

VM kernel tracing with ftrace, trace-cmd and Kernel Shark

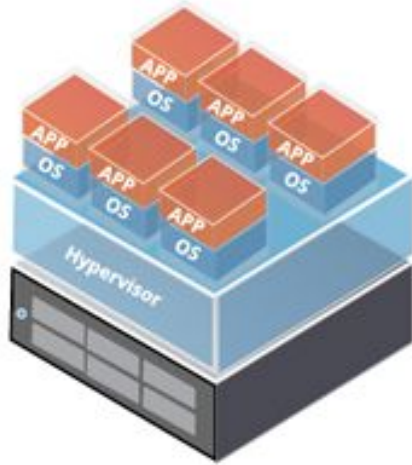
Tzvetomir Stoyanov

VMware Open Source Technology Center

Agenda

- The problem
- The general idea
- The challenges
- The solution

The problem



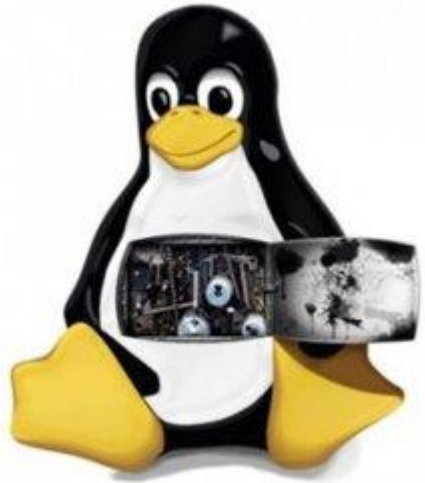
- + scaling
- + redundancy
- + hardware utilization
- nightmare for developers and administrators

The problem



According to the Linux foundation Linux OS runs **90%** of the public cloud workload

The problem



In majority of use cases
Linux is used in both sides -
as host and guest OS

ftrace

```
# tracer: function_graph
```

```
#
```

```
# CPU DURATION
```

```
# | | |
5) | | |
5) | | |
5) 0.148 us |
5) | | |
5) 0.167 us |
5) 0.415 us |
5) | | |
5) 0.216 us |
5) 0.519 us |
5) 1.931 us |
5) 0.151 us |
5) | | |
5) 0.193 us |
5) 0.529 us |
```

```
FUNCTION CALLS
```

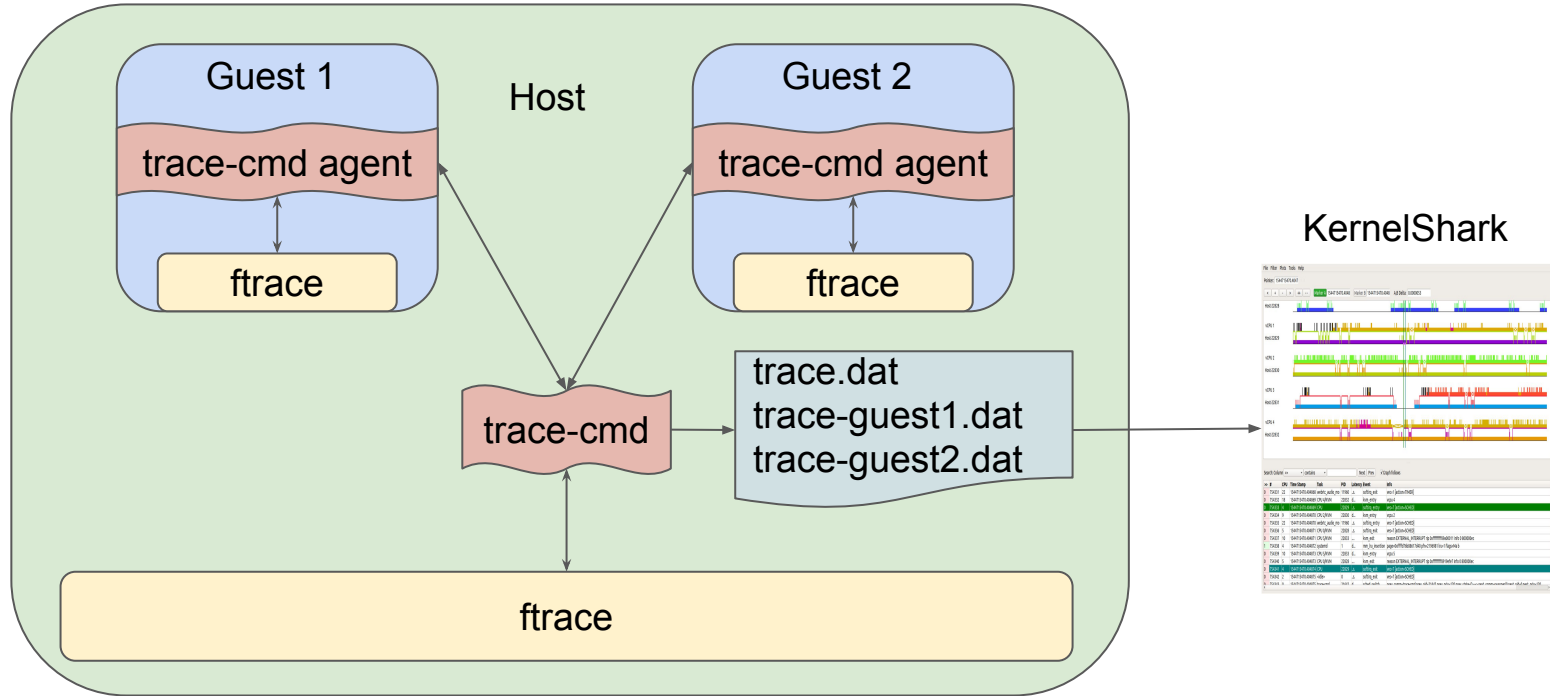
```
 | | | |
wake_up_q() {
  try_to_wake_up() {
    _raw_spin_lock_irqsave();
    select_task_rq_fair() {
      available_idle_cpu();
      update_cfs_rq_h_load();
      select_idle_sibling() {
        available_idle_cpu();
      }
    }
  }
  _raw_spin_lock();
  update_rq_clock() {
    update_irq_load_avg();
  }
}
```

