

容器环境下智能运维技术研发 与实践

才振功 博士

浙江大学SEL实验室

QCon

全球软件开发大会

10月17-19日 上海·宝华万豪酒店



扫码锁定席位

九折即将结束

团购还享更多优惠，折扣有效期至9月17日

扫描右方二维码即可查看大会信息及购票



如果在使用过程中遇到任何问题，可联系大会主办方，欢迎咨询！

微信：qcon-0410

电话：010-84782011

ArchSummit

全球架构师峰会 2017



扫码锁定席位

12月8-9日 北京·国际会议中心

七折即将截止立省2040元

使用限时优惠码AS200，

以目前最优惠价格报名ArchSummit

仅限前20名用户，优惠码有效期至9月19日，

扫描右方二维码即可使用



如果在使用过程中遇到任何问题，可联系大会主办方，欢迎咨询！

微信：aschina666

电话：15201647919

极客搜索

全站干货，一键触达，只为技术

s.geekbang.org



扫描二维码立即体验

有没有一种搜索方式，能整合 InfoQ 中文站、极客邦科技旗下12大微信公众号矩阵的全部资源？

极客搜索，这款针对极客邦科技全站内容资源的轻量级搜索引擎，做到了！

扫描上方二维码，极客搜索！

这里只有 技术领导者

EGO会员第二季招募季正式开启



E小欧

报名时间：9月1日-9月15日

扫描添加E小欧，
邀您进入EGO会员预报名群

立即报名



TABLE OF CONTENTS

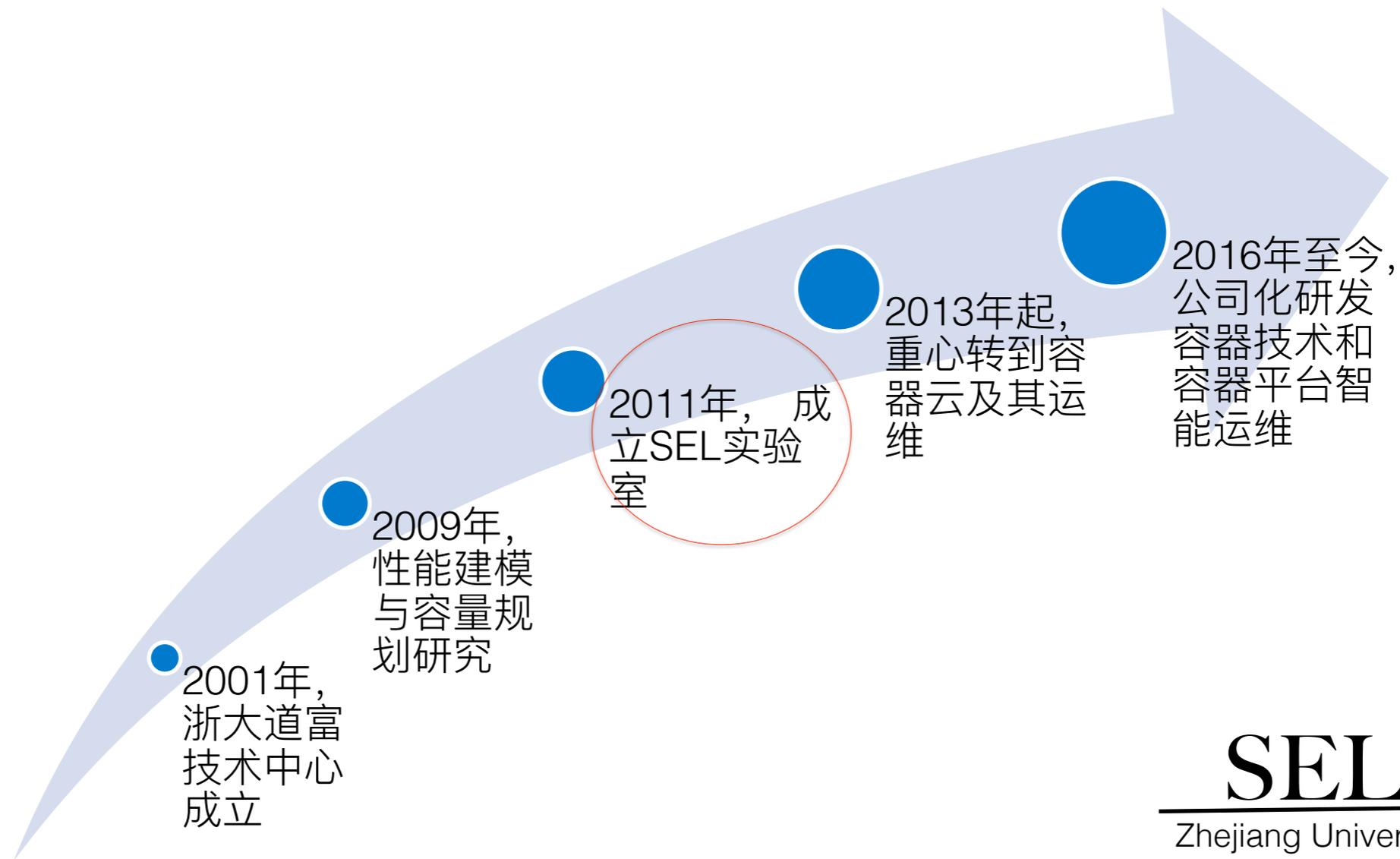
浙江大学SEL实验室

容器环境下智能运维的挑战

容器环境下智能运维研发实践

智能运维发展展望

浙江大学SEL实验室



SEL

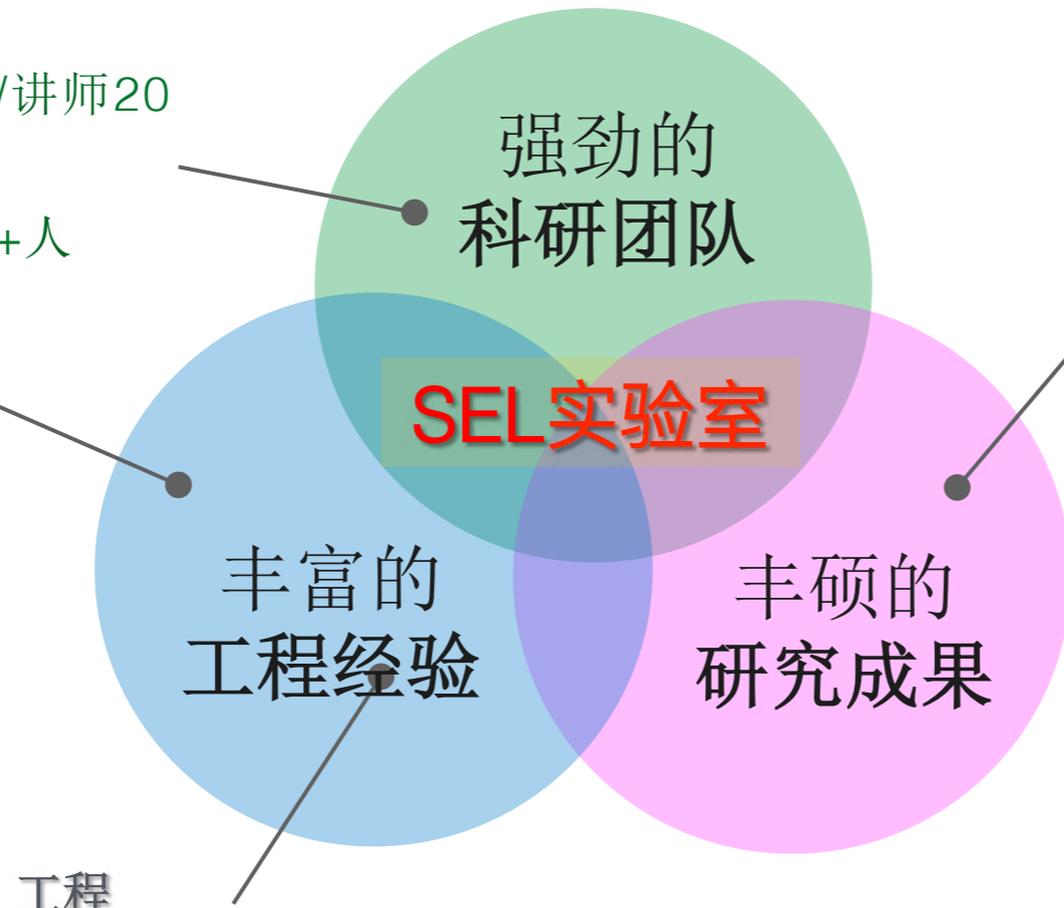
Zhejiang University

实验室团队

团队

- 教授/副教授/讲师20人
- 博士硕士100+人

媒体声音：
被美国CIO杂志誉为“一朵在中国开放的金融IT奇葩”



