Improvements
New Tools
Scripting
KVM Support
Work in Progress
That is all folks!

Perf Tools: Recent Improvements Recent developments and discussion about TODO

Arnaldo Carvalho de Melo

Red Hat Inc.

Netconf and Linux Plumbers Conference, Cambridge November, 2010



- Improvements
 - Tool Integration
 - Slang Based TUI
- New Tools
 - perf diff
 - perf archive
 - perf probe
- Scripting
 - Available Scripts
 - Generate Scripts
- 4 KVM Support
- Work in Progress
- That is all folks!



Improvements

- ① On the initial set of tools
- Tools Integration
- Slang based Text User Interface
- Use of build ids

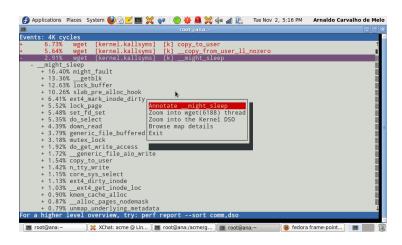
Tools Integration

- One tool doesn't do it all
- 2 Combine steps to achieve multiple results
- Allows spreading work flows over multiple machines
- Profiling fast path
- report to annotate
- Reuse perf.data parsing

Slang Based TUI

- GUIs not necessarily better
- mutt x thunderbird
- Out the changes paves the way for GUIs
- report to annotate fast path
- See Keys used: arrows + ENTER mostly, TAB sometimes
- Still don't like it? Use -stdio

perf report TUI



perf annotate

- Starts at the line with most hits
- Tabs through ordered list of hot lines

perf annotate TUI

```
🛐 Applications Places System 🥹 🔂 🗾 💥 🕡 💨 🔴 🎉 🐗 📶 🖺
                                                                  Tue Nov 2, 5:16 PM
                                                                                   Arnaldo Carvalho de Melo
 might_sleep
                          return (struct thread info *)
                                  (current_stack_pointer & ~(THREAD_SIZE - 1));
                  c042b562:
   0.80:
                                                          mov
                                                                  %esp,%eax
                                  25 00 e0 ff ff
                                                                  $0xffffe000,%eax
   0.00:
                  c042b564:
                                                          and
   0.00:
                                  89 d3
                                                          mov
                                                                  %edx.%ebx
                  #ifdef CONFIG DEBUG SPINLOCK SLEEP
                 static inline int preempt count equals(int preempt offset)
                          int nested = (preempt count() & ~PREEMPT ACTIVE) + rcu preempt depth();
   1.60 :
                  c042b56b:
                                  8b 40 14
                                                                 0x14(%eax),%eax
                                                          mov
                                  25 ff ff ff ef
  26.40
                                                                 $0xefffffff.%eax
                 void __might_sleep(const char *file, int line, int preempt offset)
                  #ifdef in atomic
                          static unsigned long prev jiffy;
                                                                 /* ratelimiting */
                          if ((preempt_count_equals(preempt_offset) && !irqs_disabled()) ||
   4.00 :
                  c042b573:
                                                                  %ecx.%eax
                                                          CMD
                                                                  c042b586 <__might_sleep+0x30>
   6.40
                                                          ine
                 #define PV IS CALLEE SAVE(func)
                          ((struct paravirt callee save) { func })
                 static inline unsigned long arch_local_save_flags(void)
<-, -> or ESC: exit, TAB/shift+TAB: cycle thru samples
 m root@ana:~

XChat: acme @ Lin... 
☐ root@ana:/acme/g... ☐ root@ana:~

                                                                          (a) fedora frame-point...
```

UI - TODO

- perf top
- Allow selecting events to record at any time
- Start with top
- Freeze == report
- Save == record
- perf probe
- Go from annotate to probe, restart top

perf top

Considers user space symbols too:

```
PerfTop: 155 irqs/sec kernel:83.9% [1000Hz cycles], (all, 2 CPUs)
```

