

# L41: Kernels and Tracing

Dr Robert N. M. Watson

15 October 2015

## Reminder: last time

1. What is an operating system?
2. Systems research
3. About the module
4. Lab reports

# This time: Tracing the kernel

1. DTrace
2. The probe effect
3. The kernel: Just a C program?
4. A little on kernel dynamics: How work happens

## Dynamic tracing with DTrace

- ▶ Bryan M. Cantrill, Michael W. Shapiro, and Adam H. Leventhal. *Dynamic Instrumentation of Production Systems*, USENIX Annual Technical Conference, USENIX, 2004.
  - ▶ “Facility for dynamic instrumentation of production systems”
  - ▶ Unified and safe instrumentation of kernel and userspace
  - ▶ Zero *probe effect* when not enabled
  - ▶ Dozens of ‘providers’ representing different trace sources
  - ▶ Tens (hundreds?) of thousands of instrumentation points
  - ▶ C-like high-level control language with predicates and actions
  - ▶ User-defined variables, thread-local variables, associative arrays
  - ▶ Data aggregation and speculative tracing
- ▶ Adopted in Solaris, Mac OS X, and FreeBSD; module for Linux
- ▶ Heavy influence on Linux SystemTap
- ▶ **Our tool of choice for this module**

## DTrace scripts

- ▶ Human-facing C-like language
- ▶ One or more  $\{probe\ name, predicate, action\}$  tuples
- ▶ Expression limited to control side effects (e.g., no loops)
- ▶ Specified on command line or via a `.d` file

```
fbt::malloc:entry /execname == "csh"/ { trace(arg0); }
```

**probe name** Identifies the probe(s) to instrument; wildcards allowed;  
identifies the *provider* and a provider-specific *probe name*

**predicate** Filters cases where action will execute

**action** Describes tracing operations

## D Intermediate Format (DIF)

```
root@beaglebone:/data # dtrace -Sn
  'fbt::malloc:entry /execname == "csh"/ { trace(arg0); }'
```

DIFO 0x0x8047d2320 returns D type (integer) (size 4)

```
OFF OPCODE      INSTRUCTION
00: 29011801    ldgs DT_VAR(280), %r1 ! DT_VAR(280) = "execname"
01: 26000102    sets DT_STRING[1], %r2 ! "csh"
02: 27010200    scmp %r1, %r2
03: 12000006    be 6
04: 0e000001    mov %r0, %r1
05: 11000007    ba 7
06: 25000001    setx DT_INTEGER[0], %r1 ! 0x1
07: 23000001    ret %r1
```

NAME	ID	KND	SCP	FLAG	TYPE
execname	118	scl	g1b	r	string (unknown) by ref (size 256)

DIFO 0x0x8047d2390 returns D type (integer) (size 8)

```
OFF OPCODE      INSTRUCTION
00: 29010601    ldgs DT_VAR(262), %r1 ! DT_VAR(262) = "arg0"
01: 23000001    ret %r1
```

NAME	ID	KND	SCP	FLAG	TYPE
arg0	106	scl	g1b	r	D type (integer) (size 8)

## Some kernel DTrace providers in FreeBSD

<b>Provider</b>	<b>Description</b>
<code>callout_execute</code>	Timer-driven callouts
<code>dtmalloc</code>	Kernel <code>malloc()/free()</code>
<code>dtrace</code>	DTrace script events ( <code>BEGIN</code> , <code>END</code> )
<code>fbt</code>	Function Boundary Tracing
<code>io</code>	Block I/O
<code>ip, udp, tcp, sctp</code>	TCP/IP
<code>lockstat</code>	Locking
<code>proc, sched</code>	Kernel process/scheduling
<code>profile</code>	Profiling timers
<code>syscall</code>	System call entry/return
<code>vfs</code>	Virtual filesystem

- ▶ Providers represent data sources – types of instrumentation
- ▶ Apparent duplication: FBT vs. event-class providers?
  - ▶ Efficiency, expressivity, interface stability, portability

## Tracing kernel `malloc()` calls

- ▶ Trace first argument to kernel `malloc()` for `csch`
- ▶ Note: captures both successful and failed allocations

```
root@beaglebone:/data # dtrace -n
'fbt::malloc:entry /execname=="csch"/ { trace(arg0); }'
```