工程

- 完成超过100个海外合作工程项目，掌握世界一流的大型信息系统建设核心工程技术
- 丰富的大型金融信息系统开发、云计算、大数据实践经验
- 具备全球化金融领域业务知识
- 拥有美国道富、DST、IFDS、路透社、中国外汇交易中心等合作伙伴

研究

- 容器云技术
- 企业级应用运维模型
- 机器学习及其应用
- 大规模信息系统软件开发模型
- 300余篇高水平论文

实验室容器云贡献

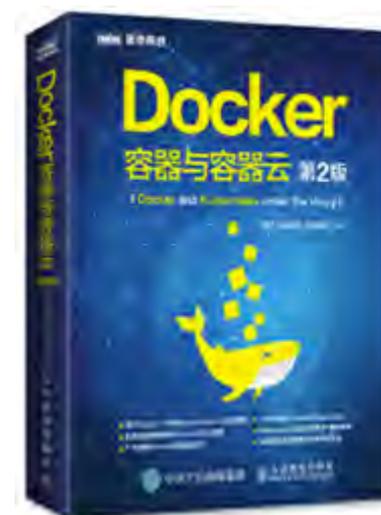
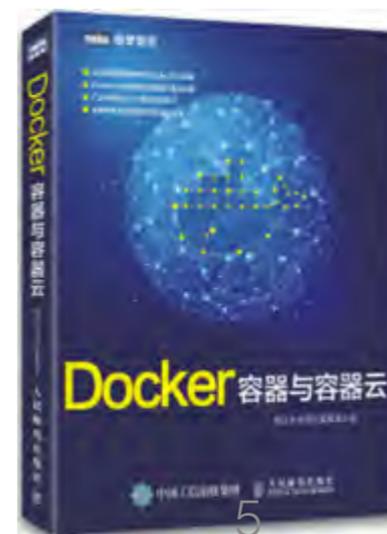
□ CNCF创始会员，在容器开源项目Kubernetes等贡献核心代码200多万行，贡献度国内第一，全球第五，并出版了国内第一本深度解析容器技术的专业书籍《Docker容器与容器云》。

□ 实验室容器管理平台2016年7月推出以来，迅速在电信，银行，电商，智能制造，智慧城市等行业服务于20多个合作伙伴。

#	Company	Lines of code
1	Google	65530604
	*independent	36866996
2	Red Hat	28110737
3	FathomDB	16332008
4	Fujitsu	2767331
5	Zhejiang University	1321654
6	Intel	1134044
7	HarmonyCloud	934708
8	IBM	829507
9	Inspur	463112

