

Tools

	<i>DTrace</i>	<i>SystemTap</i>
Tool	<code>dtrace(1M)</code>	<code>stap(1)</code>
List probes	<code># dtrace -l</code> <code># dtrace -l -P io</code>	<code># stap -l 'ioblock.*'</code> <code># stap -L 'ioblock.*'</code>
One-liner	<code># dtrace -n 'syscall::read:entry { trace(arg1); } '</code>	<code># stap -e 'probe syscall.read { println(fd); } '</code>
Script	<code># dtrace -s script.d</code> (optionally add <code>-C</code> for preprocessor, <code>-q</code> for quiet mode)	<code># stap script.stp</code>
Custom probe	<code># dtrace -P io -n start</code>	-
Integer arguments	<code># dtrace -n 'syscall::read:entry / cpu == \$1 / ' 0</code>	<code># stap -e 'probe syscall.read { if(cpu() != \$1) next; println(fd); } ' 0</code>
String arguments	<code># dtrace -n 'syscall::read:entry / execname == \$1 / ' "cat"</code>	<code># stap -e 'probe syscall.read { if(execname() == @1) println(fd); } ' cat</code>
Guru/destructive mode (!)	<code># dtrace -w ...</code>	<code># stap -g ...</code>
Redirect to file	<code># dtrace -o FILE ...</code> (appends)	<code># stap -o FILE ...</code> (rewrites)
Tracing process	<code># dtrace -n 'syscall::read:entry / pid == \$target / { ... }' -c 'cat /etc/motd'</code> (or <code>-p PID</code>)	<code># stap -e 'probe syscall.read { if(pid() == target()) ... }' -c 'cat /etc/motd'</code> (or <code>-x PID</code>)

Probe names

	<i>DTrace</i>	<i>SystemTap</i>
Begin/end	<code>dtrace:::BEGIN, dtrace:::END</code>	<code>begin, end</code>
<code>foo()</code> entry	<code>fbt::foo:entry</code>	<code>kernel.function("foo") module("mod").function("foo")</code>
<code>foo()</code> return	<code>fbt::foo:return</code>	<code>kernel.function("foo").return</code>
Wildcards	<code>fbt::foo*:entry</code>	<code>kernel.function("foo*")</code>
Static probe mark	<code>sdt:::mark</code>	<code>kernel.trace("mark")</code>
System call	<code>syscall::read:entry</code>	<code>syscall.read</code>
Timer once per second	<code>tick-1s</code>	<code>timer.s(1)</code>
Profiling	<code>profile-997hz</code>	<code>timer.profile(), perf.*</code>
<code>read()</code> from libc	<code>pid\$target:libc:read:entry</code> Traces process with <code>pid == \$target</code>	<code>process("/lib64/libc.so.6").function("read")</code> Traces any process that loads libc

In DTrace parts of probe name may be omitted: `fbt::foo:entry -> foo:entry`
 Units for timer probes: ns, us, ms, s, hz, jiffies (SystemTap), m, h, d (all three - DTrace)

Context variables

Description	DTrace	SystemTap
Thread	curthread	task_current()
Thread ID	tid	tid()
PID	pid	pid()
Parent PID	ppid	ppid()
User/group ID	uid/gid	uid()/gid() euid()/egid()
Executable name	execname curpsinfo-> ps_fname	execname()
Command line	curpsinfo-> ps_psargs	cmdline_*
CPU number	cpu	cpu()
Probe names	probeprov probemod probefunc probename	pp() pn() ppfunc() probefunc() probemod()

Time

Time source	DTrace	SystemTap
System timer	`lbolt `lbolt64	jiffies()
CPU cycles	-	get_cycles()
Monotonic time	timestamp	local_clock_unit() cpu_clock_unit(cpu)
CPU time of thread	vtimestamp	-
Real time	walltimestamp	gettimeofday_unit()

Where *unit* is one of s, ms, us, ns

Printing

	DTrace	SystemTap
Value	trace(v)	print(v)
Value + newline	-	println(v)
Delimited values	-	printf(", ", v1, v2) println(", ", v1, v2)
Memory dump	tracemem(ptr, 16)	printf("%16M", ptr)
Formatted	printf("%s", str)	
Backtrace	ustack(n) ustack()	print_ubacktrace() print_ustack(ubacktrace())
Symbol	usym(addr) ufunc(addr) uaddr(addr)	print(usymname(addr)) print(usymdata(addr))

If *u* prefix is specified, userspace symbols and backtraces are printed, if not — kernel symbols are used

String operations

Operation	DTrace	SystemTap_
Get from kernel	stringof(expr) (string) expr	kernel_string*()
Convert scalar		sprint() and sprintf()
Copy from user	copyinstr()	user_string*()
Compare	==, !=, >, >=, <, <=	
Concat	strjoin(str1, str2)	str1 . str2
Get length	strlen(str)	
Check for substring	strstr(haystack, needle)	isinstr(haystack, needle)

