

# 基于eBPF构建无侵入的可观测系统

张海彬 阿里云





张海彬

花名古琦，阿里云技术专家

负责阿里云应用监控eBPF版的产品研发工作，有多年的微服务治理、可观测实践经验，目前在阿里云原生应用平台负责eBPF的应用监控工作。



# 目录

CONTENTS



01

微服务应用可观测的挑战

02

使用eBPF进行监控数据采集

03

构建完整的可观测系统

04

未来与展望





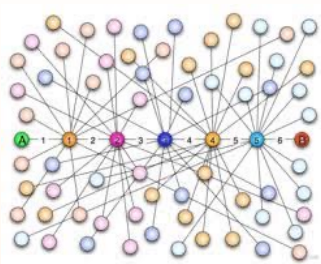
# Part 01

## 微服务应用可观测的挑战





# 微服务应用可观测的挑战



应用：微服务架构、多语言、多协议

挑战1：微服务、多语言、多协议环境下，端到端观测复杂度上升，埋点成本居高不下



如何关联？

基础设施层复杂度日益增加



挑战3：数据散落，工具多，缺少上下文，排查效率低下

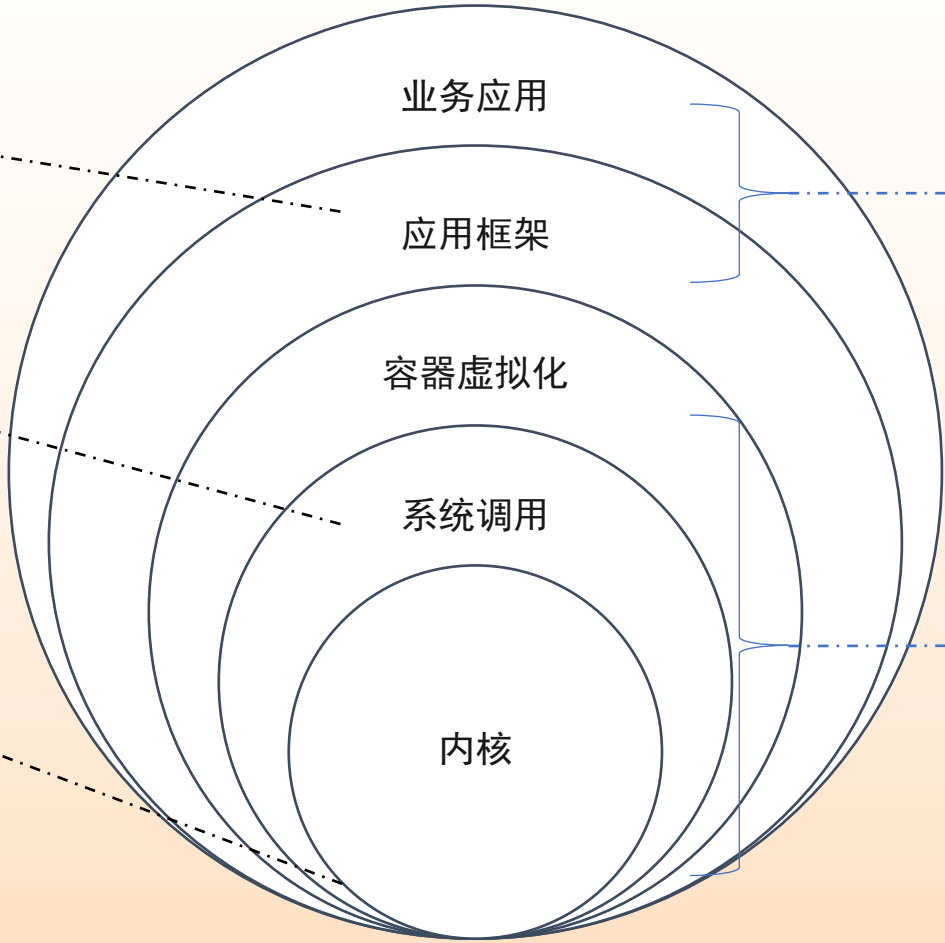
# Kubernetes下的可观测

Kubernetes组件异常:

Scheduler, KCM,  
etcd, api-server,  
coredns...

系统调用异常: 网络请  
求, 内存申请, 文件操  
作, CGroup...

内核异常: 进程调度,  
内存管理, 文件管理,  
夯机宕机, 资源异  
常...



应用组件异常: 线程池满, 数据库连接无法获取,  
OOM, 文件读取错误...

应用性能监控 (APM)

Kubernetes监控

无法自顶向下端到端  
串联导致棘手问题频  
发。

# Part 02

## 使用eBPF进行监控数据采集



# 使用eBPF进行监控数据采集

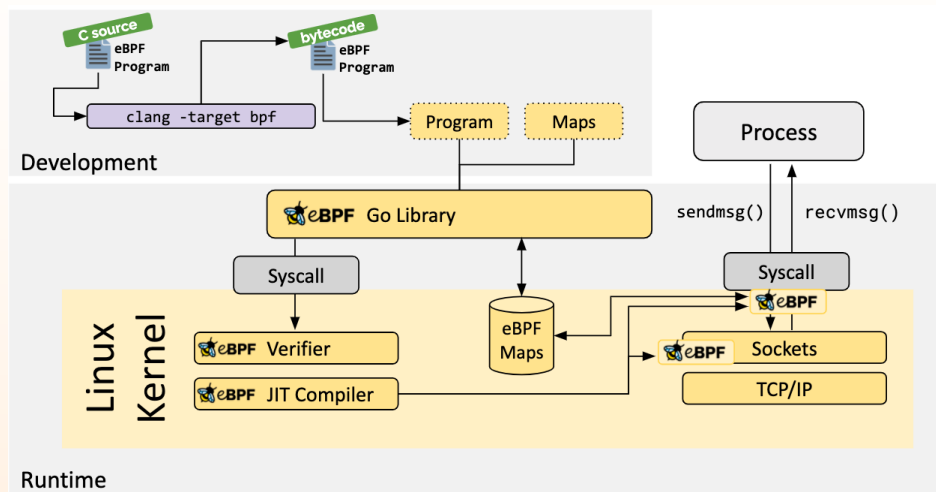
## 数据处理流程

加载、编译BPF程序

监听内核事件

识别网络协议

聚合指标、输出Trace



extend Berkeley Packet Filter，在Linux 内核中运行沙盒程序，而无需更改任何源代码或加载任何内核模块，有以下特性：

- 无侵入：成本极低。应用不需要改任何代码。
- 动态可编程：不需要重启探针，动态下发采集脚本
- 高性能：自带JIT编译成native代码
- 安全：verifier机制保障内核运行稳定