Kubernetes 项目的代码贡献, 按照贡献行数统计, 截止 2017/8/30
-- by <http://stackalytics.com/>

开源项目	国内排名	全球排名
Kubernetes	1	5
Docker	5	24
Cloud Foundry	3	26



我们为什么做容器平台的智能运维



SEL
Zhejiang University

TABLE OF CONTENTS

浙江大学SEL实验室

容器环境下智能运维的挑战

容器环境下智能运维研发实践

智能运维发展展望

传统运维挑战



事故报警以投诉为主



客服记录故障描述



故障诊断定位



现场跟踪处理

容器环境特点

- 容器化是IT架构主要发展趋势之一
- 组件和服务数量众多
- 部署模式多样化
- 自带服务治理功能

容器环境下智能运维挑战

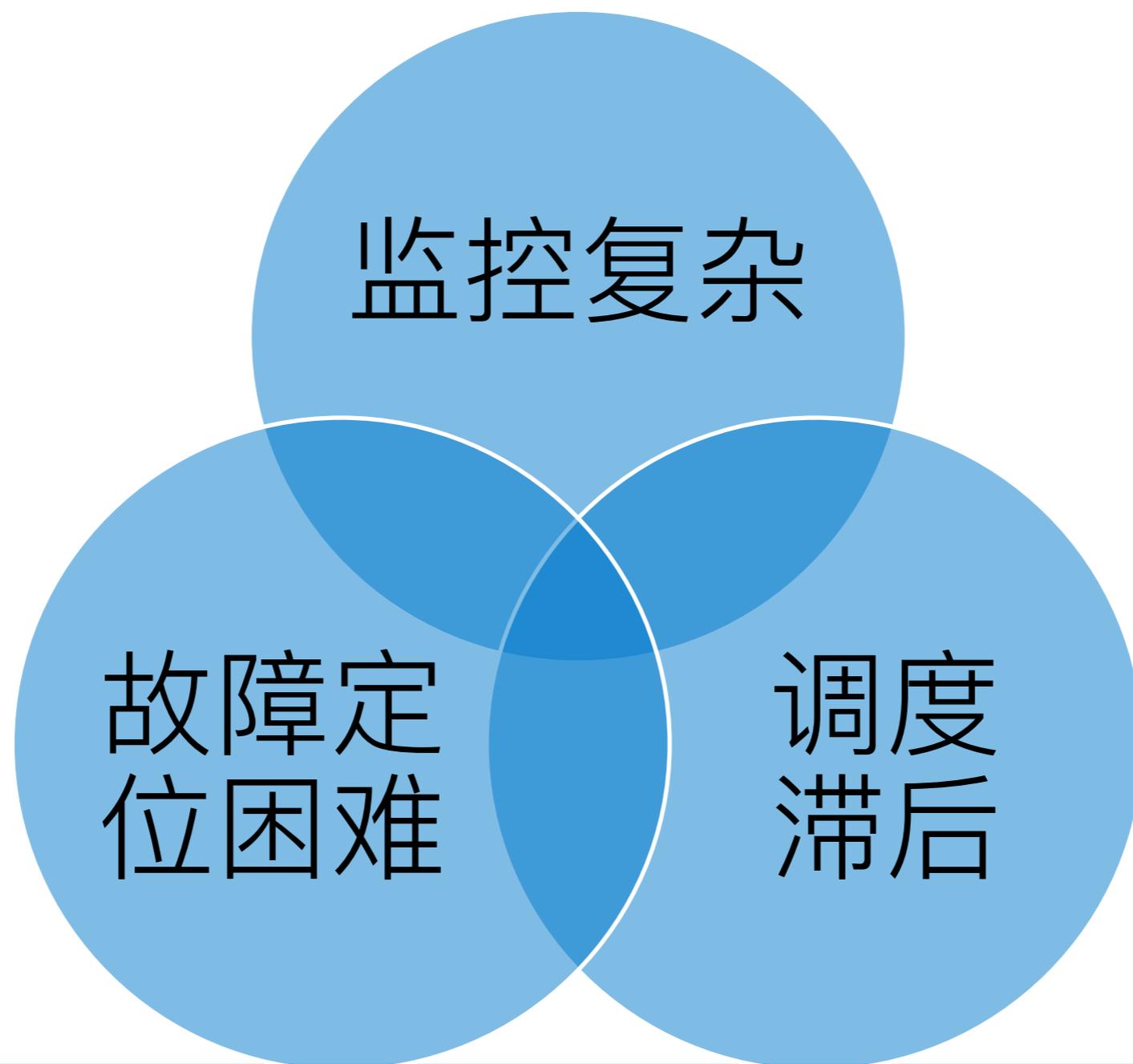


TABLE OF CONTENTS

浙江大学SEL实验室

容器环境下智能运维的挑战

容器环境下智能运维研发实践

智能运维发展展望

智能运维技术研发与实践

1. 容器统一 监控

- 容器监控
- 链路数据跟踪
- 日志数据汇聚

2. 故障根源 定位

- 动态拓扑分析
- 因果关系提取
- 变更影响分析

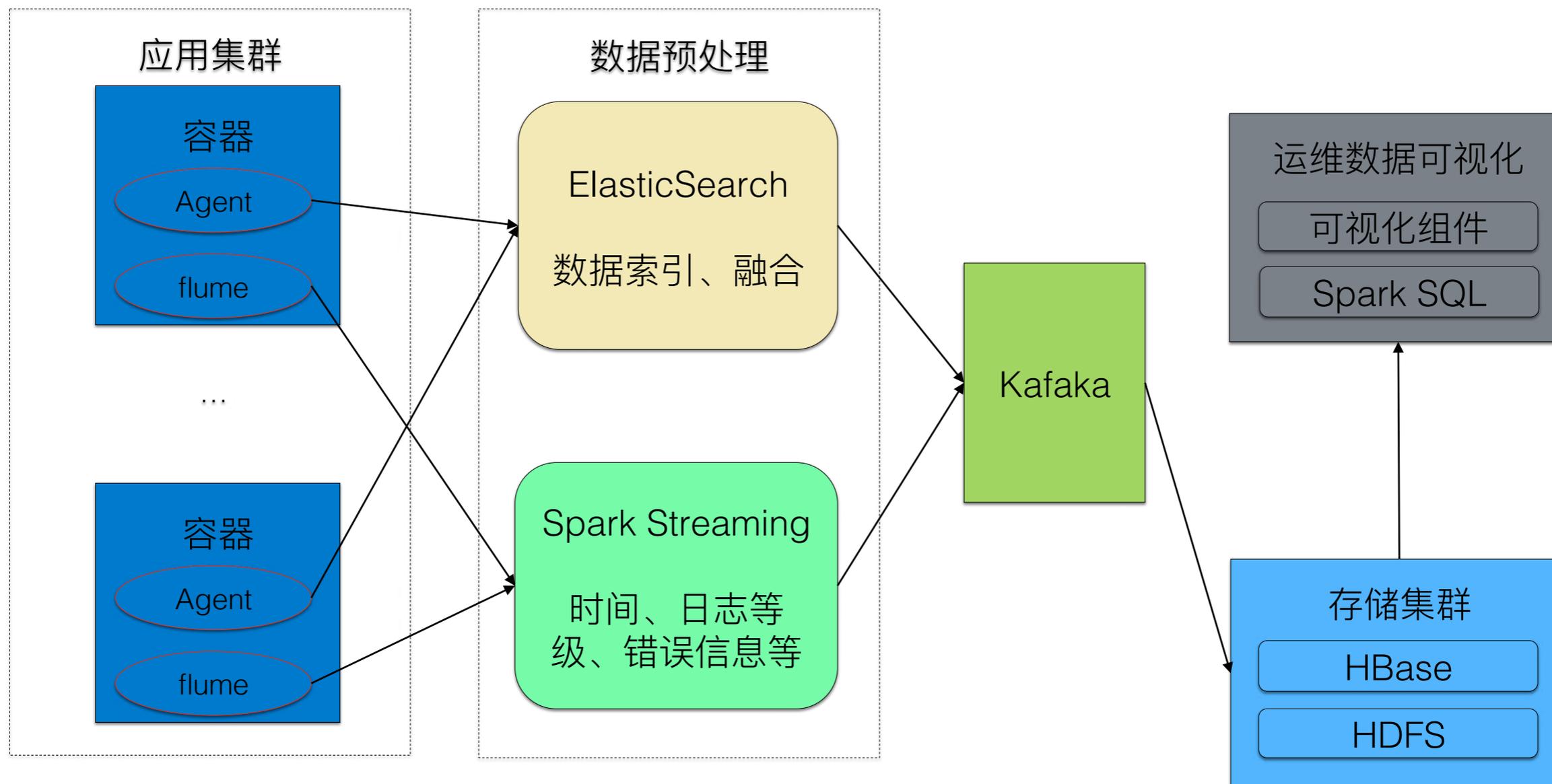
3. 智能调度

- 性能瓶颈分析
- 资源利用率估算
- 资源调度

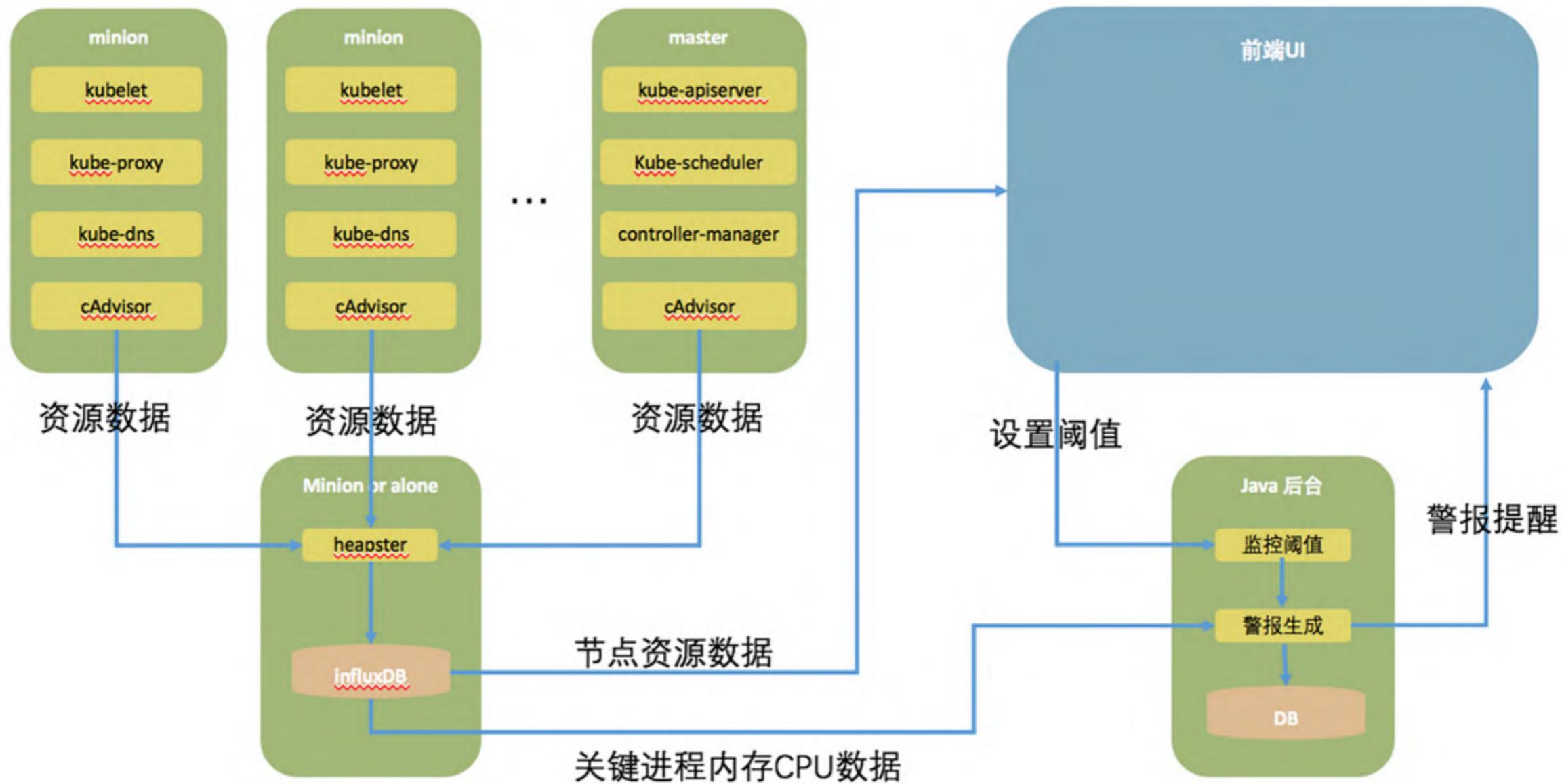
1. 容器统一监控

- 案例场景为某政务数据中心，采用IaaS+容器技术方案
- 主要业务：支撑各部门单位应用系统发布运维
- 特点：系统托管在数据中心，发现问题后，应用方请求数据中心提供相关数据进行故障诊断
- 目标：提供适合应用的监控数据，供运维运行分析

设计思路



容器监控方案



如何监控容器镜像

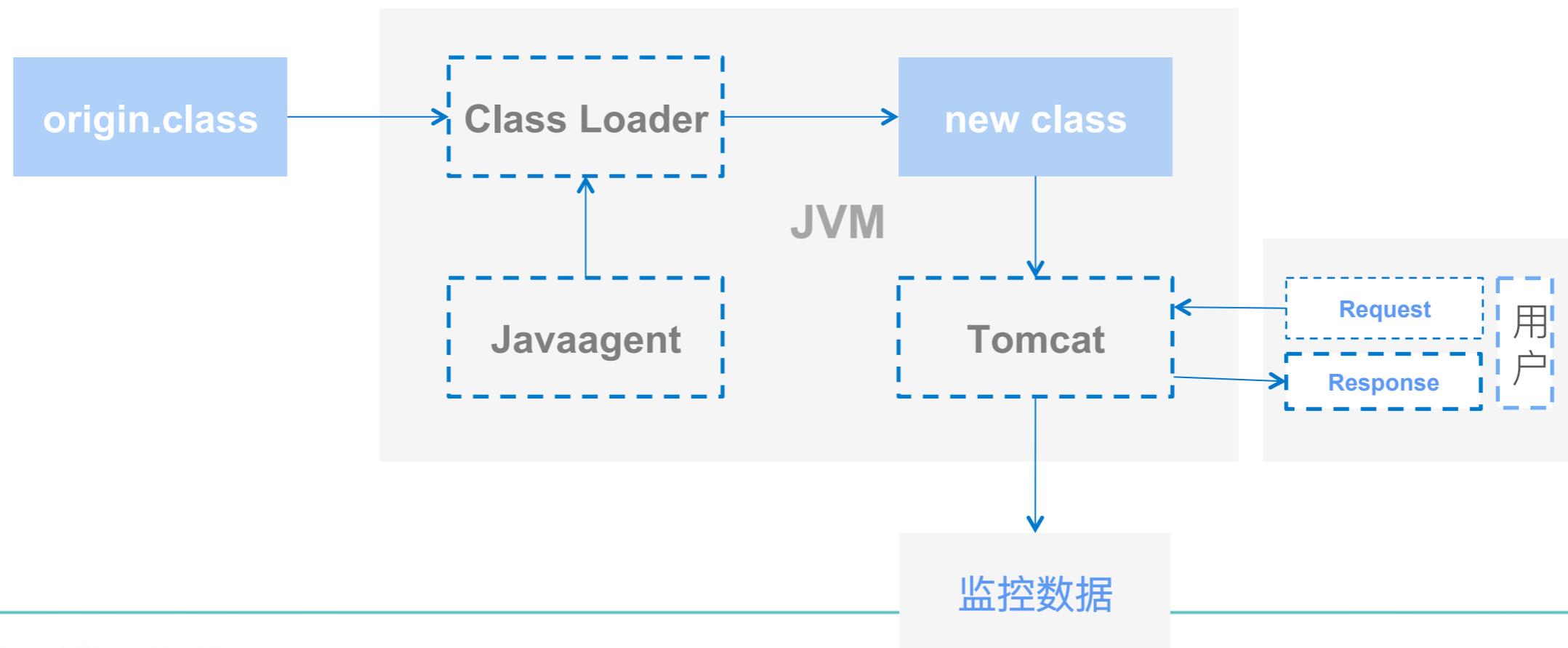
- Option 1: 镜像提前做好，tomcat，jetty等启动命令提前加好
- Option 2: 脚本Batch 批量更新镜像
- 但：根据只跟踪service ip，没有pod ip，如何对应到pod

