

## Virtio-trace

- Towards the flexible fast interconnection between guest and host for tracing

### Tracing Summit 2012

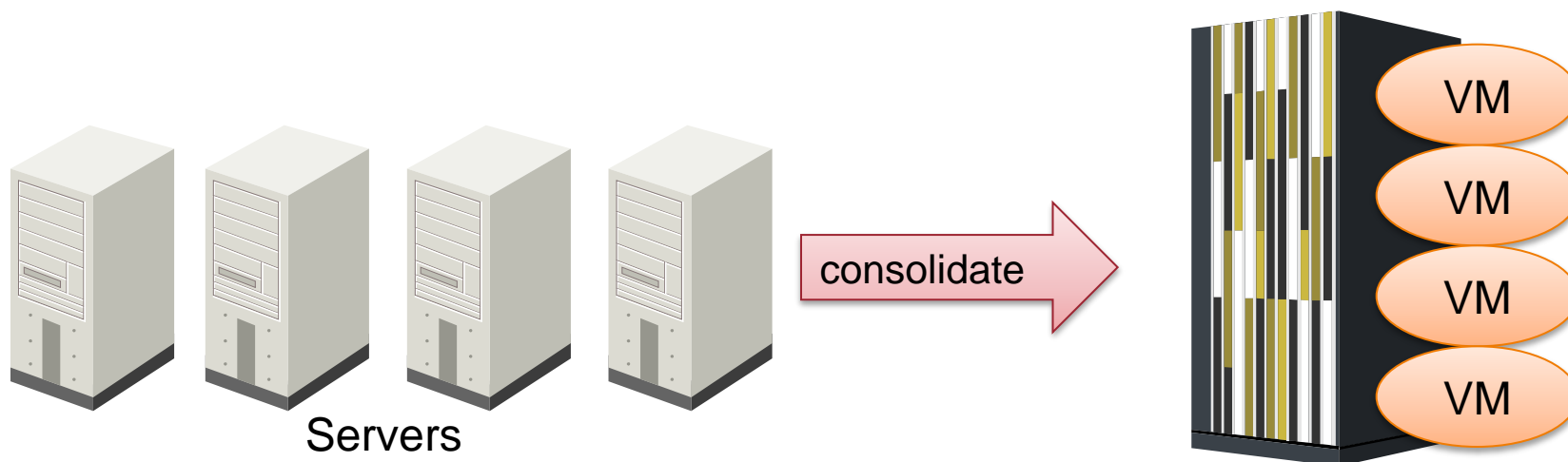
Masami Hiramatsu <masami.hiramatsu.pt@hitachi.com>

Yoshihiro YUNOMAE <yoshihiro.yunomae.ez@hitachi.com>

Linux Technology Center

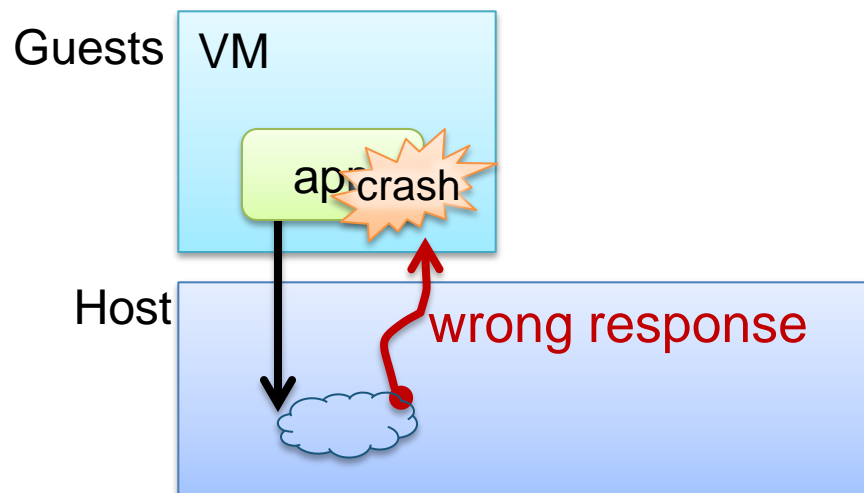
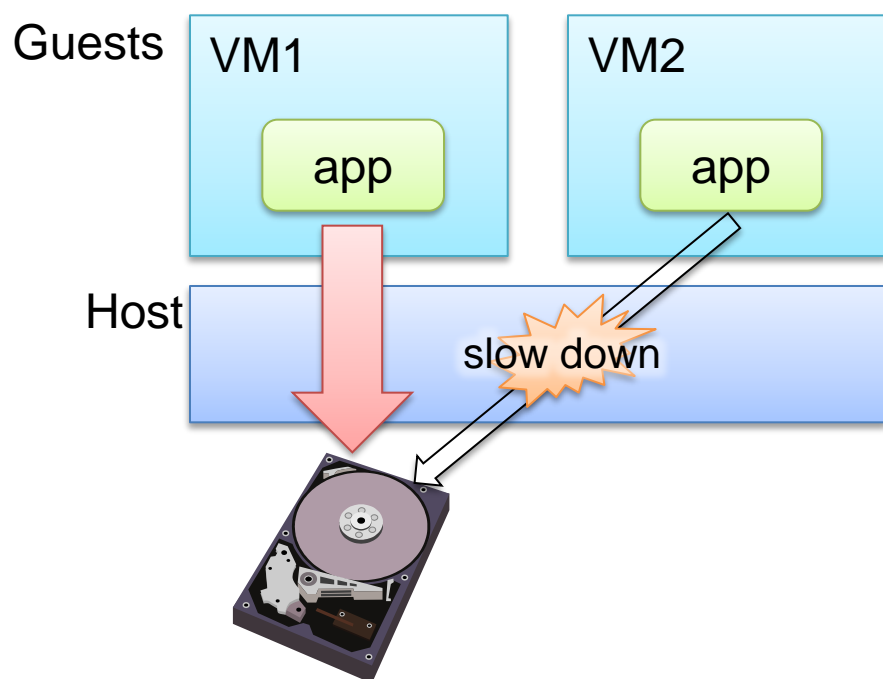
Yokohama Research Lab., Hitachi, Ltd.