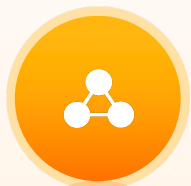




# eBPF技术在可观测领域的优势



无侵入

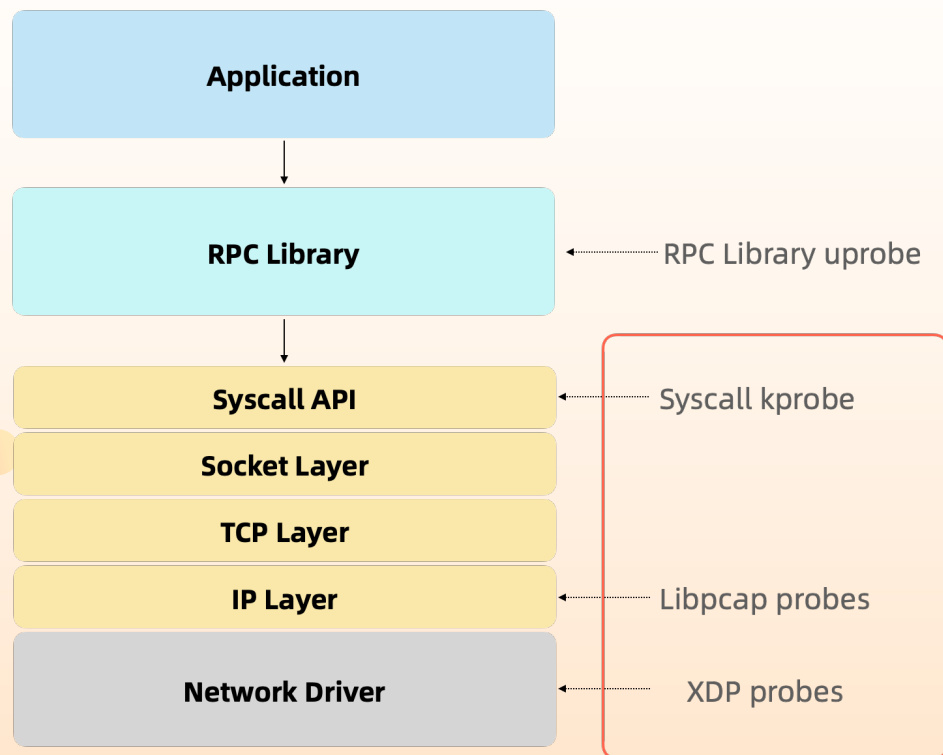


多语言/多协议/多框架



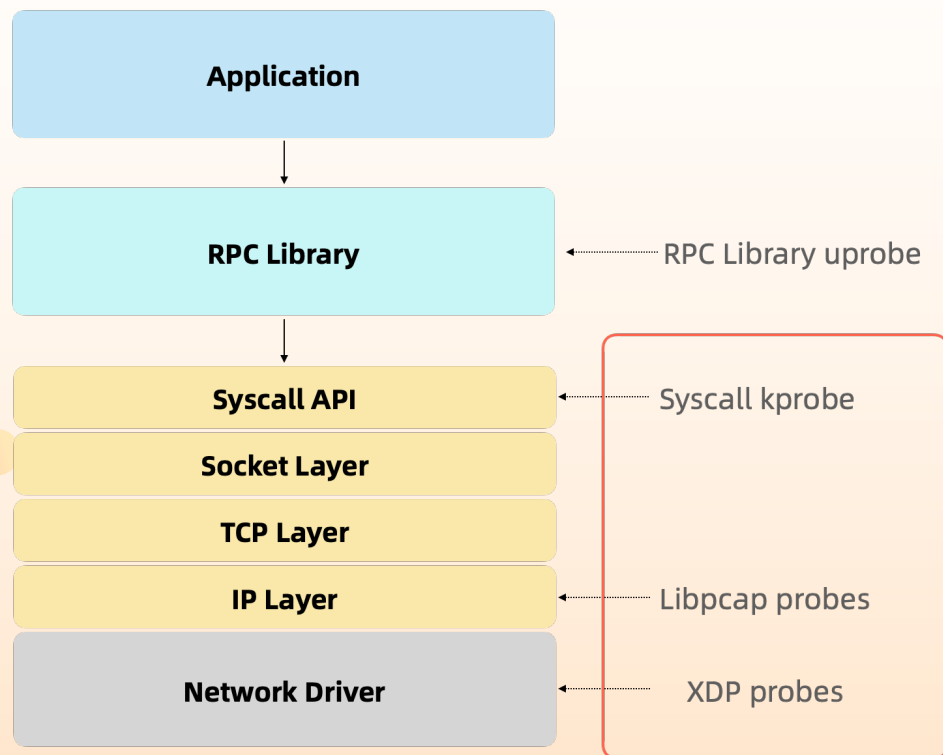
全栈覆盖

# 无侵入



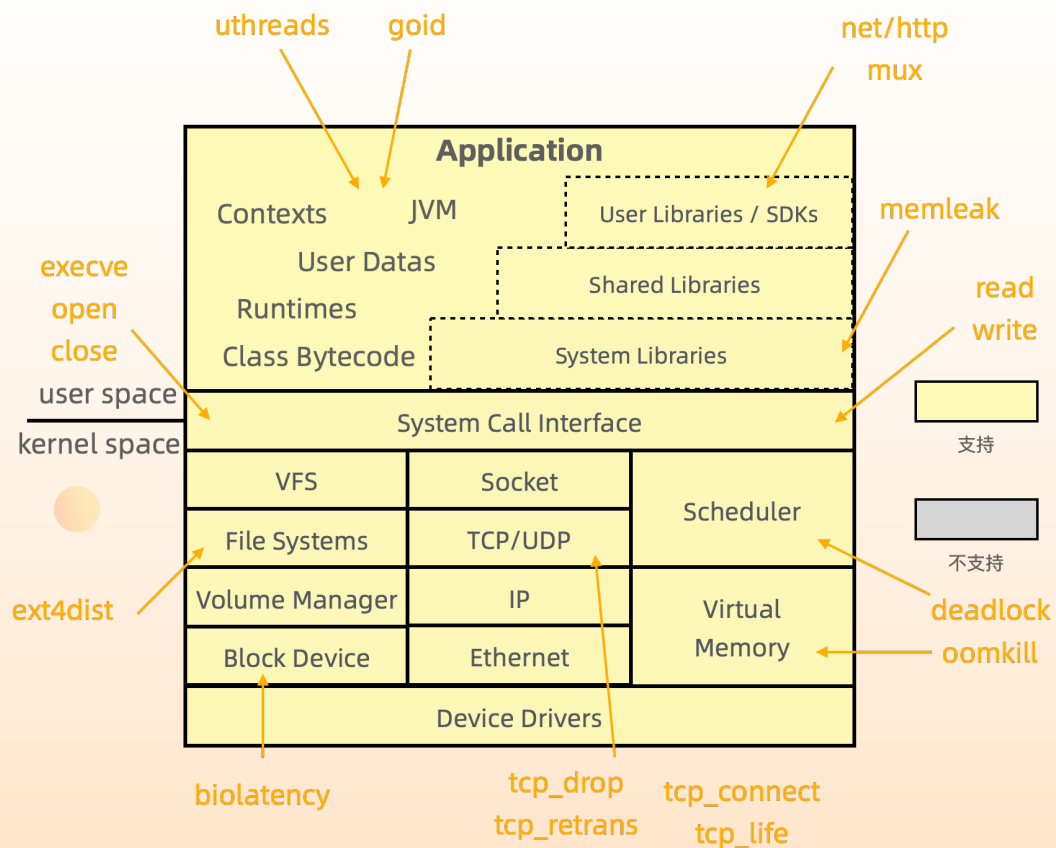
- 无需修改代码
- 无需重启应用
- Verifier保证运行安全

# 多语言、多协议、多框架



- 捕获网络字节流
- 无需适配编程语言
- 无需适配协议框架
- 同时支持用户态插桩

# 全栈覆盖



- uprobe
- kprobe
- tracepoint
- USDT
- perf
- ...



