



第三届 eBPF开发者大会

[www.ebpftravel.com](http://www.ebpftravel.com)

# eBPF - 交流研讨

中国·西安



第三届 eBPF开发者大会

[www.ebpftravel.com](http://www.ebpftravel.com)

# eBPF技术在滴滴自动驾驶场景下的相关实践

陈涛·滴滴

中国·西安



第三届 eBPF 开发者大会

[www.ebpftravel.com](http://www.ebpftravel.com)

- 1 项目整体概况
- 2 如何做到低频问题的定位/定界
- 3 问题案例
- 4 未来展望

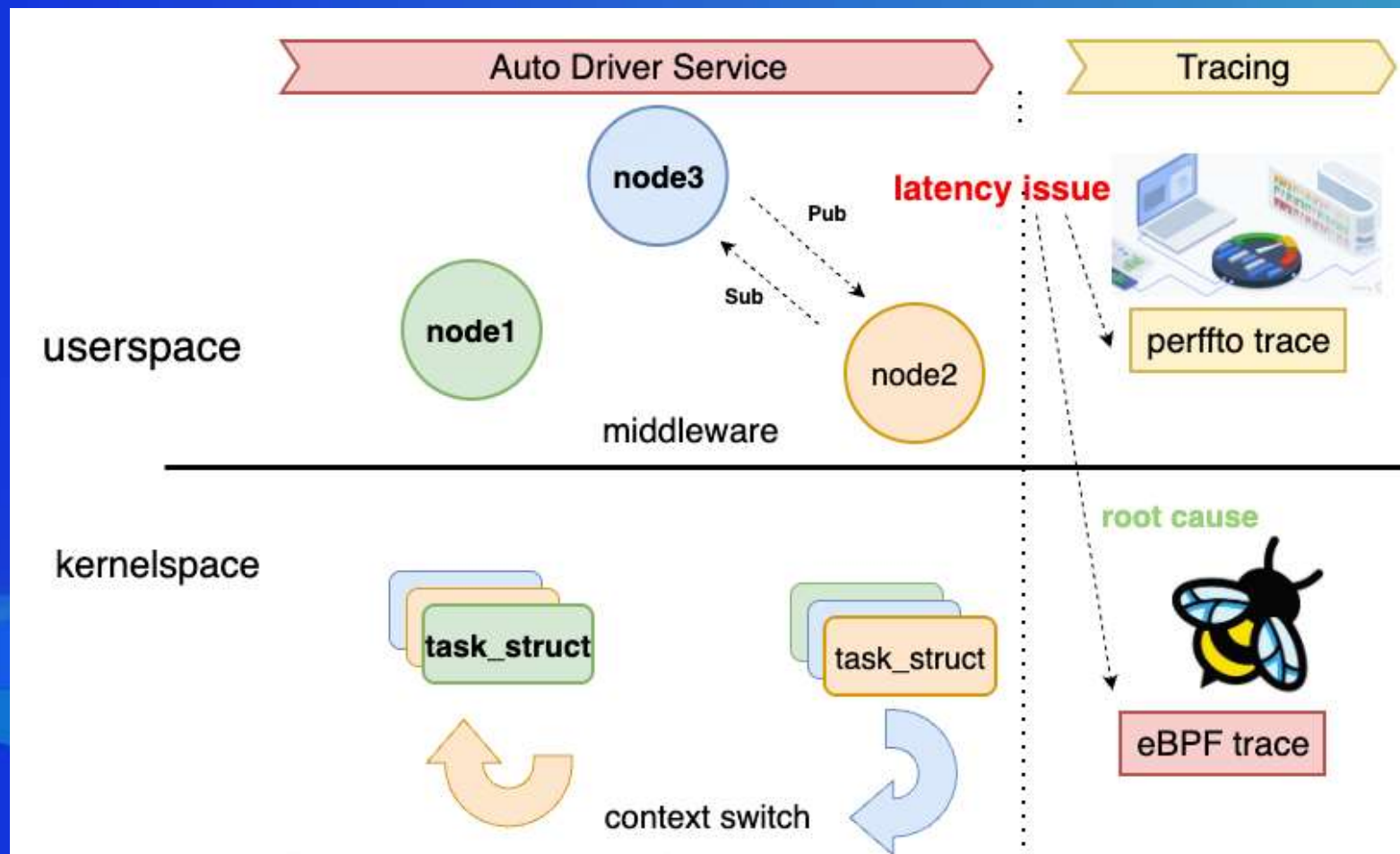


第三届 eBPF开发者大会

[www.ebpftravel.com](http://www.ebpftravel.com)