- Enterprise systems are moving on (private/public) cloud which uses virtualization technology and aim for system consolidation.
  - Multiple servers run on one physical system
  - This makes system trouble shooting harder



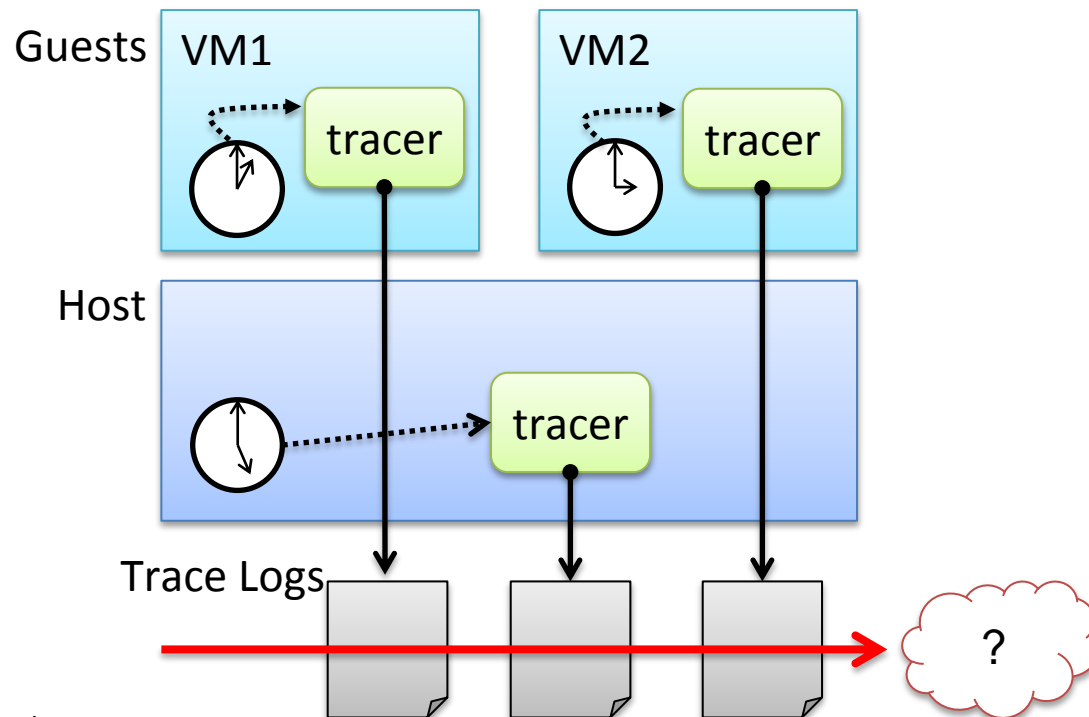
## ■ A guest VM ...

- Can be affected by other VM operation
- Can die by host OS or hypervisor's bug



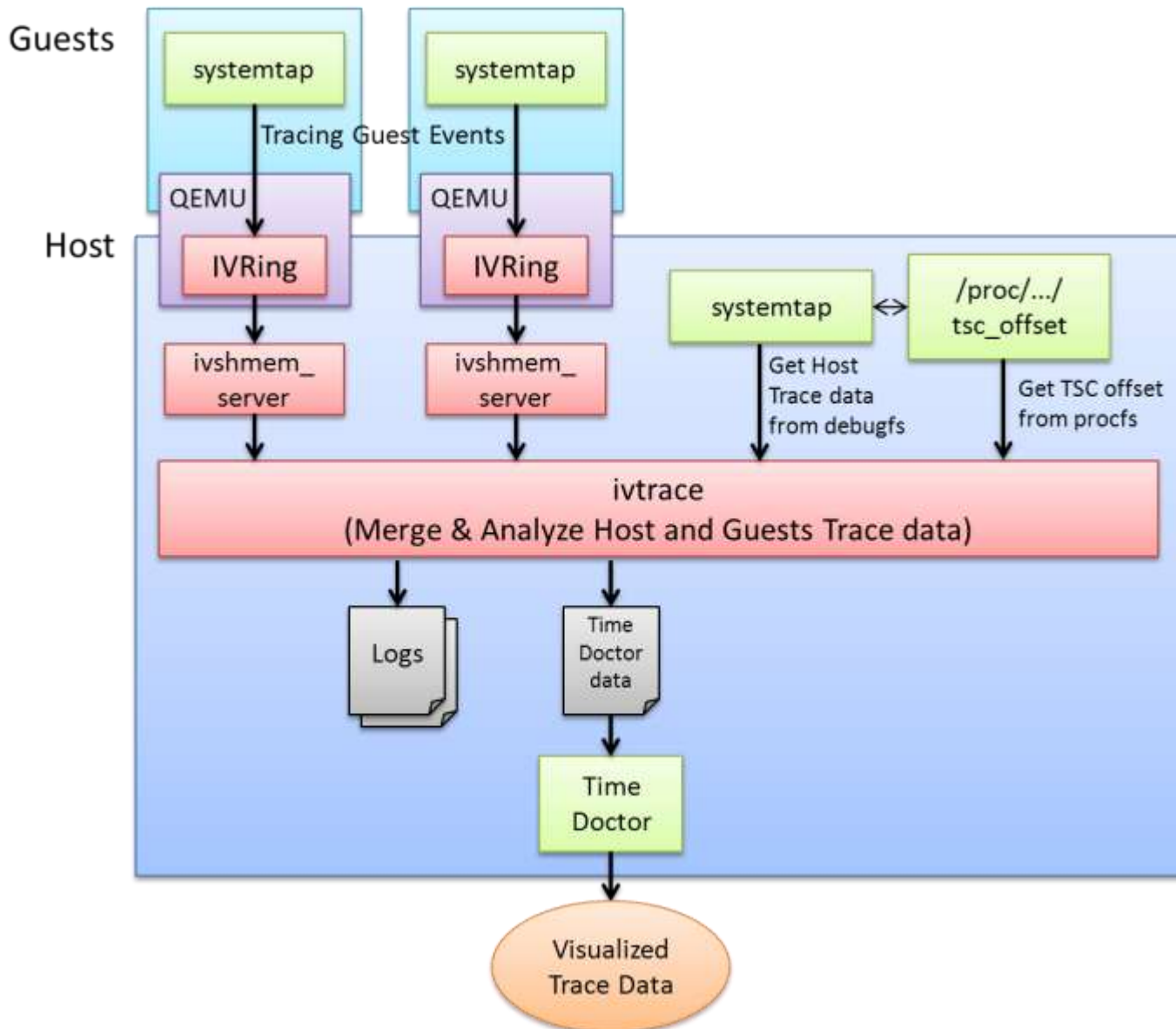
## Host-Guest inter-VM tracing helps root cause analysis