sa	mples	pcnt	function	DSO
1	19.00	12.0%	read_hpet	[kernel]
	43.00	4.4%	strstr_ia32	/lib/libc-2.12.1.so
	28.00	2.8%	system_call	[kernel]
	25.00	2.5%	unix_poll	[kernel]
	24.00	2.4%	aes_enc_blk	[aes_i586]
	21.00	2.1%	schedule	[kernel]
	21.00	2.1%	_raw_spin_lock_irqsave	[kernel]
	19.00	1.9%	_raw_spin_unlock_irqrestore	[kernel]
	19.00	1.9%	aes_dec_blk	[aes_i586]
	18.00	1.8%	probe_workqueue_insertion	[kernel]
	17.00	1.7%	hpet_next_event	[kernel]
	13.00	1.3%	fget_light	[kernel]
	13.00	1.3%	do_select	[kernel]
	12.00	1.2%	audit_syscall_entry	[kernel]
	12.00	1.2%	ktime_get	[kernel]
	11.00	1.1%	test_ti_thread_flag	[kernel]
	11.00	1.1%	std::_List_node_base::transfer(std::_L	libstdc++.so.6.0.13
	11.00	1.1%	native_sched_clock	[kernel]
	11.00	1.1%	vsnprintf	[kernel]
	11.00	1.1%	format_decode	[kernel]
	10.00	1.0%	index	/lib/libc-2.12.1.so

perf stat

- List of CPUs to monitor
- Ask for precise events(PEBS) using suffix: "-e cycles:p"
- Multiple 'p' characters == more precise
- Proof of concept patch for printing counters periodically ready
- Merge app log output sorting by timestamps

New Tools

Introduced after Plumbers'2009:

- diff
- archive
- probe
- trace
- several trace ones (timechart, etc)

perf diff

- Shows difference in symbol hits between two perf.data files
- Weyed by build-ids in the cache
- Should support more than two files
- Generating version X samples symbol plottings
- Read "Differential Profiling" paper by Paul McKenney on how to use it

perf archive

- Looks at perf.data files for DSOs with hits
- Creates tarball
- Transfer to another machine
- Populate the cache
- Use report and annotate
- Mandles endianness

perf probe

- Inserts dynamic probes
- ② Doesn't necessarily requires debuginfo
- Can collect variables
- Struct members can be specified to any level
- Works with callchains
- Together with perf trace == systemtap subset
- Example of use together with scripting later in this presentation
- Ontributed by Masami Hiramatsu



Scripting

- Use scripting languages to process events
- 2 Python and Perl
- 4 Allows tapping into tons of language libraries
- Several scripts available
- Generate scripts from perf.data
- Ontributed by Tom Zanussi

Available Scripts

```
[root@ana ~] # perf trace --list
List of available trace scripts:
  rw-by-pid
                                system-wide r/w activity
  wakeup-latency
                                system-wide min/max/avg wakeup latency
                                workqueue stats (ins/exe/create/destroy)
  workqueue-stats
  rwtop [interval]
                                system-wide r/w top
  failed-syscalls [comm]
                                system-wide failed syscalls
  rw-by-file <comm>
                                r/w activity for a program, by file
  syscall-counts-by-pid [comm]
                                system-wide syscall counts, by pid
  netdev-times [tx] [rx] [dev=]
                                display a process of packet and processing
  sctop [comm] [interval]
                                syscall top
                                futex contention measurement
  futex-contention
  sched-migration
                                sched migration overview
  failed-syscalls-by-pid [comm]
                                system-wide failed syscalls, by pid
  syscall-counts [comm]
                                system-wide syscall counts
[root@ana ~]#
```

Generate Scripts

- From the events found in perf.data file
- Quickly start writing event handling
- Oreates function skeletons for each trace event
- With a common set of parameters
- Opening Plus event specific parameters
- O Calls methods at init, exit and for unhandled events
- Comes with library of tracing specific methods

Listing Possible probe points

```
[root@ana icmp]# perf probe -L icmp_rcv
<icmp_rcv:0>
     0 int icmp_rcv(struct sk_buff *skb)
      1
        {
    59
                if (rt->rt_flags & (RTCF_BROADCAST | RTCF_MULTICAST)) {
                        /*
                         * RFC 1122: 3.2.2.6 An ICMP ECHO to broadcast MAY be
                         * silently ignored (we let user decide with a sysctl).
                         * RFC 1122: 3.2.2.8 An ICMP_TIMESTAMP MAY be silently
                         * discarded if to broadcast/multicast.
                         */
    66
                        if ((icmph->type == ICMP_ECHO ||
                             icmph->type == ICMP_TIMESTAMP) &&
                            net->ipv4.sysctl_icmp_echo_ignore_broadcasts) {
                                goto error;
                        }
    71
                        if (icmph->type != ICMP_ECHO &&
```

