# 主机安全领域 eBPF的探索与实践

巫强 – 字节跳动终端安全

中国·西安
2025/04/19

# ByteHIDS Elkeid CWPP/主机入侵检测

字节跳动终端安全团队开发的主机入侵检测系统，驱动组件的目标是跟踪和审计所有高危操作，如进程创 建、文件创建、网络连接和权限提升等



ByteHIDS：https://github.com/bytedance/elkeid

# 为什么 **用** eBPF

- 高可用场景，低侵入性
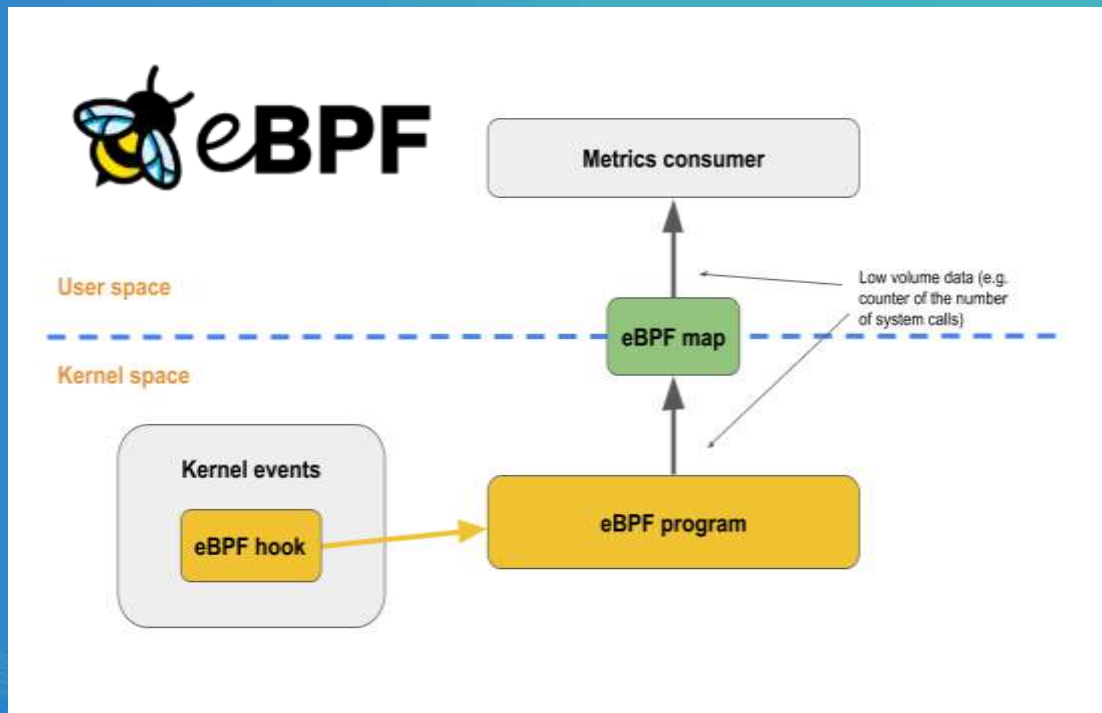- 用户需求以及对LKM的忧虑
- Linux发行版及众多内核版本的适配

# 为什么 **不得不用** eBPF

- 模块签名校验：Secure Boot、可信计算环境、云主机证书链管理
- 用户自定义内核：缺失头文件，模块支持 (CONFIG_MODULES)
- 权限限制：禁止root权限 (CAP_SYS_MODULE)，SECCOMP

