

围绕gala-gopher构建的eBPF云原生网络可观测性

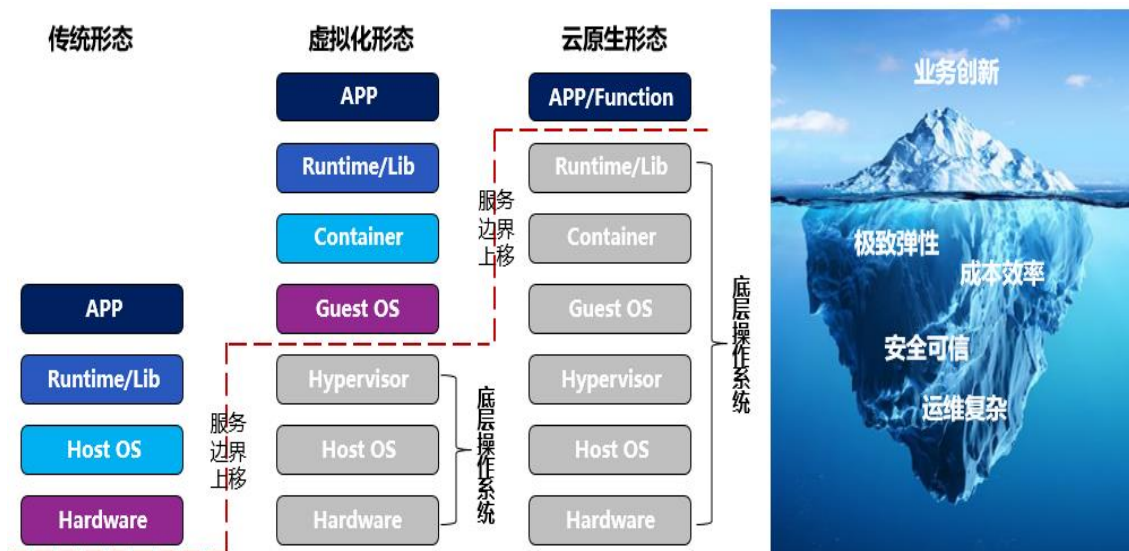
网络请求分布式链路追踪能力

汇报人：孙景浩 华为高级工程师

openEuler运维面临挑战

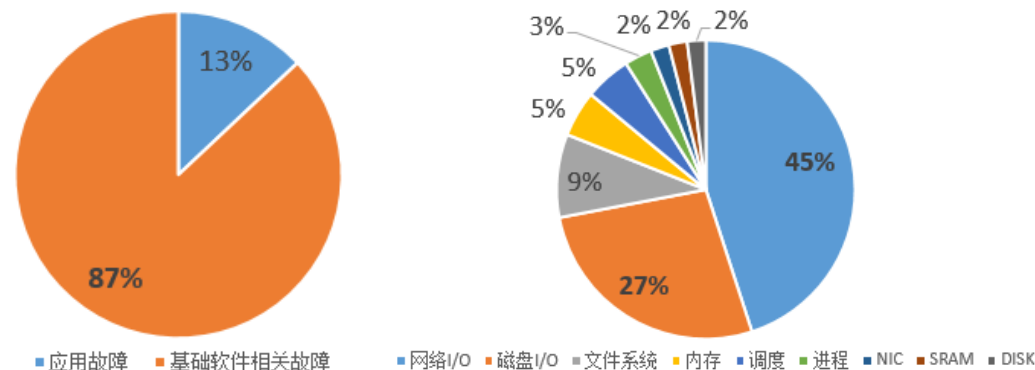
01

云基础设施的复杂性更易产生**隐性故障**，如何更快定位根因是云场景故障诊断的重要挑战。



02

XXX云故障87%定位与OS相关（**自身问题 < 5%**），**平均定时时长 > 10H**，影响客户体验

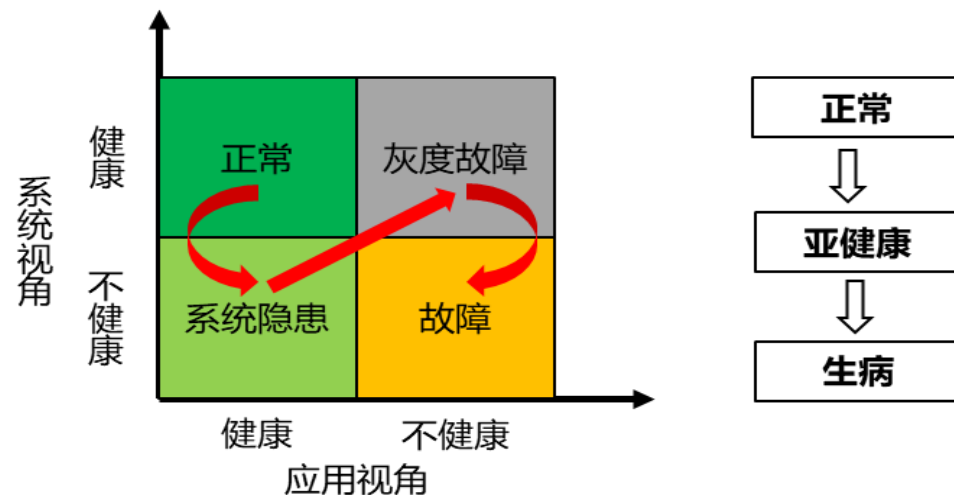


基础软件占比87%

网络，磁盘I/O分类占比最高

03

灰度故障正在成为操作系统的致命弱点，openEuler缺少灰度故障监测、定界、定位和自愈的高级运维能力。



openEuler操作系统全栈观测能力（gala-gopher）

- 基于**非侵入观测技术**（eBPF + Java agent等）技术立足操作系统，提供包括云原生在内场景的可观测能力。
- 围绕系统行为可观测、业务全流程观测、应用行为可观测、性能观测四个能力方向，构建**基础设施监控、应用性能监控、应用安全、自动化及监控四大解决方案**。