# Part 03

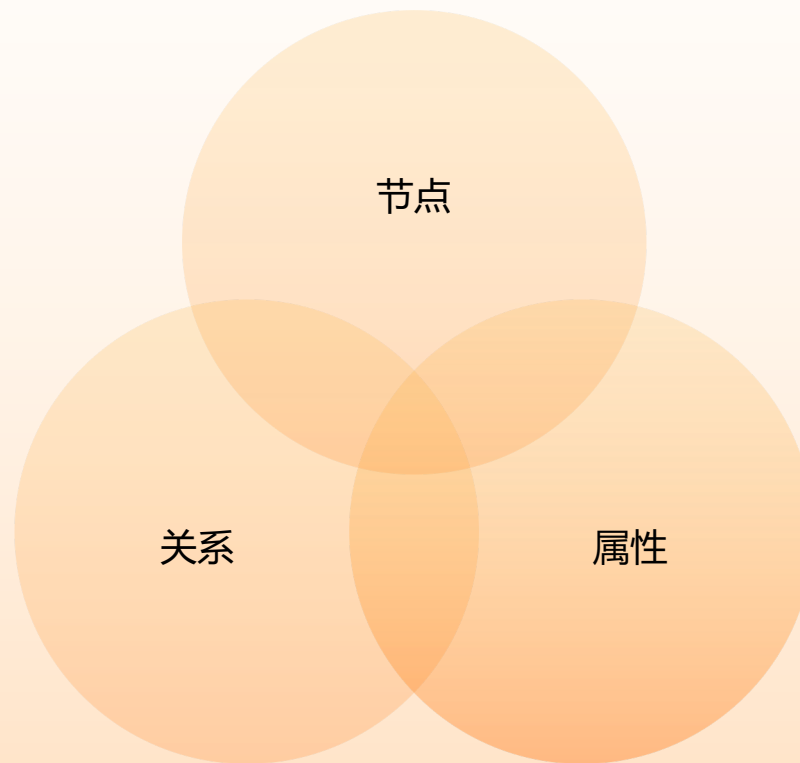
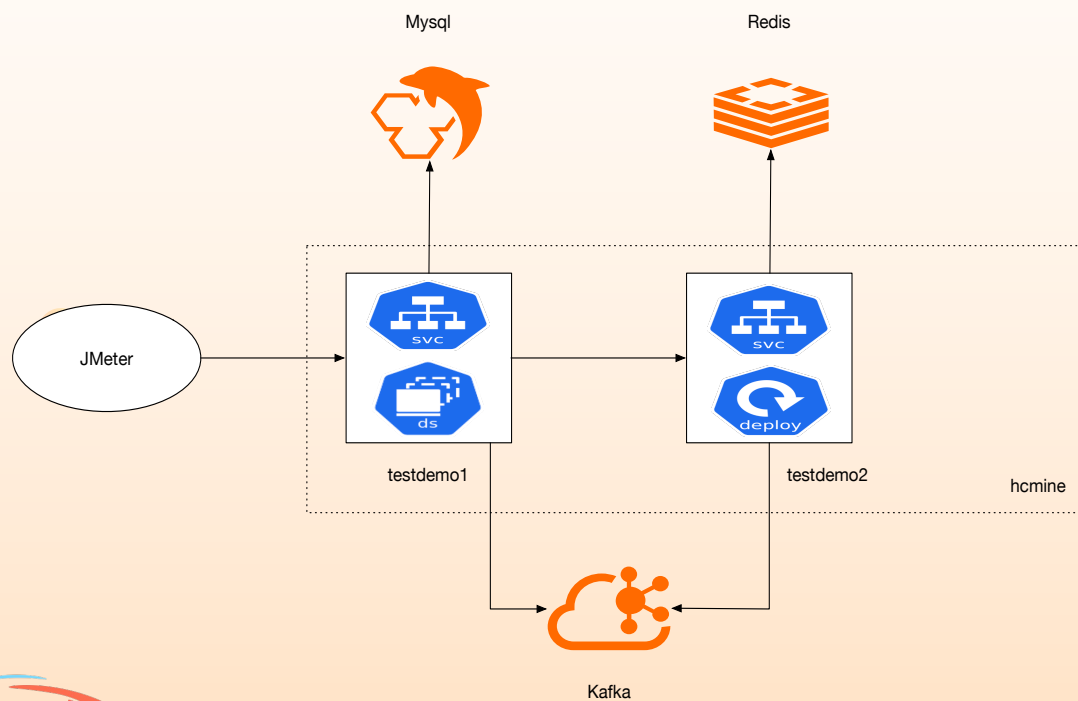
## 构建完整的可观测系统





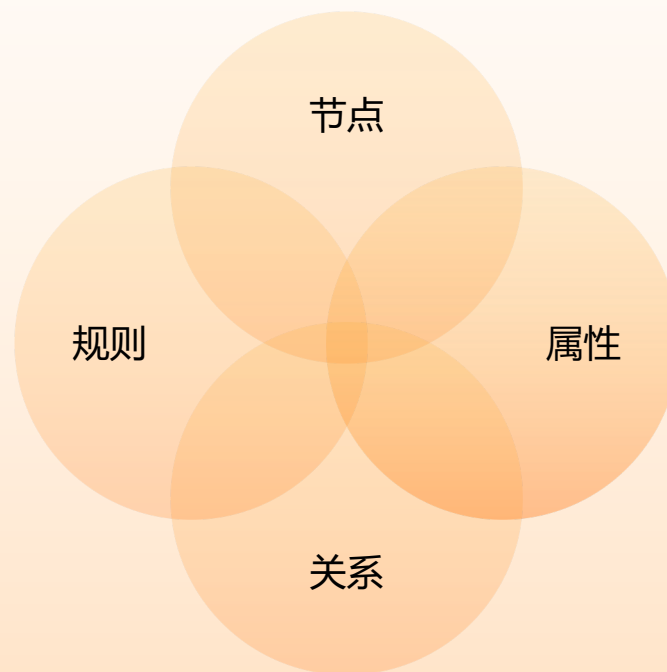
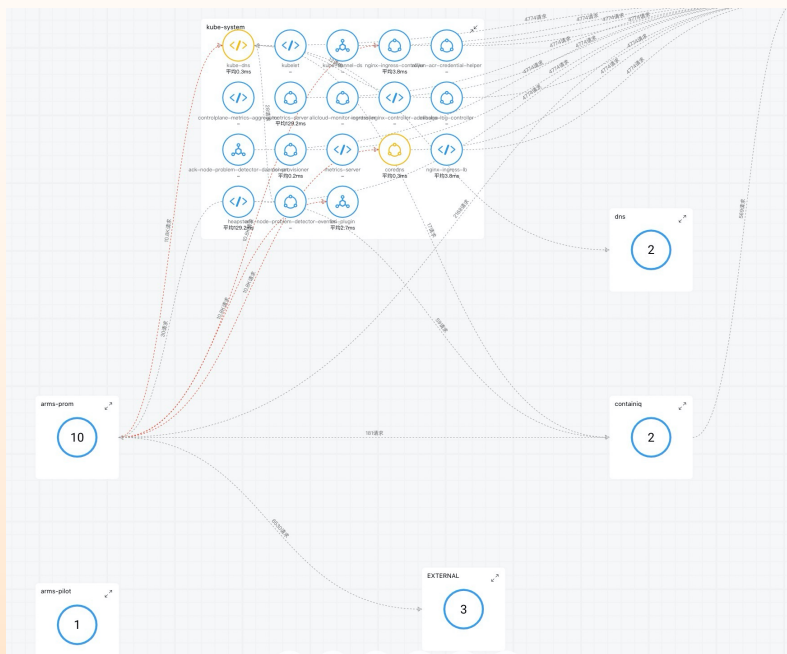
# 架构感知

架构感知，节点和关系以及他们的属性，能够正确地反应当前运行的网络关系，帮助用户感知架构，通过对比期望架构，发现问题，通常在新应用上线，新地区开服，整体链路梳理等场景使用。