# 为什么 **不用** eBPF

- 内核版本：支持与否/支持程度，如tail call可用性 4.2->5.10
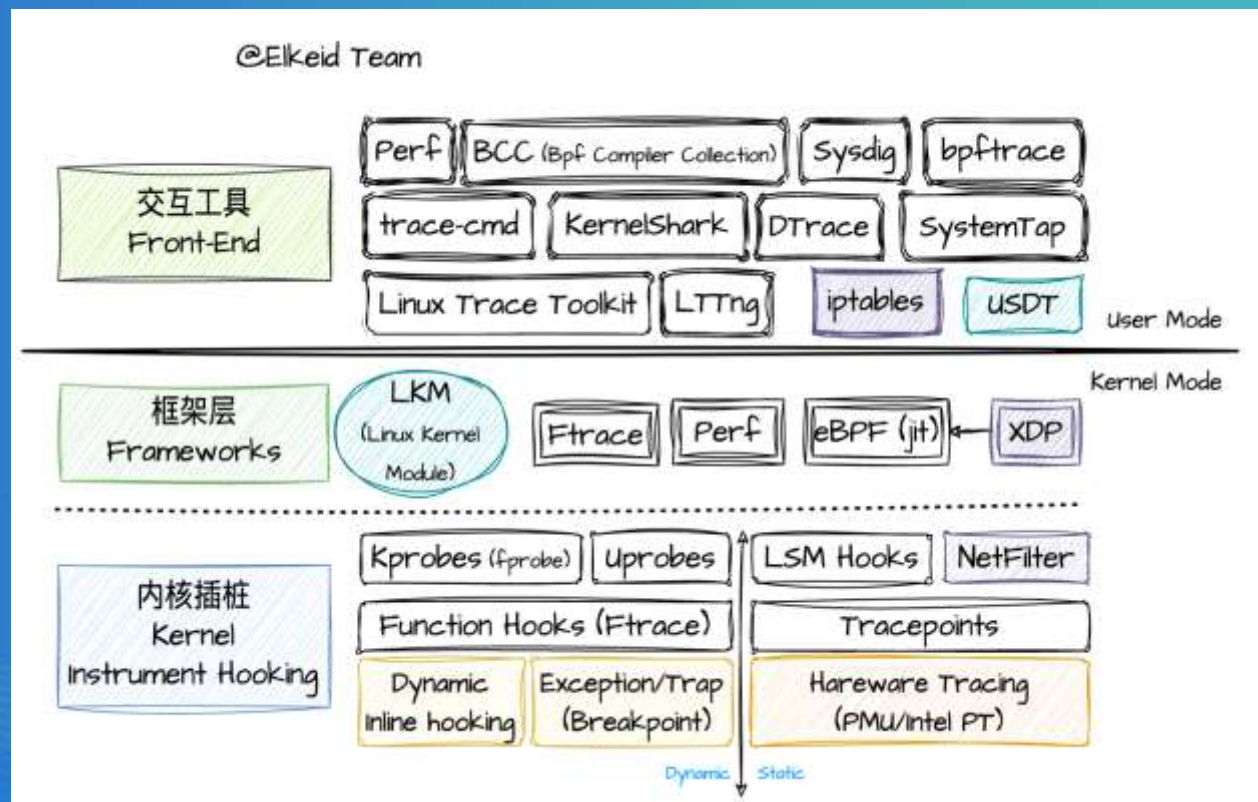- 特殊场景及性能要求: lockdown模式 (v5.4)，架构/硬件相关



