# 基于eBPF的函数调用栈观测工具

## 西安邮电大学--刘冰

中国·西安

# 个人简介

刘冰，就读于西安邮电大学计算机学院，研究生二年级，曾获得2019全国大学生数学建模竞赛省奖，2023年全国大学计算机系统能力大赛国奖，2023年CCF开源夏令营优秀项目奖

0、背景

1、高效的数据组织方式

2、PSI触发及及时采集

3、多指标关联

4、持续性观测和长期存储

中国·西安

# 0、使用不便

**！相关工具繁杂，问题不一**

- ！ profile用于on_cpu调用栈采样，只能设置采样频率，功能单一
- ！ offcputime用于采集off_cpu调用栈，无法实时显示数据
- ！ memleak用于采集内存泄露调用栈，无法正常显示用户态调用栈，无法有效执行子进程
- ！ ...

**！不便于分析**

- ！ 上述等大多数工具都会将采集的数据按调用栈进行输出，篇幅大，不便于直接查阅和分析

```
COMM: profile (pid=3203151) @ CPU 1
No Kernel Stack
Userspace:
0000586b012429c3:
blazesym::symbolize::symbolizer::Symbolizer::symbolize @
0x1239a0+0x23
0000586b011550e7: event_handler @ 0x3603a+0xad /profile.c:111
0000586b0116f70f: ringbuf_process_ring @ 0x506c0+0x4f
/ringbuf.c:261

COMM: sshd (pid=3183771) @ CPU 3
Kernel:
ffffffffc023c74b: e1000_xmit_frame @ 0xffffffffc023c130+0x61b
ffffffff81e3c5a5: dev_hard_start_xmit @ 0xffffffff81e3c540+0x65
ffffffff81eaf25b: sch_direct_xmit @ 0xffffffff81eaf150+0x10b
ffffffff81e38c3f: __dev_xmit_skb @ 0xffffffff81e38930+0x30f
ffffffff81e3cb9b: __dev_queue_xmit @ 0xffffffff81e3c810+0x38b
ffffffff81ef98b3: neigh_hh_output @ 0xffffffff81ef9820+0x93
ffffffff81efa21e: ip_finish_output2 @ 0xffffffff81efa040+0x1de
ffffffff81efb166: __ip_finish_output @ 0xffffffff81efb0b0+0xb6
ffffffff81efb269: ip_finish_output @ 0xffffffff81efb240+0x29
ffffffff81efb3c3: ip_output @ 0xffffffff81efb350+0x73
ffffffff81efca31: ip_local_out @ 0xffffffff81efc9d0+0x61
ffffffff81efcbdd: __ip_queue_xmit @ 0xffffffff81efca50+0x18d
ffffffff81efcf15: ip_queue_xmit @ 0xffffffff81efcf00+0x15
ffffffff81f269a1: __tcp_transmit_skb @ 0xffffffff81f26050+0x951
ffffffff81f282af: tcp_write_xmit @ 0xffffffff81f27e00+0x4af
ffffffff81f28907: __tcp_push_pending_frames @
0xffffffff81f288d0+0x37
ffffffff81f0d683: tcp_push @ 0xffffffff81f0d560+0x123
ffffffff81f0e82f: tcp_sendmsg_locked @ 0xffffffff81f0de70+0x9bf
ffffffff81f0ec1c: tcp_sendmsg @ 0xffffffff81f0ebf0+0x2c
ffffffff81fb3472: inet6_sendmsg @ 0xffffffff81fb3430+0x42
ffffffff81e03895: sock_write_iter @ 0xffffffff81e03770+0x125
ffffffff814b0fd4: vfs_write @ 0xffffffff814b0c40+0x394
ffffffff814b1499: ksys_write @ 0xffffffff814b13d0+0xc9
ffffffff814b14f9: __x64_sys_write @ 0xffffffff814b14e0+0x19
ffffffff8213bdb9: do_syscall_64 @ 0xffffffff8213bd60+0x59
ffffffff822000e6: entry_SYSCALL_64_after_hwframe @
0xffffffff82200078+0x6e
Userspace:
0000767b3391b294: write @ 0x11b280+0x14
0000000000000000: <no-symbol>

COMM: swapper/2 (pid=0) @ CPU 2
Kernel:
ffffffff8214301b: pv_native_safe_halt @ 0xffffffff82143010+0xb
ffffffff821459e0: acpi_idle_do_entry @ 0xffffffff821459a0+0x40
ffffffff82145e16: acpi_idle_enter @ 0xffffffff82145d60+0xb6
```

