

## Linux perf\_events Off-CPU Time Flame Graph

26 Feb 2015

I've been asked how to do these several times, so here's a quick blog post.

[CPU Flame Graphs](#) are great for understanding CPU usage, but what about performance issues of latency, when programs are blocked and not running on-CPU? There's a generic methodology to attack these, which I've called [Off-CPU Analysis](#). This involves instrumenting call stacks and timestamps when the scheduler blocks a thread, and when it later wakes it up. There are even [Off-CPU Time Flame Graphs](#) to visualize this blocked time, which are the counterpart to on-CPU flame graphs.

Off-CPU time flame graphs may solve (say) 60% of the issues, with the remainder requiring walking the thread wakeups to find root cause. I explained off-CPU time flame graphs, this wakeup issue, and additional work, in my LISA13 talk on flame graphs ([slides](#), [youtube](#)).

Here I'll show one way to do off-CPU time flame graphs using Linux perf\_events. Example (click to zoom):

__schedule	__schedule				
schedule	schedule				
schedule_hrtimeou..	schedule_hrti..	__schedule		__sch..	
schedule_hrtimeou..	schedule_hrti..	schedule		sched..	
poll_schedule_tim..	poll_schedule..	schedule_hrti..		pipe_..	
do_sys_poll	do_select	schedule_hrti..		pipe_..	__schedule
sys_poll	core_sys_select	poll_schedule..		new_s..	schedule
system_call	sys_select	do_select	__sche..	__vfs..	do_nanosleep
__poll	system_call	core_sys_select	schedule	vfs_r..	hrtimer_nanosleep
run_builtin	__select	sys_select	rcu_gp_k..	sys_r..	sys_nanosleep
main	[unknown]	system_call	kthread	sys_e..	system_call
__libc_start_main	[unknown]	__select	ret_from..	read	nanosleep
perf	postgres		rcu_sched	readp..	sleep

I was discussing how to do these with Nicolas on [github](#), as he had found that `perf_events` could almost do this using the `perf inject -s` feature, when I noticed `perf_events` has added "-f period" to `perf script` (added in about 3.17). That simplifies things a bit, so the steps can be:

Update: on more recent kernels, use "perf script -F ..." instead of "perf script -f ...". Your kernel will also need CONFIG\_SCHEDSTATS for the tracepoints to all be present, which may be missing (eg, RHEL7).

**Warning:** scheduler events can be very frequent, and the overhead of dumping them to the file system (perf.data) may be prohibitive in production environments. Test carefully. This is also why I put a "sleep 1" in the perf record (the dummy command that sets the duration), to start with a small amount of trace data. If I had to do this in production, I'd consider other tools that could summarize data in-kernel to reduce overhead, including perf\_events once it supports more in-kernel programming (eBPF).

**Update 1:** Since Linux 4.5 you may need the following for this to work:

**Update 2:** Since Linux 4.6 you can use BPF to do this *much* more efficiently: aggregating the stacks in-kernel context (using the BPF stack trace feature in 4.6), and only passing the summary to user-level. I developed the offcputime tool [bcc](#) to do this. I also wrote a post about it, [Off-CPU eBPF Flame Graph](#), although that was written before Linux 4.6's stack trace support, so I was unwinding stacks the hard way.

You can comment here, but I can't guarantee your comment will remain here forever: I might switch comment systems at some point (eg, if Disqus add advertisements).

---

Copyright 2017 Brendan Gregg.

[About this blog](#)

[Perf Methods](#)

[USE Method](#)

[TSA Method](#)

[Off-CPU Analysis](#)

[Active Bench.](#)

[Flame Graphs](#)

[Heat Maps](#)

[Frequency Trails](#)

[Colony Graphs](#)

[perf Examples](#)

[eBPF Tools](#)

[DTrace Tools](#)

[DTraceToolkit](#)

[DtkshDemos](#)

[Guessing Game](#)

[Specials](#)

[Books](#)

[Other Sites](#)