**Probe** Use FBT to instrument `malloc()` prologue

**Predicate** Limit actions to processes executing `csch`

**Action** Trace the first argument (`arg0`)

CPU	ID	FUNCTION:NAME	
0	8408	malloc:entry	64
0	8408	malloc:entry	2748
0	8408	malloc:entry	48
0	8408	malloc:entry	392

^C



## Aggregations

- ▶ Often we want summaries of events, not detailed traces
- ▶ DTrace allows early, efficient *reduction* using aggregations

Aggregation	Description
<code>count()</code>	Number of times called
<code>sum()</code>	Sum of arguments
<code>avg()</code>	Average of arguments
<code>min()</code>	Minimum of arguments
<code>max()</code>	Maximum of arguments
<code>stddev()</code>	Standard deviation of args
<code>lquantize()</code>	Linear frequency distribution (histogram)
<code>quantize()</code>	Log frequency distribution (histogram)

- ▶ Scalable multicore implementations (i.e., commutative)
- ▶ `@variable = function(); printa()` to print

## Profiling kernel `malloc()` calls by `csch`

```
root@beaglebone:/data # dtrace -n 'fbt::malloc:entry
    /execname=="csch"/ { @traces[stack()] = count(); }'
```

**Probe** Use FBT to instrument `malloc()` prologue

**Predicate** Limit actions to processes executing `csch`

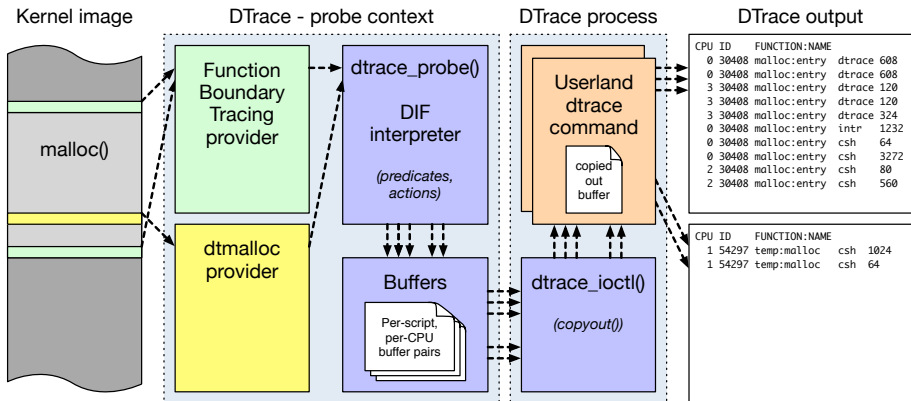
**Action** Keys of associative array are stack traces (`stack()`);  
values are aggregated counters (`count()`)

```
^C
    kernel`malloc
    kernel`fork1+0x14b4
    kernel`sys_vfork+0x2c
    kernel`swi_handler+0x6a8
    kernel`swi_exit
    kernel`swi_exit
    3
```

...

# DTrace: implementation

```
dtrace -n 'fbt::malloc:entry { trace(execname); trace(arg0); }'
```



```
dtrace -n 'dtmalloc::temp:malloc /execname="csh"/ { trace(execname); trace(arg3); }'
```

## The ‘probe effect’

- ▶ The *probe effect* is the unintended alteration in system behaviour that arises from measurement
- ▶ Why? Software instrumentation is *active*: execution is changed
- ▶ DTrace minimises *probe effect* when not being used...
  - ▶ ... but has a very significant impact when it is
  - ▶ Disproportionate effect on probed events
- ▶ Potential perturbations:
  - ▶ Execution speed relative to other cores (e.g., lock hold times)
  - ▶ Execution speed relative to external events (e.g., timer ticks)
  - ▶ Microarchitectural effects (e.g., cache footprint, branch predictor)
- ▶ What does this mean for us?
  - ▶ Don't benchmark while running DTrace ...
  - ▶ ... unless benchmarking DTrace
  - ▶ Be aware that traced application may behave differently
  - ▶ E.g., more timer ticks will fire, I/O will “seem faster”

## Probe effect example: `dd` execution time

- ▶ Simple (naive) microbenchmark
- ▶ `dd(1)` copies blocks from an input to an output
- ▶ Copy one 10M buffer from `/dev/zero` to `/dev/null`
- ▶ Execution time measured with `/usr/bin/time`

```
# dd if=/dev/zero of=/dev/null bs=10m count=1 status=none
```