- Tracing virtualized system has following challenges
  1. Synchronize time-stamp for each VM's log
    - Each VM has own clock
  2. Collect guests trace-log from host without overhead
    - Too huge trace logs (# of VM times of logs)

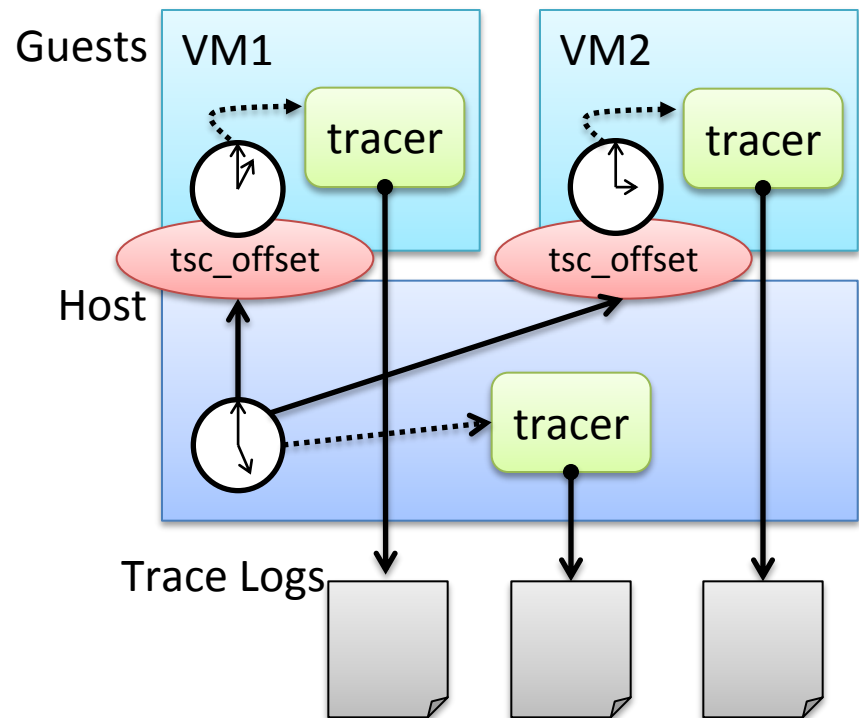


- Use **systemtap** as a tracer
  - Trace both of kernel and user applications
- Use **tsc\_offset** to adjust guests' time-stamp-counter(tsc)
  - This prototype use TSC for timestamp
  - Also, a VM is pinned on a CPU
- Use **IVRing** to pass the trace logs from guests to host
  - IVRing is an implementation of Inter-VM Ring buffer using IVShmem
- Use **TimeDoctor** to visualize trace data

