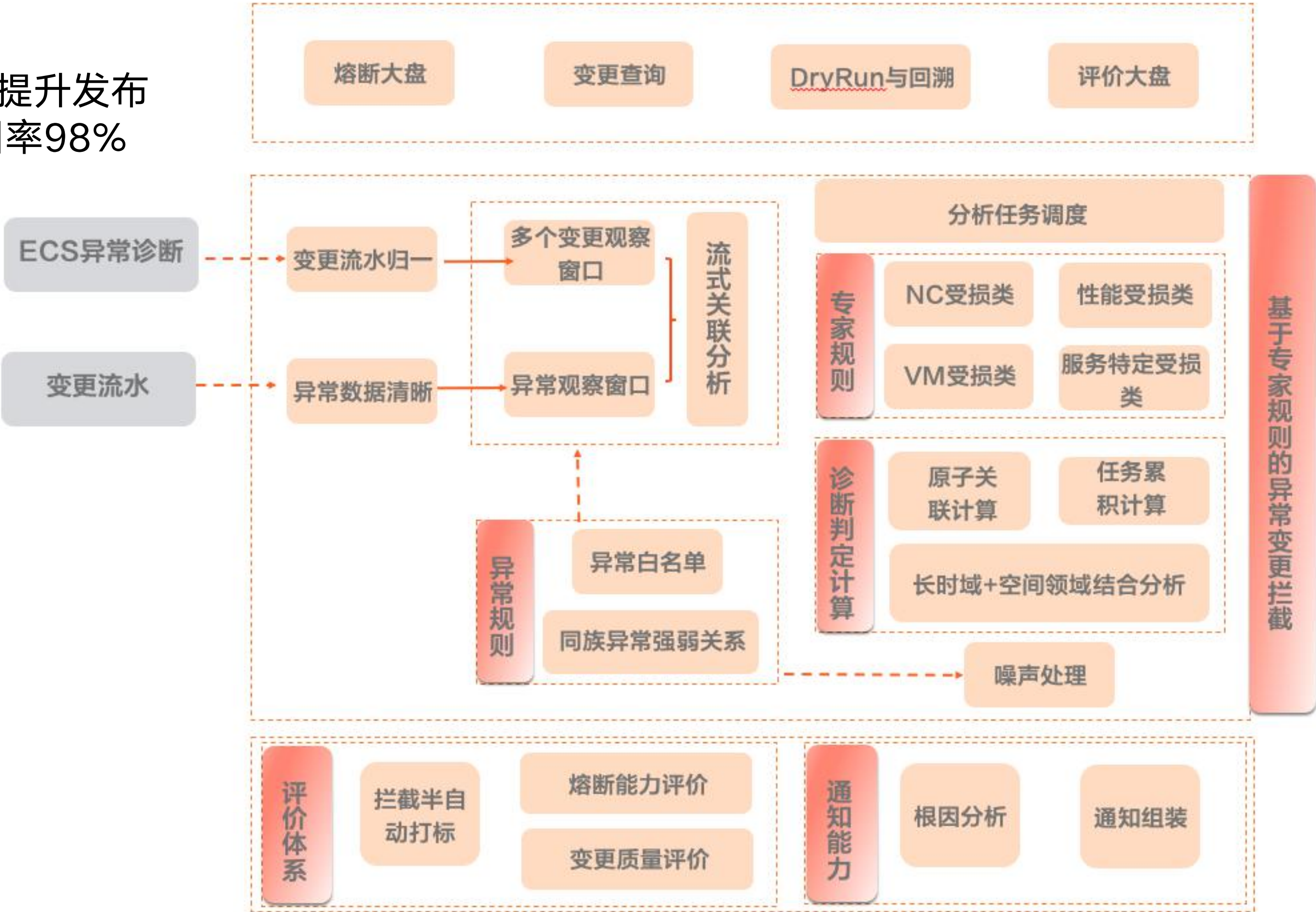
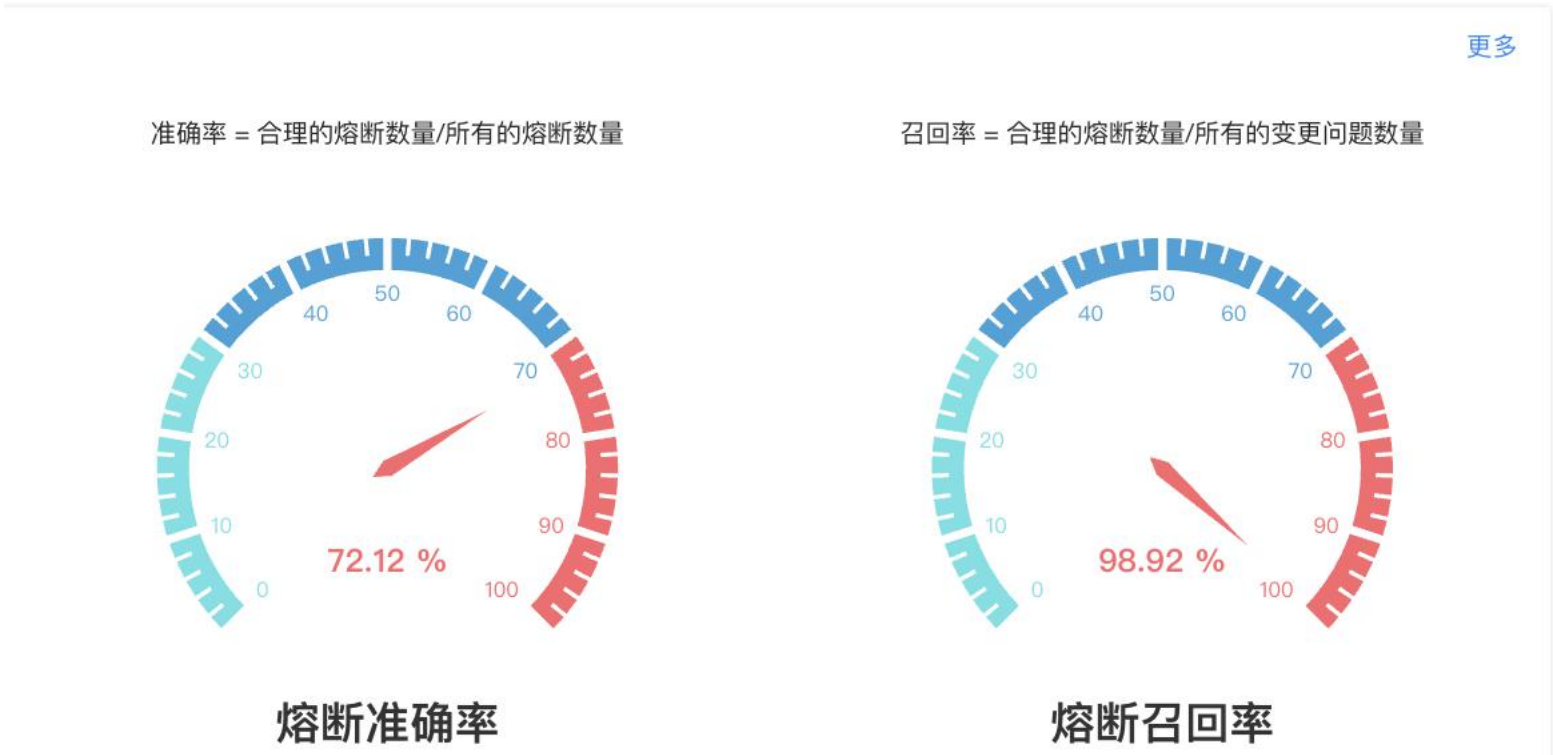
# 我们的方案 - 灰度发布&熔断

