



第二届 eBPF开发者大会

www.ebpftravel.com

eBPF交流研讨

中国·西安



第二届 eBPF开发者大会

www.ebpftravel.com

gala-gopher: openEuler基于 eBPF的全栈可观测方案及其实践

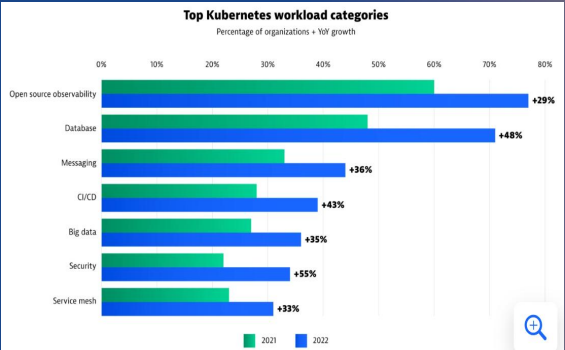
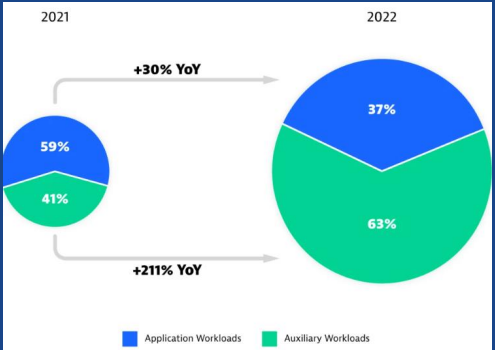
中国·西安

① 背景介绍

- Gartner将**应用可观测性**作为十大技术趋势之一，长期看好可观测性驱动**企业运营最佳决策**。
- 同时指出云原生场景中基础设施观测能力不足，为下一代**云原生可观测**提供机会与挑战。

Top Strategic Technology Trends for 2023		
Optimize	Scale	Pioneer
<ul style="list-style-type: none"> Digital Immune System Applied Observability AI Trust, Risk and Security Management 	<ul style="list-style-type: none"> Industry Cloud Platforms Platform Engineering Wireless Value Realization 	<ul style="list-style-type: none"> Superapps Adaptive AI Metaverse

- 2023年云原生报告：云原生集群内辅助类应用工作负载上升至63%，其中**近80%**的企业部署可观测性方案（同比**增长29%**）
- 变化背后意义：企业在云原生技术实施过程中逐渐**意识到可观测性的重要性**



云原生给可观测带来的变化与挑战



- 变化1**：虚拟化单一架构中“一刀切”分层运维（基础设施、应用分层）向云原生场景融合式运维发展，需要**提供全栈观测、运维能力**；
- 变化2**：云原生多技术体系（Linux、CNCF等）、快速演进等特点，要求可观测性解决方案与其应用/基础设施技术栈解耦，提供**非侵入观测能力**；
- 变化3**：云原生高密度、分布式部署方式，要求具备**集群运维**视角，从业务集群视角逐层/级定界、定位至具体问题根因；