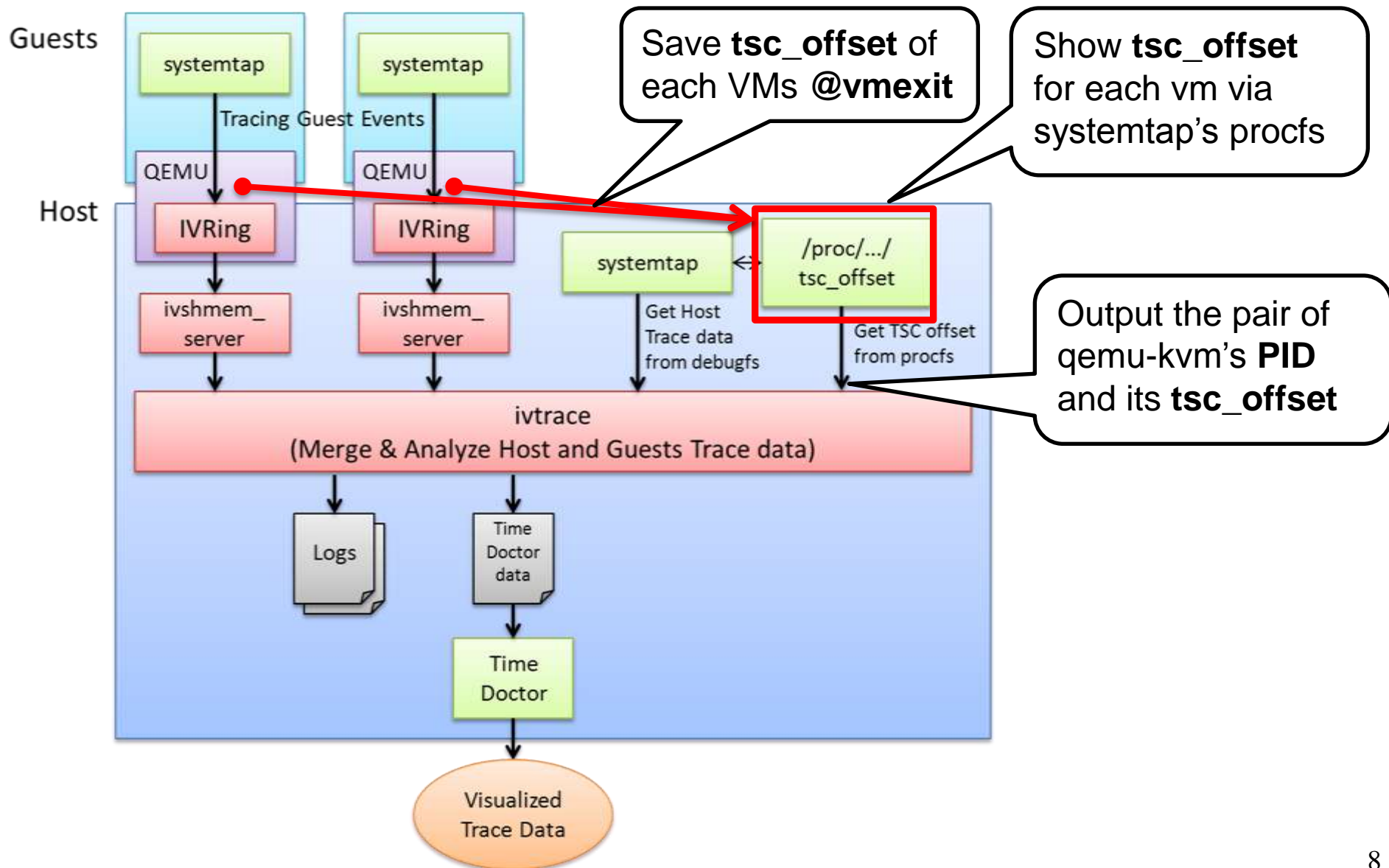
# First Prototype Overview



- Each guests has virtualized independent tsc
  - The substitution of host's tsc and guests' one is called **tsc\_offset**
  - Each tracers in guests uses own tsc to record time stamp, and we can adjusting this tsc using **tsc\_offset** to merge logs.
  - Of course, this depends on “**constant\_tsc**”