## 熔断能力

发布过程，通过诊断识别异常或隐患，主动阻断发布，是提升发布稳定性的利器，XX 期间，发布有效熔断近百次，有效召回率98%

召回率 =  $\frac{\text{熔断次数}}{\text{熔断次数} + \text{漏召回次数}} * 100\%$

准确率 =  $\frac{\text{有效熔断}}{\text{熔断次数} + \text{漏召回次数}} * 100\%$



[\*] Li Z , et al. Gandalf: An Intelligent, End-To-End Analytics Service for Safe Deployment in Large-Scale Cloud Infrastructure[C]// Networked Systems Design and Implementation. 2020.



# 总结

## □ 概述

- 客户上云需求&稳定性述求
- 解决问题的难点

## □ 自动化运维解决方案

- 监控采集：解耦、可降级、单元化、灰度
- 运维策略：全链路 Trace
- 运维评价：过度运维
- 稳定性利器：DryRun、A/B Test，流控、灰度、熔断

## □ 业界方案

- 监控告警体系演进

## □ 推荐更多内容

- Predictive and Adaptive Failure Mitigation to Avert Production Cloud VM Interruptions
- Gandalf: An Intelligent, End-To-End Analytics Service for Safe Deployment in Large-Scale Cloud Infrastructure
- Localizing Failure Root Causes in a Microservice through Causality Inference
- Predicting Node Failures in an Ultra-large-scale Cloud Computing Platform- an AIOps Solution



欢迎交流



# 精彩继续！ 更多一线大厂前沿技术案例

📍 北京站

## GITC

全球大前端技术大会

时间：10月30-31日

地点：北京·国际会议中心

扫码查看大会  
详情>>



📍 北京站

## QCon

全球软件开发大会

时间：10月30-11月1日

地点：北京·国际会议中心

扫码查看大会  
详情>>



📍 上海站

## QCon

全球软件开发大会

时间：11月25-26日

地点：上海·宏安瑞士大酒店

扫码查看大会  
详情>>





想一想，我该如何把这些  
技术应用在工作实践中？

---

THANKS