# 0、性能消耗大

! 数据量大
! IO损耗高

例如perf工具或者基于eBPF的profile工具，其会输出每个调用栈样本，会造成频繁输出，并产生大量数据。

本项目设计了一个统一的框架，主要针对程序性能瓶颈问题，来进行基于eBPF的调用栈采集

Agents sample the stack trace and send profiling data to the Pyroscope serve r

0、背景

1、友好的数据组织方式

2、PSI触发及及时采集

3、多指标关联

4、持续性观测和长期存储

中国·西安

# 1、本地拆分展示

## 默认49hz采集样本 输出Top10调用栈

- ✓ **第一个表为pid，用户栈id，内核栈id及其关联指标 按指标值进行排序，便于找出高消耗进程和其瓶颈 调用栈的id**

- ✓ **第二个表为栈id和栈，可根据栈id找到调用栈， 查看函数调用关系**

- ✓ **第三个表为pid、命名空间pid、命令名、 命名空间tgid和cgroup id，由此可看出pid所属的 组，容器等信息，有利于进一步回溯**



```
$ sudo ./stack_analyzer on_cpu -u -k -t 5

StackAnalyzer

Attach collecotor1 OnCPUStackCollector.
Running for 5s or Hit Ctrl-C to end.
time:20240410_08_53_55
counts:
pid      usid    ksid    OnCPUTime/20408163nanoseconds
3264294  21430   29713   1
3264295  3504    97919   1
3264301  35038   -14     1
3264303  121564  9044    1
3264305  29980   106847  1
3264306  57728   128280  1
3264307  127282  36663   1
3264308  -14     213     1
28612    124966  80747   2
3230503  10771   78673   2
traces:
sid      trace
213      entry_SYSCALL_64_after_hwframe+0x6e;do_syscall_64+0x59;__x64_sys_exit_group+0x18;do_gr
xa8;down_write+0x25;
3504     0x3639203231343220;0x72df6591b294;
9044     asm_exc_page_fault+0x27;exc_page_fault+0x94;irqentry_exit+0x43;irqentry_exit_to_user_r
10771    _Fork+0x27;
21430    0x7f483d8ea2f7;
29713    entry_SYSCALL_64_after_hwframe+0x6e;do_syscall_64+0x59;__x64_sys_clone+0x25;__do_sys_c
.isra.0+0x110;mas_wr_modify+0x19e;mas_update_gap.part.0+0xd6;mas_leaf_max_gap+0xba;