业务全流程跟踪

- ✓ 提供云原生业务全流程跟踪能力，包括分布式业务流Trace，DNS访问过程。
- ✓ 云原生网络监控等能力，实现应用、基础设施之间的网络问题定界能力。

系统行为观测

- ✓ 提供系统关键事件、性能事件、错误信息观测能力，包括Linux、Container、K8S基础软件范围。
- ✓ 提供集群动态拓扑能力，包括Service Map、集群资源状态、应用视角软件栈拓扑等。



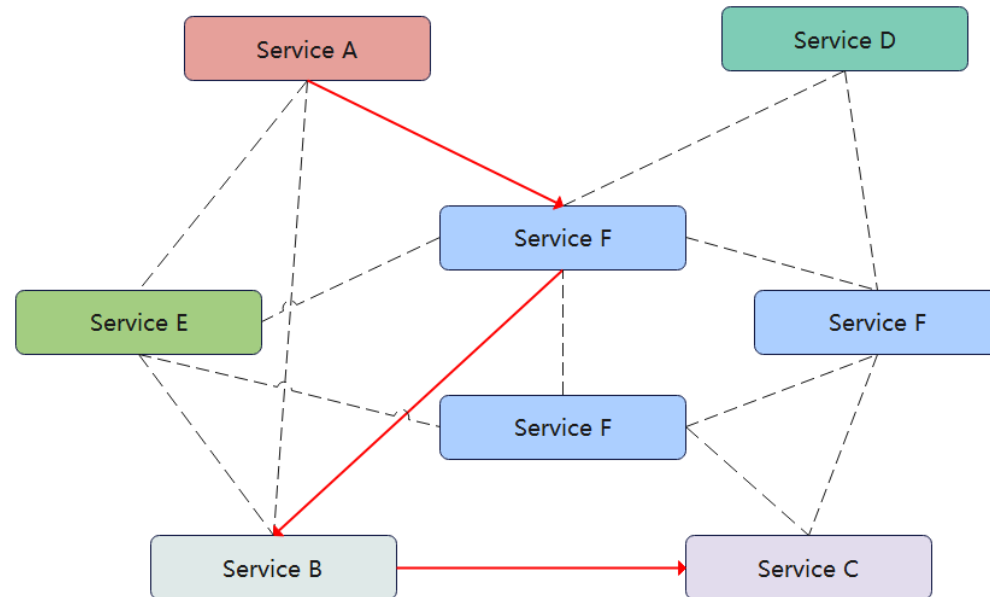
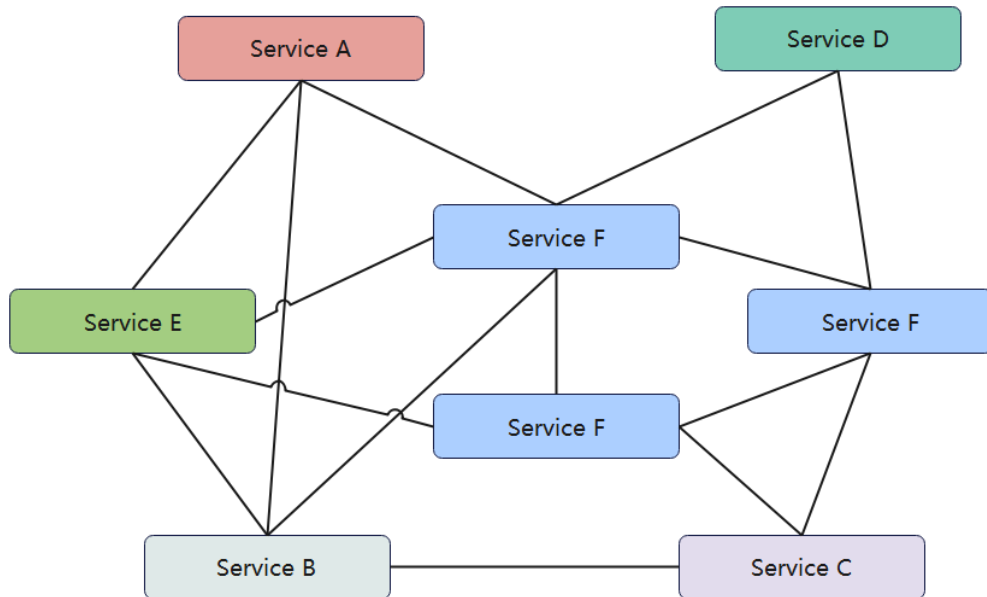
应用行为观测

- ✓ 提供应用视角下钻式观测能力（包括应用粒度的CPU&内存、TCP/IP数据，I/O数据等）。
- ✓ 应用性能观测能力（覆盖云原生常见协议gRPC/HTTP/RPC等），基础中间件观测能力（包括redis、kafka、DNS、PG等）。