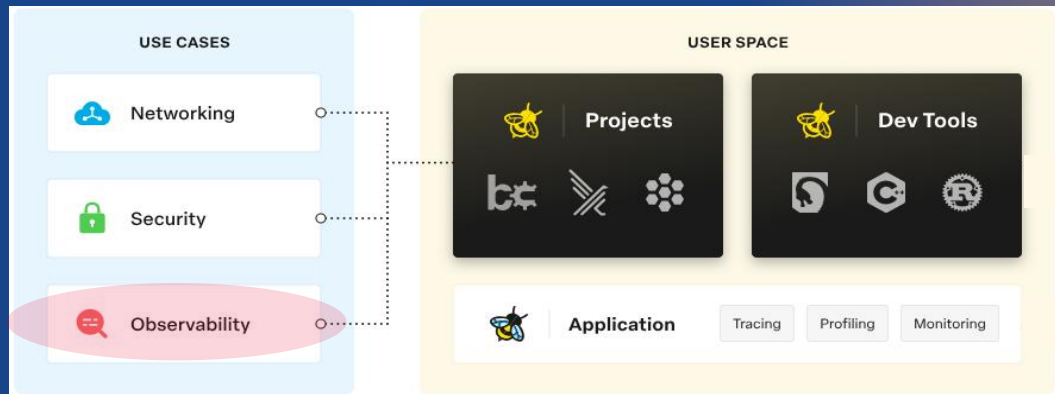
② 技术洞察：eBPF已成为新一代全栈观测技术趋势

为什么选择eBPF

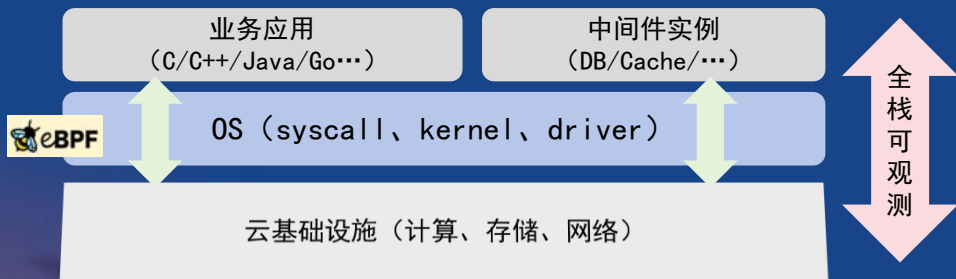
eBPF 是一个能够在**内核运行沙箱程序**的技术，通过安全注入代码的机制，使得安全的访问、控制内核状态、行为，主流应用场景有**可观测、安全、网络**。

为什么云原生场景适合eBPF可观测能力实施：

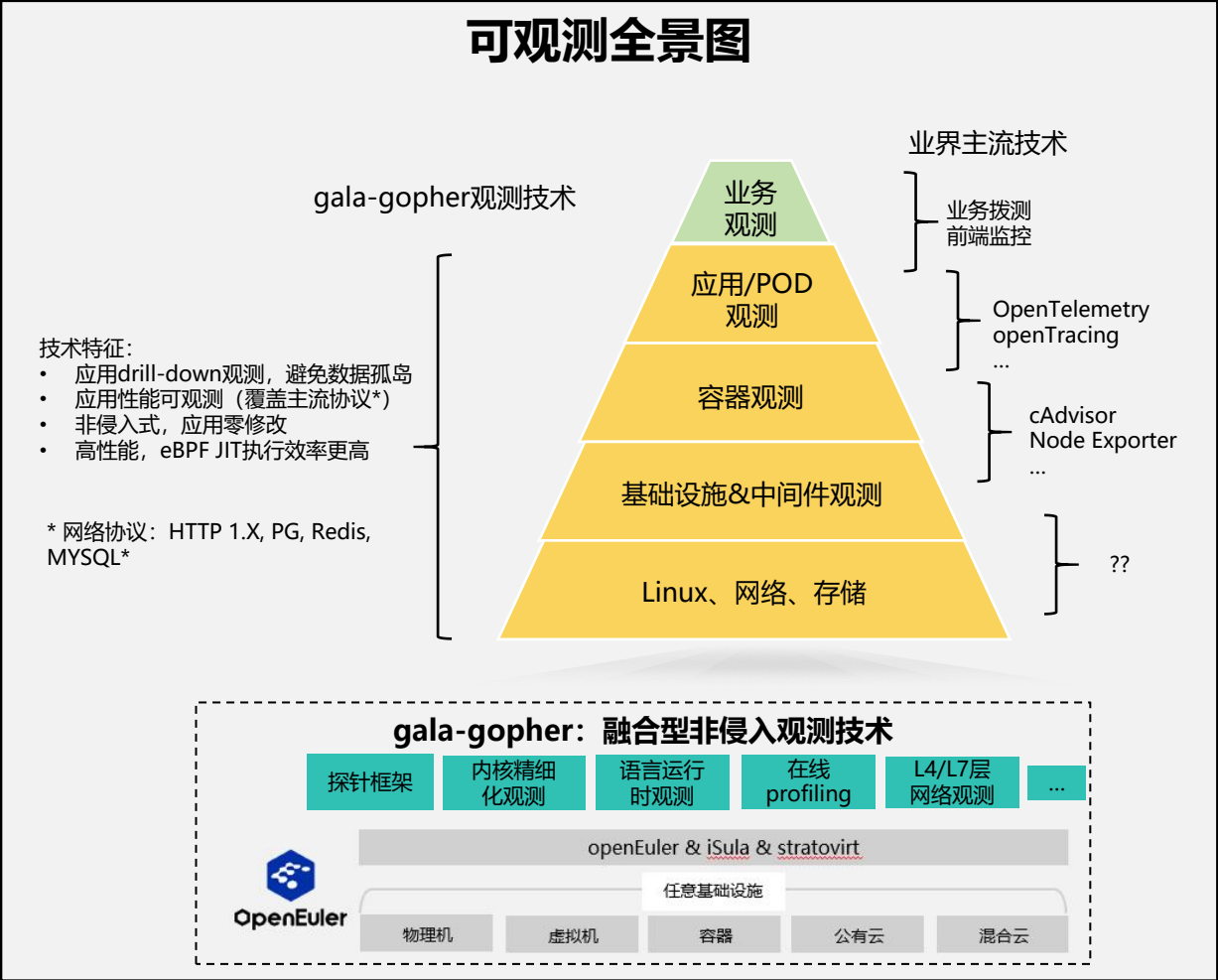
- **无侵入**：通过eBPF字节注入技术可以快速的进行可编程无侵入式观测逻辑注入，轻松应对云原生场景快速迭代的场景特征。
- **可移植&跨平台**：通过标准eBPF ISA、CO-RE等技术可以自适应适配云原生集群内不同Linux版本、不同ISA架构平台场景。
- **全栈**：通过eBPF + USDT、eBPF + Tracepoint、eBPF + kprobe等技术，可以覆盖内核、运行时、基础库等大部分基础软件，轻松应对云原生多语言、多网络协议、厚重软件栈的场景特征。



