

Linux Tracing Embedded Systems

Marcel Ochsendorf, Muhammed Parlak

27.06.2022

Inhaltsverzeichnis

- **Einleitung**
- **Grundlagen**
- **Tools**
- **Live Reverse Engineering**

Einleitung

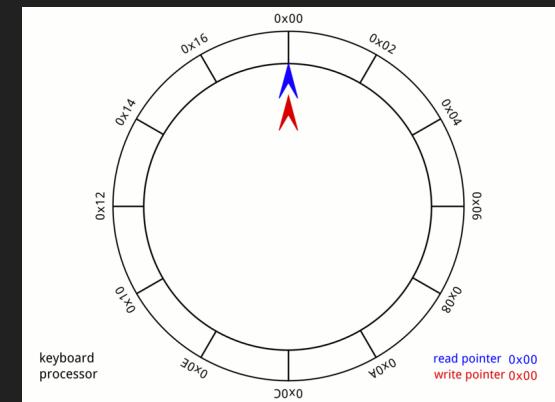
- Tracing?
- Tracing Systeme
- Linux Tracing
- Warum sollte man Tracen?

Grundlagen

- Ringbuffer
- Debug-Filesystem
- Tracing
 - Tracer
 - Events
- Auffangen von Events
- Probes
- Ressourcen

Grundlagen | Ringbuffer

- Basis für das Debug-FS
- Asynchrones Lesen und Schreiben
- Puffer mit Lesezeiger und Schreibzeiger
- Für jeden Core



Grundlagen | Debug-Filesystem

- Eingeführt mit Kernel 2.6.10-rc3
- Notwendig für die Tracing-Funktionalität
- Zugriff auf Diagnose und Debug-Informationen

```
1 # ENABLE DEBUG FS
2 $ sudo mount -t debugfs debugfs /sys/kernel/debug
3 # PRINT FOLDERS OF DEBUG FS
4 $ ls -lah /sys/kernel/debug | awk '{print $9}'
5 hid
6 usb
7 Tracing
8 [...]
```

Grundlagen | Debug-Filesystem

- Gezieltere Aggregierung von Events (z.B. Filterung)

```
1 # GET TRACERS
2 $ cat /sys/kernel/debug/tracing/available_tracers
3 hwlat blk mmiotrace function_graph wakeup_dl wakeup_rt wakeup
   function nop
4 # USE SPECIFIC TRACER
5 $ echo function_graph > /sys/kernel/debug/tracing/
   current_tracer
6 # DISABLE TRACER USAGE
7 $ echo nop > /sys/kernel/debug/tracing/current_tracer
```

Grundlagen | Debug-Filesystem

- Gezieltere Aggregierung von Events (z.B. Filterung)

```
1 # CALL STACK USING FUNCTION_GRAPH TRACER
2 $ echo function_graph > /sys/kernel/debug/tracing/
   current_tracer
3 $ cat /sys/kernel/debug/tracing/trace
4 # tracer: function_graph
5 # CPU- DURATION          FUNCTION CALLS
6 # |      |      |          |      |
7 0)           |          |      |      |      |
8 0) 4.442 us  |          |      |      |      |
9 0) 1.382 us  |          |      |      |      |
10 0) 2.478 us |          |      |      |      |
11 [...]
```

Grundlagen | Tracing Events

- Auslösung sehr unterschiedlich
 - Beispiel Lesen und Schreiben von I2C Bus
 - Scheduler-Switching
 - Disk-Events (z.B. Dateioperationen)
 - IO-Events
- Nach Aktivierung
 - Informationen zu Event werden geloggt

Grundlagen | Tracing Events

```
1 # GET AVAILABLE EVENT LIST
2 $ cd /sys/kernel/debug/tracing/events
3 $ ls -lah | awk '{print $9}'
4 alarmtimer
5 drm
6 exceptions
7 ext4 #READPAGE, WRITEPAGE, ERROR, FREE_BLOCKS
8 filelock
9 filemap
10 gpio #GPIO_DIRECTION, GPIO_VALUE
11 hda
12 i2c
13 irq
14 net
15 smbus #READ, WRITE, REPLY
16 sock #STATE_CHANGED, EXCEED_BUFFER_LIMIT, REC_QUEUE_FULL
17 spi
18 tcp
19 timer #TIMER_STOP, TIMER_INIT, TIMER_EXPIRED
20 [...]
```

```
1 $ cd /sys/kernel/debug/tracing/events/ext4
2 $ ls -lah | awk '{print $9}'
3 # EVENTS FOR EXT4
4 ext4_write_end
5 ext4_writepage
6 ext4_readpage
7 ext4_error
8 [...]
9 # INTERFACE FOR EVENT SETUP
10 enable
11 filter
12 format
```