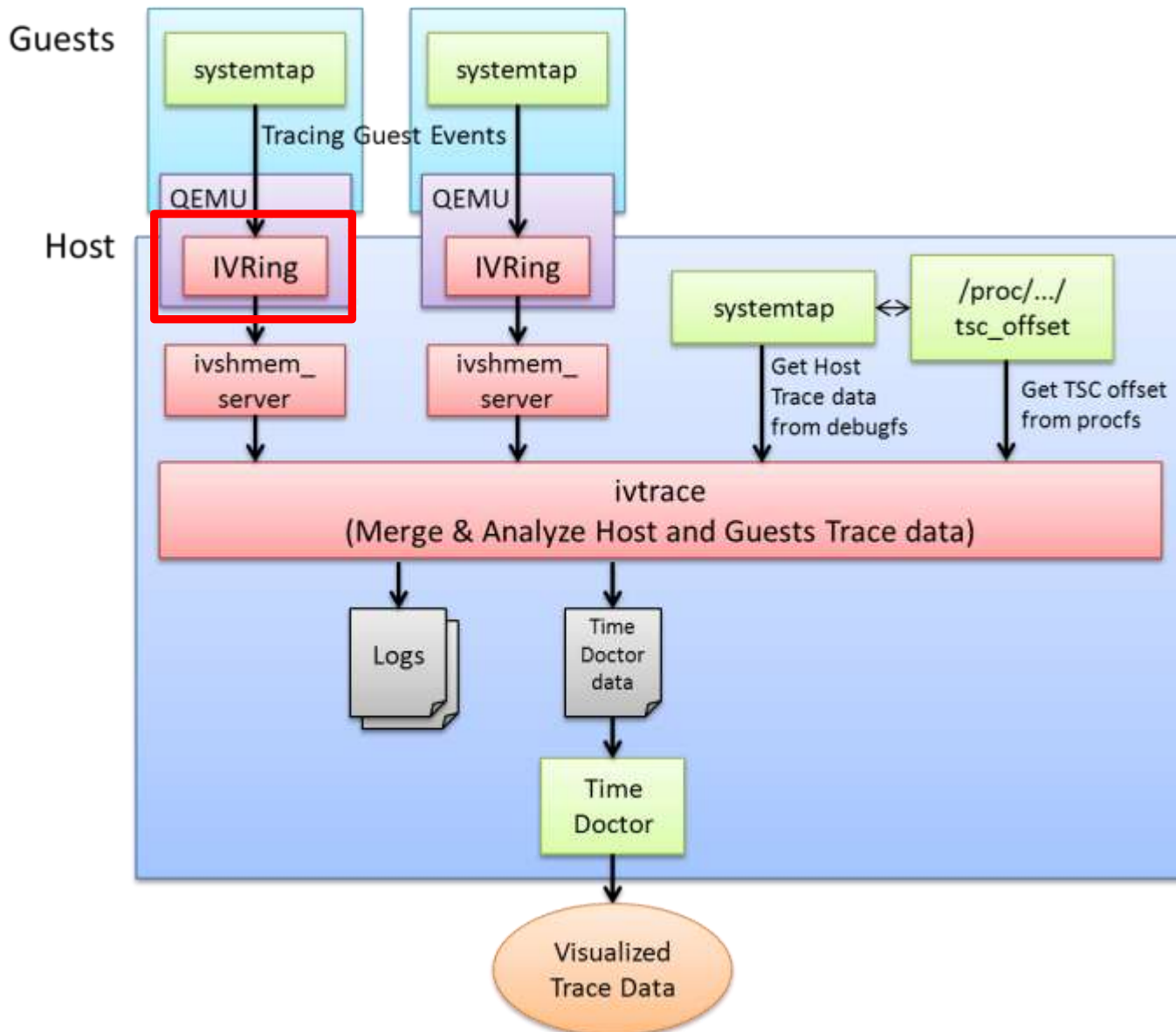


# Adjusting TSC on Prototype

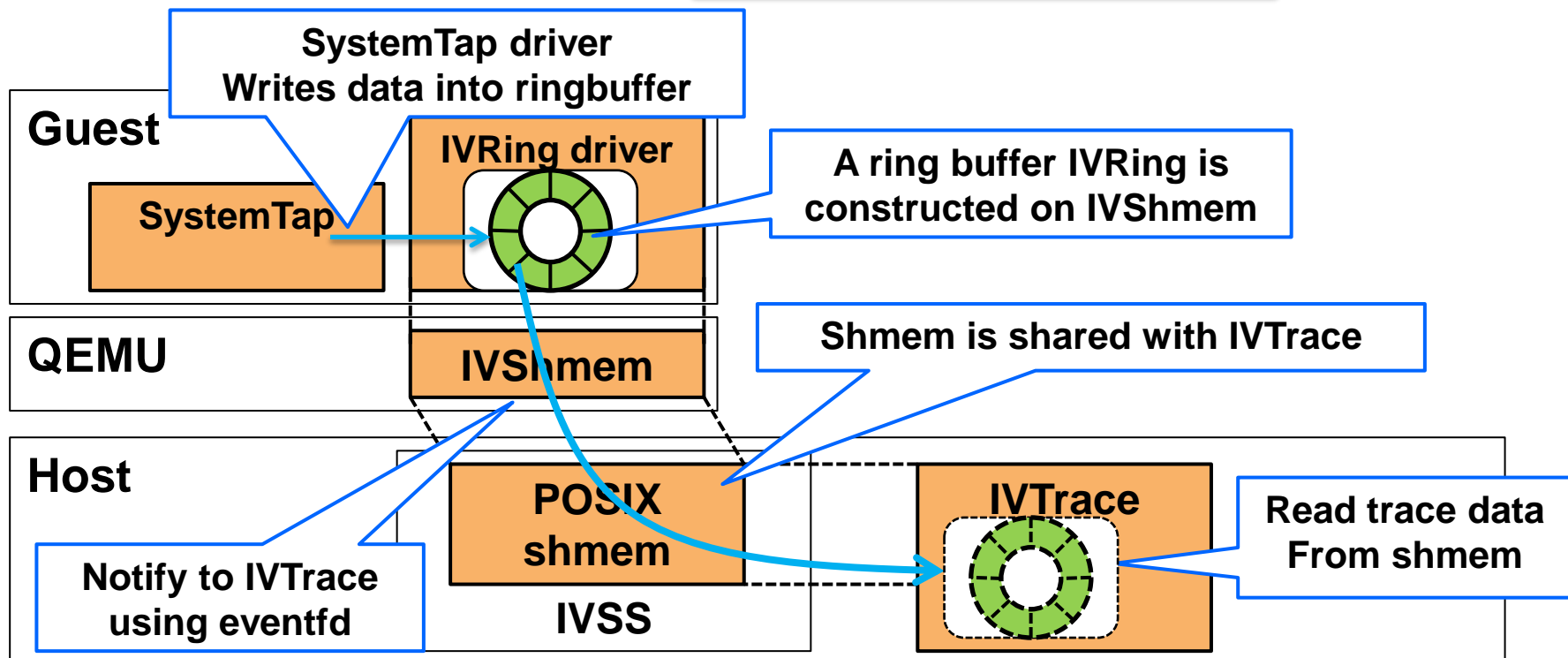