OS提供应用视角的全栈观测能力



③ openEuler gala-gopher整体介绍

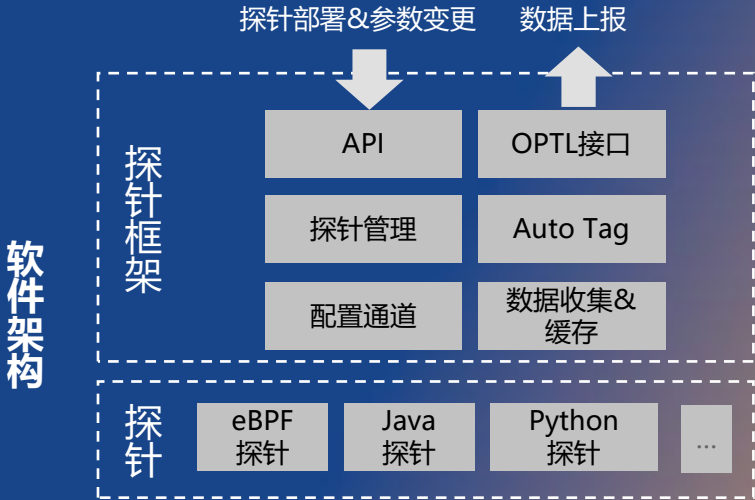


软件开销: 观测底噪<5%, 应用性能干扰<3%。

关键能力

- **基础设施**: I/O时延、错误率、I/O分类统计*、进程I/O、OOM*;
- **网络**: 进程TCP流量、进程TCP建链、进程TCP状态;
- **应用性能**: L7网络RED性能 (HTTP(S), Redis, PGSQL, MySQL*等)
- **性能Profiling**: OS runtime Profiling, 系统关键事件Profiling, 系统关键资源Profiling*;