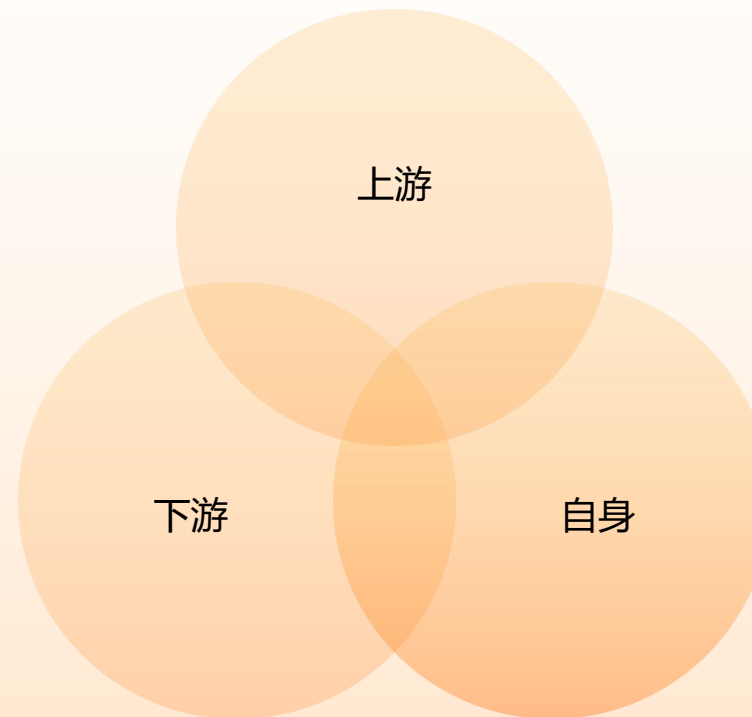
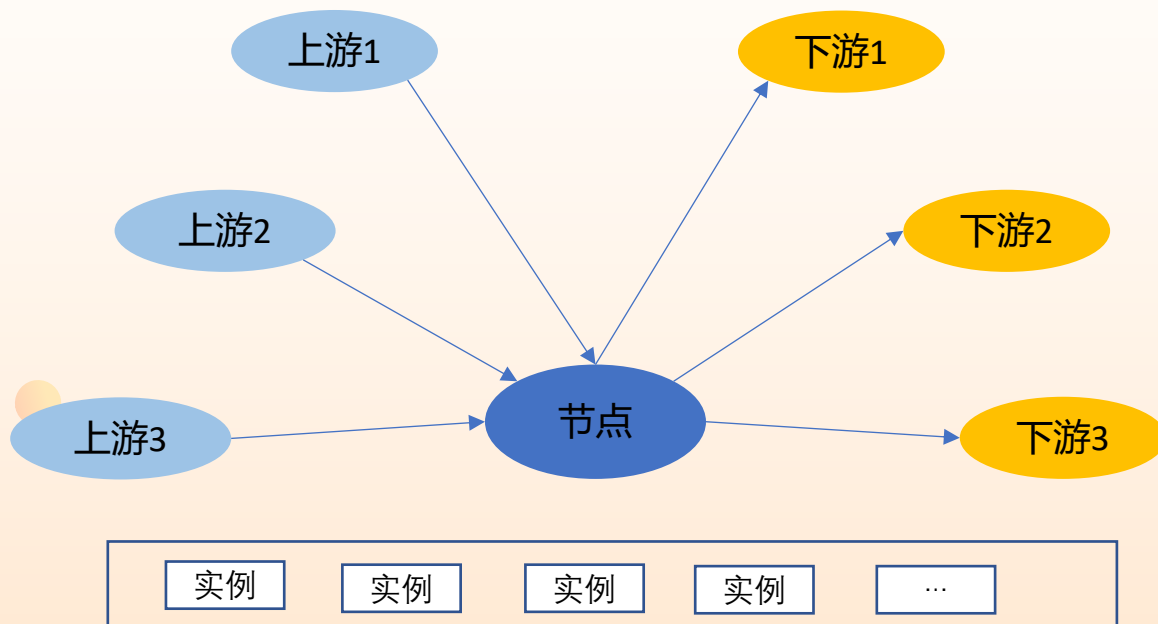
# 异常发现

异常发现，通过节点和关系颜色表达，能够快速地发现特点的节点和关系异常，进一步提升问题发现和定位的效率，通常在应用运行时整体链路梳理和特定问题节点上下游分析等场景使用。