Agent与K8S

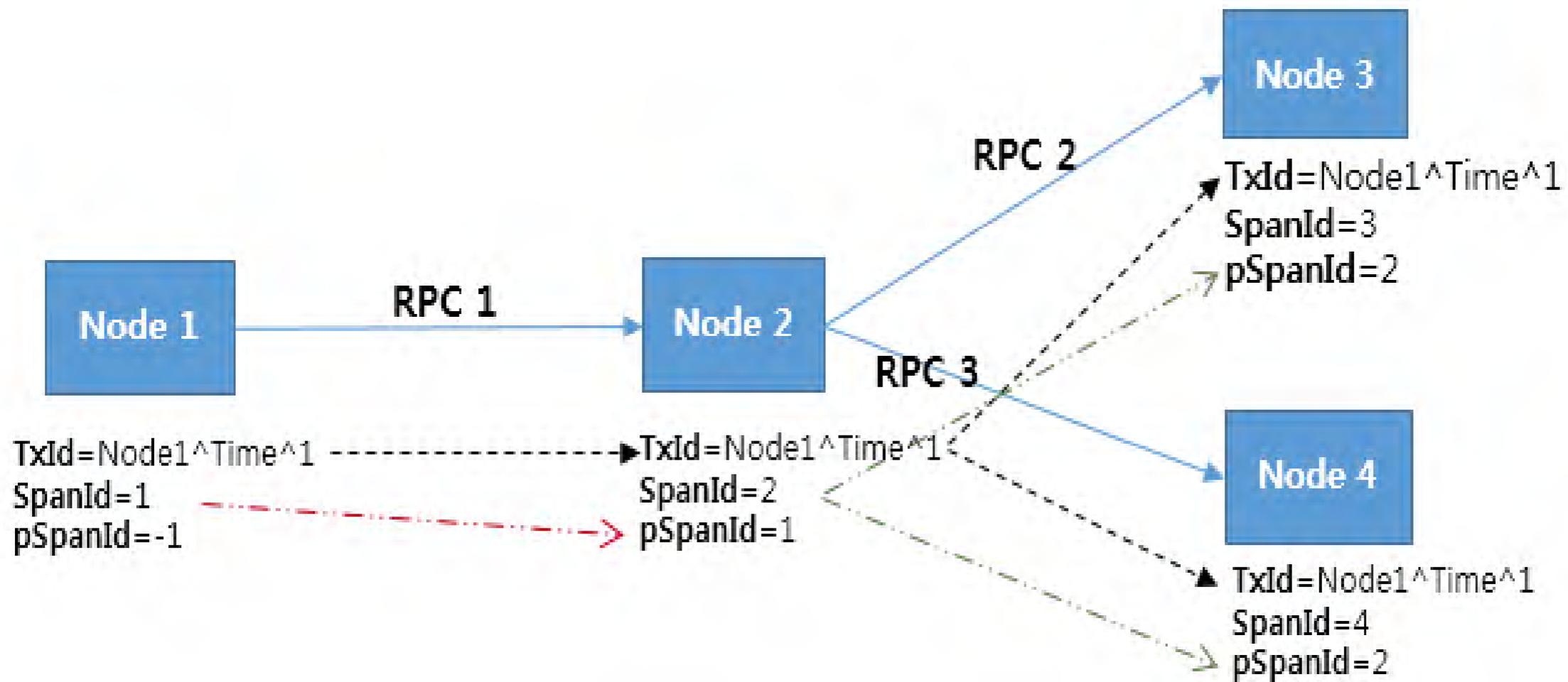
- Agent如何识别自己是哪个service？哪个POD？
 - ETCD维护着service与pod关系
 - Agent启动时查询自己所属的service

如何获取链路

- 对代码掌控能力较强采用OpenTracing
- 对代码进行插装实现Google Dapper的原理



Google Dapper



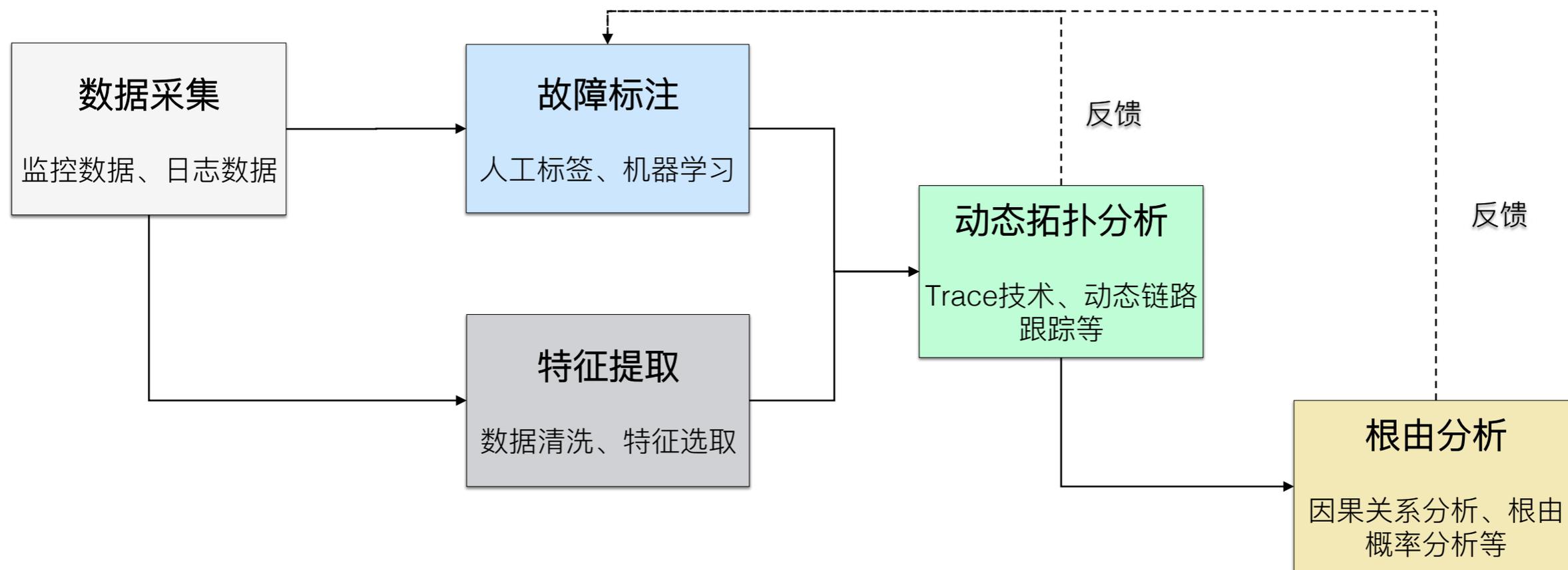
链路跟踪

- TraceID：识别用户一次请求，所有全链路上的节点共用一个TraceID
- SpanID：正在处理用户请求的节点
- ParentSpanID：正在处理用户请求节点的上一个节点