29980    0x53464900202b0053;0x6451d19cc280;
35038    0x75b66da4377a;
36663    entry_SYSCALL_64_after_hwframe+0x6e;do_syscall_64+0x68;syscall_exit_to_user_mode+0x32
57728    0x75fc00f7d096;
78673    entry_SYSCALL_64_after_hwframe+0x6e;do_syscall_64+0x59;__x64_sys_clone+0x25;__do_sys_c
nge+0x3cd;copy_pte_range+0x142;copy_present_pte+0x26b;
80747    entry_SYSCALL_64_after_hwframe+0x6e;do_syscall_64+0x68;syscall_exit_to_user_mode+0x29
er_8+0x10;
97919    entry_SYSCALL_64_after_hwframe+0x6e;do_syscall_64+0x59;__x64_sys_write+0x19;ksys_write
106847   asm_exc_page_fault+0x27;exc_page_fault+0x83;do_user_addr_fault+0x212;handle_mm_fault+0
121564   0x7f483d9157b0;
124966   __nptl_death_event+0x186;
127282   0x8bfffcc801e808ec;0x0;0x7f483d842a63;
128280   entry_SYSCALL_64_after_hwframe+0x6e;do_syscall_64+0x59;__x64_sys_mmap+0x33;ksys_mmap_p
x239;vma_complete+0x26b;mas_store_prealloc+0x6c;mas_destroy+0x71;kmem_cache_free_bulk+0x13;kme
info:
pid      NSpid    comm    tgid     cgroup
3264292  3264292  ps      3264292  session-409.scope
3264291  3264291  node    3264291  session-409.scope
3264268  3264268  ps      3264268  session-409.scope
3264289  3264289  node    3264289  session-409.scope
3264276  3264276  cat     3264276  session-409.scope
3257031  3257031  pyroscope        3257012  session-409.scope
3264303  3264303  cpuUsage.sh      3264303  session-409.scope
3204135  3204135  node    3204135  session-409.scope
```
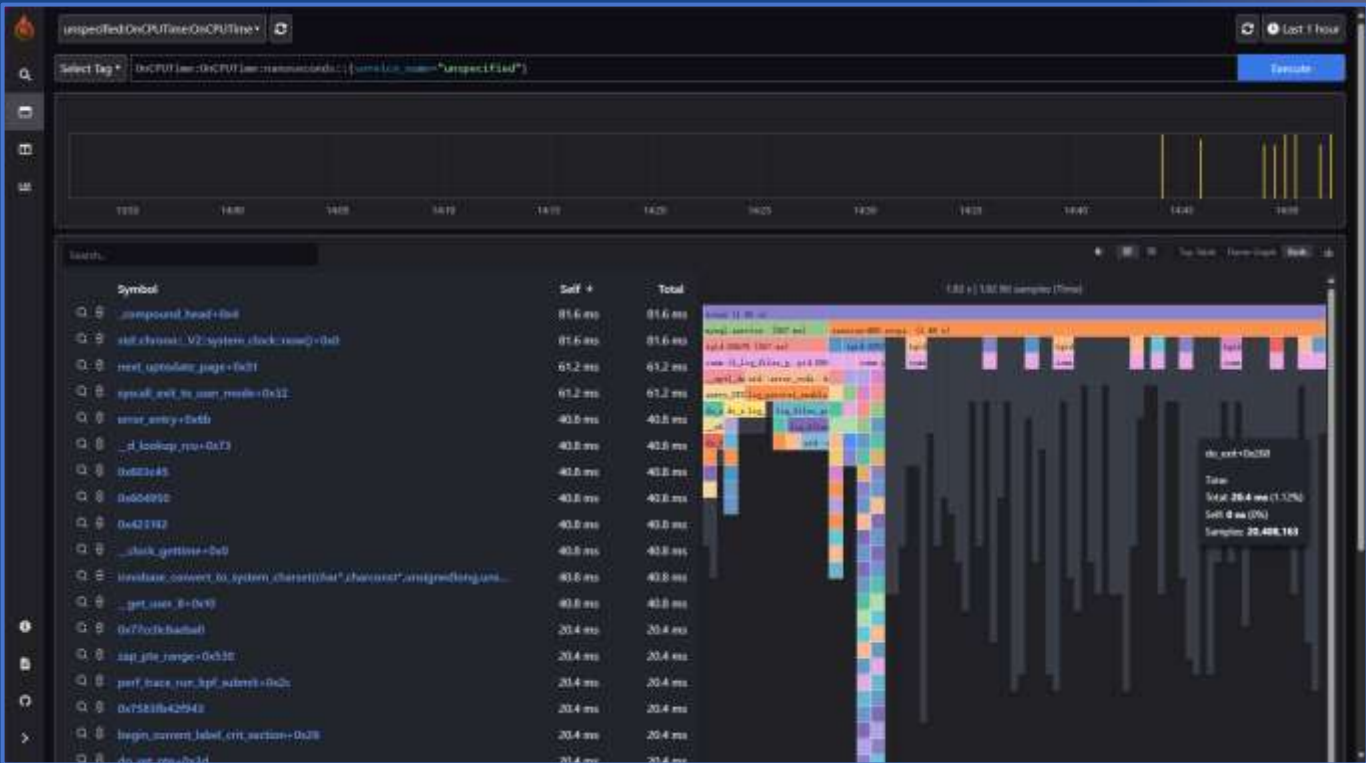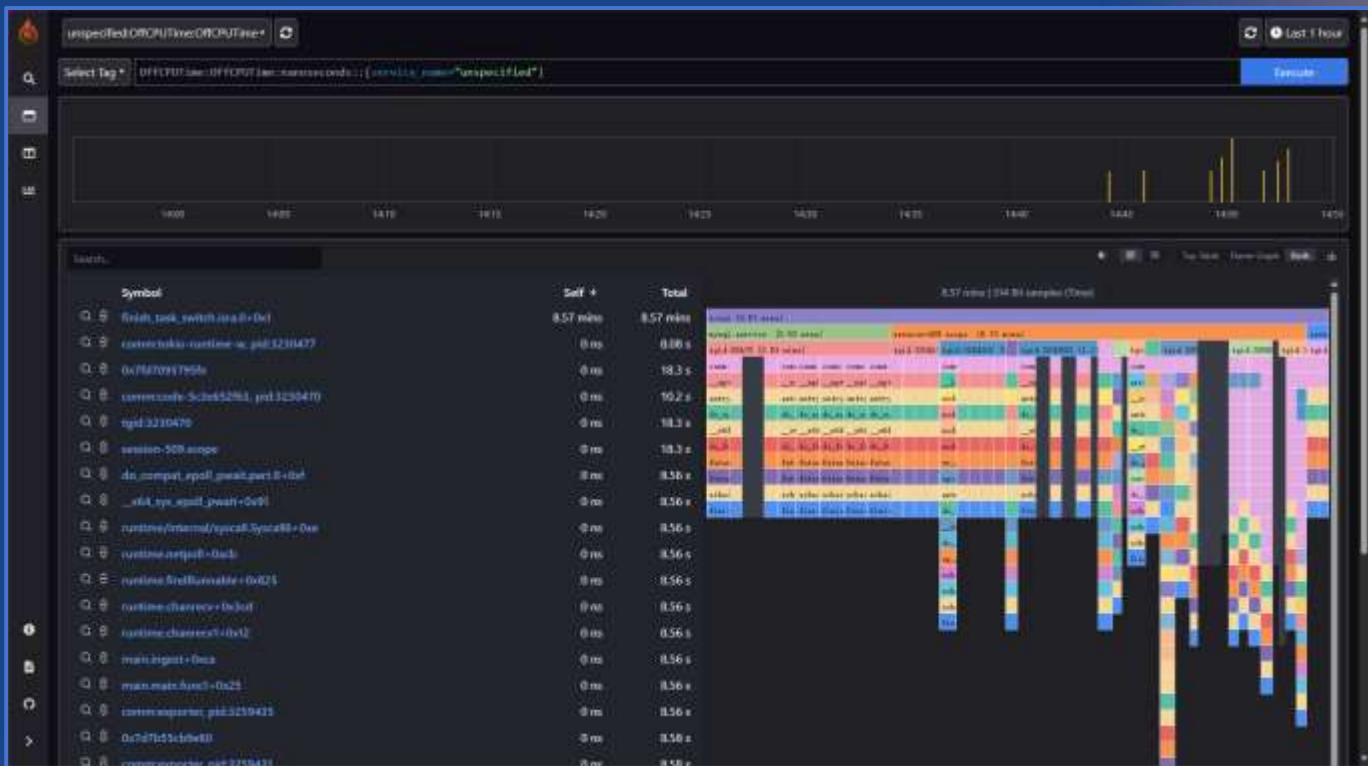
# 1、精准数据采集