# 1 项目整体概况

# 1.eBPF使用场景



以调度场景为列

## perffto

- 用户态函数trace, 能看到时延不知道底层原因

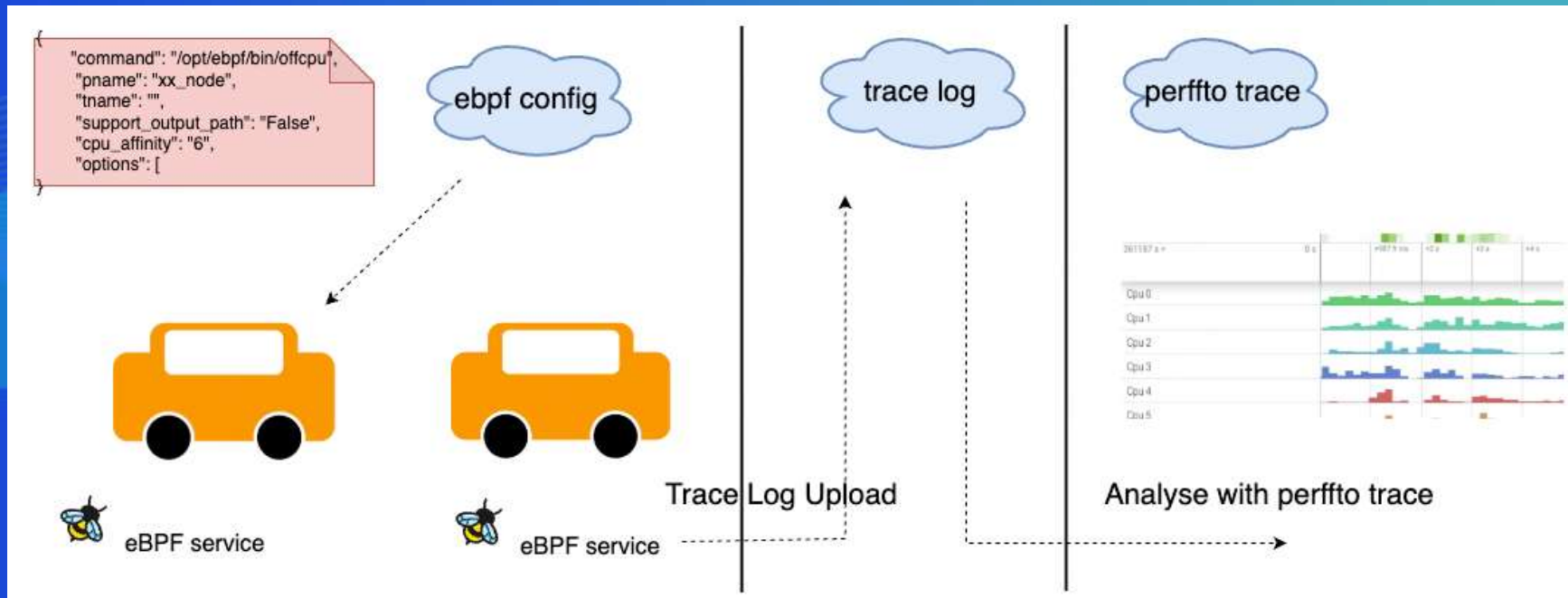
## eBPF

- 内核级trace, 善于定位根因
- 安全、高效、随时可上可下, 不需要修改内核代码

## 2.覆盖的子系统

 调度子系统	 net子系统	 others
offcpu 任务执行时延	netrack 网络时延	top_cache 查看cache占用
syscall 陷入内核时延	tcp_alloc skb分配时延	direct_reclaim 直接内存回收造成的系统hang
 oncpu 实时 profiling	 traffic_stat 进程级流量统计	rcetrace
irq_delay 查看中断时延	tcprrt 网络健康度衡量	systctl_watch 查看配置更改