图片来源：https://sysdig.com/blog/sysdig-and-falco-now-powered-by-ebpf/

# eBPF ≠ 高枕无忧

## eBPF常用监控机制:

- Kprobe (v4.1)
- Kretprobe ❌
- Tracepoint (v4.7)
  - ➢ sys_enter_xxx
  - ➢ sys_exit_xxx
- Raw-tracepoint (v4.17)
  - ➢ sys_enter
  - ➢ sys_exit
- Fprobe (x86: v5.5 arm64: 6.0)
  - ➢ fentry
  - ➢ fexit
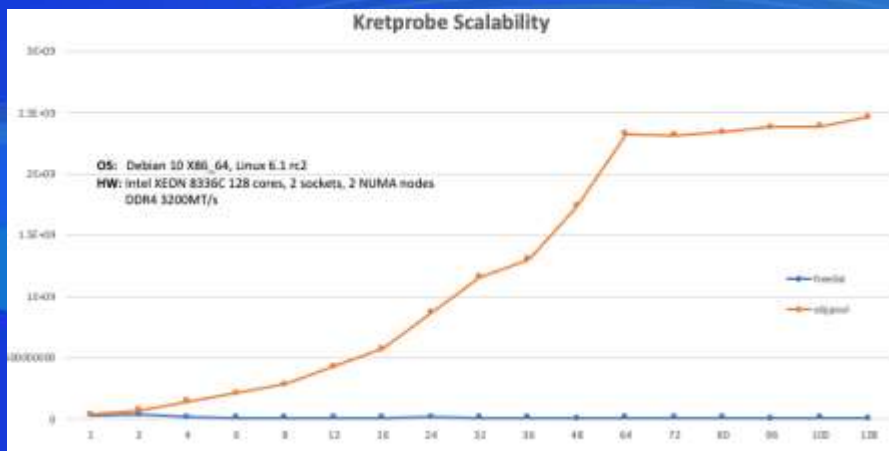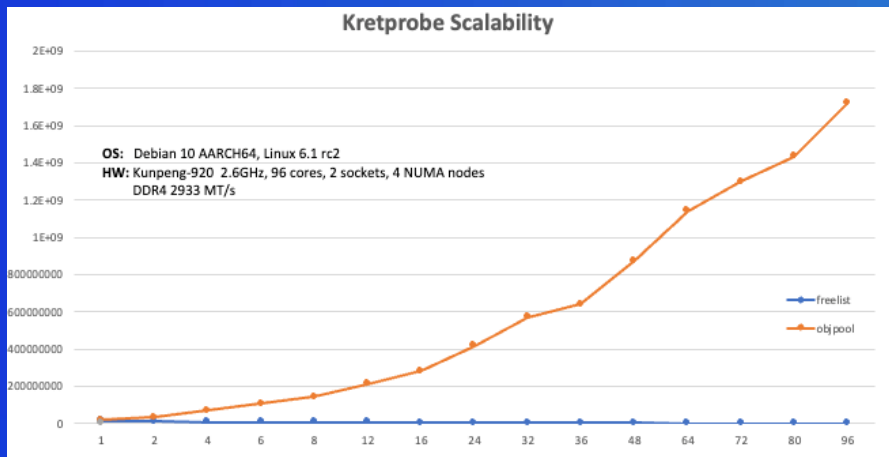- BPF LSM (v5.7: Kernel Runtime Security Instrumentation)



