

bpftrace 速查表

smallest

1 基本语法

概念	描述
探针	格式: provider:probe
Action	代码在 { ... } 中
内建变量	\$1, \$2 (参数), pid, tid, uid, gid, comm, cpu, curtask, rand args, arg0, arg1, ..., argN (函数参数)
Maps	@name[key] = value;
过滤探针	/filter/ { actions }

2 常见探针

探针	描述
kprobe:function	内核函数入口
kretprobe:function	内核函数返回
uprobe:/path/to/binary:function	用户空间函数入口
uretprobe:/path/to/binary:function	用户空间函数返回
tracepoint:subsystem:event	静态跟踪点
profile:hz:rate	周期采样
interval:ms:rate	定时器
software:event:count	软件事件 (例如, cpu-clock, task-clock)
hardware:event:count	硬件事件 (例如, cache-misses, cpu-cycles)

3 列出探针

命令	描述
bpftrace -l	列出所有探针
bpftrace -l "tracepoint:*"	列出所有 tracepoint 探针
bpftrace -l "kprobe:vfs_*"	列出所有 vfs 相关的 kprobe 探针

`bpftrace -lv "tracepoint:syscalls:sys_enter_*` 列出所有 `sys_enter` 相关的 `tracepoint` 探针及其参数

4 内置函数

函数	描述
<code>printf(format, args...)</code>	打印格式化字符串
<code>time(format)</code>	当前时间，格式化输出
<code>count()</code>	计数器
<code>sum(int n)</code>	求和
<code>avg(int n)</code>	平均值
<code>min(int n)</code>	最小值
<code>max(int n)</code>	最大值
<code>hist(int n)</code>	直方图
<code>lhist(int n, int min, int max, int step)</code>	线性直方图
<code>kstack()</code> , <code>ustack()</code>	内核和用户栈跟踪
<code>ntop(int n)</code>	转换 IP 地址为字符串
<code>reg(char *name)</code>	读取 CPU 寄存器值

5 高级特性

特性	描述
Wildcards	在探针定义中使用 <code>*</code> (例如, <code>kprobe:vfs_*</code>)
Frequency Counting	<code>count()</code> 可以与映射结合进行频率分析
Associative Arrays	<code>@map[key1, key2, ...] = value;</code>
Strings	使用 <code>str()</code> 转换为字符串, <code>strcmp()</code> 进行比较
Aggregations	<code>clear(@map)</code> , <code>print(@map)</code> , <code>zero(@map)</code>
Join	<code>join(char *arr[])</code> 连接数组元素
Timestamps	<code>nsecs</code> , <code>elapsed</code>

6 示例

6.1 系统调用计数

```
#!/usr/bin/env bpftrace

tracepoint:raw_syscalls:sys_enter
{
```

```

    @syscalls[comm] = count();
}

interval:s:5
{
    print(@syscalls);
    clear(@syscalls);
}

```

6.2 函数延迟测量

```

#!/usr/bin/env bpftool

kprobe:vfs_read
{
    @start[tid] = nsecs;
}

kretprobe:vfs_read
/ @start[tid] /
{
    @duration = hist(nsecs - @start[tid]);
    delete(@start[tid]);
}

interval:s:10
{
    print(@duration);
    clear(@duration);
}

```

6.3 内存分配跟踪

```

#!/usr/bin/env bpftool

uprobe:/lib/x86_64-linux-gnu/libc.so.6:malloc
/comm == "target_process"/
{
    @bytes[ustack] = sum(arg0);
}

interval:s:30
{
    print(@bytes);
    clear(@bytes);
}

```

7 一行代码示例

bpfttrace 一行代码	解释
<code>bpfttrace -e 'tracepoint:syscalls:sys_enter_* { @[probe] = count(); }'</code>	统计所有系统调用的次数
<code>bpfttrace -e 'kprobe:vfs_read { @bytes = sum(arg2); }'</code>	累计读取的字节数
<code>bpfttrace -e 'kprobe:vfs_write { @bytes[comm] = sum(arg2); }'</code>	按进程名统计写入的字节数
<code>bpfttrace -e 'tracepoint:block:block_rq_issue { @[args->dev] = count(); }'</code>	统计每个块设备的 I/O 请求数
<code>bpfttrace -e 'software:major-faults:1 { @[comm] = count(); }'</code>	统计每个进程的主缺页次数
<code>bpfttrace -e 'kprobe:do_nanosleep { @[comm] = count(); }'</code>	统计每个进程的睡眠次数
<code>bpfttrace -e 'tracepoint:sched:sched_process_exec { @[args->filename] = count(); }'</code>	统计每个可执行文件的执行次数
<code>bpfttrace -e 'kprobe:tcp_retransmit_skb { @[sock_ip(args->sk)] = count(); }'</code>	统计每个 IP 地址的 TCP 重传次数
<code>bpfttrace -e 'tracepoint:irq:irq_handler_entry { @[args->name] = count(); }'</code>	统计每个中断处理程序的调用次数
<code>bpfttrace -e 'kprobe:kmalloc { @bytes = hist(arg1); }'</code>	创建内核内存分配大小的直方图
<code>bpfttrace -e 'kprobe:do_sys_open { @[comm] = count(); }'</code>	统计每个进程的文件打开次数
<code>bpfttrace -e 'uprobe:/lib/x86_64-linux-gnu/libc.so.6:malloc { @bytes = hist(arg0); }'</code>	创建用户空间内存分配大小的直方图
<code>bpfttrace -e 'profile:hz:99 @[comm] = count(); '</code>	每秒采样 99 次，并统计每个进程名的调用次数
<code>bpfttrace -e 'interval:s:5 printf("Every 5 seconds"); '</code>	每 5 秒打印一次“Every 5 seconds”

8 最佳实践

- 使用过滤器减少开销并关注相关事件
- 结合多个探针以关联事件
- 使用映射和聚合进行高效数据收集
- 小心使用栈跟踪和字符串操作（高开销）
- 首先在非生产系统上测试脚本
- 使用 interval 探针定期输出和清除数据