性能观测

- ✓ 提供全栈性能热点分析能力（包括CPU、内存、I/O、Lock等）。
- ✓ 应用微基准性能分析能力（包括网络访问、文件操作、锁操作等系统性能事件），
- ✓ 资源占用区位分析能力（包括CPU、内存等）。

网络静态拓扑VS 动态链路追踪



网络静态拓扑

依赖五元组实现

关注整体网络拓扑，结构相对稳定

仅采集数据，包括应用粒度的CPU&内存、TCP/IP数据，I/O数据等

网络故障时仅能定界，无法定位

动态链路追踪

依赖五元组作为基础数据，需要补齐其他上下文信息以对整条链路进行追踪

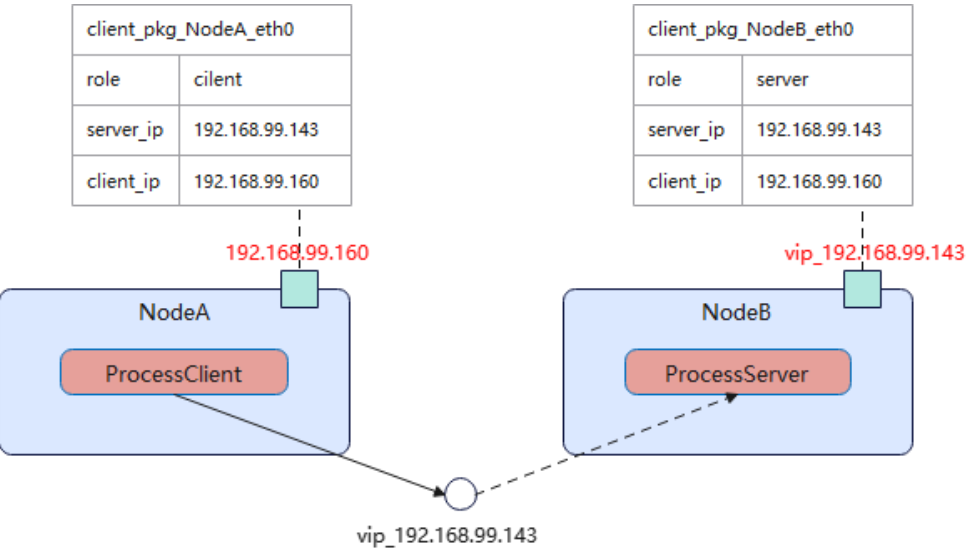
关注单个网络请求，单次追踪无法获取整体网络拓扑结构，链路变化较大

关注数据间联系

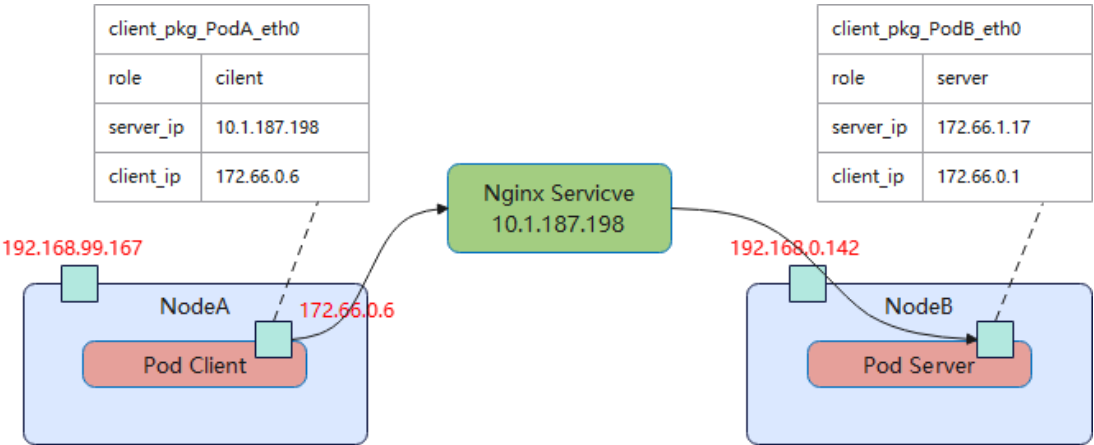