Grundlagen| Abfangen von Events **EVENT ACTIVATION**

```
1 $ cd /sys/kernel/debug/tracing/events/ext4
2 # ENABLE ALL EVENTS FROM THIS GROUP
3 $ echo 1 > ./enable
4 # DISBALE ALL EVENTS
5 $ echo 0 > ./enable
6 # ENBABLE SPECIFIC EVENT
7 $ echo 1 > ./ext4_readpage/enable
8 $ echo 1 > ./ext4_writepage/enable
9 # GET EVENT LOG
10 $ cat /sys/kernel/debug/tracing/trace
```

Grundlagen| Auffangen von Events **EVENT OUTPUT**

```
# tracer: nop
#
# entries-in-buffer/entries-written: 5394/5394    #P:4
#
#                                     ----=> irqs-off
#                                     /----=> need-resched
#                                     | /----=> hardirq/softirq
#                                     || /----=> preempt-depth
#                                     ||| /---=> delay
# TASK-PID   CPU#  ||||  TIMESTAMP      FUNCTION
#                 |   |  |  |  |
bash-13401 [001] d... 1928821.587704: sched_waking: comm=kworker/u16:1 pid=13198 prio=120 target_cpu=002
bash-13401 [001] d... 1928821.587712: sched_wake_idle_without_ipi: cpu=2
bash-13401 [001] d... 1928821.587713: sched_wakeup: comm=kworker/u16:1 pid=13198 prio=120 target_cpu=002
```

The diagram illustrates the structure of the event output shown above. Five vertical arrows point from labels to specific fields in the log lines:

- An arrow points from "Event name" to the first column of the log, which contains process names like "bash-13401".
- An arrow points from "Timestamp" to the second column, which shows timestamps such as "1928821.587704".
- An arrow points from "CPU on which the event arrived" to the third column, which shows CPU numbers like "[001]".
- An arrow points from "Active process when the event arrived" to the fourth column, which shows the full command-line arguments of the processes.
- An arrow points from "Event specific information" to the fifth column, which contains function names and additional parameters like "comm=kworker/u16:1 pid=13198 prio=120 target_cpu=002".

Grundlagen Probes **Kprobes**

- Sammeln von Laufzeit und Performance-Daten des Kernels
- Vorteil:
 1. Daten ohne Unterbrechung der Ausführung auf Prozessor-Instruktion-Ebene aggregieren
 2. Registrierung dynamisch zur Laufzeit und ohne Änderungen des Programmcodes

Grundlagen | Probes **User-Level-Probes**

- Weiterentwicklung von kprobes
- Registrierung von Laufzeit Events in Applikation
- Adressoffset wegen eigenem virtuellen Adressraum

Grundlagen | Probes User-Level-Probes **FUNCTION TRACING**

```
1 //test.c
2 #include <stdio.h>
3 int main(void)
4 {
5     int i;
6     for (i = 0; i < 5; i++)
7         printf("Hello uprobe\n");
8     return 0;
9 }
```

```
1 # BUILD APPLICATION
2 $ gcc ./test.c -o /tmp/test
3 # GET OFFSET
4 $ objdump -F -S -D /tmp/test | less | grep main
5 0000000000001149 <main> (File Offset: 0x1149):
6 # REGISTER uprobe_event
7 $ echo "p:my_uprobe /tmp/test:0x1149" > /sys/kernel/debug/
       tracing/uprobe_events
8 # ACTIVATE UPROBE EVENTS
9 $ echo 1 > /sys/kernel/tracing/events/uprobes/enable
10 # EXECUTE PROGRAM
11 $ /tmp/test
12 Hello uprobe
13 [...]
14 # PRINT TRACED EVENTS
15 $ cat /sys/kernel/debug/tracing/trace
16 # tracer: nop
17 # TASK-PID  CPU#  TIMESTAMP  FUNCTION
18 # |          |          |          |
19 test-24842 [012] 258544.995456: printf: [...]
20 [...]
```

Grundlagen | Ressourcen

- Zusätzliche Ressourcen werden benötigt
- Unerhebliche Menge an Speicherplatz
- Bandbreite vom genutzten Medium
- Betrieb und Aufzeichnung wird beeinflusst

Was sollte getan werden?

