

# Tracen mit dem Linux-Kernel

Jürgen Stuber

2017-04-17



ftrace



perf\_events



eBPF



SystemTap



LTTng



ktap



dtrace4linux



OEL DTrace



sysdig

- System im laufenden Betrieb beobachten
- Debuggen von Performance-Problemen
- Debuggen allgemein

ohne zu viel Overhead an die relevanten Informationen kommen

- über > 15 Jahre historisch entwickelt
- ein Etikett auf mehreren Dingen aus einer Quelle
  - ftrace
    - function tracer
    - Infrastruktur (Ringpuffer, debugfs-Anbindung)
- mehrere Namen für das Gleiche im Laufe der Entwicklung
  - perf\_events
  - perf
  - Performance Counters for Linux (PCL)
  - Linux Perf Events (LPE)
- Noch ziemlich neu, größere Erweiterungen in Linux 4.x

- Einsprungpunkte für das Tracen (Probes)
- mögliche Aktionen (Events)
- Ausgabe und Analyse von Daten

## Ziele

- geringer Overhead
- flexibel

## Hauptvarianten

- ftrace
- tracepoints, \*probes

# ftrace

ftrace: ein Tracer für viele Funktionen gleichzeitig  
alle Funktionen möglich (außer durch inlining verschwundene)

Tracer sind vorgegeben:

nop: abgeschaltet

function: Loggt Einsprung in Funktion, mit Aufrufer

function\_graph: Loggt Einsprung und Rücksprung

gut für die Beobachtung des Programmablaufs

Implementierung:

GCC-Profiler-Aufrufe, wird bei Boot durch NOPs gepatcht

deaktiviert  $\rightsquigarrow$  Overhead vernachlässigbar

vor Aktivieren erst Filter für Funktionen setzen

Demo

- Event erzeugen (Event-Infrastruktur, s.u.)
- eigenen Code aufrufen
  - registrieren durch Kernel-Modul
  - Live-Patching möglich

ca 500 feste Tracepoints im Kernel-Code

erzeugen Events

NOPs wenn deaktiviert ("static keys")

werden über Macros im Kernel-Code definiert:

TRACE\_EVENT: definiert Tracepoint & Event (Normalfall im Kernel) DECLARE\_TRACE: nur Tracepoint

eigentlicher Tracepoint ruft Funktion auf:

```
static inline void trace_##name()
```



kprobe auf jeder Maschineninstruktion im Kernel möglich

ersetzt Instruktion durch Breakpoint

single-step durch Originalinstruktion anderswo ( $\sim 1 \mu\text{s}$ )

Interrupts abgeschaltet

Optimierung durch Jump wo möglich ( $\sim 100 \text{ ns}$ )

außer Tracen auch Live-Patching möglich

in Rohform schwer zu bedienen  $\rightsquigarrow$  Tools wie SystemTap

uprobe: Wie kprobe aber im User-Space

kretprobe, uretprobe: Rücksprung zu gegebener Funktion

Implementierung:

ersetzt Rücksprungadresse durch Trampolin-Adresse trampoline macht Probe  
Sprung zu eigentlicher Rücksprungadresse

jprobe: einfacher Zugriff auf Funktionsargumente

```
mount -t debugfs debugfs /sys/kernel/debug Tracing in /sys/kernel/debug/tracing Ausgabe  
/sys/kernel/debug/tracing/trace /sys/kernel/debug/tracing/trace_pipe
```

- debugfs: `/sys/kernel/debug/tracing/events`
- Events im Kernel:
  - ca 500 tracepoints vorgegeben
  - ca 500 System-Call Enter/Exit
  - kprobe-Events mit debugfs definierbar
  - einzeln aktivier- und konfigurierbar
  - debugfs: `events/<subsystem>/<event>/`

- diverse aktionen konfigurierbar in `events/<subsystem>/<event>/trigger`
  - default: Daten in Ringpuffer speichern
  - `traceon/traceoff`: Trace ein/aus
  - `enable_event/disable_event`: anderen Event ein/aus
  - `stacktrace`
  - `snapshot`: bei Problemsituation Schnappschuss, dann Trace analysieren
  - `hist` (ab 4.7): Histogramme von Daten bei Events
  - `enable_hist/disable_hist`: Histogramm für anderen Event ein/aus
- Filter
  - Beispiel: Filtern nach PID
- Anzahl der Ausführungen begrenzen
- `Documentation/trace/events.txt`

- Kommandozeilentool: `trace-cmd`  
dünnere Wrapper um `/sys/kernel/debug/tracing`
- `perf`: Tool, um Performance Counter der CPU zu lesen  
Profiling:  
Welcher Code verbraucht die Rechenzeit?  
Effekte von Caching/Branch-Prediction usw. analysieren  
Beispiel: `$ perf stat <command>`
- eBPF: Berkeley Packet Filter (BPF)  $\rightsquigarrow$  VM im Kernel  
Datenanalyse schon im Kernel  
spart Kontext-Switches  
sicherer als C-Code  
ab Linux 4.9 vollständig (Funktionsumfang wie Sun DTrace)

- Documentation/trace
- Brendan Gregg (Beispiele, Marketing)  
<http://www.brendangregg.com/blog/index.html>  
  
Ponies by Deirdré Straughan and General Zoi's pony creator

# Ende