2. 故障根源分析

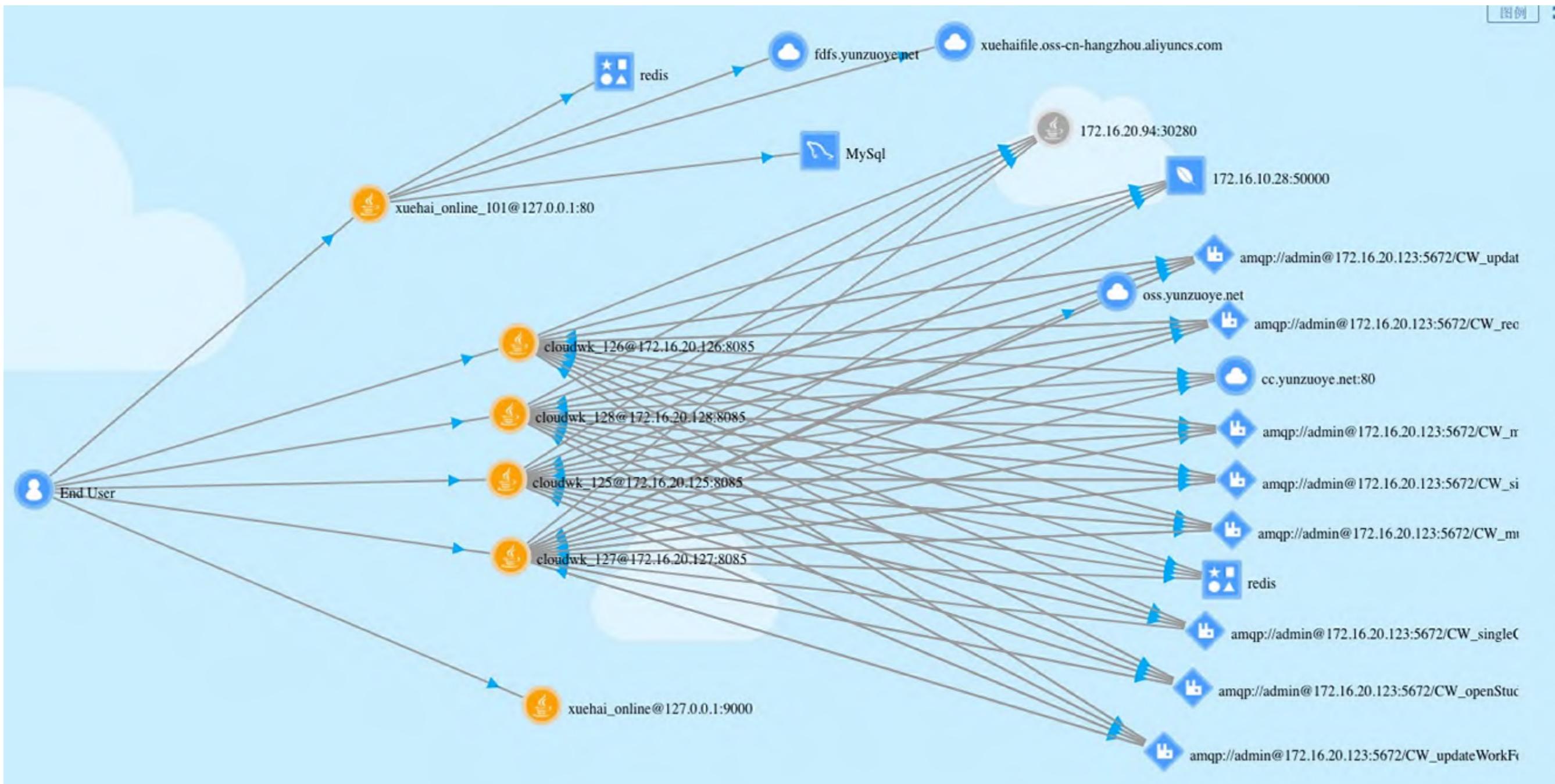
- 案例场景为某在线教育平台
- 提供各类教育资源分享，用户广泛，IT技术能力差异显著
- 客服经常接到比较奇葩的投诉，很难给出合理建议
- 发生故障后，IT团队修复问题需要花费较长时间定位问题根源
- 目标：减轻客服和运维工作量，快速定位故障或投诉根由

分析流程

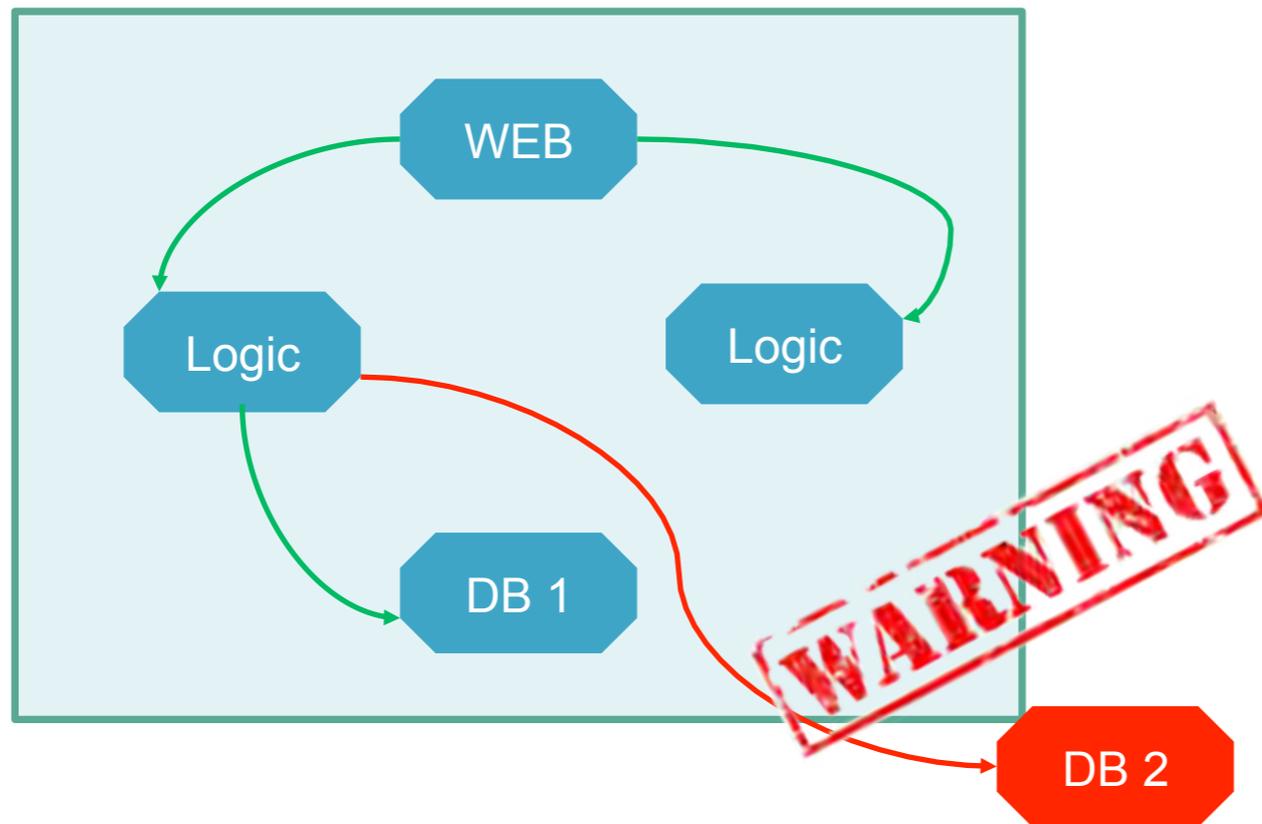


动态拓扑结构

图例 5



基于拓扑的异常检测

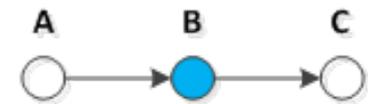


- 构建执行链正常模型
 - 资源访问
 - 性能基线
- 异常行为检测
 - 行为异常
 - 性能异常
 - 访问量异常

因果关系分析

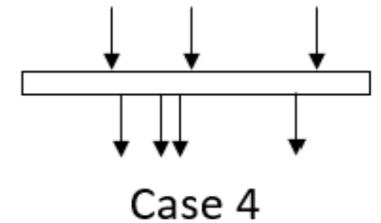
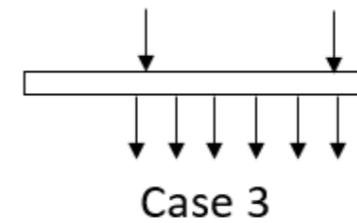
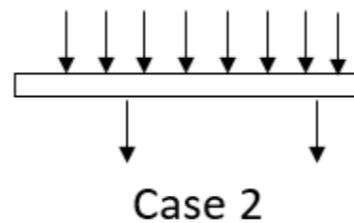
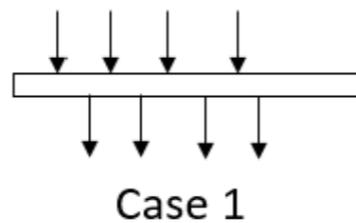
- 因果关系分析定义