网络故障时可定位到服务级别

云原生场景下分布式链路追踪的挑战

云原生场景中在容器网络下如何获取pod真实IP

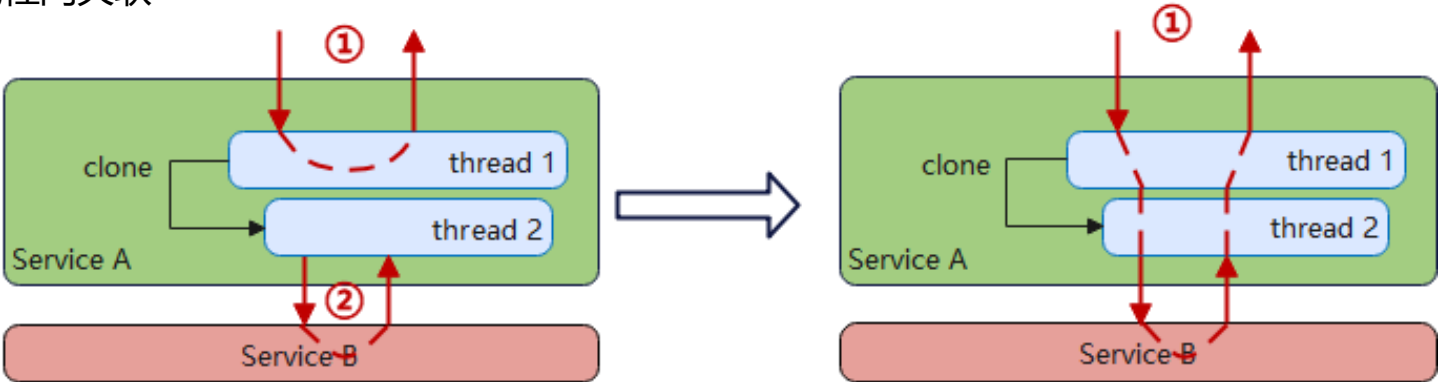


如果存在浮动IP，则采集的IP数据为浮动IP，无法直接关联到对应主机



容器化场景中，通过kubernetes service访问，主对端地址不一致无法匹配

多线程异步场景如何获取线程间关联



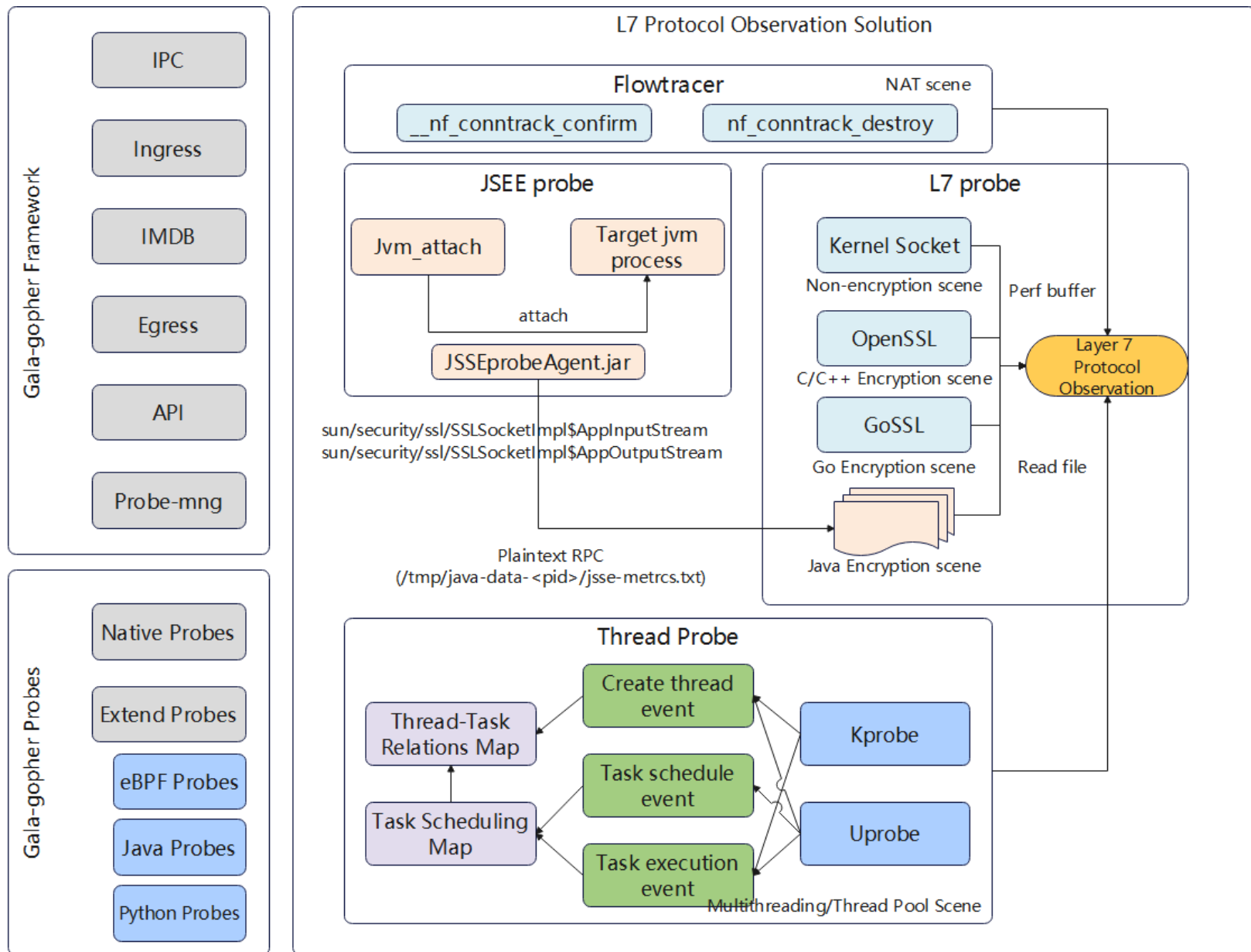
基于gala-gopher的云原生网络可观测性增强

L7 Probe: 7层协议基础观测能力增强

Flowtracer: NAT转换场景增强，如k8s容器网络中，将采集到的TCP link的IP转换为k8s容器网络对应的pod IP。