- The official tracer of the Linux kernel
- Developed by Steven Rostedt more than 10 years ago
- Part of the kernel, compiled by default in most popular Linux distros.
- Allows you to look inside every corner of a live running kernel.

The general idea



The general idea



Challenges

- Fast transfer of huge tracing data between guest and host
- Time stamps synchronisation
- User friendly visualisation of huge host and guest trace data

VM Data Transfer: TCP/IP sockets

● Pros

- Hypervisor independant
- Allows bi-directional communication
- Can be used for nested virtualisation tracing

● Cons

- Slow, very slow
- Pass through the whole kernel's network stack
- Communication can happen between any host or VM which is undesirable

VM Data Transfer: FIFOs

● Pros

- Fastest
- Has no security concerns, as the TCP / IP solution

● Cons

- works only for KVM
- 1to1 unidirectional communication only
- a FIFO for each direction is needed
- Cannot be used for nested virtualization use case

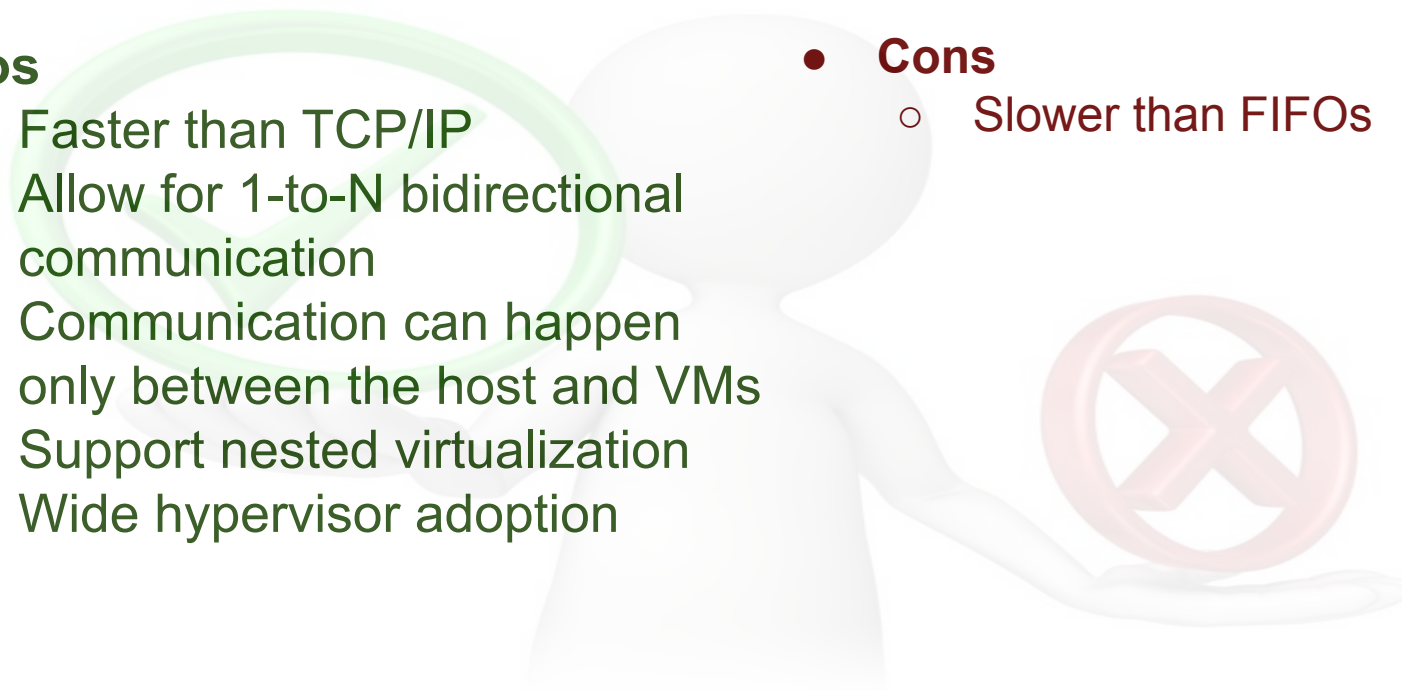
VM Data Transfer: vSockets

● Pros

- Faster than TCP/IP
- Allow for 1-to-N bidirectional communication
- Communication can happen only between the host and VMs
- Support nested virtualization
- Wide hypervisor adoption