- 构建数据中心服务拓扑结构图，如右图A调用B触发B调用C，A->B称为B的因边，B->C称为B的果边



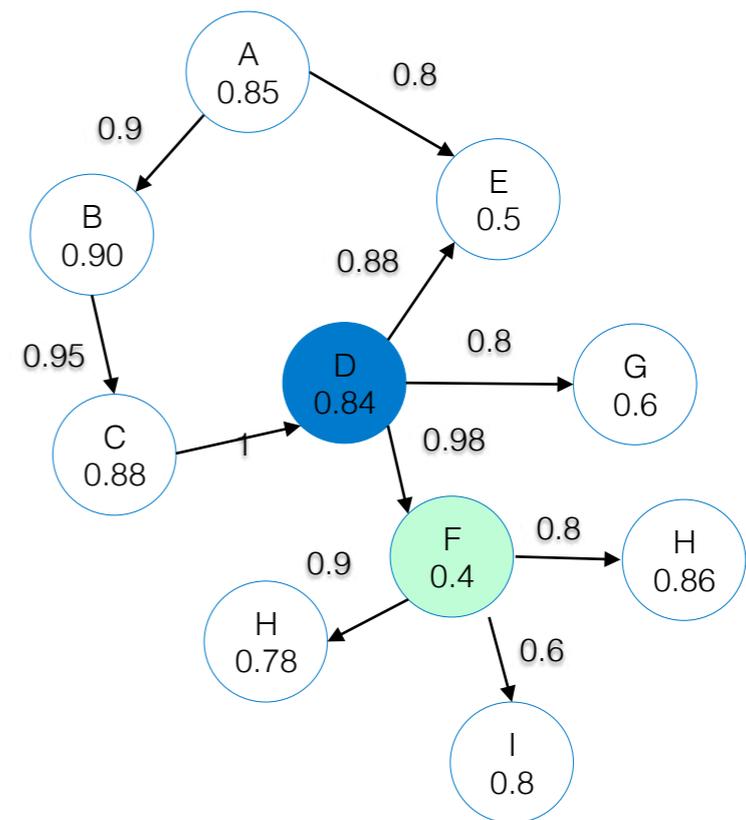
- 如何判定因果关系

- 果边应发生在因边之后的一个特定的时间段，即服务延时
- 服务延时不能超过超时时间也不能低于服务处理时间
- 通过关联分析，判断潜在因果关系



故障根由定位

- 对故障相关容器节点的可能性进行量化分析
 - 形成故障因果关系链，执行链的历史执行频度将作为该链路上节点权重
 - 计算因果链上每个节点的影响
 - 以右图D为例，影响D的节点为E，F等，受D影响的节点为C，B，A
 - 量化单个节点对其他节点影响
 - $P(D|F) = P(DF) * w_F / P(F)$
 - $P(C|F) = P(C|D) * P(D|F)$ ，其中权重由节点历史执行频率决定
 - 异常发生时，可快速计算出其他节点的嫌疑程度