### 3.车端工具部署

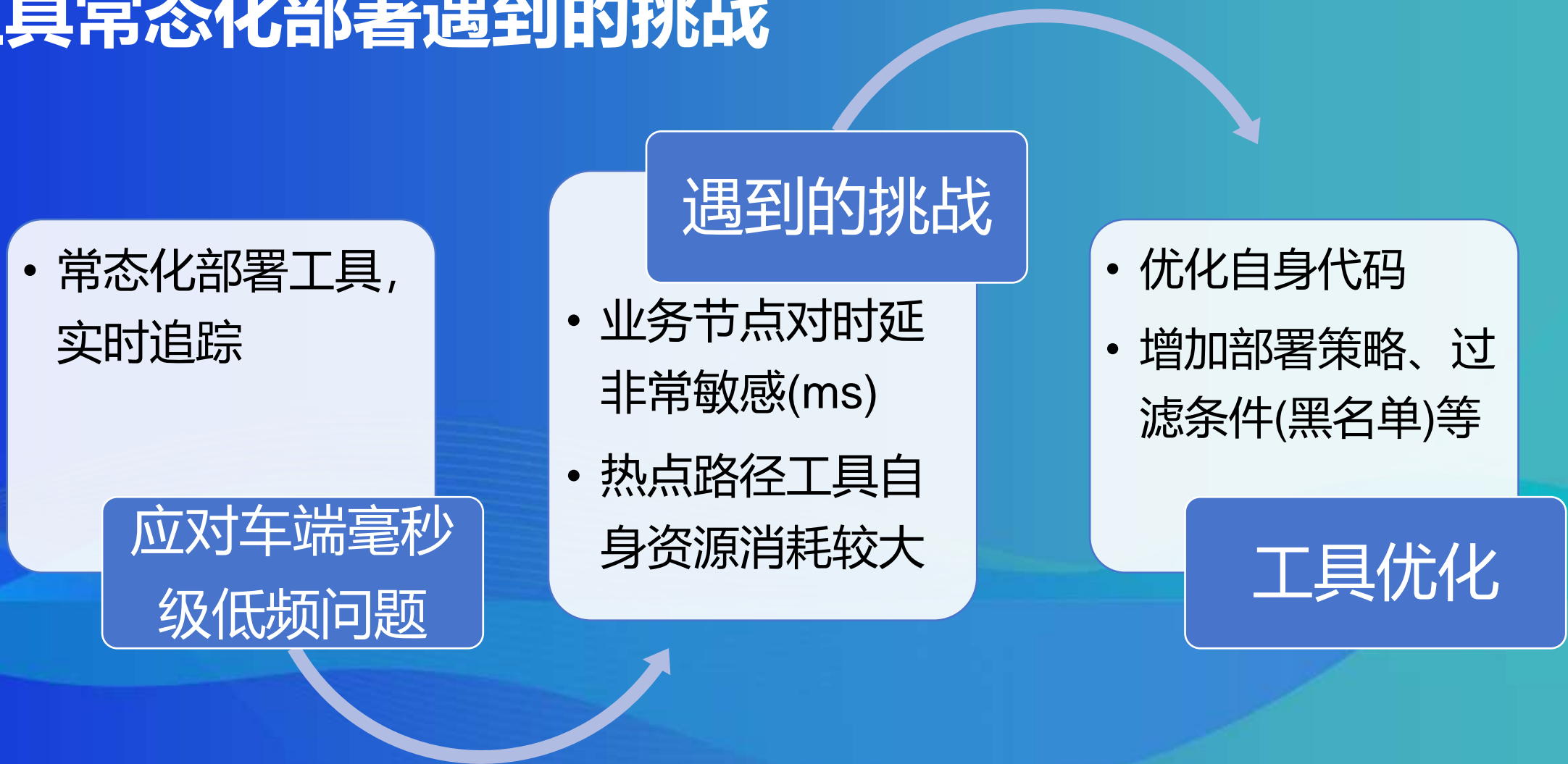




## ② 如何做到低频问题的定位/定界



# 工具常态化部署遇到的挑战



# eBPF编程范式

## 用户态项目框架

### libbpf-bootstrap

- 增加静态编译

### libbpf

- 增加自定义help func  
bpf\_preempt\_ccount等

### bazelsym 栈解析库

- cpu消耗较大替换成bcc库

## 内核程序

- 1.hook目标函数
- 2.将数据存在map
- 3.抓取栈信息(root cause)