● Cons

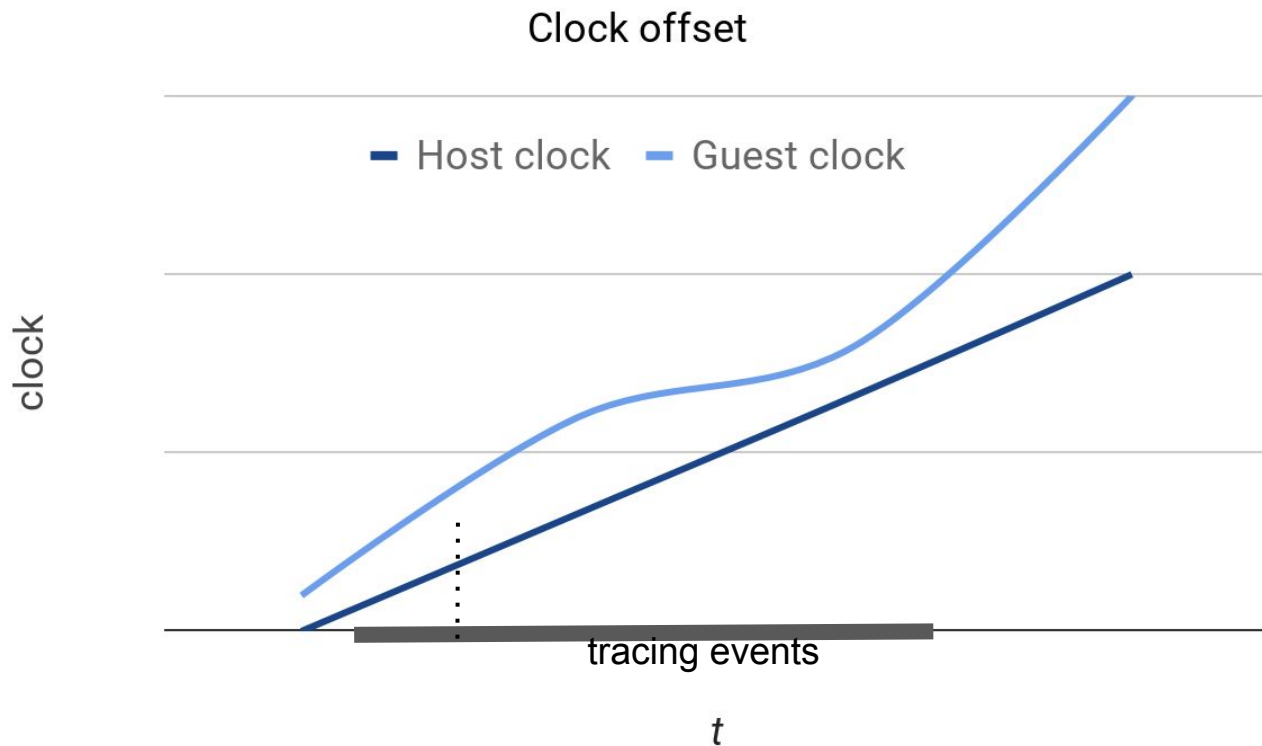
- Slower than FIFOs



KVM Channels Throughput

Channel	Throughput
FIFOs	1000 MB/s
vsockets	900 MB/s
TCP/IP sockets	275 MB/s

Synchronisation of time stamps



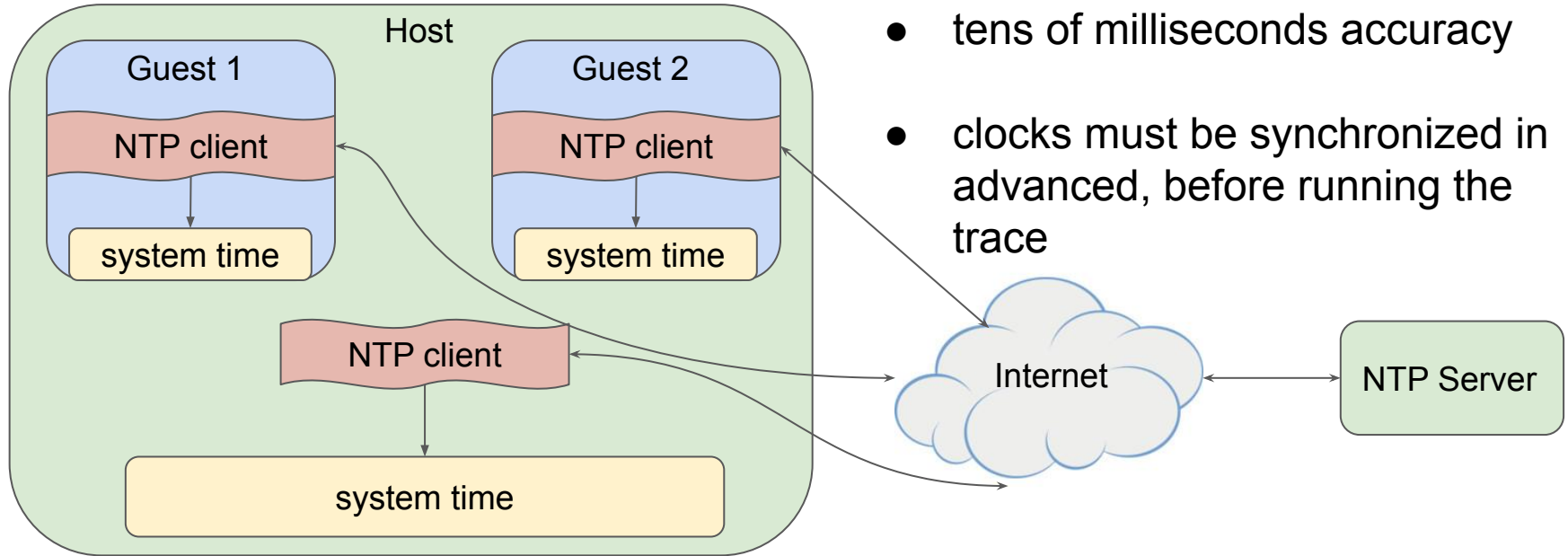
Synchronisation of time stamps