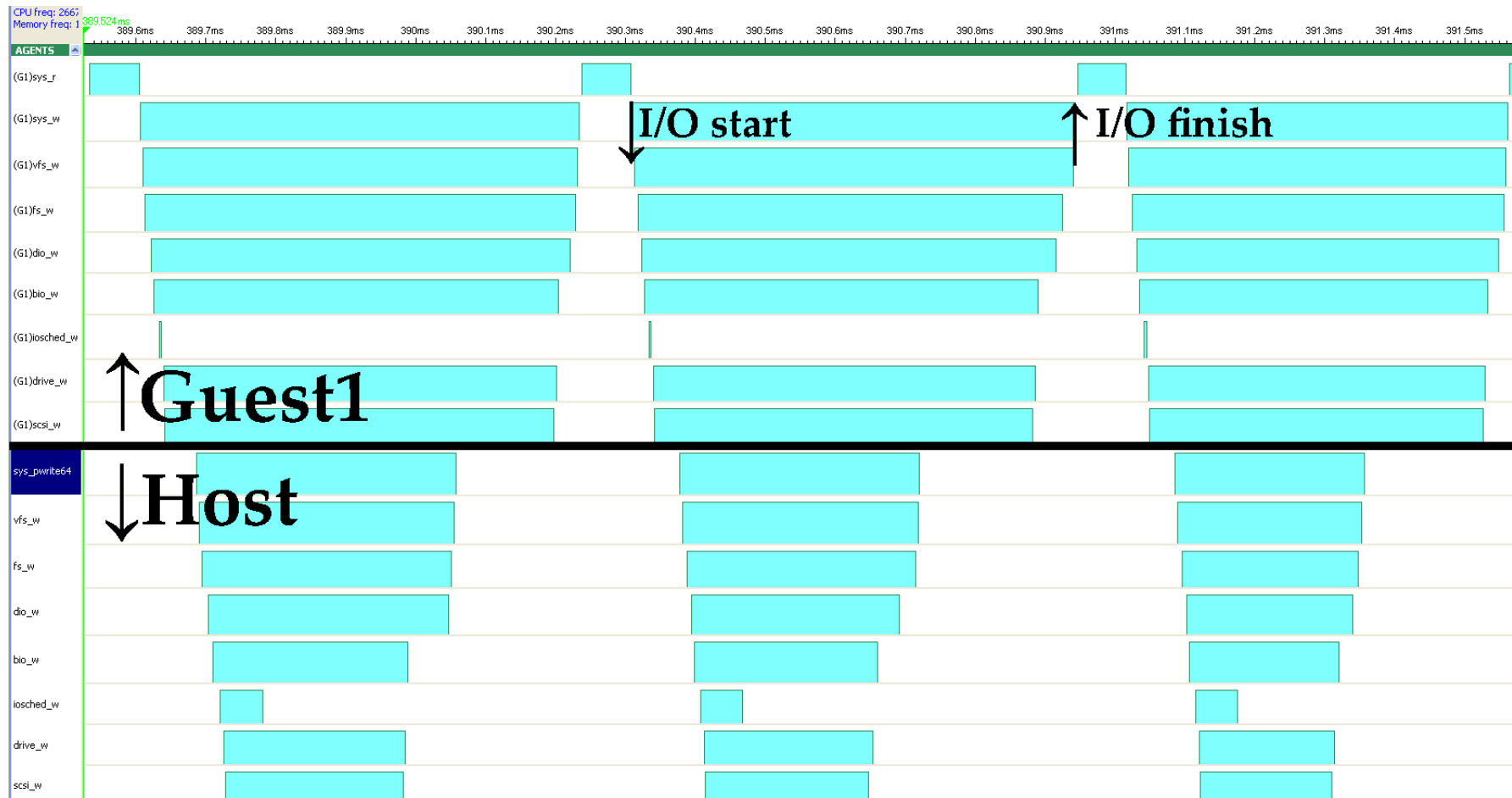
# IVRing - Fast Interconnect for Tracing