- 4.数据输出

```
SEC(kprob/tracepoint/raw_tp/fentry xx)
int prog_handler()
{
    ...
    bpf_map_update()
    bpf_map_lookup()
    ....
    bpf_get_stack
    bpf_perf_event_output
}
```

1

2

3

4

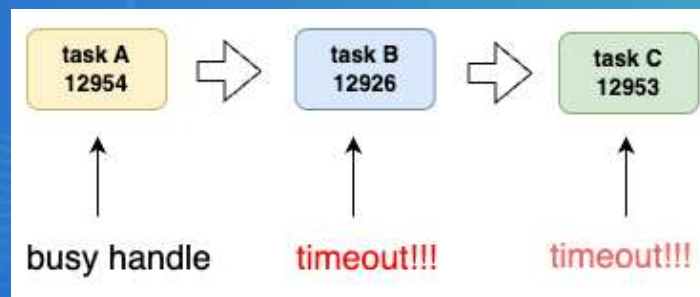
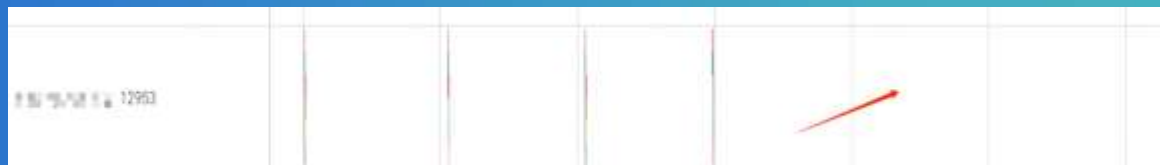
# 栈解析的重要性