Ftrace clock sources

- local
- global
- counter
- uptime
- x86-tsc
- ppc-tb
- mono
- mono_raw
- boot

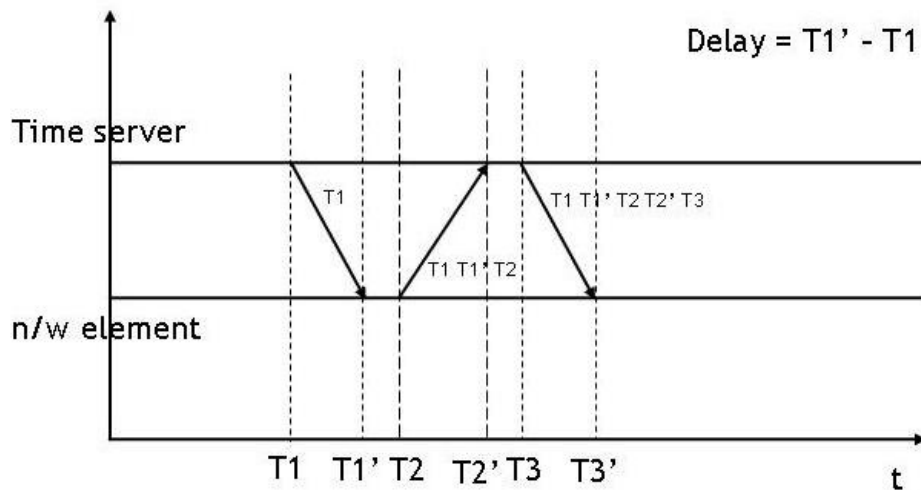
Synchronisation of time stamps

NTP approach



Synchronisation of time stamps

PTP approach

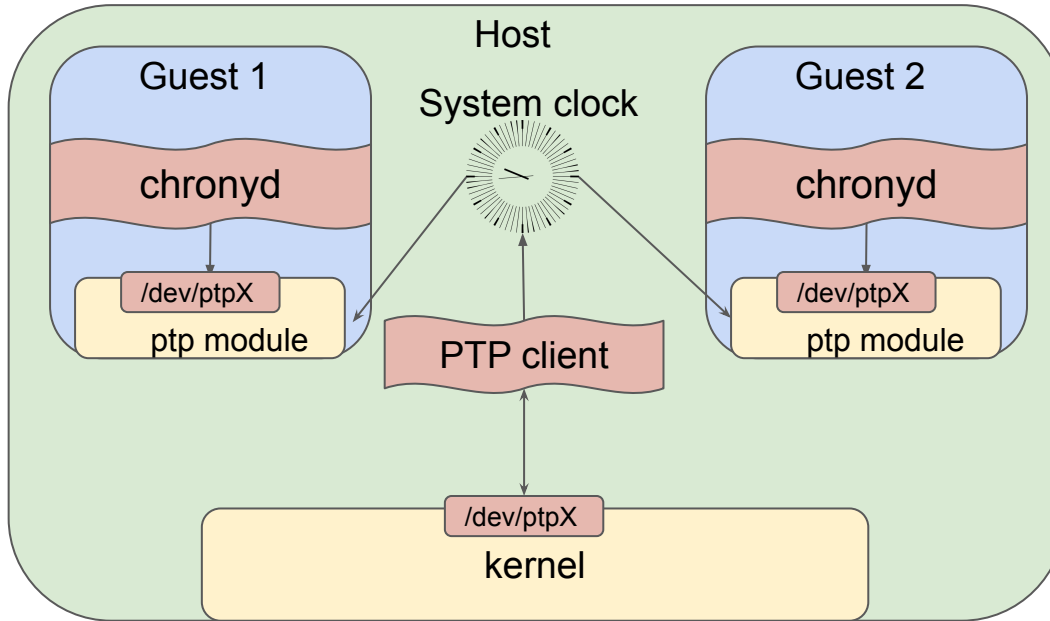


Clock offset