Aggregations

Time source	DTrace	SystemTap
Add value	@aggr[keys] = func(value);	aggr[keys] <<< value;
Printing	printa(@aggr); printa("format string", @aggr);	foreach([keys] in aggr) { print(keys, @func(aggr[keys])); }
Clear	clear(@aggr); or trunc(@aggr);	delete aggr;
Normalization by 1000	normalize(@aggr, 1000); denormalize(@aggr);	@func(aggr) / 1000 in printing
Select 20 values	trunc(@aggr, 20);	foreach([keys] in aggr limit 20) { print(keys, @func(aggr[keys])); }
Histograms (linear in [10;100] with step 5 and logarithmical)	@lin = lquantize(value, 10, 100, 5); @log = quantize(value); ... printa(@lin); printa(@log);	aggr <<< value; ... print(@hist_linear(aggr, 10, 100, 5)); print(@hist_log(aggr));

Where *func* is one of count, sum, min, max, avg, stddev

Process management

SystemTap

Getting task_struct pointers:

- `task_current()` – current task_struct
- `task_parent(t)` – parent of task t
- `pid2task(pid)` – task_struct by pid

Working with task_struct pointers:

- `task_pid(t)` & `task_tid(t)`
- `task_state(t)` – 0 (running), 1-2 (blocked)
- `task_execname(t)`

DTrace

kthread_t* curthread fields:

- `t_tid`, `t_pri`, `t_start`, `t_pctcpu`

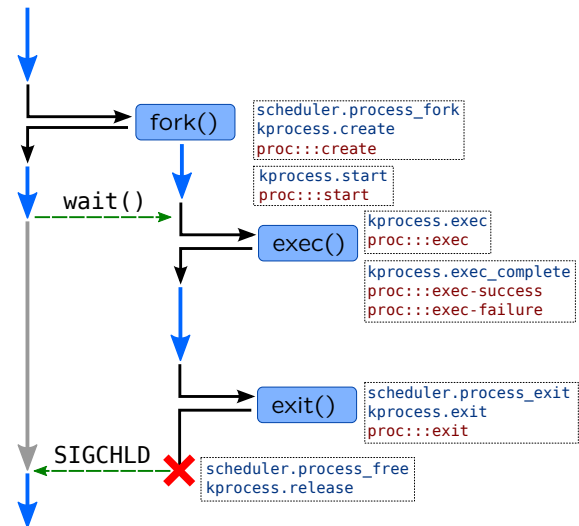
psinfo_t* curpsinfo fields:

- `pr_pid`, `pr_uid`, `pr_gid`, `pr_fname`, `pr_psargs`, `pr_start`

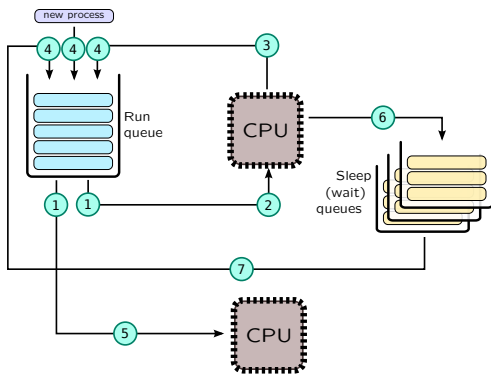
lwpsinfo_t* curlwpsinfo fields:

- `pr_lwpid`, `pr_state/pr_sname`

`psinfo_t*` and `lwpsinfo_t*` are passed to some `proc::` probes



Scheduler



	<i>DTrace</i>	<i>SystemTap</i>
1	<code>sched::dequeue</code>	<code>kernel.function("dequeue_task")</code>
2	<code>sched::on-cpu</code>	<code>scheduler.cpu_on</code>
3	<code>sched::off-cpu</code>	<code>scheduler.cpu_off</code>
4	<code>sched::enqueue</code>	<code>kernel.function("enqueue_task")</code>
5	-	<code>scheduler.migrate</code>
6	<code>sched::sleep</code>	-
7	<code>sched::wakeup</code>	<code>scheduler.wakeup</code>

Virtual memory

Probes

SystemTap

- `vm.brk` – allocating heap
- `vm.mmap` – allocating anon memory
- `vm.munmap` – freeing anon memory