## 主要考量因素:

- 功能实现（高危操作入口或返回点， 64位系统中32位进程、 …）
- 工程便利（参数获取与信息流关联）
- 性能影响（内核版本与架构）

# 为什么 **不用** kretprobe
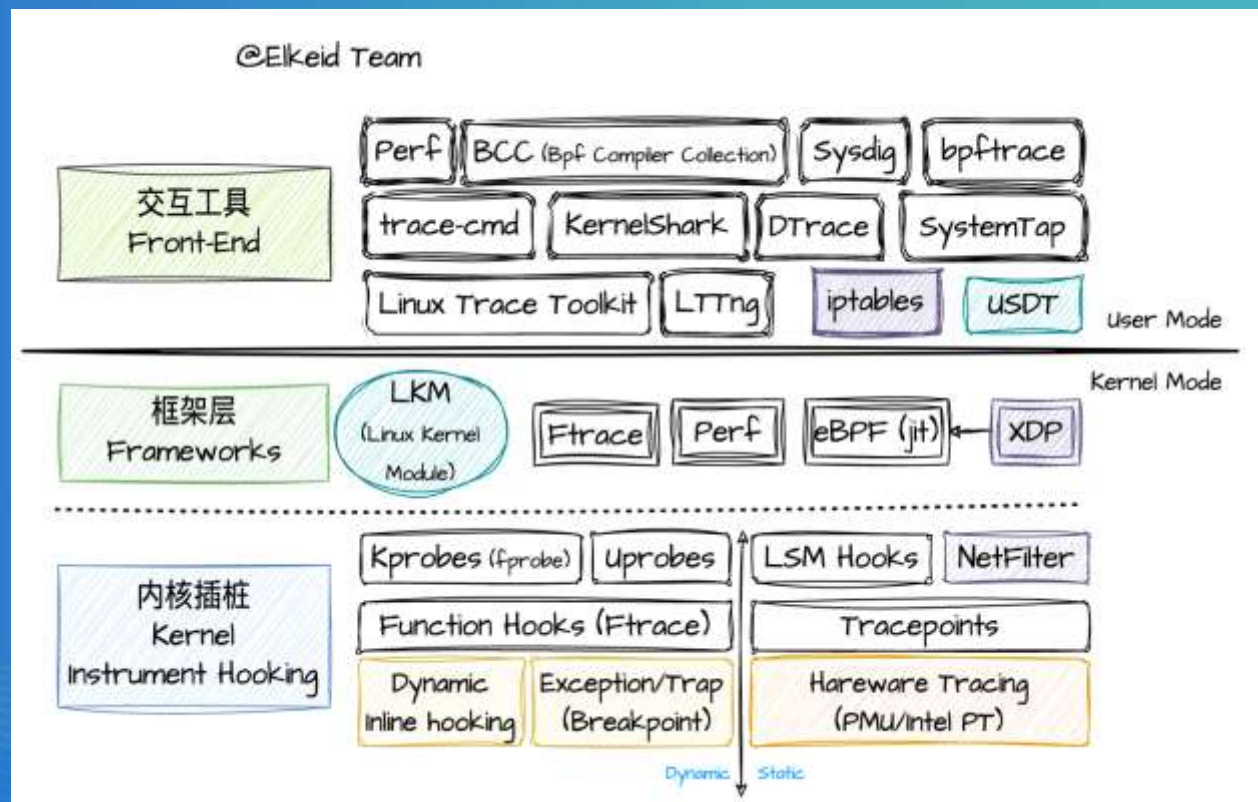
性能问题 - 可扩展性 (v6.7通过objpool解决)：

死锁风险 (5.12移除hash lock后才解决)：

# eBPF常用监控机制:

- Kprobe (v4.1)
- Kretprobe ❌
- Tracepoint (v4.7)
  - ➢ sys_enter_xxx
  - ➢ sys_exit_xxx
- Raw-tracepoint (v4.17)
  - ➢ sys_enter
  - ➢ sys_exit
- Fprobe (x86: v5.5 arm64: 6.0)
  - ➢ fentry
  - ➢ fexit
- BPF LSM (v5.7: Kernel Runtime Security Instrumentation)