Thread Probe: 多线程/线程池场景功能增强，提供多线程/线程池相关上下文信息，以支持该场景下链路追踪功能。

JSEE Probe: Java加密场景增强



分布式链路追踪能力实现

1、数据采集

基于eBPF采集TCP流数据

2、协议判断

内核态对tcp流量进行简单分析，判断所属应用层协议，并将数据发送至用户态进行进一步解析。

3、协议解析

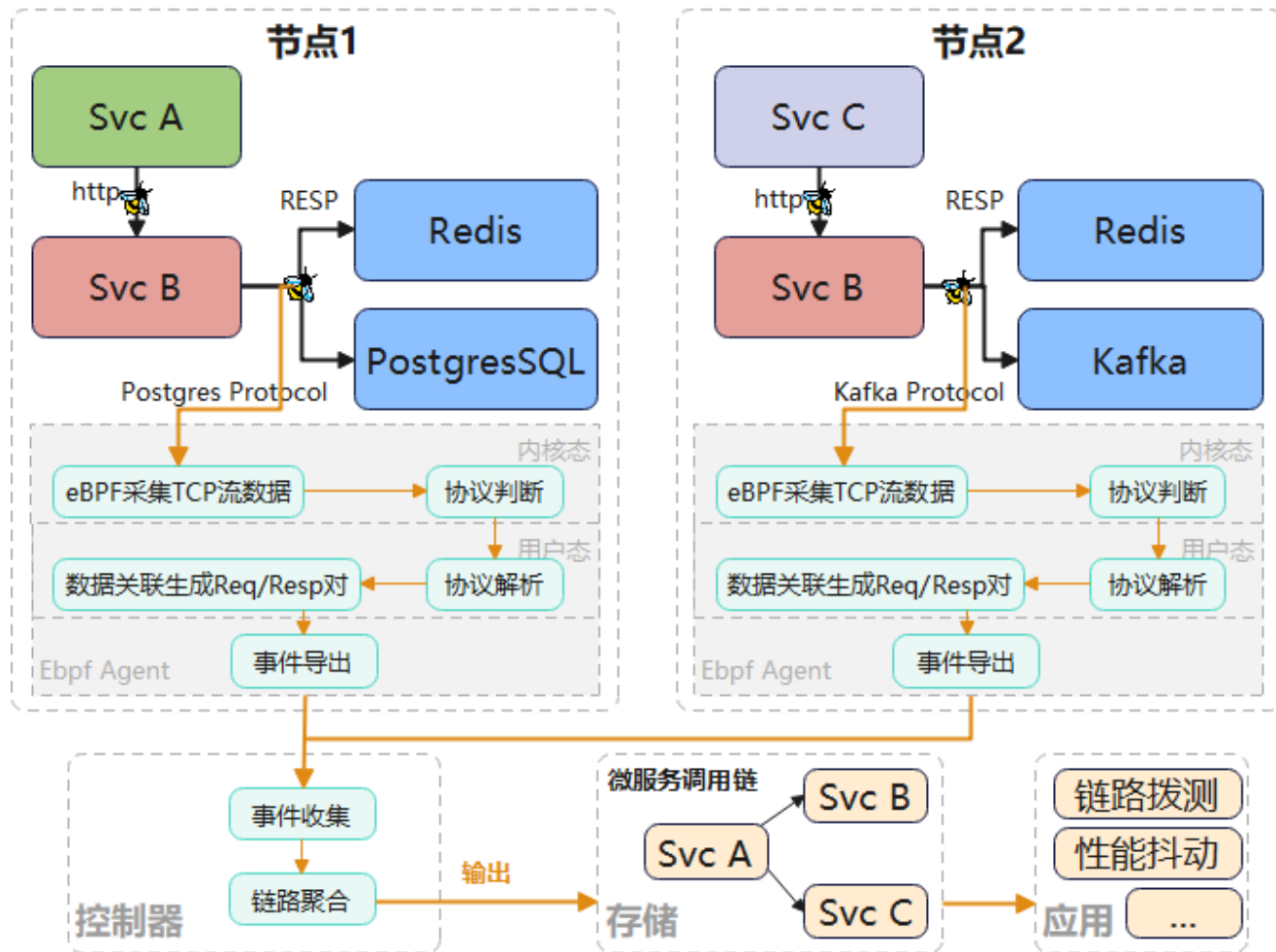
用户态程序对接受到的tcp流量信息进行解析，根据不同协议提取关键信息并存储。

4、数据匹配

对步骤3中存储的信息进行匹配，匹配基于**时间戳**和**Thread ID**和**协议报文数据**，将网络请求Req/Resp进行关联重建，并上传至Prometheus。