- Aufzeichnung separat
- Nur Events die im Fokus stehen

Tools

- trace-cmd
- bpftrace
- Kernelshark

Tools **trace-cmd**

- Meistgenutztes Tool zum aufzeichnen von Tracelogs
- Aufzeichnung von Events spezifizierter Anwendungen

Tools | trace-cmd USAGE

```
1 # CHECK IF TRACING IS ENABLED
2 $ sudo mount | grep tracefs
3 none on /sys/kernel/tracing type tracefs (rw,relatime,seclabel
    )
4 ## ONLY SCHEDULER EVENTS
5 $ echo sched_wakeup >> /sys/kernel/debug/tracing/set_event
6 ## ALL EVENTS USING set_event
7 $ echo *:* > /sys/kernel/debug/tracing/set_event
8 # RECORD
9 $ trace-cmd record ./program_executable
10 # RECORD SPECIFIC EVENT
11 $ trace-cmd record -e sched ./program_executable
12 # USING A ADDITIONAL TRACER
13 $ trace-cmd -t function ./program_executable
```

Tools | trace-cmd OUTPUT

```
cpus=4
sleep-19405 [000] 915.755341: sched_stat_runtime:      comm=trace-cmd pid=19405 runtime=168772 [ns] vruntime=124673167878 [ns]
sleep-19405 [000] 915.755344: sched_switch:           trace-cmd:19405 [120] D ==> swapper/0:0 [120]
<idle>-0   [003] 915.759205: sched_waking:          comm=rcu_sched pid=10 prio=120 target_cpu=003
<idle>-0   [003] 915.759210: sched_wakeup:          rcu_sched:10 [120] success=1 CPU:003
<idle>-0   [003] 915.759212: sched_waking:          comm=kworker/3:1 pid=2850 prio=120 target_cpu=003
<idle>-0   [003] 915.759214: sched_wakeup:          kworker/3:1:2850 [120] success=1 CPU:003
<idle>-0   [003] 915.759226: sched_switch:          swapper/3:0 [120] R ==> kworker/3:1:2850 [120]
kworker/3:1-2850 [003] 915.759231: sched_stat_runtime: comm=kworker/3:1 pid=2850 runtime=14167 [ns] vruntime=60565680765 [ns]
kworker/3:1-2850 [003] 915.759233: sched_switch:       kworker/3:1:2850 [120] W ==> rcu_sched:10 [120]
rcu_sched-10 [003] 915.759238: sched_stat_runtime:    comm=rcu_sched pid=10 runtime=7532 [ns] vruntime=60565679935 [ns]
rcu_sched-10 [003] 915.759246: sched_switch:          rcu_sched:10 [120] W ==> swapper/3:0 [120]
<idle>-0   [002] 915.762699: sched_waking:          comm=apt-get pid=2554 prio=120 target_cpu=002
<idle>-0   [002] 915.762706: sched_wakeup:          apt-get:2554 [120] success=1 CPU:002
<idle>-0   [002] 915.762709: sched_waking:          comm=kworker/2:3 pid=389 prio=120 target_cpu=002
```

Tools | **bpftrace**

- Kernelversion >4.x
- Eigene Skriptsprache
 - Aggregation
 - Eventfilter
 - Verarbeitung der Ergebnisse