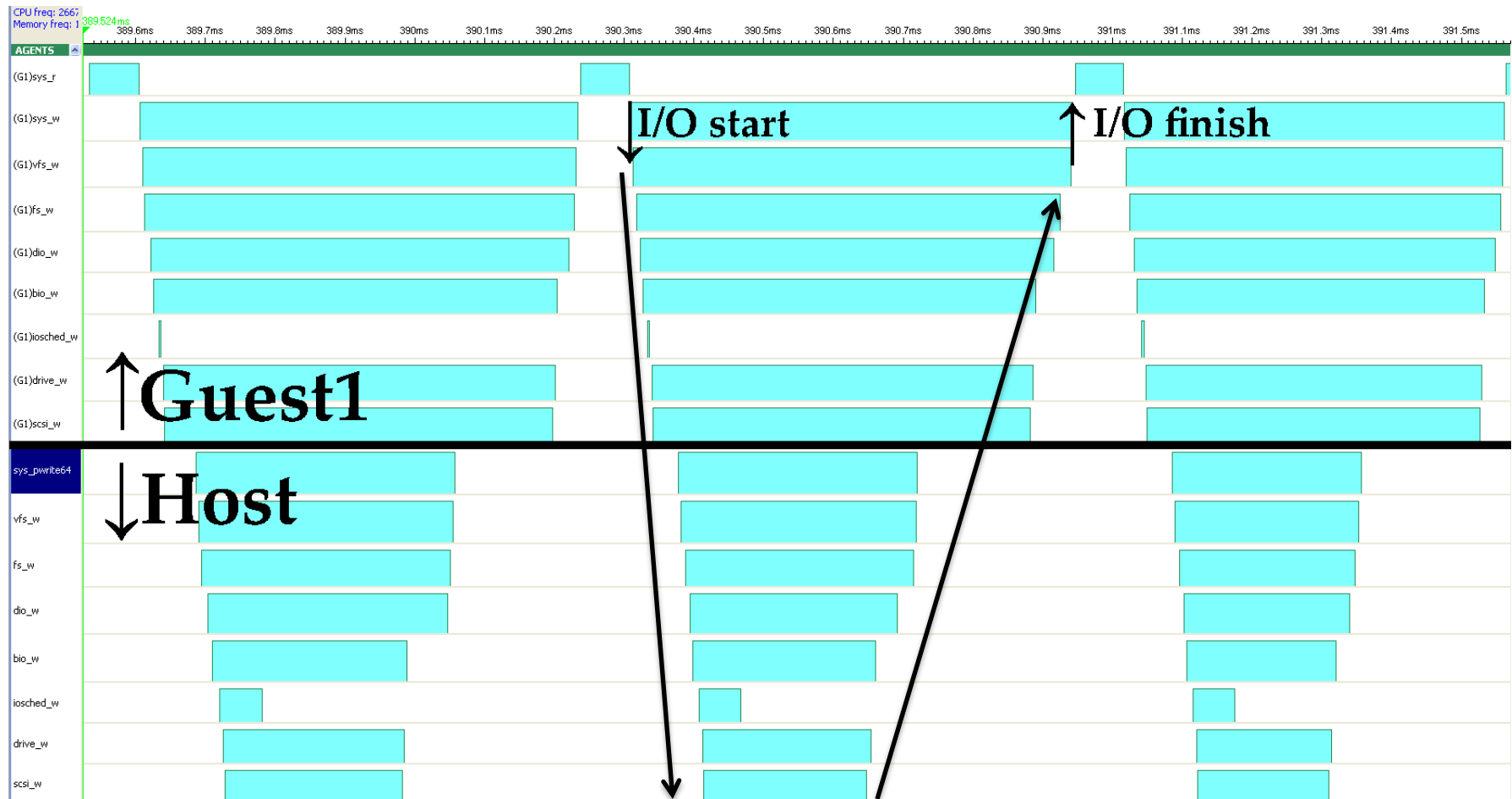
- A ring-buffer IVRing is constructed on IVShmem as a data path for trace data of a guest.
- IVShmem is a memory-PCI device
  - Backend memory is a posix shm.
- IVTrace can read the data without memory copying.