# 关联分析

关联分析，通过关联关系的切换，可以快速查看上游请求和下游依赖，以及自身服务实例的运行情况，进一步提升问题定位能力，通常在已经定位到某个异常节点后使用。

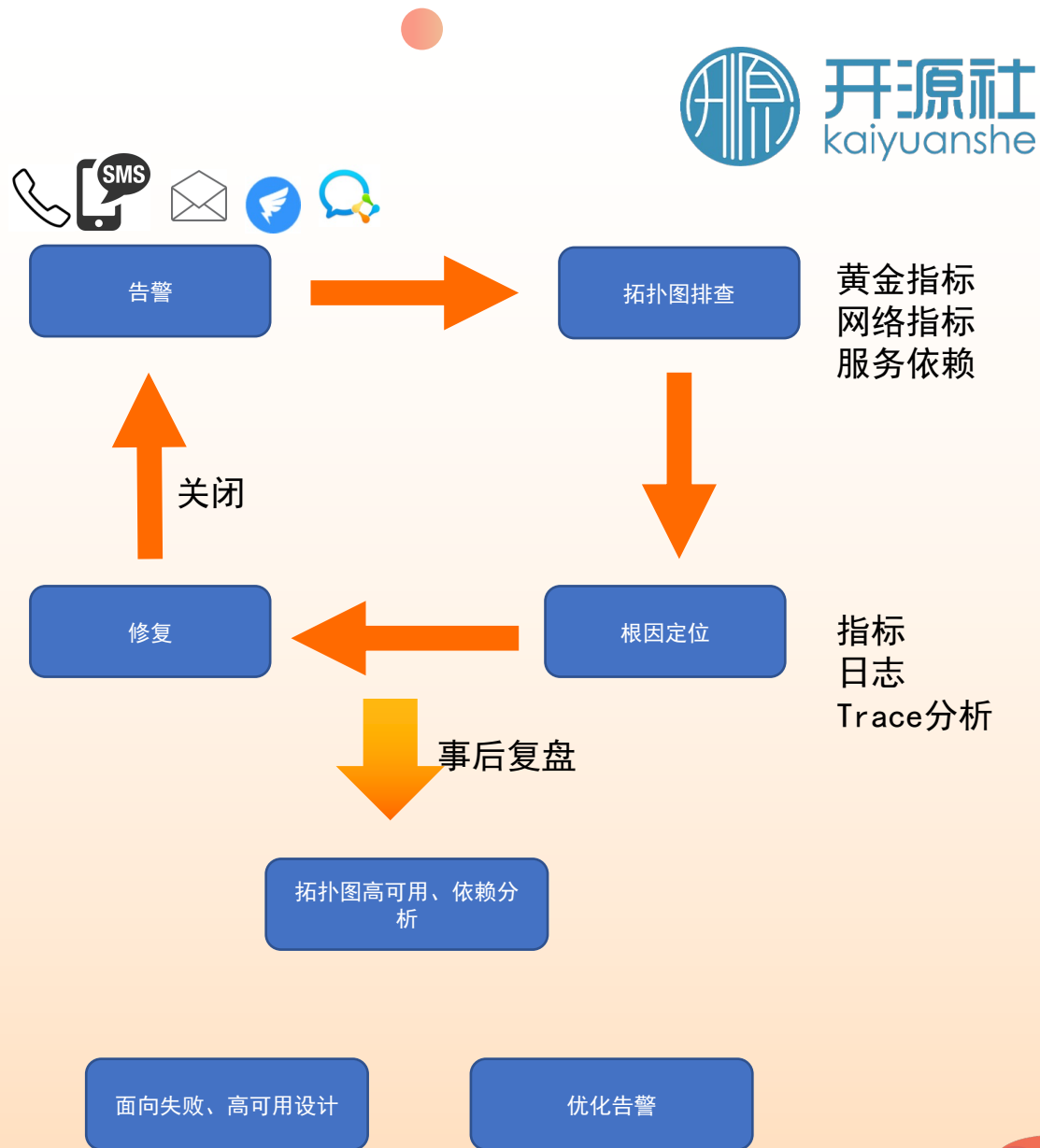


# 智能告警

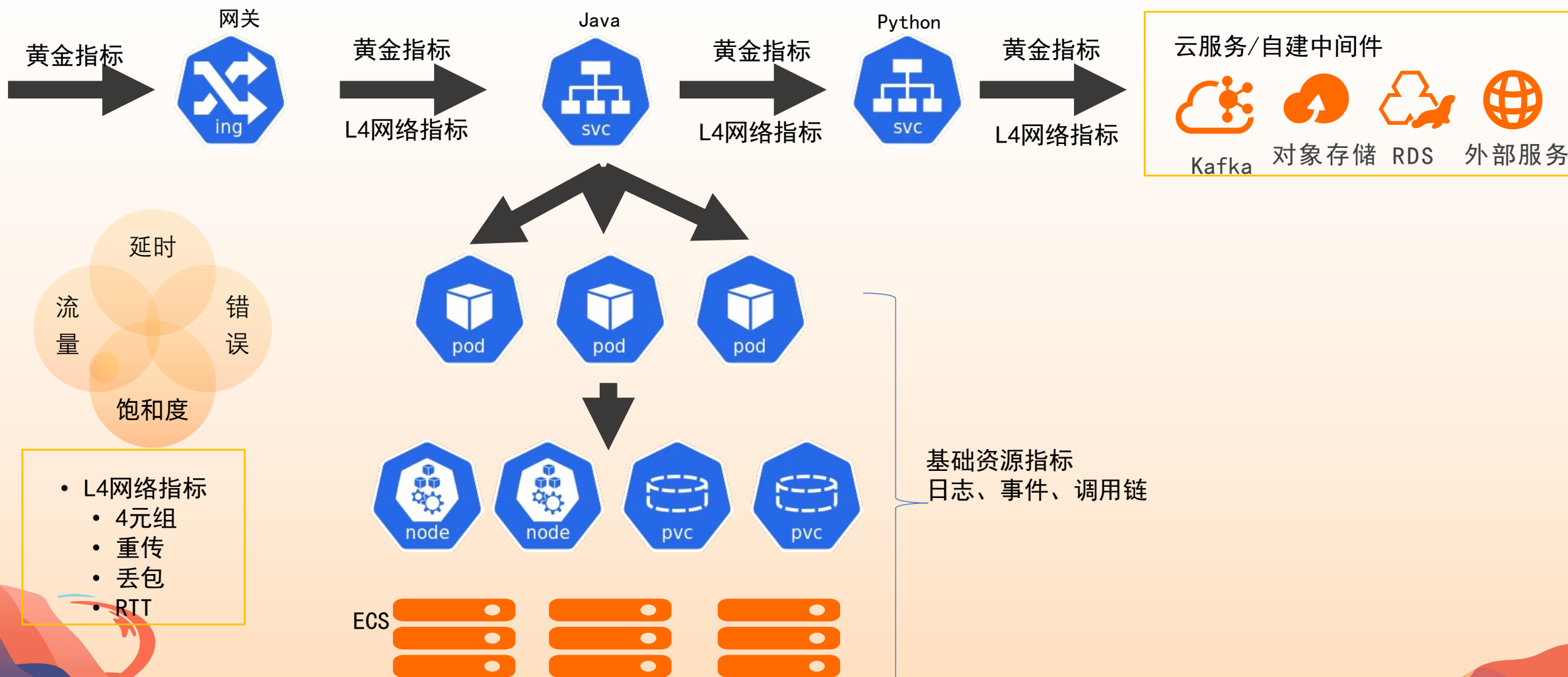
全栈数据源，70+个告警模板开箱即用：  
应用级别：Pod/Service/Deployment  
K8S控制面：apiserver/ETCD/Scheduler  
基础设施：节点、网络、存储  
云服务界别：Kafka/MySQL/Redis/