$$(T1' - T1 - T2' + T2) / 2$$

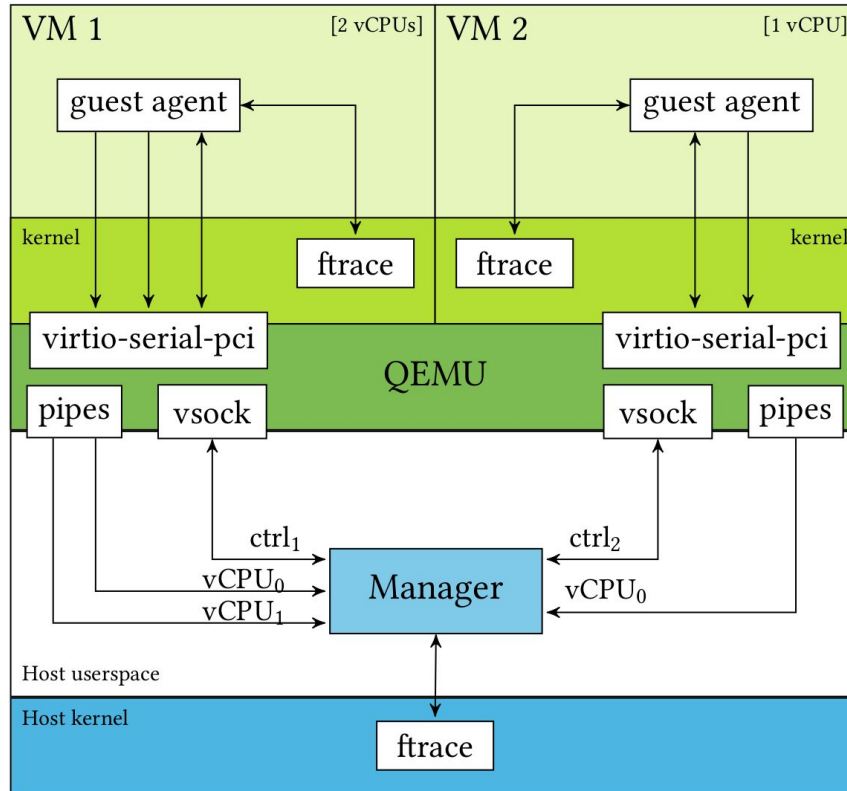
Synchronisation of time stamps

PTP in VM environment



- sub-microsecond accuracy
- clocks must be synchronized in advanced, before running the trace

VM Tracing Overview



Trace data transfer

- FIFOs, in case of KVM hypervisor