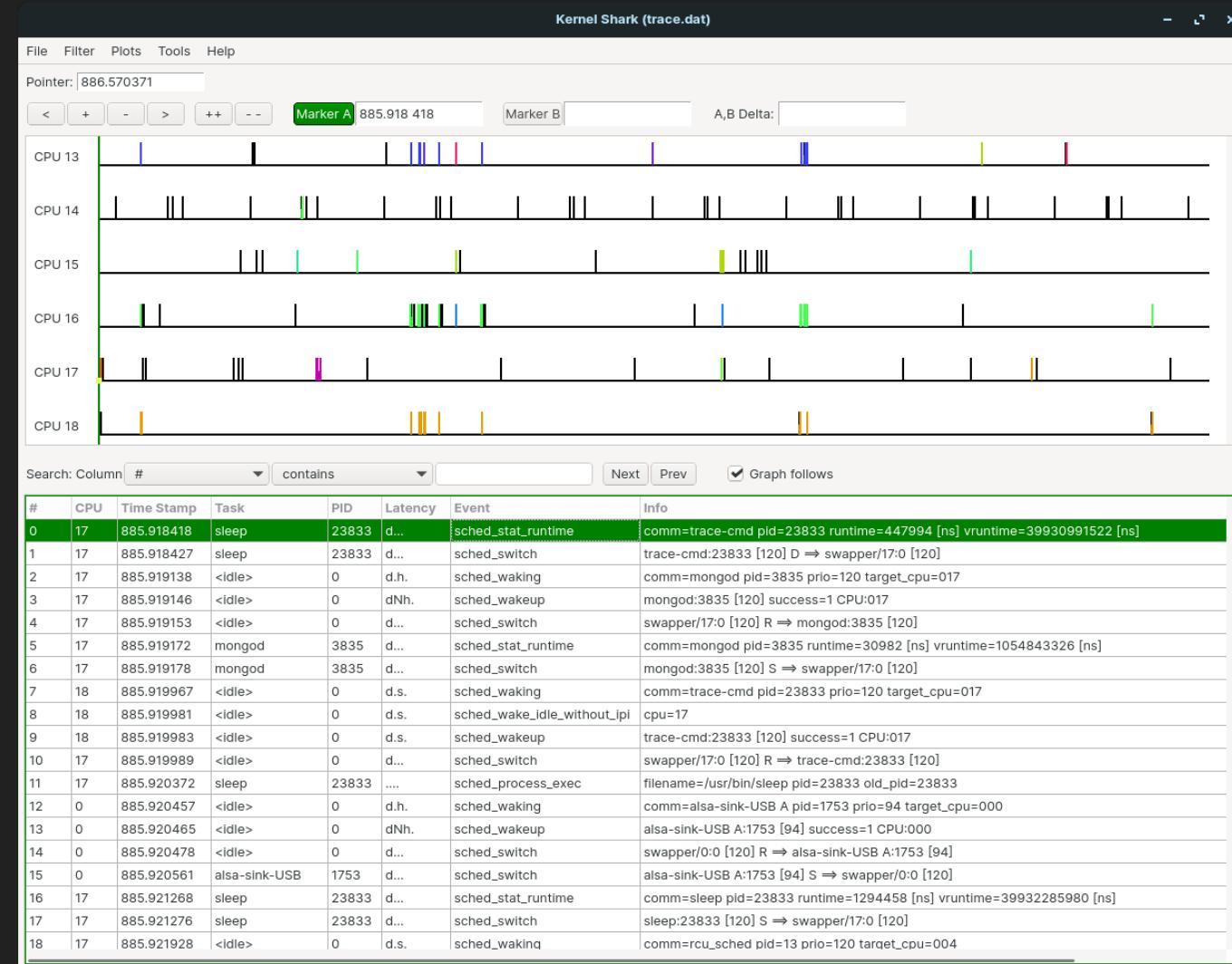
Tools | bpftrace SCRIPTING

```
50
51 kprobe:tcp_drop
52 {
53     $sk = ((struct sock *) arg0);
54     $inet_family = $sk->__sk_common.skc_family;
55
56     if ($inet_family == AF_INET || $inet_family == AF_INET6) {
57         if ($inet_family == AF_INET) {
58             $daddr = ntop($sk->__sk_common.skc_daddr);
59             $saddr = ntop($sk->__sk_common.skc_rcv_saddr);
60         } else {
61             $daddr = ntop($sk->__sk_common.skc_v6_daddr.in6_u.u6_addr8);
62             $saddr = ntop($sk->__sk_common.skc_v6_rcv_saddr.in6_u.u6_addr8);
63         }
64         $lport = $sk->__sk_common.skc_num;
65         $dport = $sk->__sk_common.skc_dport;
66
67         // Destination port is big endian, it must be flipped
68         $dport = bswap($dport);
69
70         $state = $sk->__sk_common.skc_state;
71         $statestr = @tcp_states[$state];
72
73         time("%H:%M:%S ");
74         printf("%-8d %-16s ", pid, comm);
75         printf("%39s:-%6d %39s:-%6d %-10s\n", $saddr, $lport, $daddr, $dport, $statestr);
76         printf("%s\n", kstack);
77     }
78 }
79
80 END
81 {
82     clear(@tcp_states);
83 }
```

Tools | Kernelshark

- Input von Trace-cmd
- Grafische Visualisierung

Tools | Kernelshark VISUALISATION



Reverse Engineering **DEMO**

LIVE DEMO