# 主要考量因素:

- 功能实现（高危操作入口或返回点， 64位系统中32位进程、 …)
- 工程便利（参数获取与信息流关联）
- 性能影响（内核版本与架构）

# Tracepoint：32位程序的syscall？

```c
static void perf_syscall_enter(void *ignore, struct pt_regs *regs, long id)
{
        struct syscall_metadata *sys_data;
        struct syscall_trace_enter *rec;
        struct pt_regs *fake_regs;
        struct hlist_head *head;
        unsigned long args[6];
        bool valid_prog_array;
        int syscall_nr;
        int rctx;
        int size;

        syscall_nr = trace_get_syscall_nr(current, regs);
        if (syscall_nr < 0 || syscall_nr >= NR_syscalls)
                return;
        if (!test_bit(syscall_nr, enabled_perf_enter_syscalls))
                return;
```

```c
        rec->nr = syscall_nr;
        syscall_get_arguments(current, regs, args);
        memcpy(&rec->args, args, sizeof(unsigned long) * sys_data->nb_args);

        if ((valid_prog_array &&
             !perf_call_bpf_enter(sys_data->enter_event, fake_regs, sys_data, rec)) ||
            hlist_empty(head)) {
                perf_swevent_put_recursion_context(rctx);
                return;
        }

        perf_trace_buf_submit(rec, size, rctx,
                              sys_data->enter_event->event.type, 1, regs,
                              head, NULL);
}
```

```c
#ifdef ARCH_TRACE_IGNORE_COMPAT_SYSCALLS
/*
 * Some architectures that allow for 32bit applications
 * to run on a 64bit kernel, do not map the syscalls for
 * the 32bit tasks the same as they do for 64bit tasks.
 *
 *      *cough*x86*cough*
 *
 * In such a case, instead of reporting the wrong syscalls,
 * simply ignore them.
 *
 * For an arch to ignore the compat syscalls it needs to
 * define ARCH_TRACE_IGNORE_COMPAT_SYSCALLS as well as
 * define the function arch_trace_is_compat_syscall() to let
 * the tracing system know that it should ignore it.
 */
static int
trace_get_syscall_nr(struct task_struct *task, struct pt_regs *regs)
{
        if (unlikely(arch_trace_is_compat_syscall(regs)))
                return -1;

        return syscall_get_nr(task, regs);
}
```
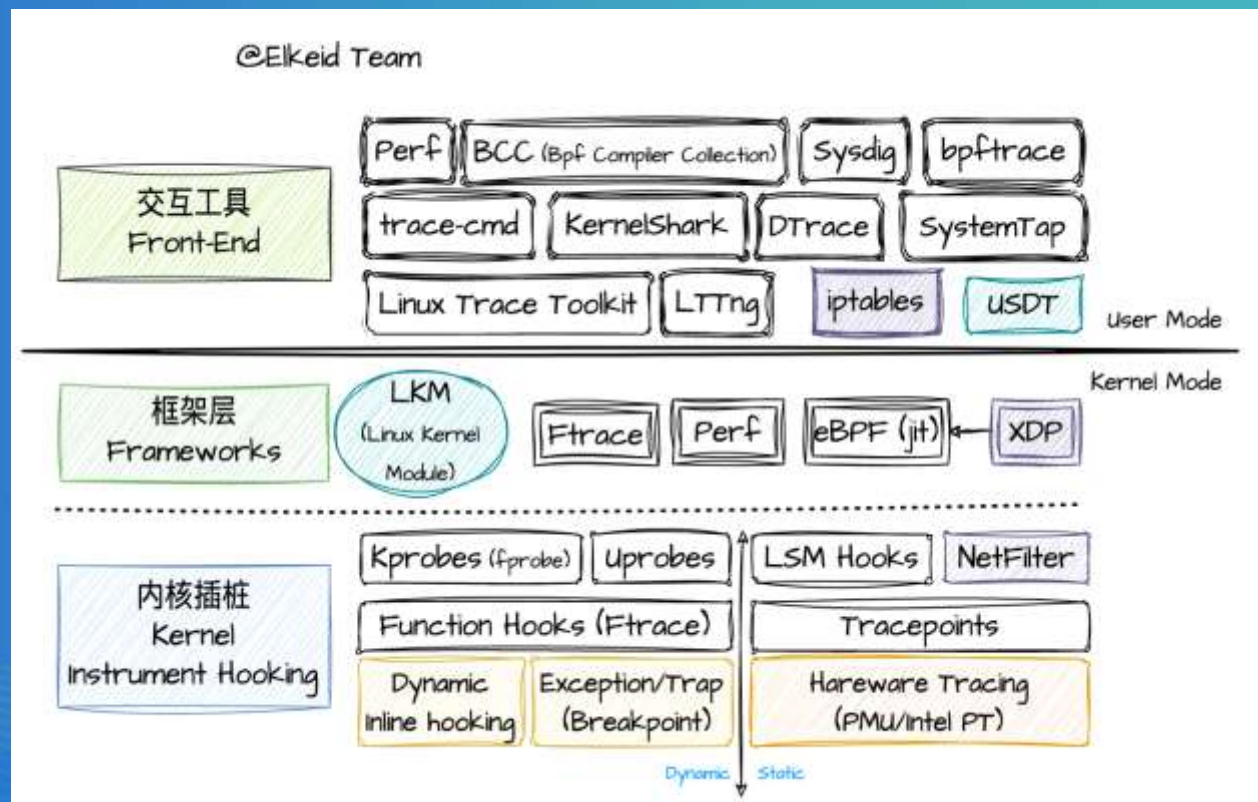
```
root@P22:/BUILD/linux-6.6-op# grep -rn IGNORE_COMPAT_SYSCALLS *
arch/x86/include/asm/ftrace.h:144:#define ARCH_TRACE_IGNORE_COM
arch/riscv/include/asm/ftrace.h:37:#define ARCH_TRACE_IGNORE_CO
arch/arm64/include/asm/ftrace.h:175:#define ARCH_TRACE_IGNORE_C
arch/s390/include/asm/ftrace.h:122:#define ARCH_TRACE_IGNORE_CO
```

## eBPF常用监控机制:

- Kprobe (v4.1)
- Kretprobe ✗
- Tracepoint (v4.7)
  - ➢ sys_enter_xxx
  - ➢ sys_exit_xxx
- Raw-tracepoint (v4.17)
  - ➢ sys_enter
  - ➢ sys_exit
- Fprobe (x86: v5.5 arm64: 6.0)
  - ➢ fentry
  - ➢ fexit
- BPF LSM (v5.7: Kernel Runtime Security Instrumentation)



## 主要考量因素:

- 功能实现（高危操作入口或返回点， 64位系统中32位进程、 …）
- 工程便利（参数获取与信息流关联）
- 性能影响（内核版本与架构）

# 问题：架构无关的悖论

**底层设计：**
ISA: eBPF 字节码、 Just-In-Time编译器、虚拟机 (BPF VM)

**细节依赖：**
Tracepoint是对Raw Tracepoint的封装，前者禁止了对pt_regs的访问，牺牲灵便性的代价来换取跨架构的兼容性

**实际需求：**
1. 64位系统下32位进程的判断（ARM64/X86_64/X64_32）
2. Rootkit检测需要对IDT表/CRx寄存器的审计（X86: SMEP/SMAP ARM64: PAN/PXN）



X86_64

ARM64

ARM64: 32b or 64b ELF

11111 

111

# 问题：两种ringbuf类型的选择

BPF ringbuf:

5.8及之后的内核才支持，此外bpf_ringbuf_reserve的输入size只能
为常量，不支持动态长度

Perf ringbuf:

通用性强，但需要先将数据准备好且需要整体做一次内存拷贝

```
57  SEC("raw_tracepoint/sys_enter")
58  int bpf_prog1(void *ctx)
59  {
60      int max_len, max_buildid_len, total_size;
61      struct stack_trace_t *data;
62      long usize, ksize;
63      void *raw_data;
64      __u32 key = 0;
65
66      data = bpf_map_lookup_elem(&stackdata_map, &key);
67      if (!data)
68          return 0;
69
70      max_len = MAX_STACK_RAWTP * sizeof(__u64);
71      max_buildid_len = MAX_STACK_RAWTP * sizeof(struct bpf_stack_build_id);
72      data->pid = bpf_get_current_pid_tgid();
73      data->kern_stack_size = bpf_get_stack(ctx, data->kern_stack,
74                                            max_len, 0);
75      data->user_stack_size = bpf_get_stack(ctx, data->user_stack, max_len,
76                                            BPF_F_USER_STACK);
77      data->user_stack_buildid_size = bpf_get_stack(
78          ctx, data->user_stack_buildid, max_buildid_len,
79          BPF_F_USER_STACK | BPF_F_USER_BUILD_ID);
80      bpf_perf_event_output(ctx, &perfmap, 0, data, sizeof(*data));
81
```

# 问题：map动态调整及内存占用

**不可动态调整：**

只允许在map创建前调整map大小，针对不同规模的系统
需要在eBPF加载前进行max entries的调整

**内存占用 (过大?):**

在开启CGROUP KMEM（CONFIG_MEMCG_KMEM=y）的情况
下eBPF程序的内存占用会被统计至加载器所在的CGROUP

**性能问题 (过小?):**
Hash table类型的map底层用到了pcpu_freelist (基于spinlock分锁),
在数组设置过小的情况下会导致性能问题

```
513  int bpf_mem_alloc_init(struct bpf_mem_alloc *ma, int size, bool percpu)
514  {
515      struct bpf_mem_caches *cc, __percpu *pcc;
516      struct bpf_mem_cache *c, __percpu *pc;
517      struct obj_cgroup *objcg = NULL;
518      int cpu, i, unit_size, percpu_size = 0;
519
520      if (percpu && size == 0)
521          return -EINVAL;
522
523      /* room for llist_node and per-cpu pointer */
524      if (percpu)
525          percpu_size = LLIST_NODE_SZ + sizeof(void *);
526      ma->percpu = percpu;
527
528      if (size) {
529          pc = __alloc_percpu_gfp(sizeof(*pc), 8, GFP_KERNEL);
530          if (!pc)
531              return -ENOMEM;
532
533          if (!percpu)
534              size += LLIST_NODE_SZ; /* room for llist_node */
535          unit_size = size;
536
537  #ifdef CONFIG_MEMCG_KMEM
538          if (memcg_bpf_enabled())
539              objcg = get_obj_cgroup_from_current();
540  #endif
541          ma->objcg = objcg;
542
543          for_each_possible_cpu(cpu) {
544              c = per_cpu_ptr(pc, cpu);
545              c->unit_size = unit_size;
546              c->objcg = objcg;
547              c->percpu_size = percpu_size;
548              c->tgt = c;
549              init_refill_work(c);
550              prefill_mem_cache(c, cpu);
551          }
552          ma->cache = pc;
553          return 0;
554      }
```
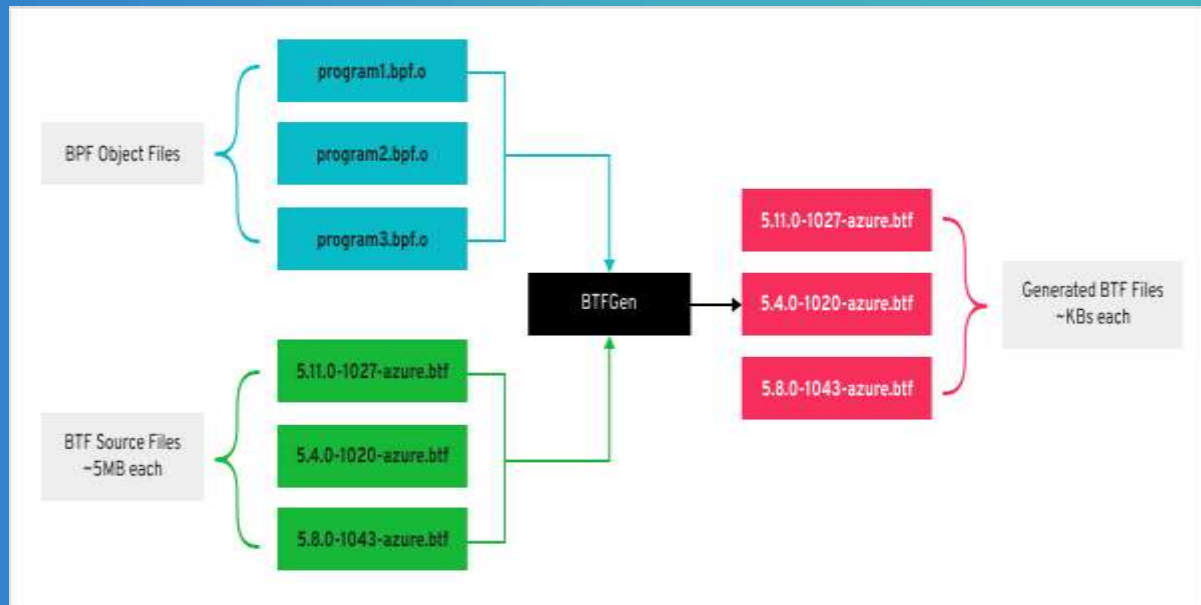
# 问题：部署与内核适配