3. 智能调度

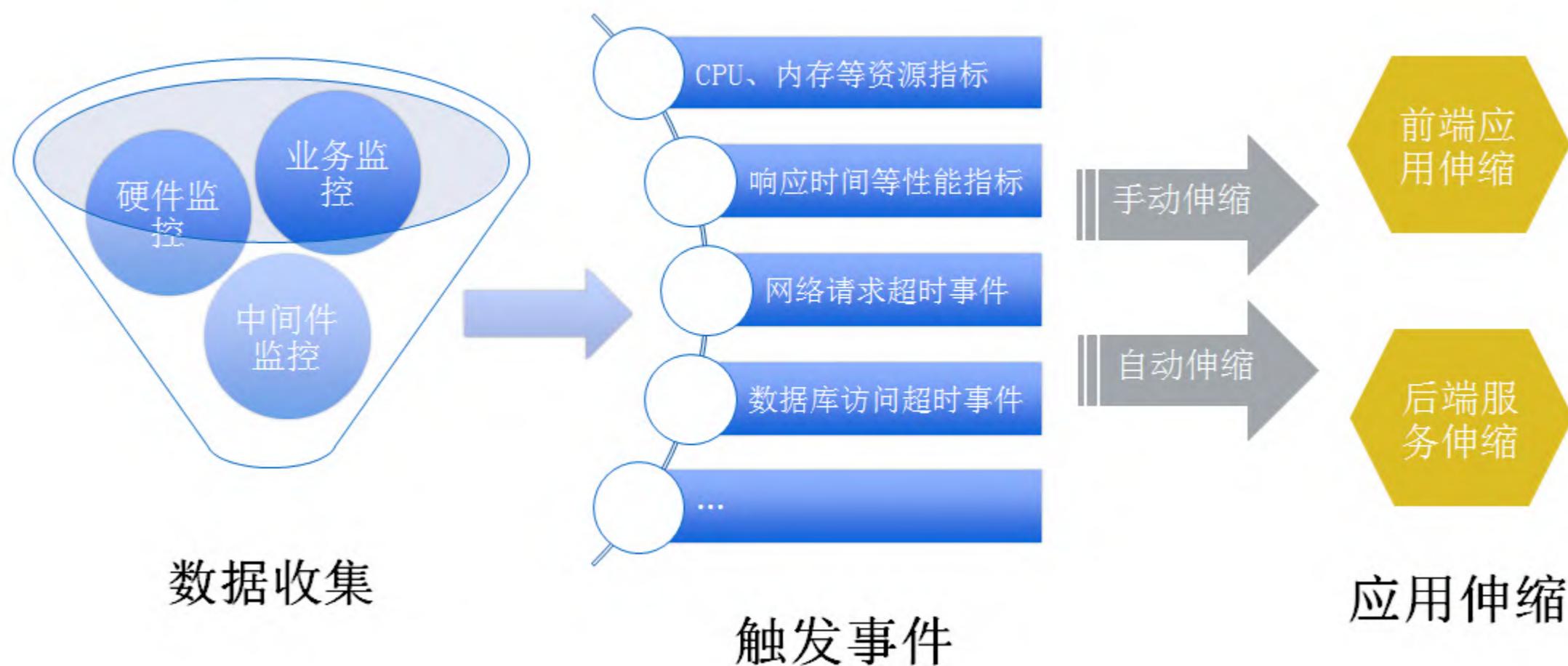
案例场景为智慧物流信息共享和监控管理平台

支持一个地级市范围内，10+万物流车辆，
数百家物流快递公司

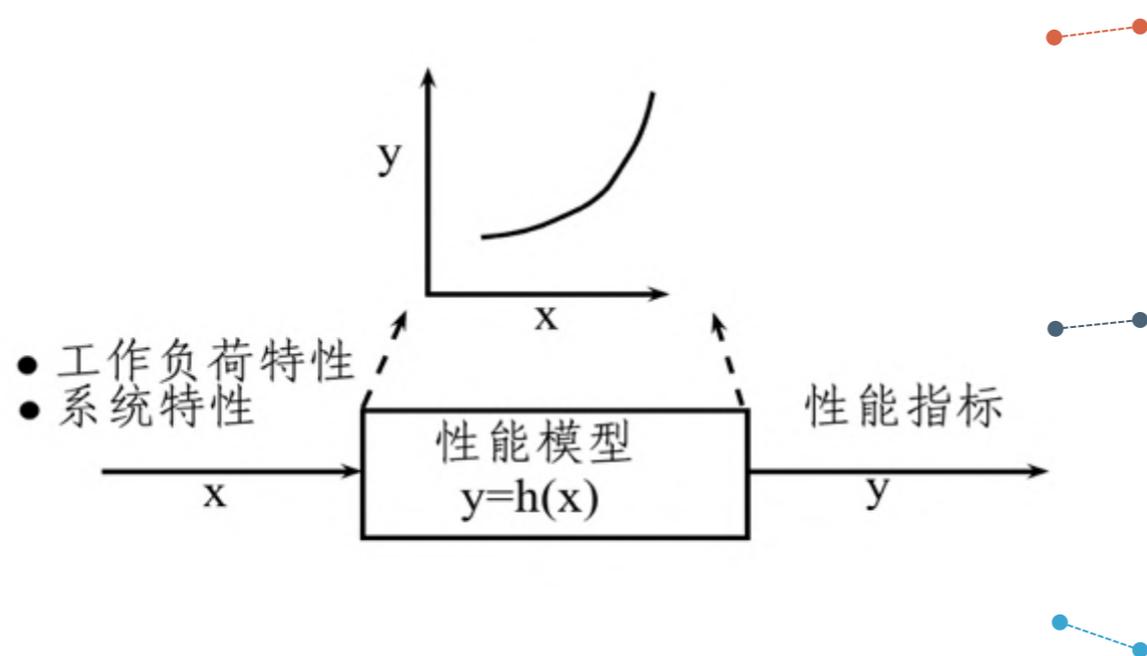
秒级实时信息更新和发布

采用云平台增强其弹性伸缩能力

设计思路



性能瓶颈定位



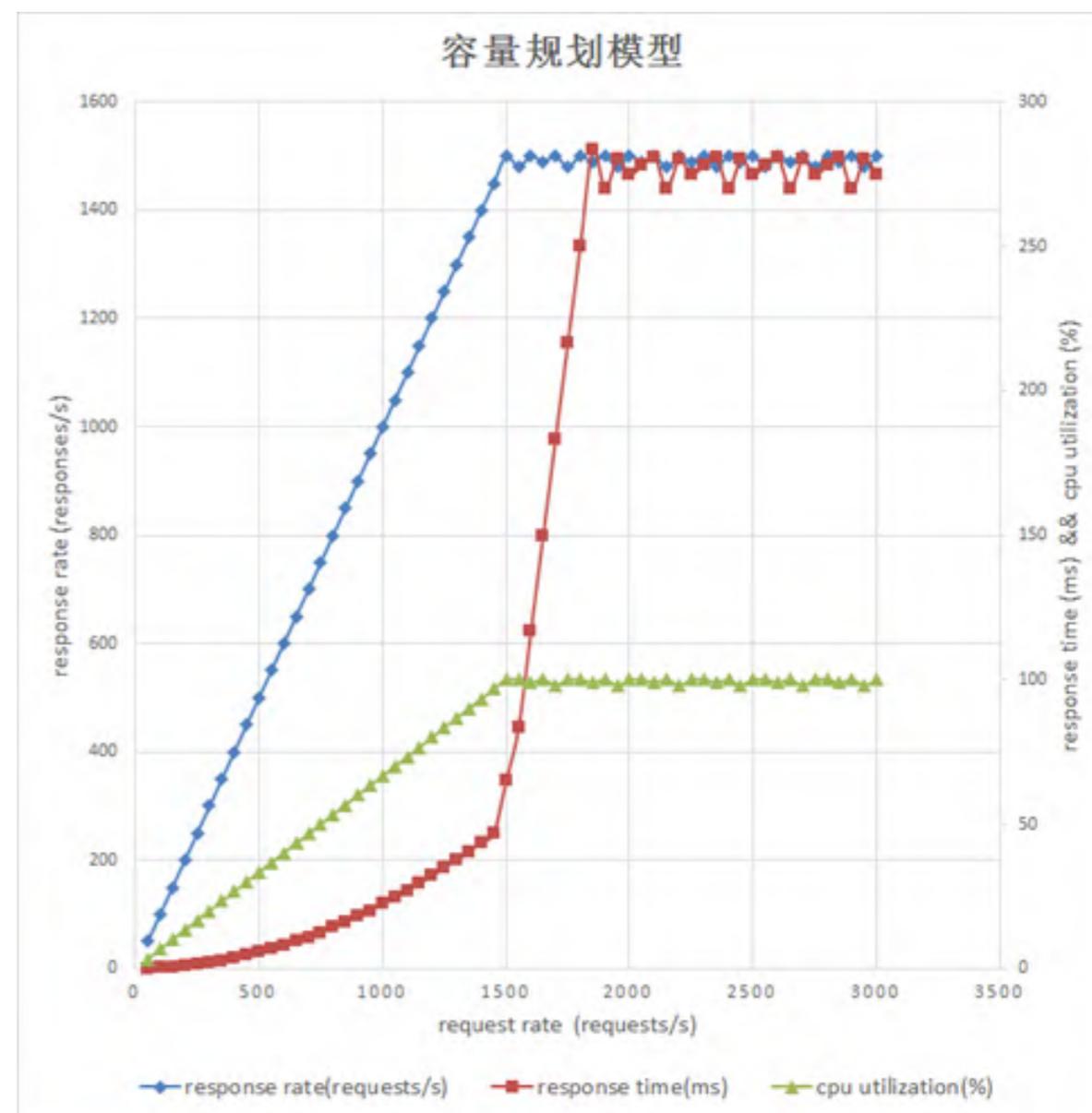
工作负荷特性：到达率、
到达间隔分布、分类

系统特性：调度策略、负载均衡策略

性能指标：响应时间、吞吐量和资源
的利用率

性能瓶颈定位

- 资源利用率
 - 服务时间估算（多元线性回归）
 - 资源利用率预测
- 响应时间预测
 - MVA队列模型
 - BP神经网络 + 遗传算法



应用容量估算相关方法

负载模型

- 基于资源需求的事务分类，CPU密集型、IO密集型等

流量分析

- 隐马尔科夫模型
- 自适应平滑 Holt-Winters

资源利用率估算

- Regression Analysis

响应时间预测

- Neural network
- Queue network

容量估算

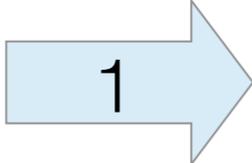
- 基于性能和资源利用率的估算，资源满足应用需求的同时性能不会明显下降

Access log

```

111.111.111.111 - [2/May/08 - 00:23:48] "GET /reportAAA.HTML"
123.123.123.123 - [2/May/08 - 00:23:48] "GET /reportBBB.HTML"
222.222.333.444 - [2/May/08 - 00:25:48] "GET /formAAA.HTML"
333.333.333.333 - [2/May/08 - 00:25:55] "GET /formBBB.HTML"
444.444.444.444 - [2/May/08 - 00:26:01] "GET /consolidationAAA.HTML"
555.555.555.555 - [2/May/08 - 00:26:03] "GET /login.HTML"
    
```

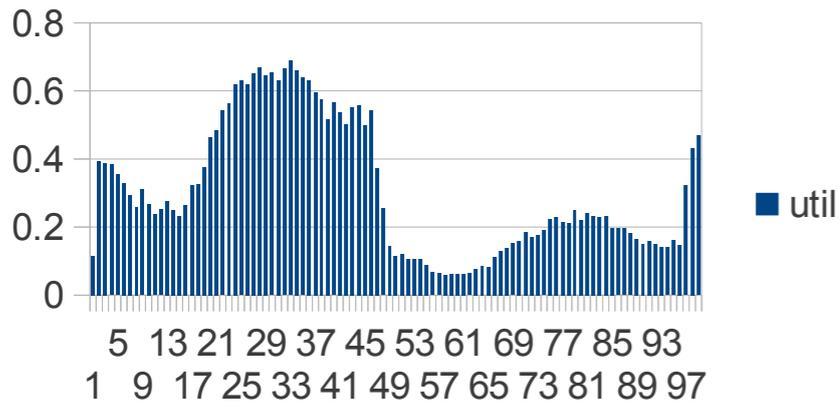
ETL



Transaction profile

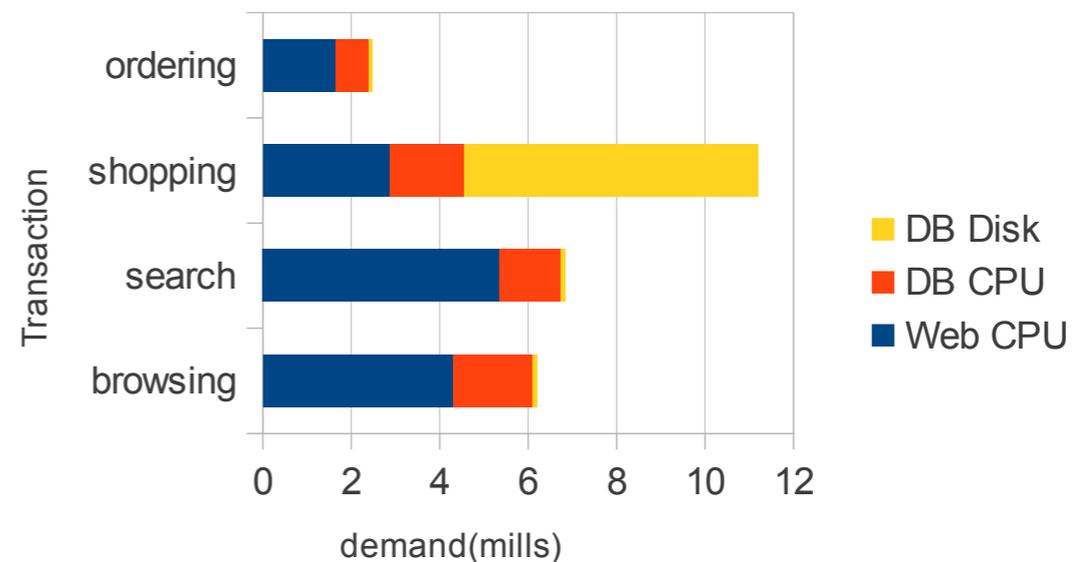
Time (hour)	N_1	N_2	N_3	N_4	...	N_{756}	$U_{CPU}(\%)$
1	21	15	21	16	...	0	13.3201
2	24	6	8	5	...	0	8.4306
3	18	2	5	4	...	0	7.4107
4	22	2	4	7	...	0	6.4274
5	38	5	6	7	...	0	7.5458
...							