- ✓ **按cgroup过滤**
- ✓ **按pid过滤**
- ✓ **按tid过滤**
- ✓ **创建指定命令进程并采集**

```
$ ./stack_analyzer -h
DESCRIPTION
   Count the function call stack associated with some metric.

                    StackAnalyzer

SYNOPSIS
   ./stack_analyzer [on_cpu] [off_cpu] [memleak [-W]] [io] [readahead] [probe <probe>] [llc_stat [-P <period>]] ([-g <cgroup
                    path>] | [-p <pid>] | [-t <tid>] | [-c <command>]) [-o <top>] [-f <freq>] [-i <interval>] [-d <duration>]
                    [-u] [-k] [-T (cpu|memory|io) <event>] [-v] [-h]

OPTIONS
   on_cpu              Collector for on-cpu trace
   off_cpu             Collector for off-cpu trace
   memleak             Collector for memleak trace
   -W                  Free when missing in kernel to alleviate misjudgments
   io                  Collector for io trace
   readahead           Collector for readahead trace
   probe               Collector for probe trace
   <probe>             Set the probe string
   llc_stat            Collector for llc_stat trace
   -P <period>         Set sampling period; default is 100

   Some overall options
   -g <cgroup path>    Set the cgroup of the process to be tracked; default is -1, which keeps track of all cgroups
   -p <pid>            Set the pid of the process to be tracked; default is -1, which keeps track of all processes
   -t <tid>            Set the tid of the thread to be tracked; default is -1, which keeps track of all threads
   -c <command>        Set the command to be run and sampled; defaults is none
   -o <top>            Set the top number; default is 10
   -f <freq>           Set sampling frequency, 0 for close; default is 49
   -i <interval>       Set the output delay time (seconds); default is 5
   -d <duration>       Set the total sampling time; default is __INT_MAX__
   -u                  Sample user stacks
   -k                  Sample kernel stacks
   -T (cpu|memory|io) <event>
                       Set a trigger for monitoring. For example, -T cpu "some 150000 100000" means triggers when
                       cpu partial stall with 1s tracking window size * and 150ms threshold.

   Information of the application
   -v, --version       Show version
   -h, --help          Show man page

LICENSE
Apache Licence 2.0
```