## 支持BTF的内核版本的适配：

尽可能复用同一个eBPF二进制程序，比如只为 5.4.x、5.10.x内核编译两个不同的eBPF目标文件，不同的发行版及架构间不通用

## 不支持BTF的内核版本适配：

利用bpftool/btfge具定制生成裁剪版本的BTF文件，同时保持eBPF程序的可复用性

```
bpftool gen min_core_btf /sys/kernel/btf/vmlinux mini.btf elkeid.btf
ls -l /sys/kernel/btf/vmlinux ./mini.btf
   4283 Apr 17 15:44 ./mini.btf
4791512 Apr 17 15:42 /sys/kernel/btf/vmlinux
```



图片来源：https://inspektor-gadget.io/blog/2022/03/btfgen-one-step-closer-to-truly-portable-ebpf-programs/

# 附录： tracepoint: sys_enter_xxx

调用栈:



参数:
struct syscall_tp_t *args; 前8字节 (struct trace_entry) 实际存放的是pt_regs，但ebpf程序不可访问

# tracepoint: sys_exit_xxx

调用栈:

```
[pwndbg> bt
#0   0xffffffffc001c054 in ?? ()
#1   0xffffffff81256c7b in bpf_dispatcher_nop_func (bpf_func=0xffffffffc001c054, insnsi=0xffffc90000911048, ctx=0xffffc90000cfbe88) at ./include/linux/bpf.h:1181
#2   __bpf_prog_run (dfunc=<optimized out>, ctx=0xffffc90000cfbe88, prog=0xffffc90000911000) at ./include/linux/filter.h:609
#3   bpf_prog_run (ctx=0xffffc90000cfbe88, prog=0xffffc90000911000) at ./include/linux/filter.h:616
#4   bpf_prog_run_array (run_prog=<optimized out>, ctx=0xffffc90000cfbe88, array=<optimized out>) at ./include/linux/bpf.h:1926
#5   trace_call_bpf (call=<optimized out>, ctx=ctx@entry=0xffffc90000cfbe88) at kernel/trace/bpf_trace.c:140
#6   0xffffffff8123a8ae in perf_call_bpf_exit (rec=0xffffe8ffffdba9b0, regs=0xffffc90000cfbf58, call=<optimized out>) at kernel/trace/trace_syscalls.c:672
#7   perf_syscall_exit (ignore=<optimized out>, regs=0xffffc90000cfbf58, ret=<optimized out>) at kernel/trace/trace_syscalls.c:712
#8   0xffffffff81187a03 in __traceiter_sys_exit (__data=0xffffc90000cfbe88, regs=0xffffc90000cfbf58, ret=9) at ./include/trace/events/syscalls.h:44
#9   0xffffffff8118823d in trace_sys_exit (ret=<optimized out>, regs=0xffffc90000cfbf58) at ./include/trace/events/syscalls.h:44
#10  syscall_exit_work (regs=regs@entry=0xffffc90000cfbf58, work=18) at kernel/entry/common.c:247
#11  0xffffffff81f32673 in syscall_exit_to_user_mode_prepare (regs=regs@entry=0xffffc90000cfbf58) at kernel/entry/common.c:278
#12  __syscall_exit_to_user_mode_work (regs=regs@entry=0xffffc90000cfbf58) at kernel/entry/common.c:283
#13  syscall_exit_to_user_mode (regs=regs@entry=0xffffc90000cfbf58) at kernel/entry/common.c:296
#14  0xffffffff81f2af36 in do_syscall_64 (regs=0xffffc90000cfbf58, nr=<optimized out>) at arch/x86/entry/common.c:86
#15  0xffffffff820000e6 in entry_SYSCALL_64 () at arch/x86/entry/entry_64.S:120
#16  0x0000000000000000 in ?? ()
```