Resource log

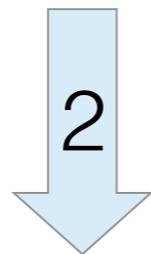


Transaction

resource demand

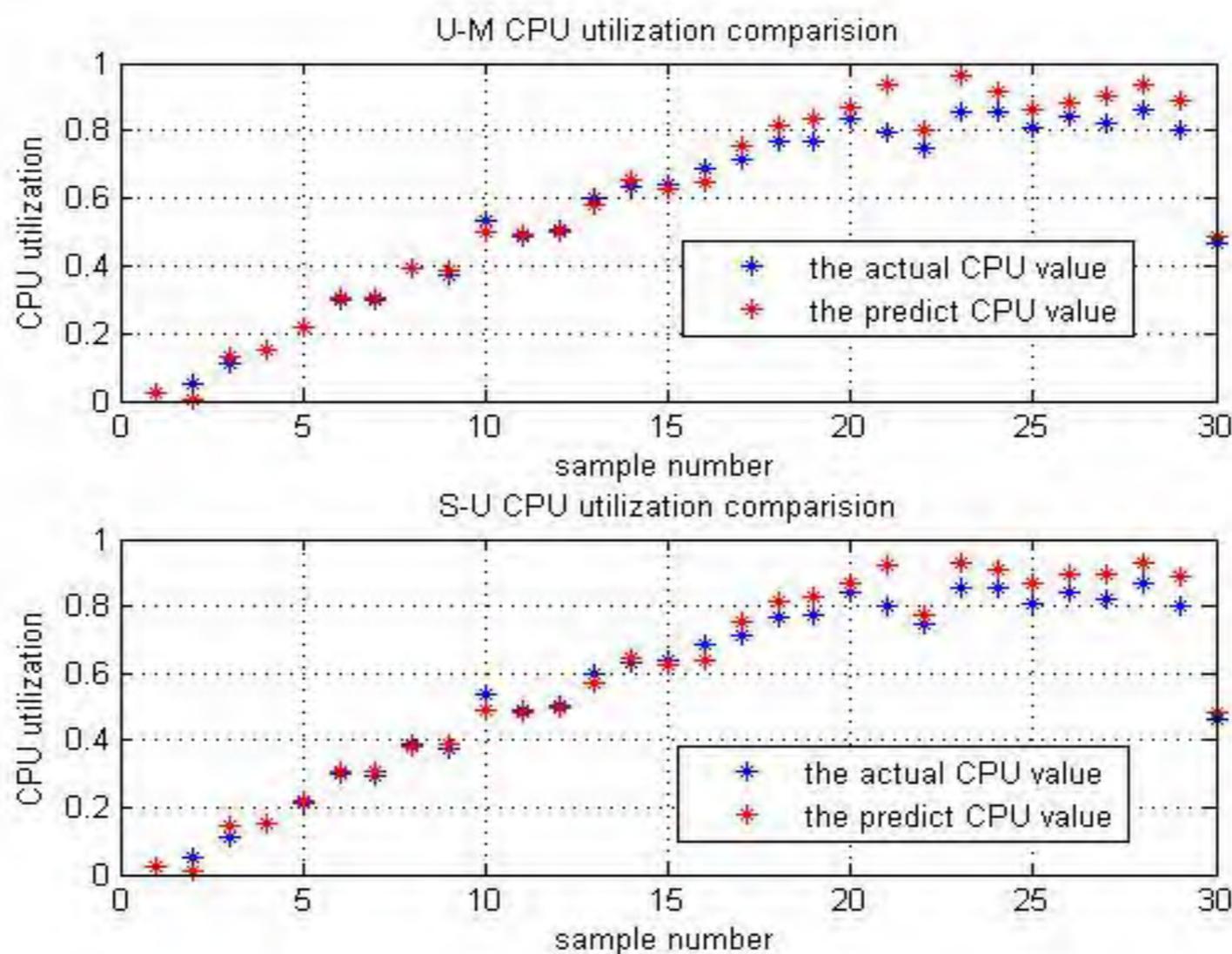


Regression Analysis



应用资源需求估算

- MySQL 节点 CPU利用率 vs 负载



容器动态资源调度

- Kubernetes容器资源调度：Predicates和Priority
- 业务指标驱动资源调度
 - 容器业务访问量预测 →
 - 容器资源需求 / 应用性能分析 →
 - 如可能过载，容器弹性伸缩或迁移

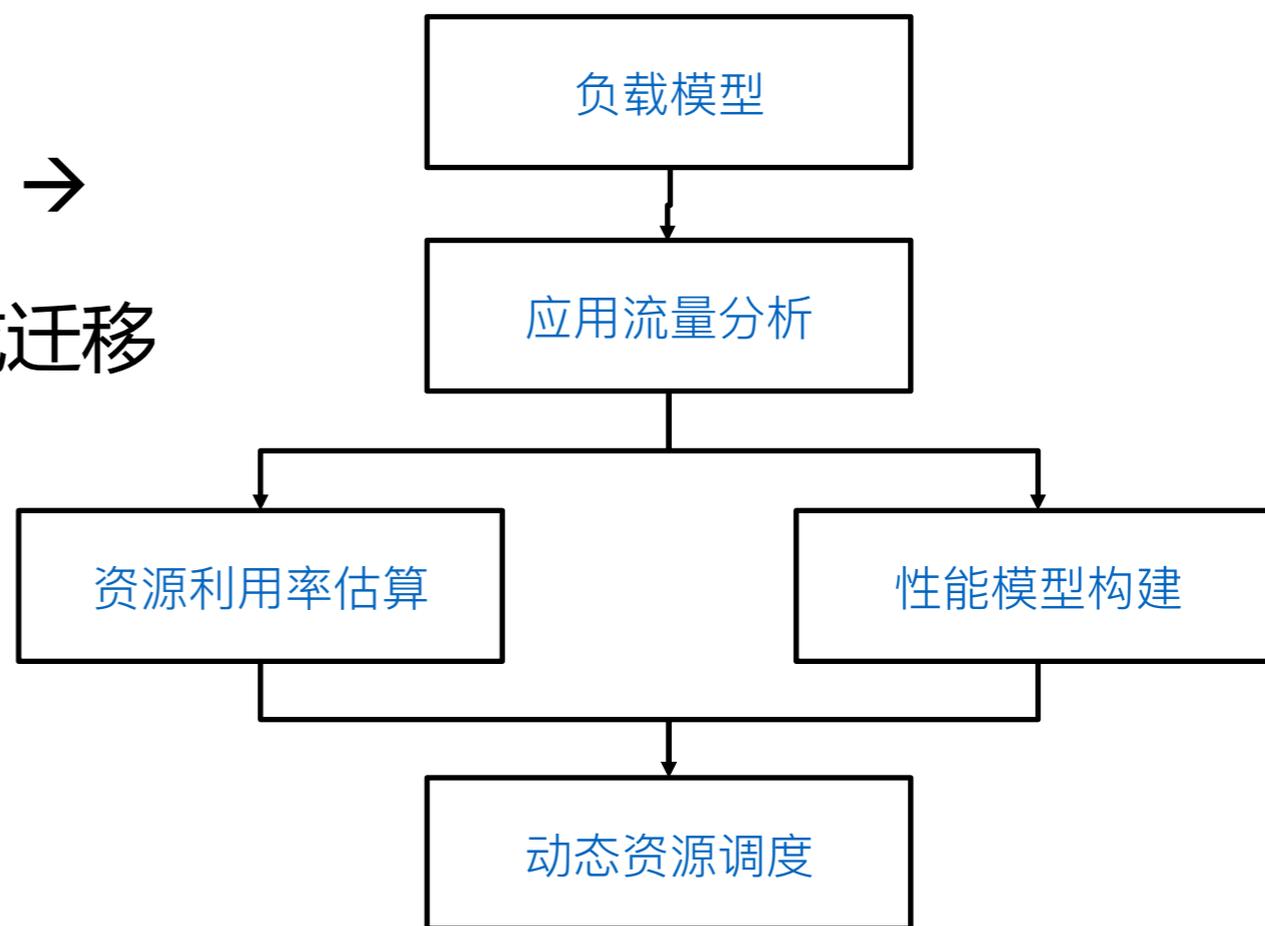


TABLE OF CONTENTS

浙江大学SEL实验室

容器环境下智能运维的挑战

容器环境下智能运维实践案例

智能运维发展展望

智能运维展望

被动故障告警到 主动预防

- 用预测代替传统的单纯监控
- 实时高性能运维数据处理

用户行为分析到 服务行为分析

- 服务资源访问
- 服务性能基线
- ...

人工故障恢复到 自动故障恢复

- 改变传统人工 or 规则修复故障的模式
- 通过机器学习, 不断完善故障恢复知识库

进一步交流请加微信



THANKS!

智能时代的新运维

CNUTCon 2017