利

- 直观呈现问题第一现场，利于问题分析

尚友

- 在bpf\_prog中耗时占大头



## 通过栈信息分析用户态程序等锁唤醒

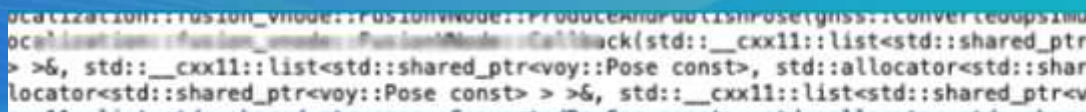


```

[12]w... do_futex+0x8: x64_sys_futex+0xf4: syscall_64+0x1entry_575CALL_64... pthread_c... waker: se...
p+0x8;00_futex+0x8: x64_sys_futex+0xf4: syscall_64+0x1entry_575CALL_64... pthread_c... waker: se...
i:11654,t_tid:12926,waker:11654,w_kid:14474,w_wid:63344,w_oid:11654,w_tid:12954,w_cpu:33,t_w_cpu:20,state:1,ts:
Bus

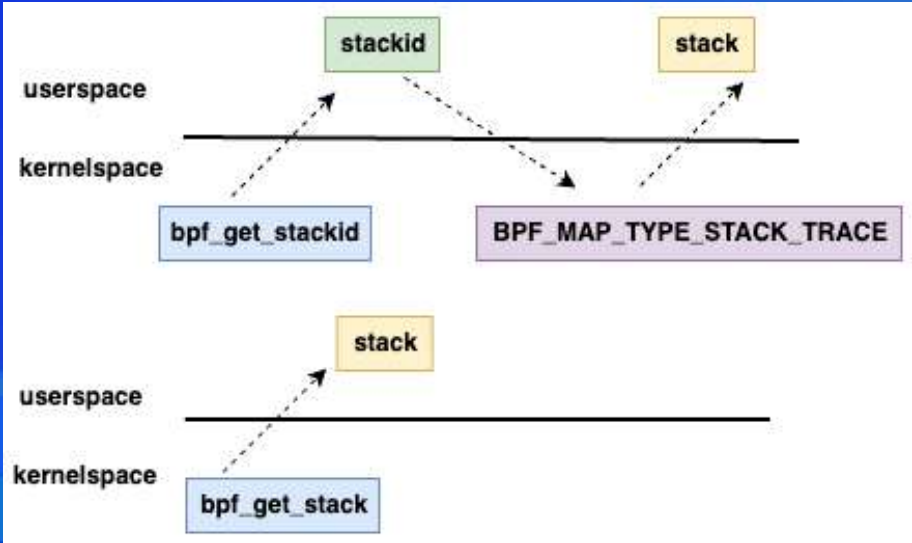
[12]w... stack id:14474,w_w stack id:63344: schedule+0x... do_futex+0x8: x64_sys_fu...
p+0x8... futex+0xf4: syscall_64+0x1entry_575CALL_64... pthread_c... waker: se...
i:11654,t... x62,w_kid:14474,w_wid:63344,w_oid:11654,w_tid:12954,w_cpu:33,state:1,ts:
d:474,t_tid:0000... w_kid:14474,w_wid:63344,w_oid:11654,w_tid:12954,w_cpu:33,state:1,ts:
id:11654,t_tid:12945,w_wid:63344,w_oid:11654,w_tid:12932,w_cpu:32,t_w_cpu:31,state:1,ts:
id:11654,t_tid:12953,waker:11654,w_kid:14474,w_wid:63344,w_oid:11654,w_tid:12926,w_cpu:20,t_w_cpu:30,state:1,ts:
i:11654,t_tid:12932... w_kid:14474,w_wid:63344,w_oid:11654,w_tid:12926,w_cpu:20,t_w_cpu:30,state:1,ts:

```



# 1.采栈优化-采栈时机

bpf\_get\_stack\_id vs bpf\_get\_stack



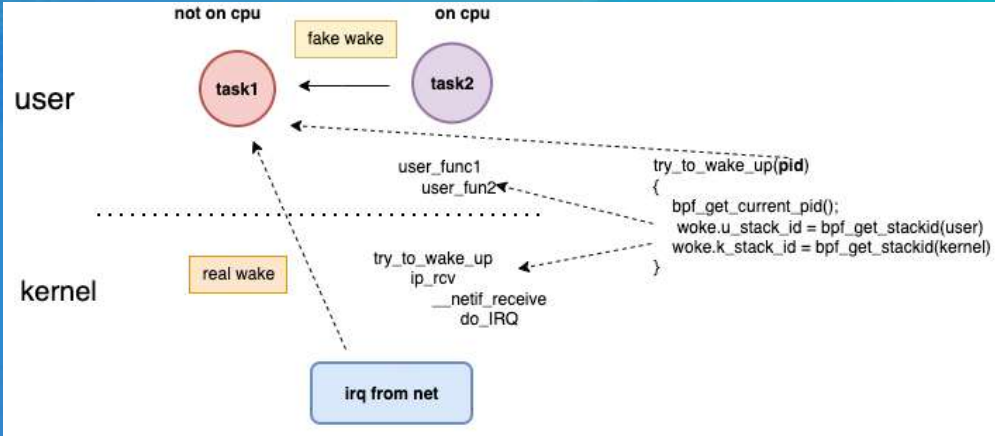
id hash conflict

USER\_STACK vs KERNEL\_STACK

## 1.内核线程上下文不采集用户态栈

```
,waker:kworker/u256:2,w_k_id:4363,w_u_id:-14,w_pid:64889,w_u_len:0
ack_id:4363;w_u_stack_id:-14;finish_task_switch+0x1;schedu
```