告警收敛，幸福感UP



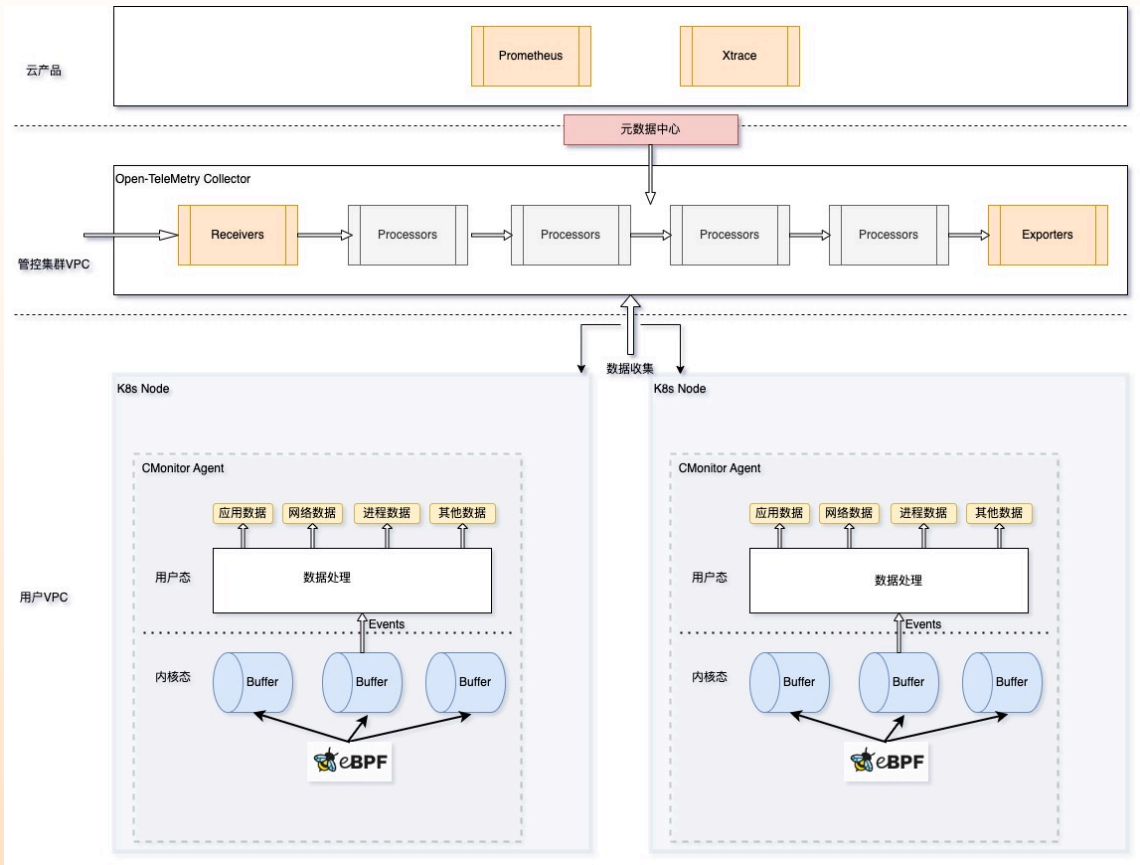
# 完整的应用可观测



- L4网络指标
  - 4元组
  - 重传
  - 丢包
  - RTT



# eBPF在阿里云可观测的实践



# 应用详情概览

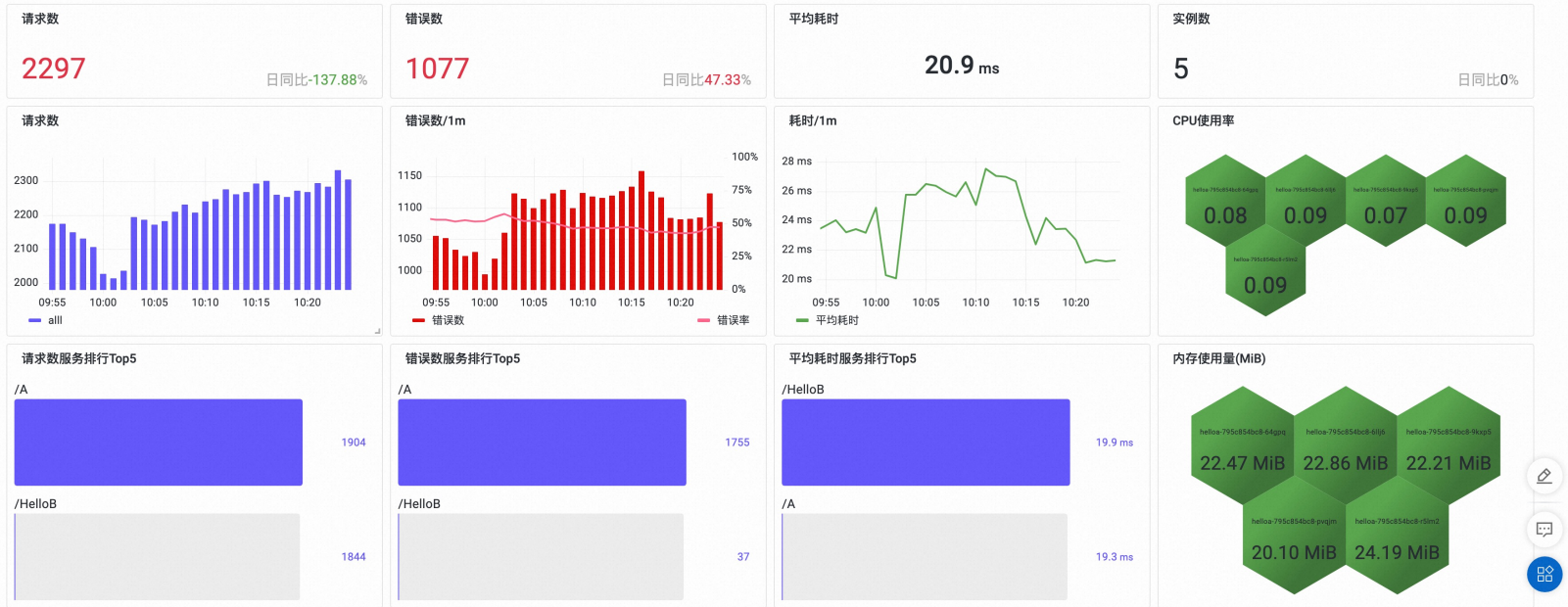
eBPF应用列表 / helloa

helloa

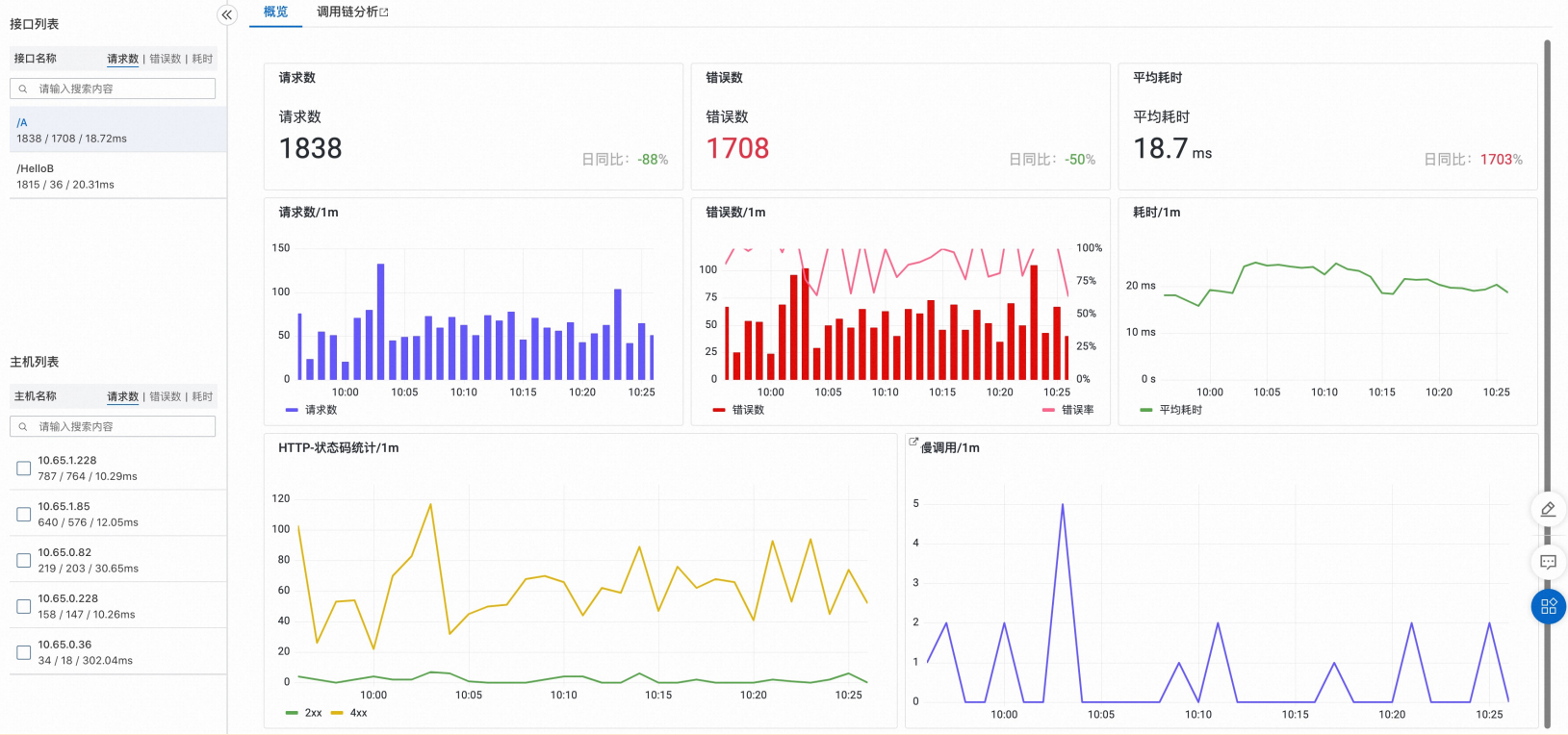
30min 最近30分钟

刷新率: off

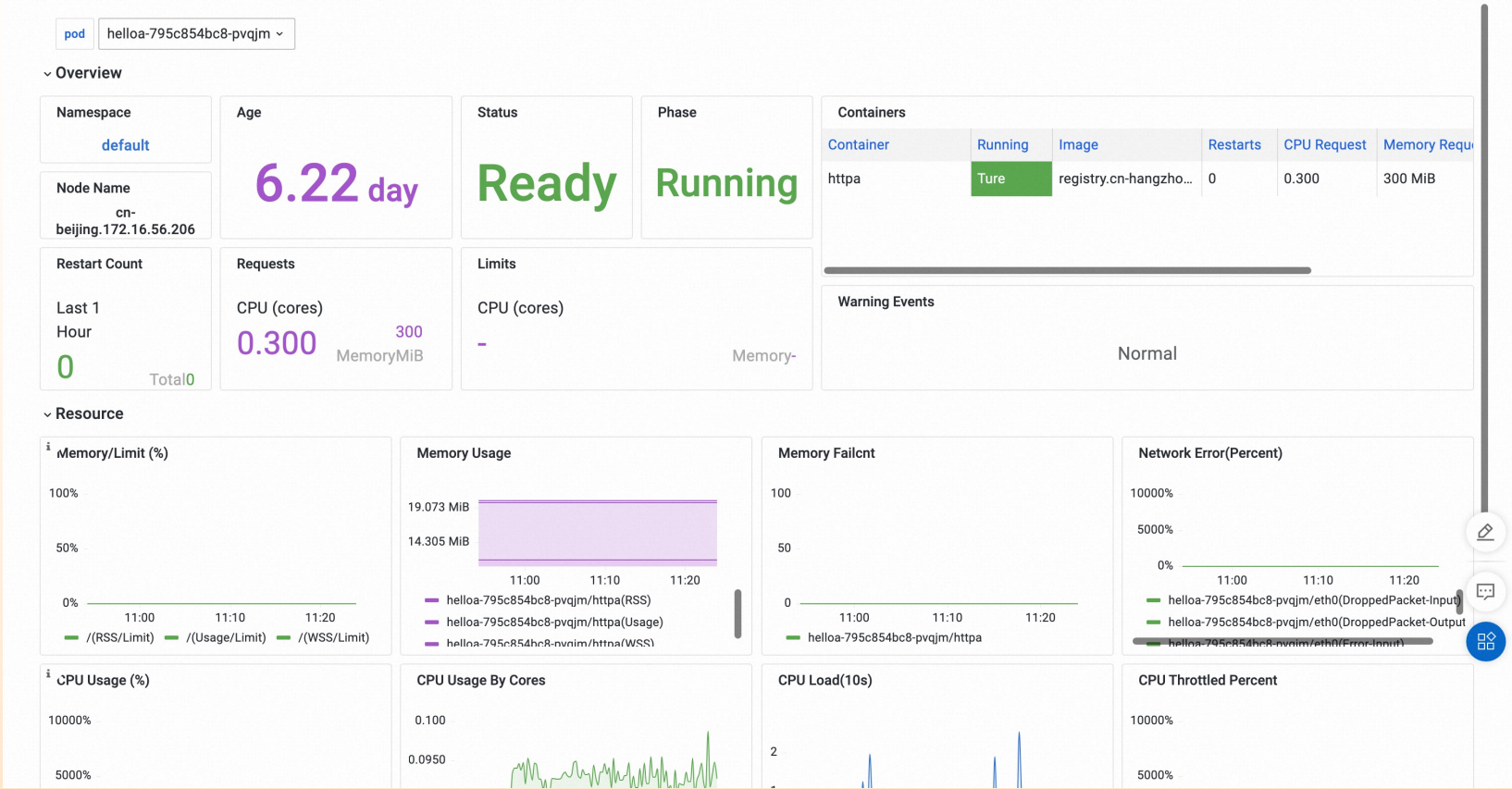