备注: *表示暂未开源



④ openEuler可观测能力介绍



技术能力

集群中微服务间调用访问PDB（无业务时延，供

误率

支持

PGS

支持

1.1.0

快速

使用

使用

使用

使用

使用

使用

使用

使用

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

技术能力

• 应用维度的TCP性能观测：提供TCP窗口、RTT、

SR

技术能力

• 应用

基于提供 L4层网络流、负载分担流、L7层网络流、

软件部署等信息，构建系统2D 拓扑

技术能力

• I/O性能：提供进程维度的 I/O操作字节数统计

Block性能指标，磁盘指标，包括磁盘读写速率、

使用率、吞吐量等指标，以及block层驱动、设备

的时延、错误统计

•

•

使用

使用

使用

使用

使用

使用

使用

使用

使用

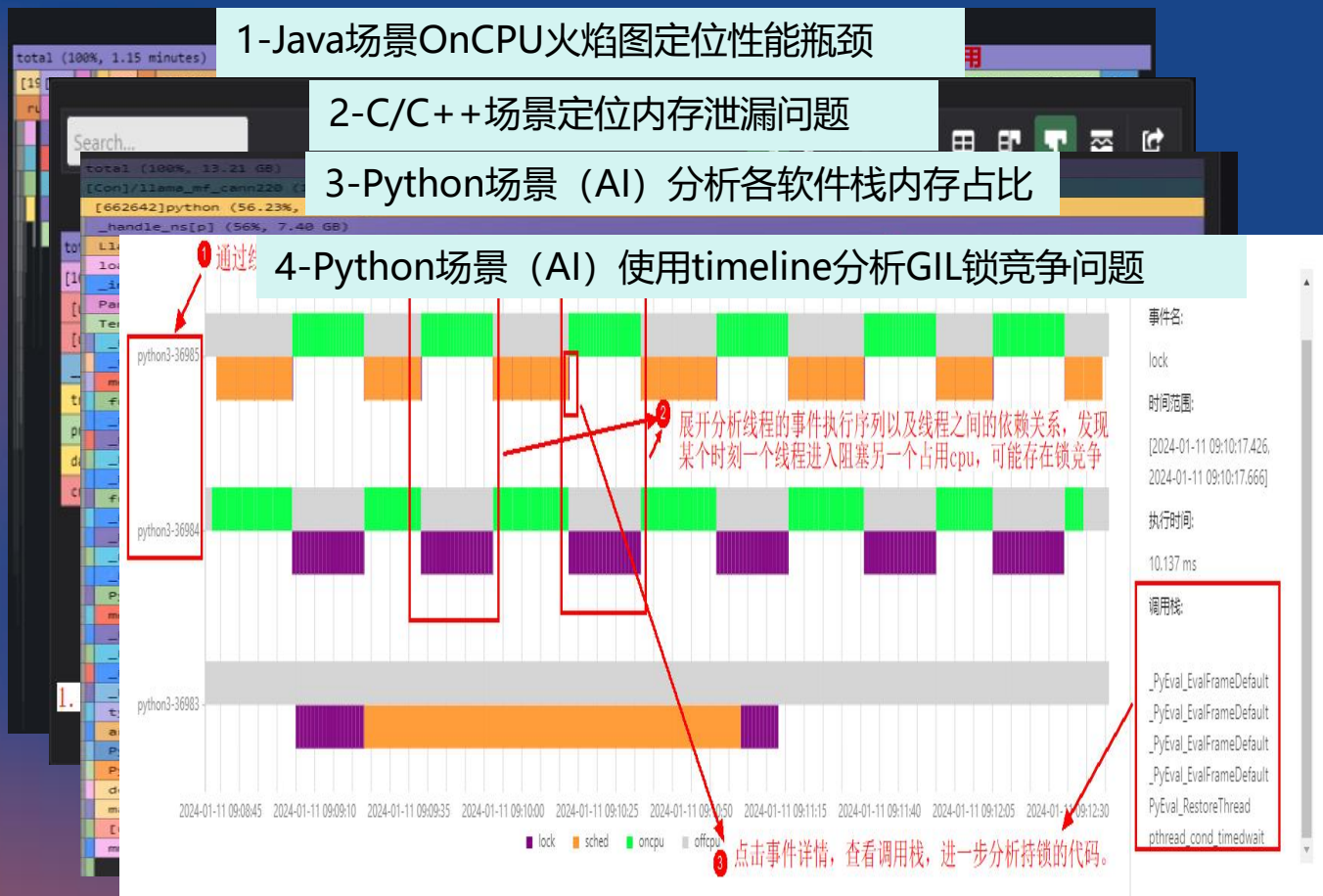
⑤ openEuler 可观测能力介绍

1-Java场景OnCPU火焰图定位性能瓶颈

2-C/C++场景定位内存泄漏问题

3-Python场景 (AI) 分析各软件栈内存占比

4-Python场景 (AI) 使用timeline分析GIL锁竞争问题



技术能力

技术能力

ebpf 技术观测线程的关键系统性能事件, 并关联丰富的事件内容, 从而实时地记录线程的运行状态和关键行为, 并在前端界面以时间线的方式进行展示, 支持观测的线程事件:

- 文件操作 (file)
 - read/write: 读写磁盘文件或网络, 可能会耗时、阻塞。
 - sync/fsync: 对文件进行同步刷盘操作, 完成前线程会阻塞。
- 网络操作 (net)
 - send/recv: 读写网络, 可能会耗时、阻塞。
- 锁操作 (lock)
 - futex: 用户态锁实现相关的系统调用, 触发 futex 往往意味出现锁竞争, 线程可能进入阻塞状态。
- 调度操作 (sched): 这里泛指那些可能会引起线程状态变化的系统调用事件, 如线程让出 cpu、睡眠、或等待其他线程等。
 - nanosleep: 线程进入睡眠状态。
 - epoll_wait: 等待 I/O 事件到达, 事件到达之前线程会阻塞。

使用场景

代码级别定位线程间由于资源竞争导致的性能问题。例如:

- 文件 I/O 耗时、阻塞问题
- 网络 I/O 耗时、阻塞问题
- 锁竞争问题
- 死锁问题

⑥ 版本发布节奏&规划

eBPF全栈可观测：应用级下钻全栈观测能力，提供应用协议性能、应用粒度的网络、I/O、CPU、MEM观测能力

容器场景eBPF全栈观测，应用/OS/容器网络分钟级定界

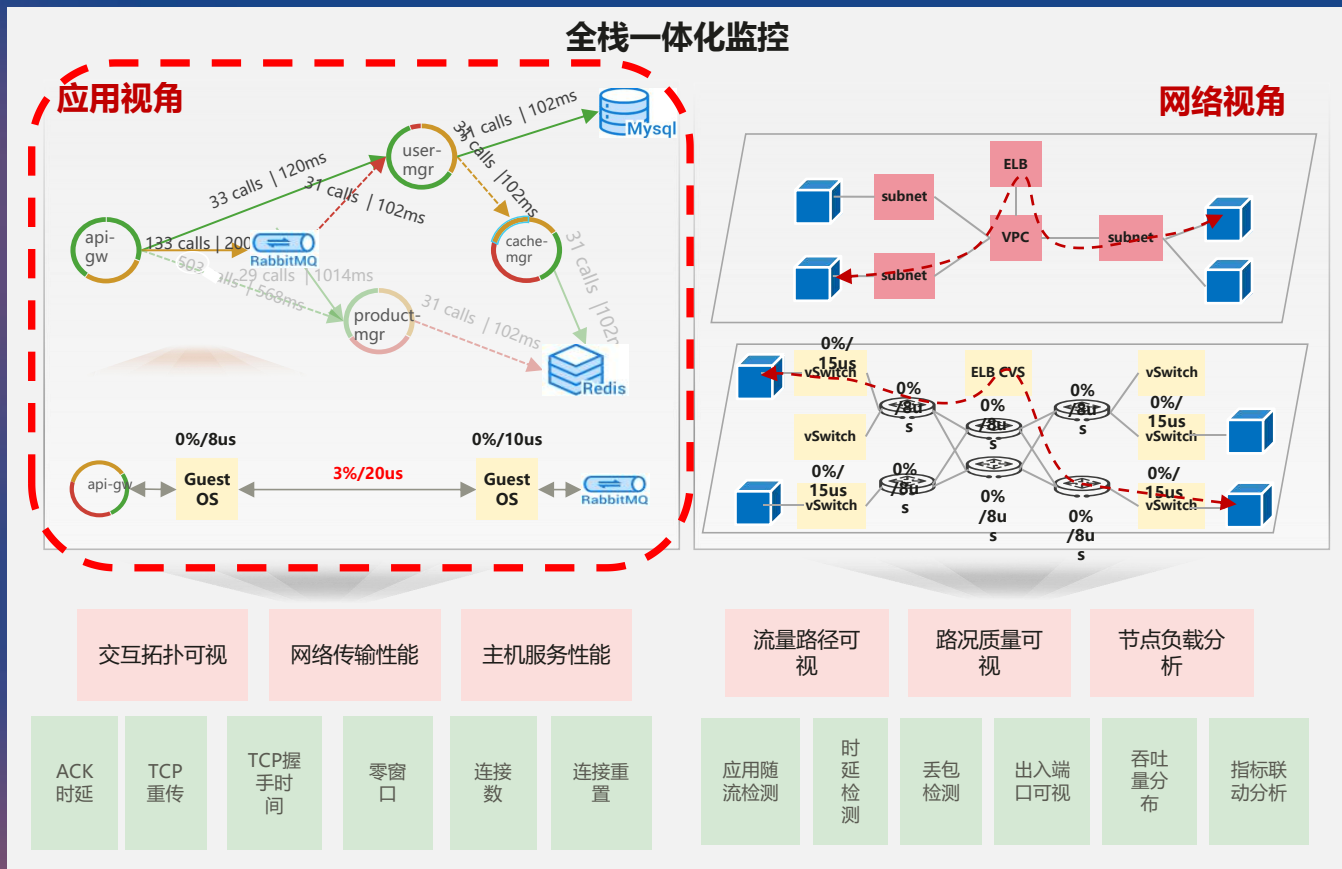
容器干扰检测，分钟级完成业务干扰源（CPU/IO）识别与干扰源发现。



规划：

1. 云原生场景：继续补齐基础设施观测能力，包括容器干扰观测、应用/网络定界观测等方面；
2. AI场景：提供CPU/NPU全栈观测能力，包括AI训练集群慢节点观测能力，NPU关键资源Profiling能力等。

⑦ gala-gopher在华为云Stack网络中的实践与应用



CloudNetDebug运维工具

故障诊断

流拨测

流抓包

主动链路监控