- vsockets
- splice()

Timestamps synchronisation

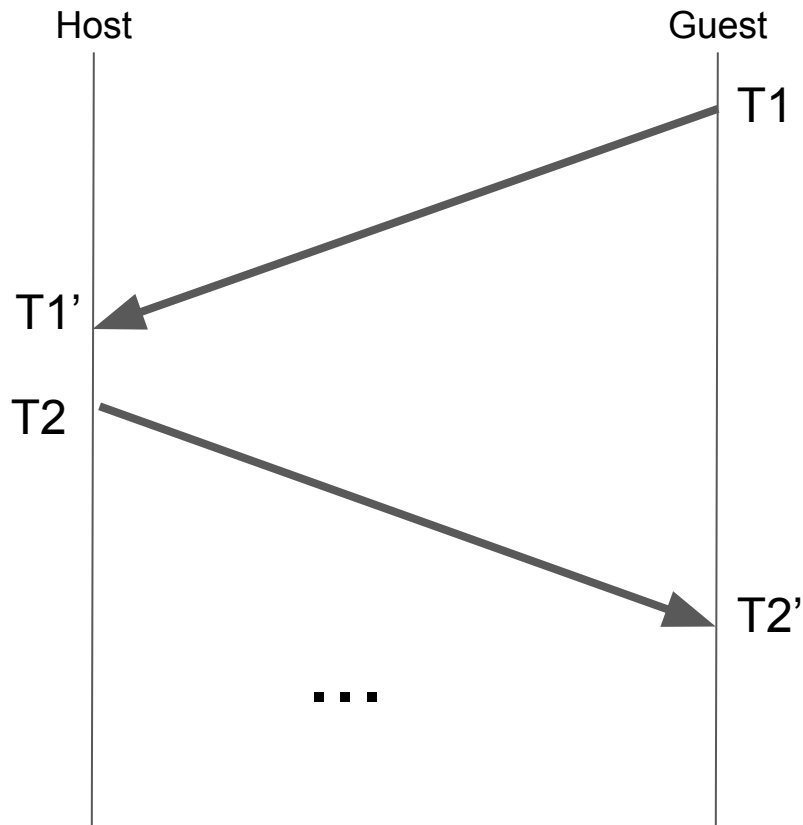
using system time

`trace-cmd record --date ...`

- converts trace timestamps to system time
- all clocks must be synchronized in advanced
- accuracy depends on clock synchronization