## 2.中断上下文不采集用户态栈 30%占比!!!





# 2.函数hook方式优化

kprobe	tracepoint	raw_tracepoint	fentry
1.787M/s	1.873M/s	3.217M/s	3.795M/s

以实际内核版本实测为准

test enviroment

```
Architecture: x86_64
CPU op-mode(s): 32-bit, 64-bit
Address sizes: 48 bits physical, 48 bits virtual
Byte Order: Little Endian
CPU(s): 4
On-line CPU(s) list: 0-3
Vendor ID: AuthenticAMD
BIOS Vendor ID: QEMU
Model name: QEMU Virtual CPU version 2.5+
```



# 3.引入trigger机制

由频繁系统调用引起的sys打高

只在监控指标burst的情况下触发工具使用

节省CPU资源消耗

节省存储资源消耗



```
99 oncpu,k_id:2731,u_id:9938,pid:39784,tid:0,comm:voy_trace_contr,ts:2024-12-04T16:34:16.612890,cpu:3,dur:10.000,prio:120,static_prio:120,normal_prio:120,rt_prio:0,exec_runtime:0,vruntime:5539415789025,nr_running:1
76 <stack>;path_lookupat+0xa5;filename_lookup+0xb6;kern_path+0x2b;unix_find_other+0x45;unix_stream_connect+0xf2;__sys_connect+0xa3;__x64_sys_connect+0xa1a;do_syscall_64+0x5a;entry_SYSCALL_64_after_hwframe+0x44;
71 oncpu,k_id:2731,u_id:9938,pid:39784,tid:0,comm:connect,ts:2024-12-04T16:34:16.612890,cpu:3,dur:10.000,prio:120,static_prio:120,normal_prio:120,rt_prio:0,exec_runtime:0,vruntime:5539415789025,nr_running:1
72 <stack>;connect+0x0;0x676e74746c2f6e75;
73 oncpu,k_id:2731,u_id:9938,pid:39784,tid:0,comm:connect,ts:2024-12-04T16:34:16.612890,cpu:3,dur:10.000,prio:120,static_prio:120,normal_prio:120,rt_prio:0,exec_runtime:0,vruntime:5539415789025,nr_running:1
```



# 3 案例分享

# 实际问题分享

内核中睡眠锁争抢引起时延(mmap\_sem、rtnl\_lock)

CFS公平调度引起的调度时延

高优先级任务抢占低优先级任务(网络包)

内存直接回收(同步等待)

用户程序逻辑问题(老大难!!!)

# 1.mmap\_sem锁争抢

eBPF 工具 1

监控进程 2

业务节点 pidxx 3

/proc/pidxx/maps

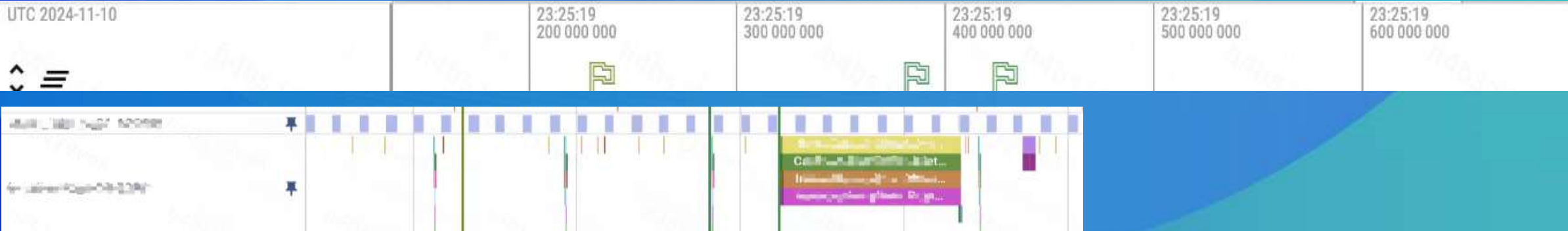
/proc/pidxx/schedstat

Libc系统调用进入

```
static const struct seq_operations proc_pid_maps_op = {
    .start = m_start,
    .next  = m_next,
    .stop  = m_stop,
    .show  = show_map
};
```