应用概览 应用拓扑 提供服务 依赖服务 调用链分析 实例监控



# 异常接口查询



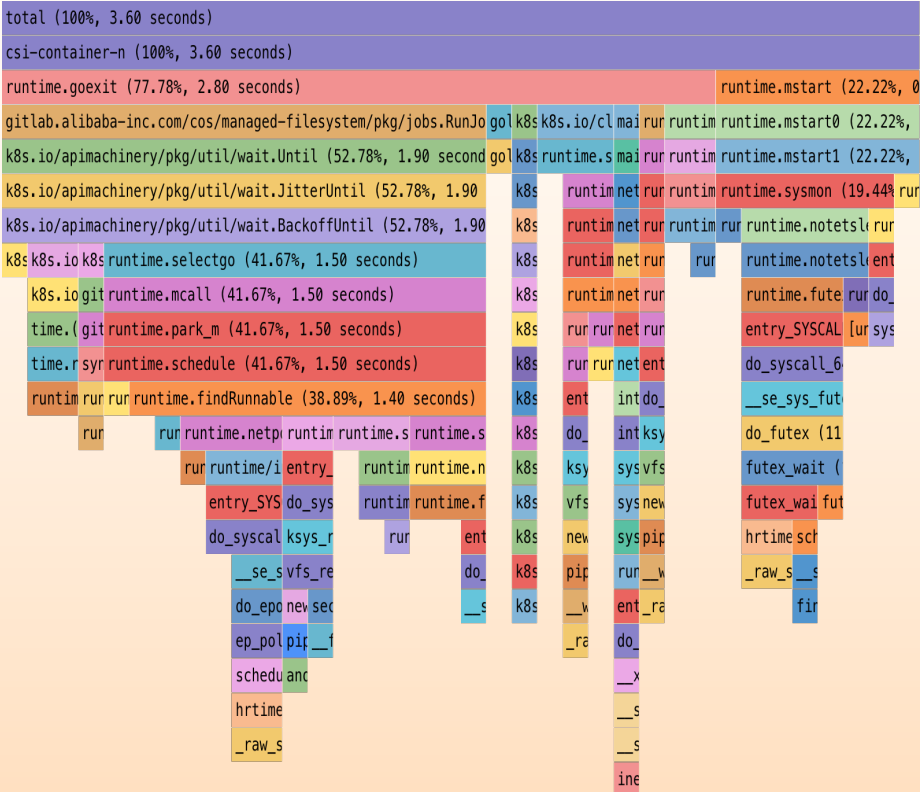
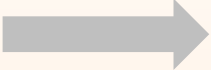
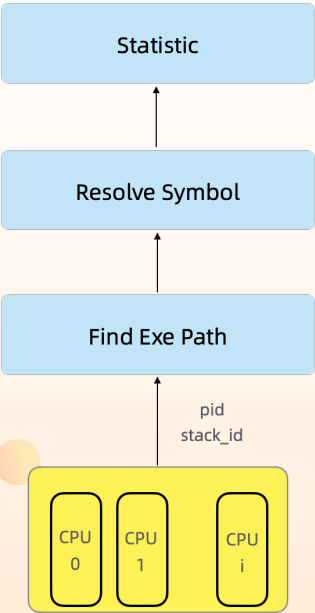
# 应用实例监控





# 持续剖析

提供基于eBPF的应用Continuous Profiling能力，基于微服务应用符号表，自动标记微服务调用栈，轻松获取整个容器集群应用Profiling数据，快速查询应用CPU热点