- ▶ Simultaneously, run various DTrace scripts
- ▶ Compare resulting execution times using `ministat`

## Probe effect example 1: Memory allocation

- ▶ Using the `dtmalloc` provider, count kernel memory allocations

```
dtmalloc::: { @count = count(); }
```

```
% minitstat no-dtrace dtmalloc-count
x no-dtrace
+ dtmalloc-count
```

```
+-----+
|                                     |
|                                     |
|                                     |
|x                                     |
|x                                     |
|*                                     |
|_____A_____M_____A_____|
+-----+
      N      Min      Max      Median      Avg      Stddev
x  11      0.2      0.22      0.21      0.20818182  0.0060302269
+  11      0.2      0.22      0.21      0.21272727  0.0064666979
No difference proven at 95.0% confidence
```

- ▶ No statistically significant overhead at 95% confidence level

## Probe effect example 2: Locking operations

- ▶ lockstat provider tracks lock acquire, release, ...
- ▶ 6 calls to malloc, but 170K locking operations!

```
lockstat::: { @count = count(); }
```

```
ministat no-dtrace lockstat-count
x no-dtrace
+ lockstat-count
```

```
+-----+
|  x                                     +|
|  x                                     +|
|  x                                     +  +|
|x  x                                     +  +|
|x  x                                     +  +|
|x  x  x                                 +  +  +|
| |_A_|                                  |_A_M|
+-----+
```

	N	Min	Max	Median	Avg	Stddev
x	11	0.2	0.22	0.21	0.20818182	0.0060302269
+	11	0.42	0.44	0.44	0.43454545	0.0068755165

Difference at 95.0% confidence

0.226364 +/- 0.00575196

108.734% +/- 2.76295%

(Student's t, pooled s = 0.0064667)

- ▶ 109% overhead to count all lockstat probes!

## Probe effect example 2: Limiting to dd?

- ▶ Limit action to processes with name dd

```
lockstat::: /execname == "dd"/ { @count = count(); }
```

```
ministat no-dtrace lockstat-count-dd
```

```
x no-dtrace
```

```
+ lockstat-count-dd
```

```
+-----+
|                                             + |
| x                                             + |
| x                                             + |
| x                                             + |
| x                                             + |
|x x                                             + |
|x x                                             + |
|x x x                                           + + + |
|_A_|                                           |_A_| |
+-----+
```

	N	Min	Max	Median	Avg	Stddev
x	11	0.2	0.22	0.21	0.20818182	0.0060302269
+	11	0.54	0.57	0.56	0.55818182	0.0075075719

Difference at 95.0% confidence

0.35 +/- 0.0060565

168.122% +/- 2.90924%

(Student's t, pooled s = 0.00680908)

- ▶ Well, crumbs. Now the overhead is 168%!



## Probe effect example 3: Locking stack traces

- ▶ Gather more information in action: capture call stacks

```
lockstat::: { @stacks[stack()] = count(); }
lockstat::: /execname == "dd"/ { @stacks[stack()] =
  count(); }
```

```
ministat no-dtrace lockstat-stack*
x no-dtrace
+ lockstat-stack
* lockstat-stack-dd
```

```
+-----+-----+
|                                         * |
|                                         * |
|                                         * |
|                                         * |
|xx                                     ++ ** |
|xx                                     ++ ** |
|xx                                     +  +++ ** **|
|AM                                     |_MA_|A| |
+-----+-----+

```

	N	Min	Max	Median	Avg	Stddev
x	11	0.2	0.22	0.21	0.20818182	0.0060302269
+	11	1.38	1.57	1.44	1.4618182	0.058449668
		1.25364 +/- 0.0369572				
		602.183% +/- 17.7524%				
*	11	1.5	1.55	1.51	1.5127273	0.014206273
		1.30455 +/- 0.00970671				
		626.638% +/- 4.66261%				

## The kernel: “Just a C program”?

I claimed that the kernel was mostly “just a C program”.

This is mostly true, especially if you look at high-level subsystems.

<b>Userspace</b>	<b>Kernel</b>
crt/csu	locore
rtld	Kernel linker
Shared objects	Kernel modules
main()	main(), platform_start
libc	libkern
POSIX threads API	kthread KPI
POSIX filesystem API	VFS KPI
POSIX socket API	socket KPI
DTrace	DTrace
...	...