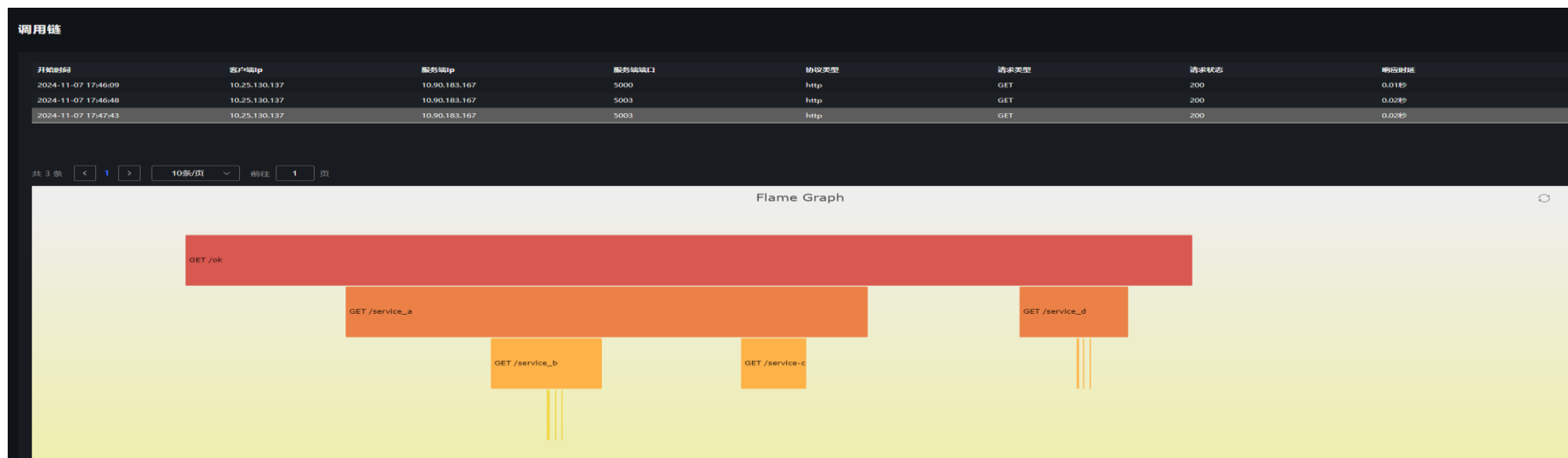
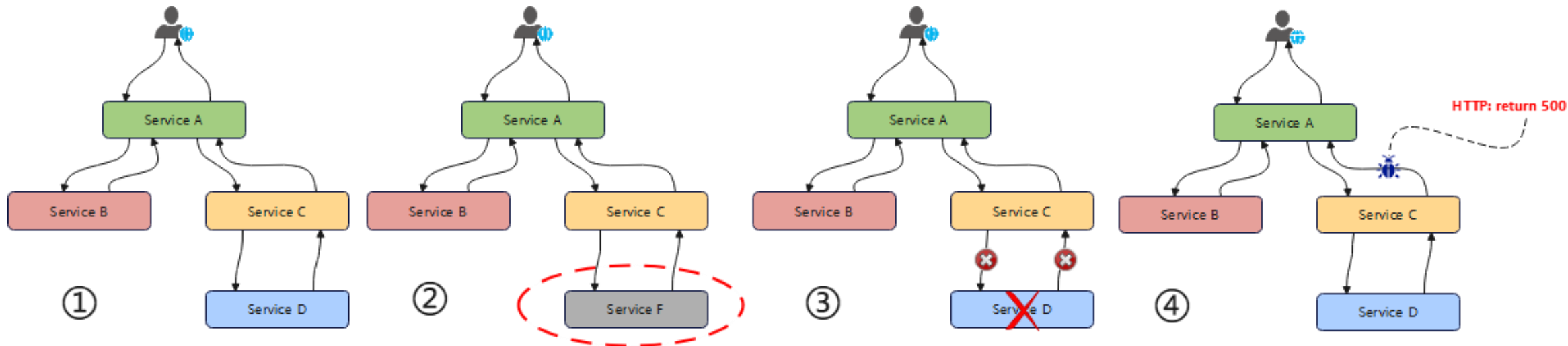
5、链路聚合