# Part 04

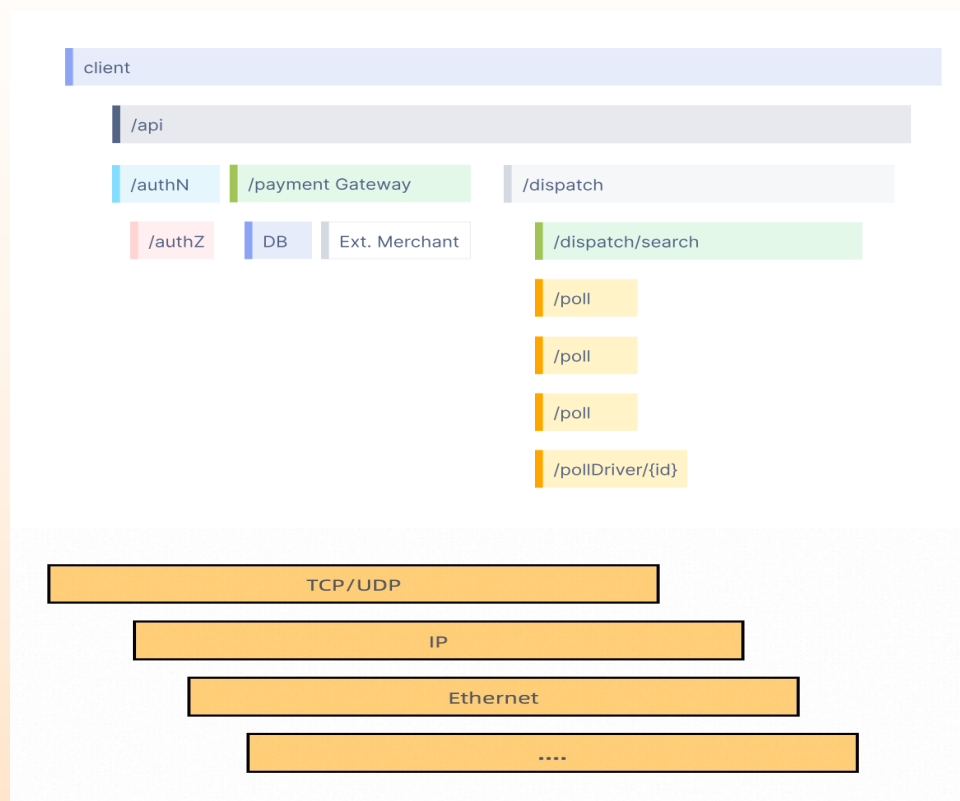
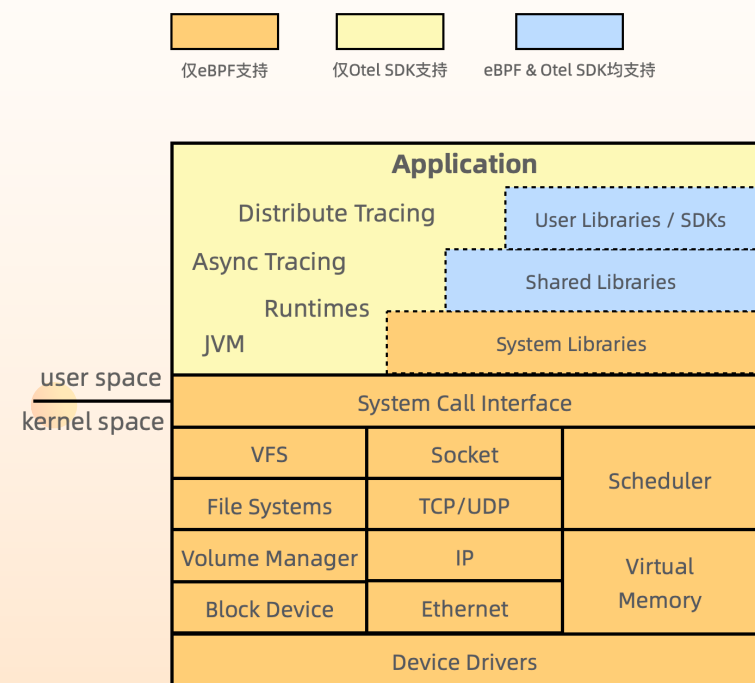
## 未来与展望



# 结合Otel SDK构建全链路Trace能力

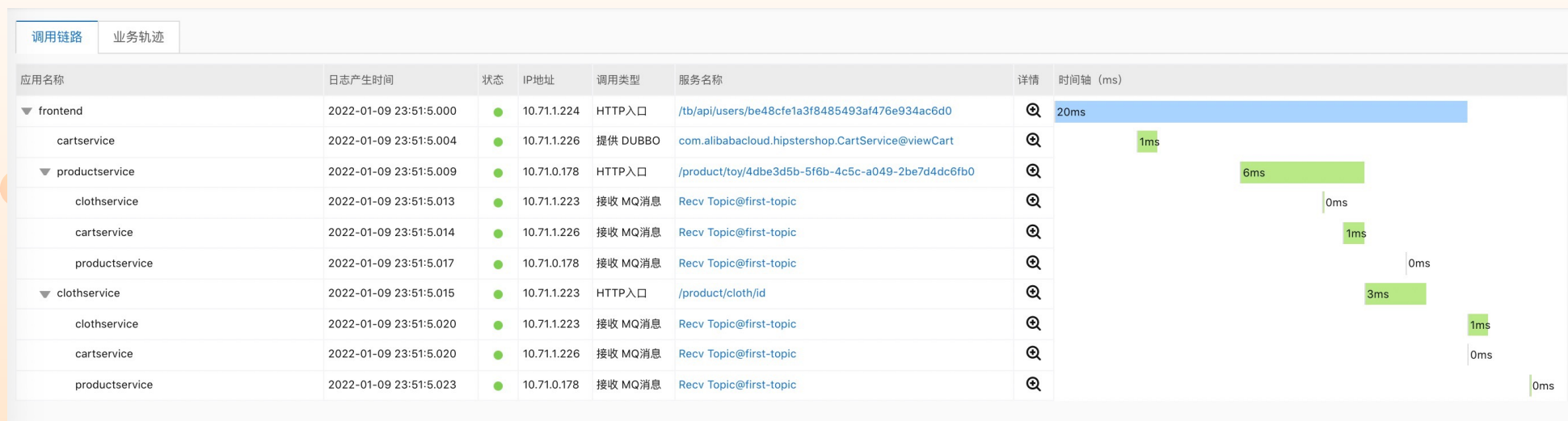


Otel SDK提供应用Trace能力，eBPF提供网络层Trace能力，打通这两种实现全栈的Trace



# 无侵入的Trace Id透传

eBPF因为它本身的无侵入性，在应用内部的Trace Id透传上无法满足要求，特别是异步场景下的Trace Id透传，当前同步场景下的Trace Id透传已经解决，类似golang、java这些高级语言的异步场景还无法满足客户的生产可用



# 拥抱开源、贡献社区



预计今年12月份左右相关工作将贡献给alibaba开源项目KubeSkoop，构建Kubernetete下的全栈可观测能力，尽情期待。

<https://github.com/alibaba/kubeskoop>

# THANK YOU

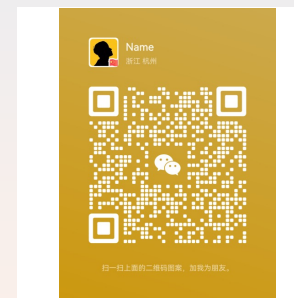
QUESTIONS?



欢迎扫码打卡  
积分可兑换对应礼品哟!



扫码关注开源社公众号



扫码添加讲师联系方式

微信公众号：开源社KAIYUANSHE

视频号：开源社KAIYUANSHE

新浪微博：开源社

B站：开源社KAIYUANSHE

简书：开源社

头条：开源社

Facebook: KaiyuansheChina

Twitter: 开源社KAIYUANSHE