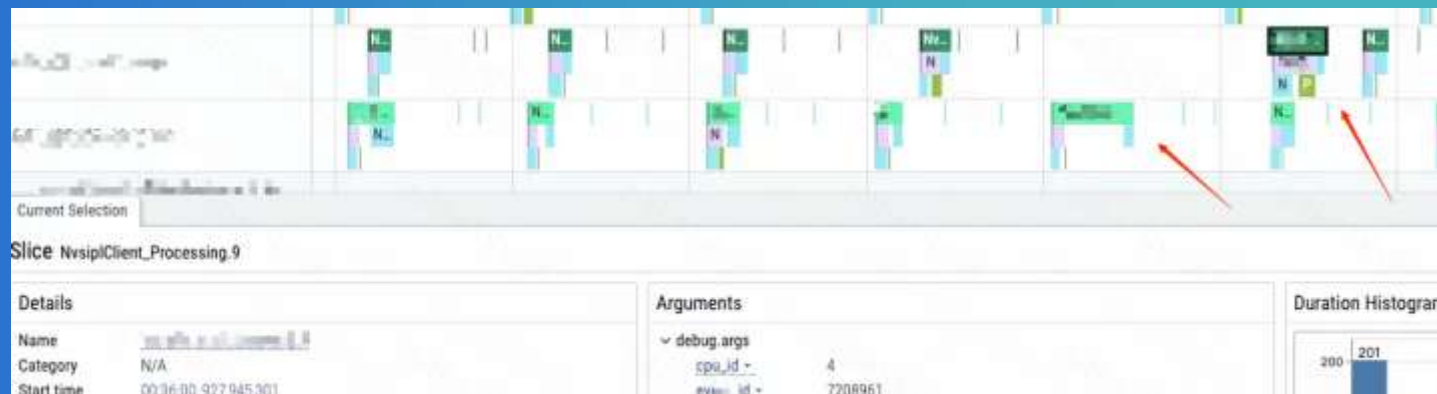
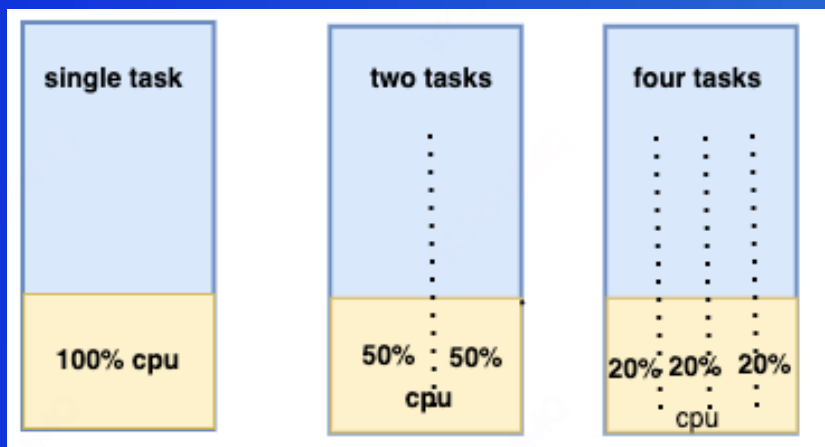
down\_read\_killable(&mm->mmap\_sem);

```
sys_mprotect
do_mprotect_pkey
down_write_killable(&current->
mm->mmap_sem)
```