Listing variables that can be collected

Adding a probe

```
[root@ana icmp]# perf probe icmp_rcv:66 'type=icmph->type'
Add new event:
  probe:icmp_rcv
                    (on icmp_rcv:66 with type=icmph->type)
You can now use it on all perf tools, such as:
perf record -e probe:icmp_rcv -aR sleep 1
[root@ana ~] # perf probe --list
  probe:icmp_rcv (on icmp_rcv:66@net/ipv4/icmp.c with type)
[root@ana icmp]# perf record -a -g -e probe:icmp_rcv
^C[ perf record: Woken up 1 times to write data ]
[ perf record: Captured and wrote 0.324 MB perf.data ]
```

Generating a python script from perf.data

```
[root@ana icmp] # perf trace -g python
generated Python script: perf-trace.py
[root@ana icmp]# cat perf-trace.py
def trace_begin():
        print "in trace_begin"
def trace end():
        print "in trace_end"
def probe__icmp_rcv(evname, cpu, secs, nsecs, pid, comm,
                    probe_ip, type):
        print "%s %u.%u type=%u" % (evname, secs, nsecs, type)
```

Running python script

```
[root@ana icmp]# perf trace -s perf-trace.py
in trace_begin
probe__icmp_rcv 71171.964568380 type=8
probe__icmp_rcv 71177.792382154 type=8
probe__icmp_rcv 71178.792236953 type=8
in trace_end
[root@ana icmp]#
```

Backtraces from probes

```
[root@ana ~]# perf report --stdio
# Events: 2
# Overhead Command Shared Object
                                       Symbol
   100.00%
              ping [kernel.kallsyms] [k] icmp_rcv
               --- icmp_rcv
                  ip_local_deliver_finish
                  NF_HOOK.clone.1
                  ip_local_deliver
                  ip_rcv_finish
                  NF HOOK.clone.1
                  ip_rcv
                  __netif_receive_skb
                  process_backlog
                  net rx action
                  __do_softirg
                  0xb7707424
```

Scripting TODO List

- Convert trace builtins to scripts (sched, kmem, etc)
- Convert net/ipv4/tcp_probe.c
- SCTP and DCCP variants too
- Write more scripts for showing where IO is happening
- Improve passing data from record to trace
- Remove requirement on using netcat for dual machine use
- Write more scripts (you can help here!)

- Collect guest OS statistics from host side.
- 2 top, record, report, diff, buildid-list
- Need to specify guest vmlinux or kallsyms, /proc/modules
- Or –guestmount directory with sshfs mounted per pid subdirs
- Use -pid to specify specific guest
- Would be great to automatically get vmlinux by build-id
- Contributed by Zhang, Yanmin.

perf top kvm example

```
# perf kvm --host --guest --guestkallsvms=guest/kallsvms \
           --guestmodules=guest/modules top
 PerfTop: 16010 irgs/sec kernel:59.1% us: 1.5% guest
         kernel:31.9% guest us:7.5% [+1000Hz cvcles]
  samples pcnt function
                                          DSO
 38770.00 20.4% __ticket_spin_lock
                                          [guest.kernel]
 22560.00 11.9% ftrace_likely_update
                                          [kernel]
  9208.00 4.8% __lock_acquire
                                          [kernel]
 5473.00 2.9% trace_hardings_off_caller [kernel]
 5222.00 2.7% copy_user_generic_string [guest.kernel]
 4450.00 2.3% validate_chain
                                          [kernel]
 4262.00 2.2% trace_hardings_on_caller
                                          [kernel]
 4239.00 2.2% do raw spin lock
                                          [kernel]
 3548.00 1.9% do_raw_spin_unlock
                                          [kernel]
  2487.00 1.3% lock release
                                          [kernel]
  2165.00 1.1% __local_bh_disable
                                          [kernel]
  1905.00 1.0% check chain key
                                          [kernel]
```

Improvements
New Tools
Scripting
KVM Support
Work in Progress
That is all folks!

- Cgroup support
- ② utrace to probe user space

Thanks!

Arnaldo Carvalho de Melo acme@infradead.org acme@redhat.com