参数:
struct syscall_tp_t *args; 前8字节实际存放的是pt_regs，但ebpf程序不可访问

# kprobe hook

调用栈:



参数:
    struct pt_regs *regs;

# raw-tracepoint: sys_enter/sys_exit

调用栈:

```
[pwndbg> bt
#0  0xfffffffffc00138ac in ?? ()
#1  0xffffffff81254bd1 in bpf_dispatcher_nop_func (bpf_func=0xfffffffffc00138ac, insnsi=0xffffc900019b1048, ctx=0xffffc90000d33e90) at ./include/linux/bpf.h:1181
#2  __bpf_prog_run (dfunc=<optimized out>, ctx=0xffffc90000d33e90, prog=0xffffc900019b1000) at ./include/linux/filter.h:609
#3  bpf_prog_run (ctx=0xffffc90000d33e90, prog=0xffffc900019b1000) at ./include/linux/filter.h:616
#4  __bpf_trace_run (args=0xffffc90000d33e90, prog=0xffffc900019b1000) at kernel/trace/bpf_trace.c:2306
#5  bpf_trace_run2 (prog=0xffffc900019b1000, arg0=<optimized out>, arg1=<optimized out>) at kernel/trace/bpf_trace.c:2345
#6  0xffffffff81187d29 in __bpf_trace_sys_exit () at ./include/trace/events/syscalls.h:44
#7  0xffffffff81187a03 in __traceiter_sys_exit (__data=0xffffc90000d33e90, regs=0xffffc90000d33f58, ret=-110) at ./include/trace/events/syscalls.h:44
#8  0xffffffff8118823d in trace_sys_exit (ret=<optimized out>, regs=0xffffc90000d33f58) at ./include/trace/events/syscalls.h:44
#9  syscall_exit_work (regs=regs@entry=0xffffc90000d33f58, work=18) at kernel/entry/common.c:247
#10 0xffffffff81f32673 in syscall_exit_to_user_mode_prepare (regs=regs@entry=0xffffc90000d33f58) at kernel/entry/common.c:278
#11 __syscall_exit_to_user_mode_work (regs=regs@entry=0xffffc90000d33f58) at kernel/entry/common.c:283
#12 syscall_exit_to_user_mode (regs=regs@entry=0xffffc90000d33f58) at kernel/entry/common.c:296
#13 0xffffffff81f2af36 in do_syscall_64 (regs=0xffffc90000d33f58, nr=<optimized out>) at arch/x86/entry/common.c:86
#14 0xffffffff820000e6 in entry_SYSCALL_64 () at arch/x86/entry/entry_64.S:120
#15 0x00007ff4307ed000 in ?? ()
#16 0x00007fffdccd0da0 in ?? ()
#17 0x0000000000000016 in fixed_percpu_data ()
#18 0xffffffffffffff78 in ?? ()
#19 0x00007ff430fecd50 in ?? ()
#20 0x000000c00005c400 in ?? ()
#21 0x0000000000000202 in ?? ()
#22 0x00007ff430fecd40 in ?? ()
#23 0x0000000000000000 in ?? ()
```

参数:

struct bpf_raw_tracepoint_args *ctx;
其中ctx->args[0]为pt_regs指针，sys_exit的情况下ctx->args[1]为syscall返回值