# 1、服务端可视化

## 兼容grafana pyroscope

✓ **本项目的发送器将采集数据发送到 pyroscope服务端**

✓ **有指标时序图、调用栈火焰图和 函数占比列表三种可视化方案**

✓ **支持对比视图和差分视图**

# 1、性能测试



对系统性能的影响（负值代表降低性能）

对系统性能的影响（负值代表降低性能）

第二届 eBPF开发者大会
www.ebpftravel.com

0、背景

1、高效的数据组织方式

2、PSI触发及及时采集

3、多指标关联

4、持续性观测和长期存储

中国·西安

# 2、PSI触发

## ✓ PSI报警后监控

**本系统可基于psi实现动态的负载监控，着眼与高阻塞环境，节省不必要的监控开销。**

# 2、高阻塞环境下及时触发

✓ **核心思想：**
 **eBPF attach需要进行系统调用，在高阻塞环境可能会使数据采集不及时，这里使用eBPF 全局变量进行采集控制**

✓ **优点：**
 **减小性能损耗**
 **及时地数据采集**

1、　　背景

2、　　PSI触发及及时采集

3、　　多指标关联

4、　　持续性观测及长期存储

中国·西安

# 3、增加数据复用性

## ✓ 关联多指标的监测功能

按进程调用栈，统计IO操作的次数和数据量

按进程调用栈，统计缓存失配次数、缓存命中次数以及缓存命中率

按进程调用栈，统计内存泄露大小、内存分配未释放次数

。。。

## ✓ 多视图观察

为调用栈数据提供了更多的视图，便于多维度剖析调用栈性能瓶颈

# 3、便于扩展的指标

✓ **自定义指标列表**
✓ **自定义指标解析方法**

**利于框架可扩展性**

中国·西安

# 4、eBPF应用案例　-- cpu占用

# 4、eBPF应用案例 -- 阻塞

# 4、eBPF应用案例 --内存泄露

# 4、eBPF应用案例 --内存泄露

# 4、eBPF应用案例  -- 输入输出

# 4、eBPF应用案例 -- 自定义跟踪点计数：vfs_open

谢谢大家！