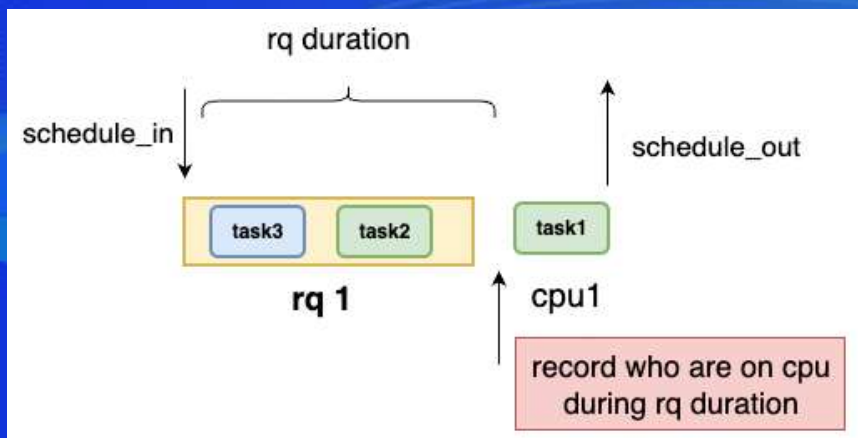


## 2.CFS公平调度引起的调度时延

### CFS虚拟时间补偿机制



在rq长的时刻工具(offcpu)抓取  
哪些任务在占用CPU



```
state:1,ts:2024-12-27T08:36:00.785,dur:113.97,q_dur:85500.00us,rq_dur:85478.00us,r
pid:11078,comm:collected,prio:120,cpu:4,state:1,id:275,ts:2024-12-27T08:36:00.879972791,ser:6179269376,vr:4519552887
pid:11077,comm:collected,prio:120,cpu:4,state:1,id:280,ts:2024-12-27T08:36:00.880273781,ser:6204805440,vr:4519554258
pid:11076,comm:collected,prio:120,cpu:4,state:1,id:285,ts:2024-12-27T08:36:00.880472179,ser:6135801760,vr:4519553740
pid:11077,comm:collected,prio:120,cpu:4,state:1,id:290,ts:2024-12-27T08:36:00.881188718,ser:6204838176,vr:4519554585
pid:52,comm:ksoftirqd/4,prio:120,cpu:4,state:1,id:295,ts:2024-12-27T08:36:00.881605195,ser:24050739616,vr:451953305
pid:11079,comm:collected,prio:120,cpu:4,state:0,id:300,ts:2024-12-27T08:36:00.881899337,ser:6227900416,vr:4519555936
pid:52,comm:ksoftirqd/4,prio:120,cpu:4,state:1,id:305,ts:2024-12-27T08:36:00.882149447,ser:24050754368,vr:451953319
pid:11079,comm:collected,prio:120,cpu:4,state:0,id:310,ts:2024-12-27T08:36:00.882404965,ser:6228023200,vr:4519557164
pid:11079,comm:collected,prio:120,cpu:4,state:1,id:315,ts:2024-12-27T08:36:00.882614084,ser:6228095136,vr:4519557883
pid:817,comm:irq/90-eth0.rx0,prio:49,cpu:4,state:2,id:320,ts:2024-12-27T08:36:00.883001921,ser:582408730080,vr:0
pid:11076,comm:collected,prio:120,cpu:4,state:1,id:325,ts:2024-12-27T08:36:00.883218527,ser:6135843392,vr:4519554157
```

```
target:...,t_k_id:30141,t_u_id:22662,t_pid:15774,t_tid:16418,waker...,ing,w_k_id:51166,w_u_id:25970,w_pid:156
8.01,q_dur:5.00us,rq_dur:0.00us,rq_len:0,ser:2257056959,vr:3700378228307,w_ser:5484873547,w_dur:38.74us,w_ldur:39.47us
```

## 4 未来展望

# 未来展望

## 进一步优化工具，覆盖更多场景

- 使用更多高版本内核特性: fentry task\_storage ringbuffer等
- 特定场景下的优化: map perfoutput

## 平台化、服务化、可视化、自动化

## 引入大模型做trace分析及性能优化





第三届 eBPF开发者大会

[www.ebpftravel.com](http://www.ebpftravel.com)

# 谢谢!

- [email: chen.dylane@linux.dev](mailto:chen.dylane@linux.dev)