# The kernel: not just *any* C program

- ▶ Core kernel:  $\approx 3.4\text{M}$  LoC in  $\approx 6,450$  files
  - ▶ Kernel foundation: Built-in linker, object model, scheduler, memory allocator, threading, debugger, tracing, I/O routines, timekeeping
  - ▶ Base kernel: VM, process model, IPC, VFS w/20+, filesystems, network stack (IPv4/IPv6, 802.11, ATM, ...), crypto framework
  - ▶ Includes roughly  $\approx 70\text{K}$  lines of assembly over  $\approx 6$  architectures
- ▶ Alternative C runtime – e.g., `SYSINIT`, `curthread`
- ▶ Highly concurrent – really, very, very concurrent
- ▶ Virtual memory makes pointers .. odd
- ▶ Debugging features such as `WITNESS` lock order verifier
- ▶ Device drivers:  $\approx 3.0\text{M}$  LoC in  $\approx 3,500$  files
  - ▶  $\approx 415$  device drivers (may support multiple devices)

# Spelunking the kernel

```

/usr/src/sys> ls
Makefile      ddb/          mips/         nfs/          sys/
amd64/        dev/          modules/      nfscclient/  teken/
arm/          fs/           net/          nfsserver/   tools/
boot/         gdb/          net80211/     nlm/         ufs/
bsm/          geom/         netgraph/     ofed/        vm/
cam/          gnu/          netinet/      opencrypto/  x86/
cddl/         i386/         netinet6/     pc98/        xdr/
compat/       isa/          netipsec/     powerpc/     xen/
conf/         kern/         netnatm/      rpc/
contrib/      kgssapi/     netpfil/      security/
crypto/       libkern/     netsmb/       sparc64/

/usr/src/sys> ls kern
Make.tags.inc  kern_racct.c      subr_prof.c
Makefile       kern_rangelock.c subr_rman.c
bus_if.m       kern_rctl.c       subr_rtc.c
capabilities.conf kern_resource.c   subr_sbuf.c
clock_if.m     kern_rmlock.c     subr_scanf.c
...

```

- ▶ Kernel source lives in `/usr/src/sys`:
  - kern/ - core kernel features
  - sys/ - core kernel headers
- ▶ Useful resource: <http://fxr.watson.org/>

# How work happens in the kernel

- ▶ Kernel code executes concurrently in multiple threads
  - ▶ User threads in kernel (e.g., system call)
  - ▶ Shared worker threads (e.g., callouts)
  - ▶ Subsystem worker threads (e.g., network-stack worker)
  - ▶ Interrupt threads (e.g., clock ticks)
  - ▶ Idle threads

```

root@beaglebone:/data # procstat -at
  PID   TID  COMM          TDNAME          CPU  PRI  STATE  WCHAN
    0   100000 kernel        swapper         -1   84  sleep  swapon
    0   100006 kernel        dtrace_taskq   -1   84  sleep  -
...
   10  100002 idle          -                -1  255  run    -
   11  100003 intr         swi3: vm         0   36  wait  -
   11  100004 intr         swi4: clock (0) -1   40  wait  -
   11  100005 intr         swil: netisr 0  -1   28  wait  -
...
   11  100018 intr         intr16: ti_adc0  0   20  wait  -
   11  100019 intr         intr91: ti_wdt0  0   20  wait  -
   11  100020 intr         swi0: uart      -1   24  wait  -
...
   739 100064 login        -                -1  108  sleep  wait
   740 100079 csh          -                -1  140  sleep  ttyin
   751 100089 procstat    -                0   140  run    -

```

## Work processing and distribution

- ▶ Many operations begin with system calls in a user thread
- ▶ But they may trigger work in many other threads; for example:
  - ▶ Triggering a callback in an interrupt thread when I/O is complete
  - ▶ Eventual write back of data to disk from the cache
  - ▶ Delayed transmission if TCP isn't able to send
- ▶ We will need to care about these things, as not all the work we are analysing will be in the user thread performing a system call.
- ▶ Several major subsystems provide this:
  - `callout` Closure called after wall-clock delay
  - `eventhandler` Closure called for key global events
  - `task` Closure called eventually
  - `SYSINIT` Function called when module loads/unloads
- ▶ (Where closure in C means: function pointer, opaque data pointer)

## For next time

- ▶ Read Ellard and Seltzer, *NFS Tricks and Benchmarking Traps*
- ▶ Skim handout, *L41: DTrace Quick Start*
- ▶ Be prepared to try out DTrace on a real system