DTrace

- `as_map:entry` – allocating proc mem
- `as_unmap:entry` – freeing proc mem

Page faults

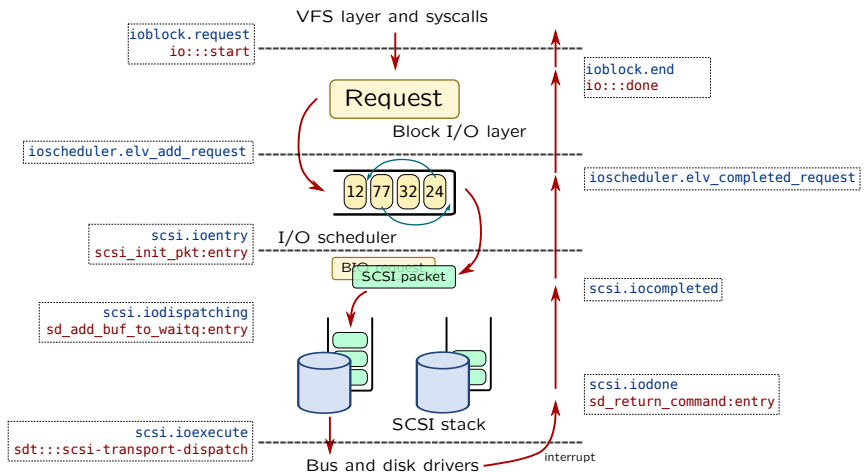
Type	<i>DTrace</i>	<i>SystemTap</i>
Any	<code>vminfo::as_fault</code>	<code>vm.pagefault</code> <code>vm.pagefault.return</code> <code>perf.sw.page_faults</code>
Minor		<code>perf.sw.page_faults_min</code>
Major	<code>vminfo::maj_fault</code>	<code>perf.sw.page_faults_maj</code>
CoW	<code>vminfo::cow_fault</code>	
Protection	<code>vminfo::prot_fault</code>	

Block Input-Output

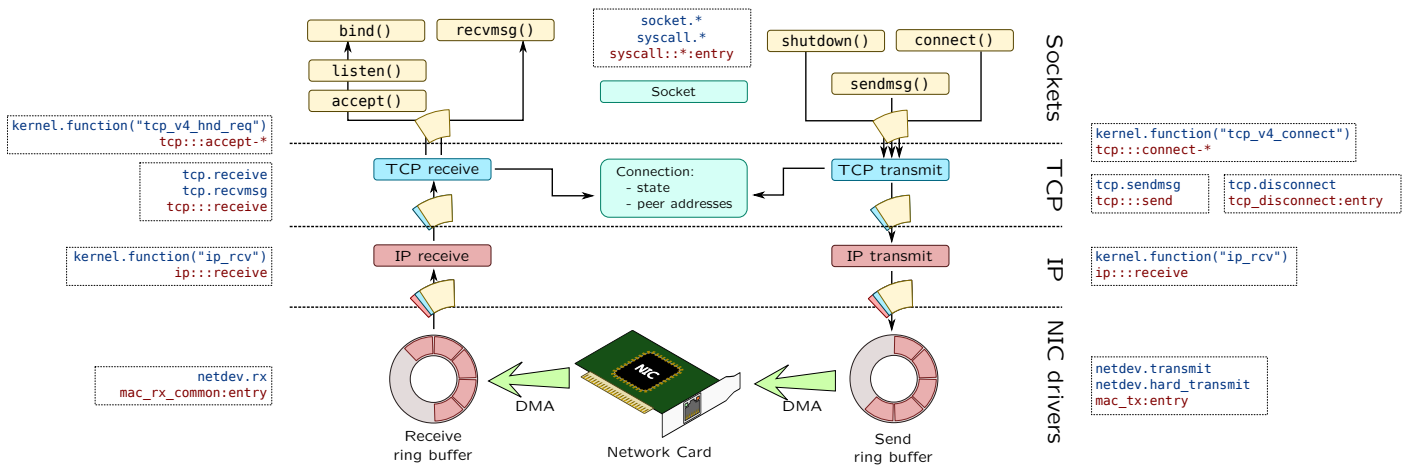
Block request structure fields:

Field	bufinfo_t struct buf	struct bio
Flags	b_flags	bi_flags
R/W	b_flags	bi_rw
Size	b_bcount	bi_size
Block	b_blkno b_lblkno	bi_sector
Callback	b_iodone	bi_end_io
Device	b_edev b_dip	bi_bdev

* flags B_WRITE, B_READ



Network stack



Non-native languages

Function call	DTrace	SystemTap
Java*	method-entry <ul style="list-style-type: none"> arg0 — internal JVM thread's identifier arg1:arg2 — class name arg3:arg4 — method name arg5:arg6 — method signature 	hotspot.method_entry <ul style="list-style-type: none"> thread_id — internal JVM thread's identifier class — class name method — method name sig — method signature
Perl	perl\$target:::sub-entry <ul style="list-style-type: none"> arg0 — subroutine name arg1 — source file name arg2 — line number 	process("...").mark("sub__entry") <ul style="list-style-type: none"> \$arg1 — subroutine name \$arg2 — source file name \$arg3 — line number
Python	python\$target:::function-entry <ul style="list-style-type: none"> arg0 — source file name arg1 — function name 	python.function.entry <ul style="list-style-type: none"> \$arg1 — source file name \$arg2 — function name
PHP	function-entry <ul style="list-style-type: none"> arg0 — function name arg1 — file name arg2 — line number arg3 — class name arg4 — scope operator :: 	process("...").mark("function__entry") <ul style="list-style-type: none"> \$arg1 — function name \$arg2 — file name \$arg3 — line number \$arg4 — class name \$arg5 — scope operator ::

*requires -XX:+DTraceMethodProbes