基于步骤4中的网络请求Req/Resp对，通过TCP的序列号进行聚合，得到一条微服务调用链。

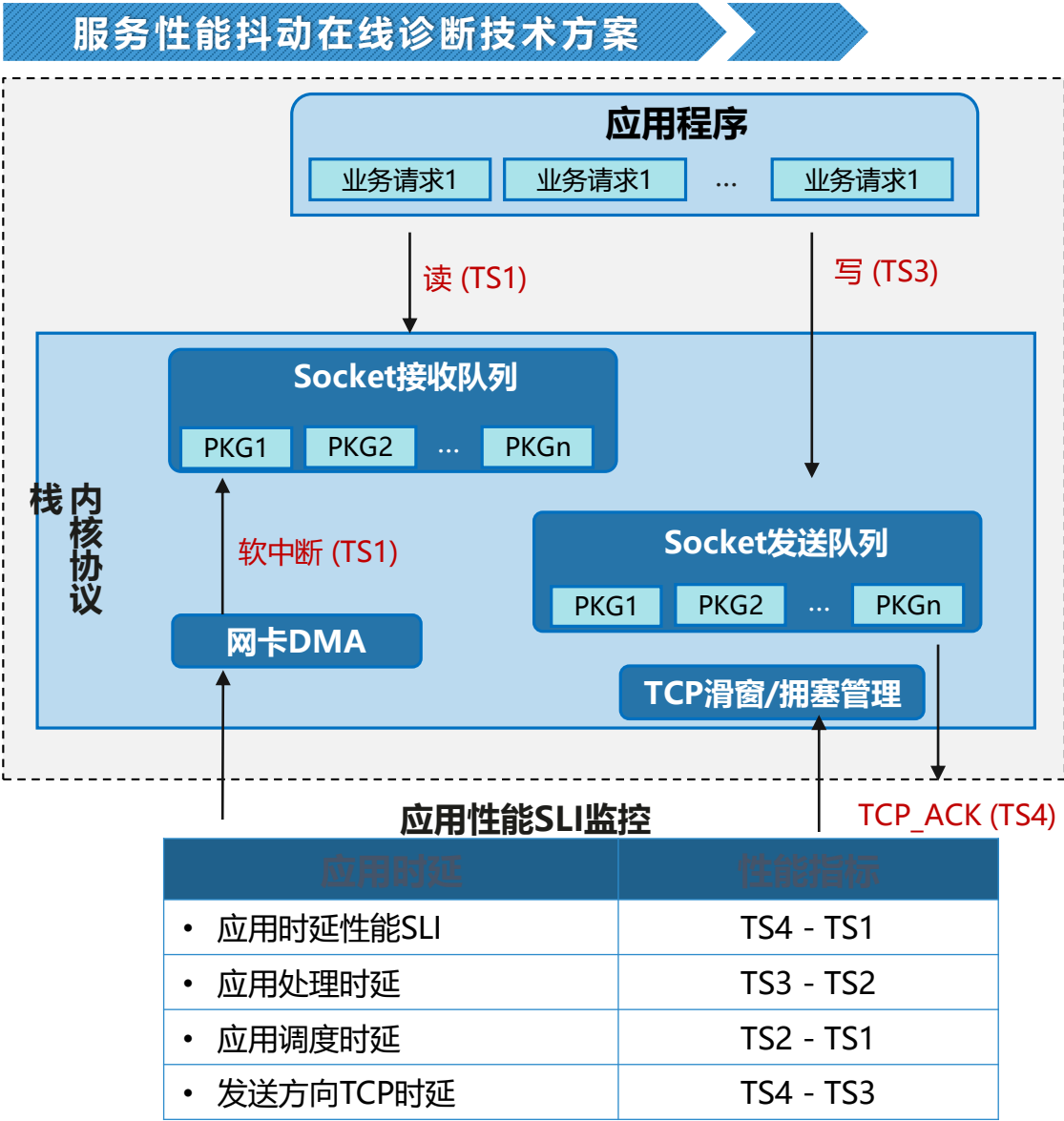
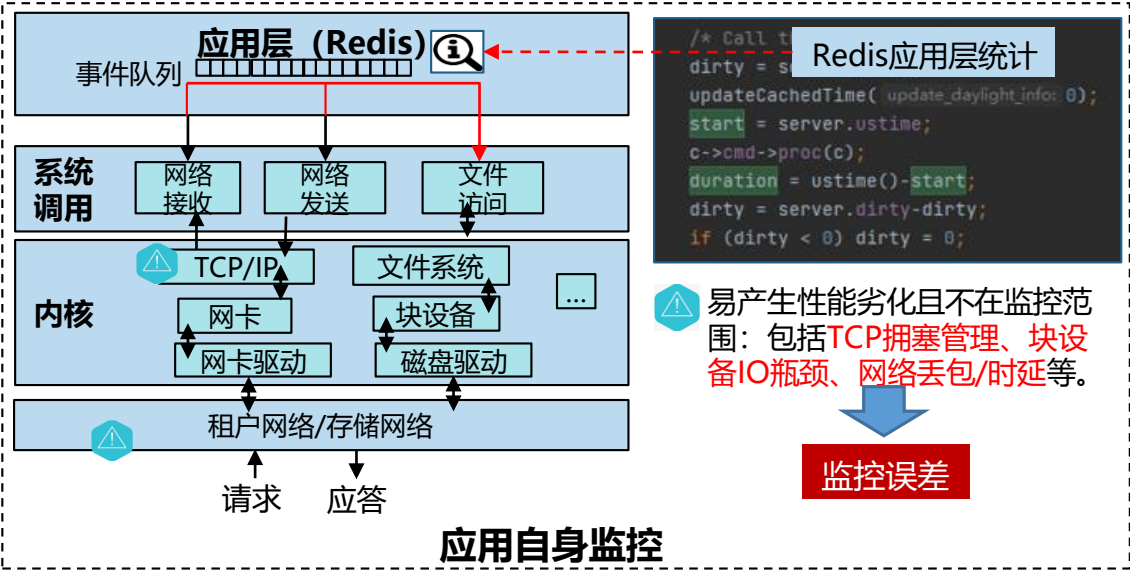
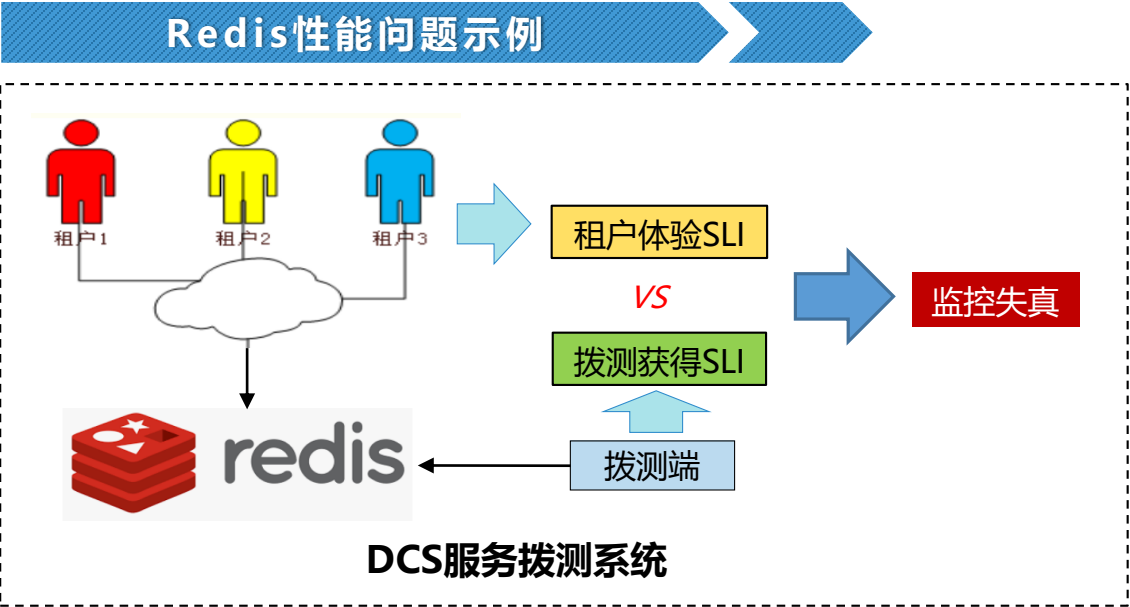