Timestamps synchronisation

using PTP-like algorithm



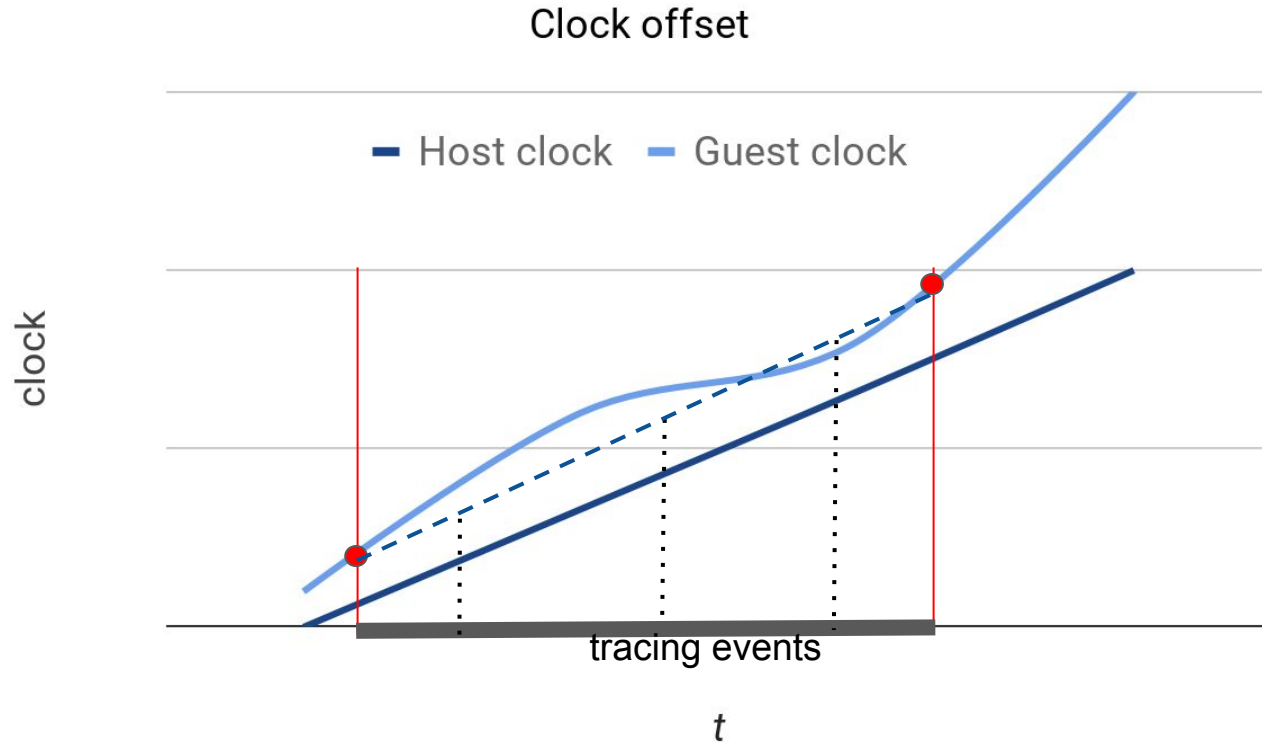
Clock offset

$$(T1' - T1 - T2' + T2) / 2$$

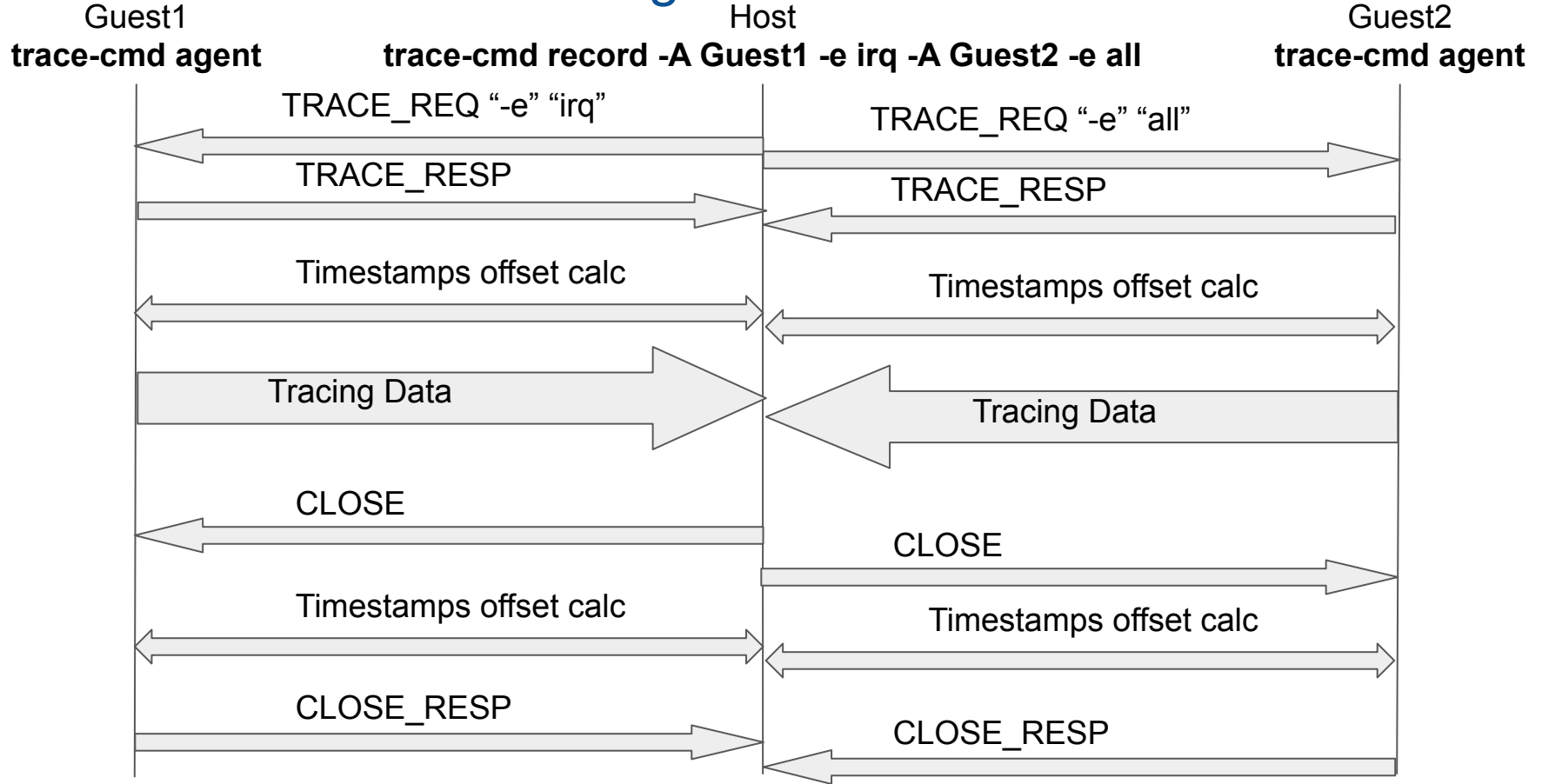
- round trip time is not symmetric
- no hardware timestamping
- Up to few hundred packets are exchanged in one clock offset measurement
- ftrace is used to get the packet times

Timestamps synchronisation

using PTP-like algorithm



Tracing Protocol Overview



trace files

trace.dat

meta data

- TraceID
- Guests trace ID array
- Host's hypervisor task <-> guest's VCPU mapping array

trace data

.....

trace-guest.dat

meta data

- TraceID
- Host trace ID
- Timestamp offsets array

trace data

.....

trace files

trace.dat

meta data

- TraceID
- Guests trace ID array
- Host's hypervisor task <-> guest's VCPU mapping array

trace data

.....

trace-guest.dat

meta data

- TraceID
- Host trace ID
- Timestamp offsets array

trace data

.....

KernelShark



Final words

- **Ftrace, trace-cmd and KernelShark are originally developed by Steven Rostedt**
- **VM tracing implementation was started by Yoshihiro Yunomae**
- **Redesigned by VMware's Open Source team last year**
- **Newly released trace-cmd 2.9 has support of VM tracing**
- **Next major release of KernelShark will include support for VM tracing visualisation**

Resources

ftrace (part of the Linux kernel)

<https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git>

trace-cmd

<https://trace-cmd.org/>

KernelShark

<https://www.kernelshark.org>

Bugzilla

<https://bugzilla.kernel.org/buglist.cgi?component=Trace-cmd%2FKernelshark>



Thank You !

tstoyanov@vmware.com

@VMWopensource
blogs.vmware.com/opensource