全链路故障诊断

审计/核查

网络随流检测

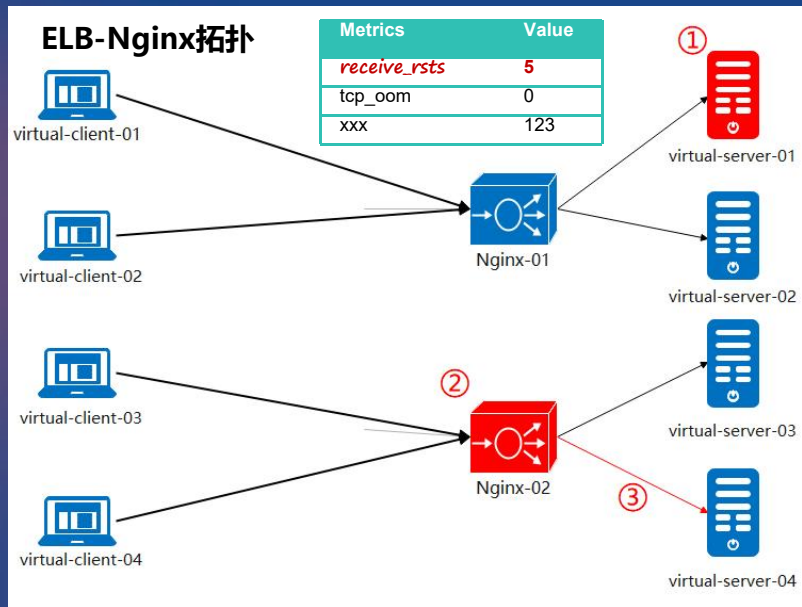
CloudNetDebug运维能力不足:

- **事后运维**: 只能在问题发生之后手动触发拨测, 非实时流量观测;
- **性能受限**: 频繁拨测造成的资源耗费和性能损耗较大, 不适合实时观测

gala-gopher带来的价值:

- 低底噪网络流量指标采集, 实时发现异常流量
- 全栈观测快速厘清应用/网络问题
- 应用资源关键指标波动可回溯

⑧ gala-gopher在华为云Stack网络中的实践与应用



问题	原因	策略
客户通过7层elb压测, 三万条有几十条报错	【后端业务问题】 后端超时配置错误导致回复reset报文。	取代抓包, 偶现故障记录: Gala-gopher可以采集到socket数据中的reset报文, 拓扑上指标可直接体现后端业务异常, 同时可生成系统告警。
apic 服务异常, 客户反馈影响某实时交易的业务	【ELB数据面问题】 Nginx进程单核卡死	关键指标波动回溯查询: 1. 采集进程CPU占用率可知Nginx进程异常; 2. Nginx和后端服务的数据量减小, 时延增大。
客户某业务经过elb达不到性能要求	【ELB性能问题】 后端服务器抓包判断elb负载合理, 最终原因是服务经过云外带宽受限	流量分布快速厘清: 拓扑可以直接体现Nginx和后端服务器的连接情况和数据量, 判断负载均衡是否合理。

监控策略: 使用eBPF监控ELB数据面高频故障组件ELB-Nginx, 采集四层网络通信状态指标数据, 并在指标异常时进行特殊标记、告警。

增强ELB现网监控和问题定位定界的能力, 补齐NGINX网元没有4层相关连接指标监控的缺陷, 同时补充健康检查离线场景定位定界的能力。