分布式链路追踪能力应用效果

调用链时空分析

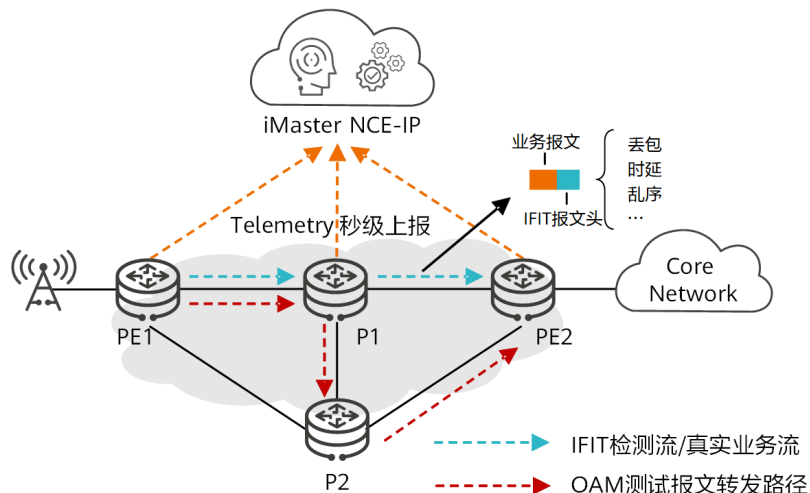


分布式链路追踪能力相关应用



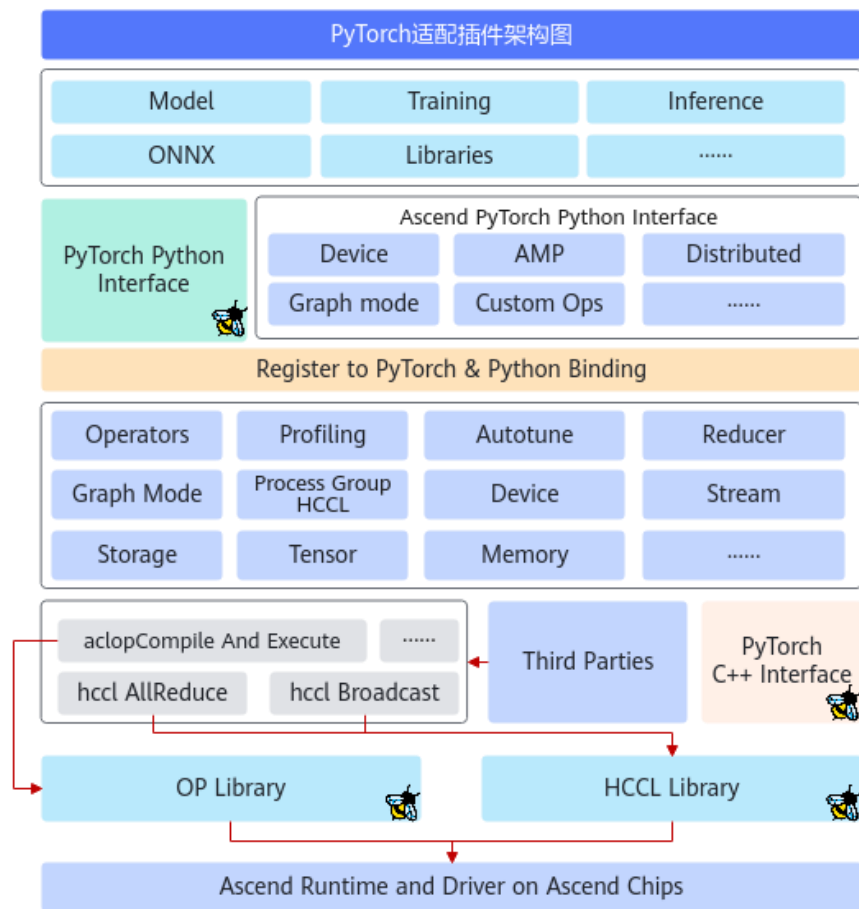
未来展望

支持采集IFIT随流检测协议



IFIT 是一种通过在业务报文中插入 IFIT 报文头标记特征，直接检测网络时延、丢包、抖动等性能指标的 IETF 标准化检测协议。通过IFIT数据获取网络流转发路径，弥补网络流转发盲区，实现报文端到端追踪。

基于eBPF uprobe的CANN/CUDA监控



基于eBPF Uprobe 可以在大模型训推全流程进行打点监控，以获取PyTorch三方库，HCCL集合通信库等重要方法的相关执行信息。

THANKS