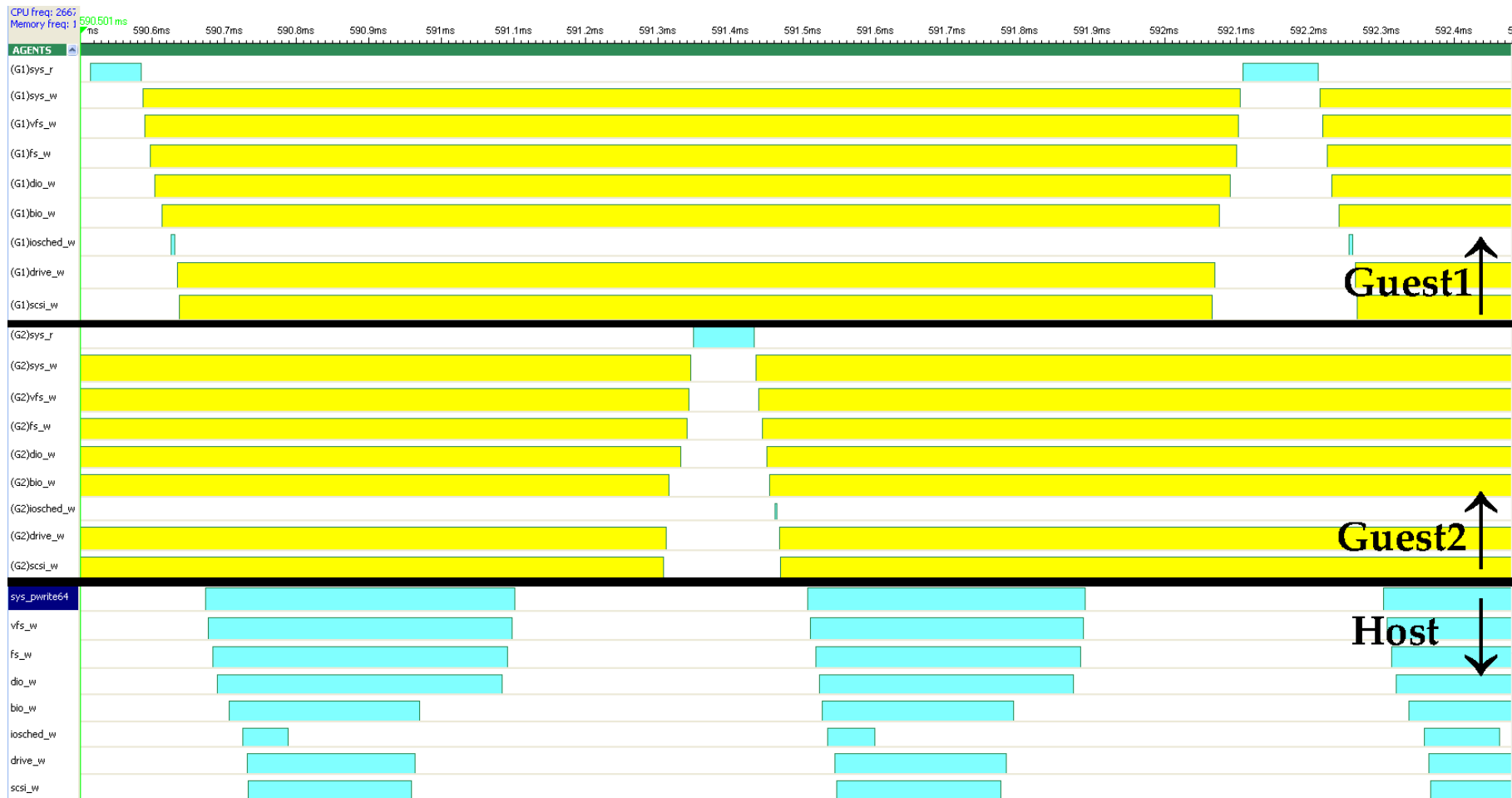
## ■ The case of no contentions



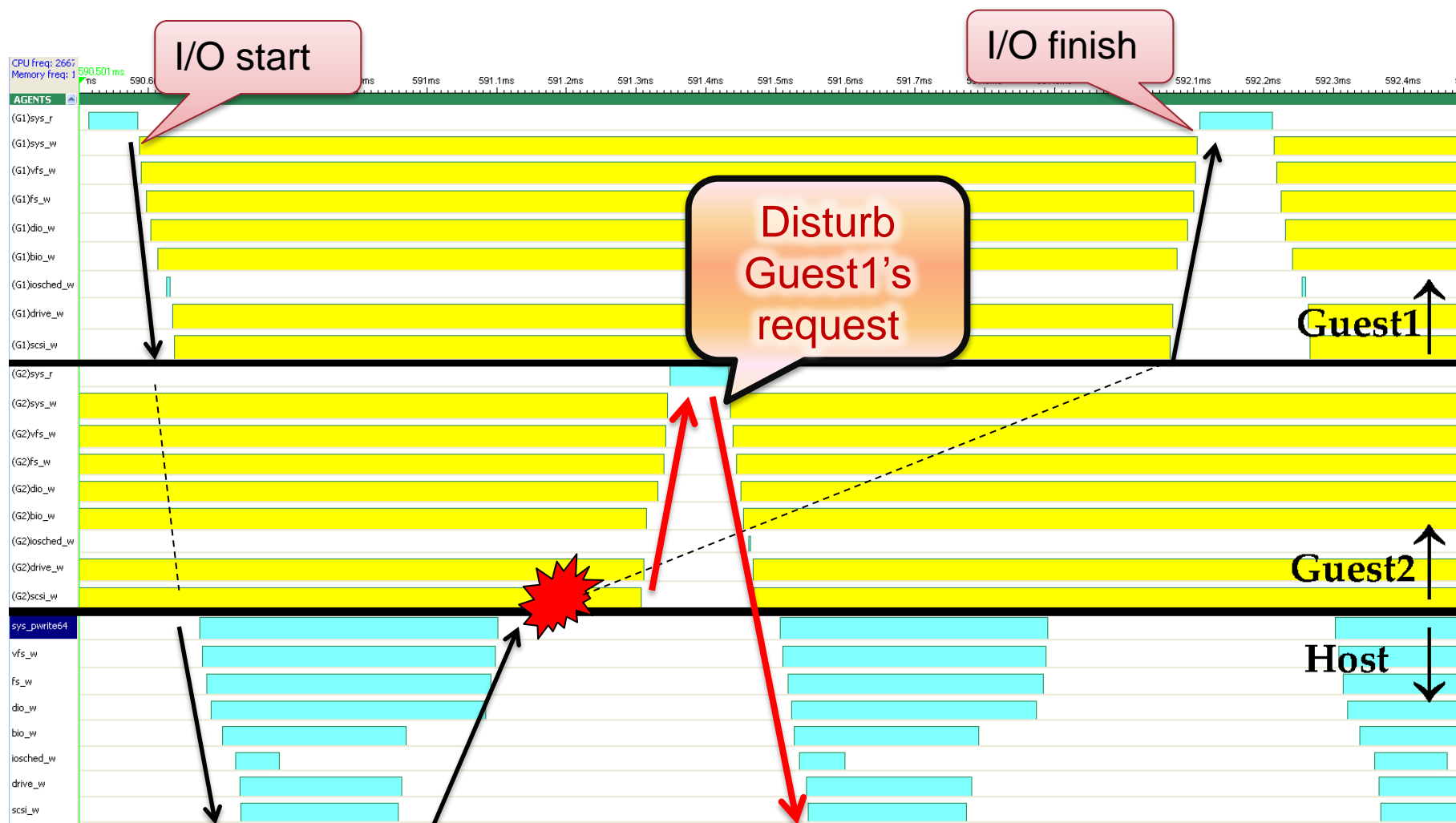
## ■ The case of no contentions



## ■ The case of contentions. I/O slowdown



## ■ The case of contentions. I/O slowdown



## ■ Maintainers Don't like it

```
(2012/06/06 8:22), Greg Kroah-Hartman wrote:  
> On Wed, Jun 06, 2012 at 07:03:06AM +0800, Anthony Liguori wrote:  
>> On 06/05/2012 09:10 PM, Borislav Petkov wrote:  
>>>  
>>> Yet another ring buffer?  
>>>  
>>> We already have an ftrace and perf ring buffer, can't you use one of those?  
>>  
>> Not to mention virtio :-)  
>>  
>> Why not just make a virtio device for this kind of thing?  
>  
> Yeah, that's exactly what I was thinking, why reinvent things again?
```

## ■ Points

- NO “yet another ring buffer” in linux kernel
  - Use ftrace and perf ring buffer for guest recording
- Use virtio instead of ivshmem

## ■ Virtio-shmem (new device)

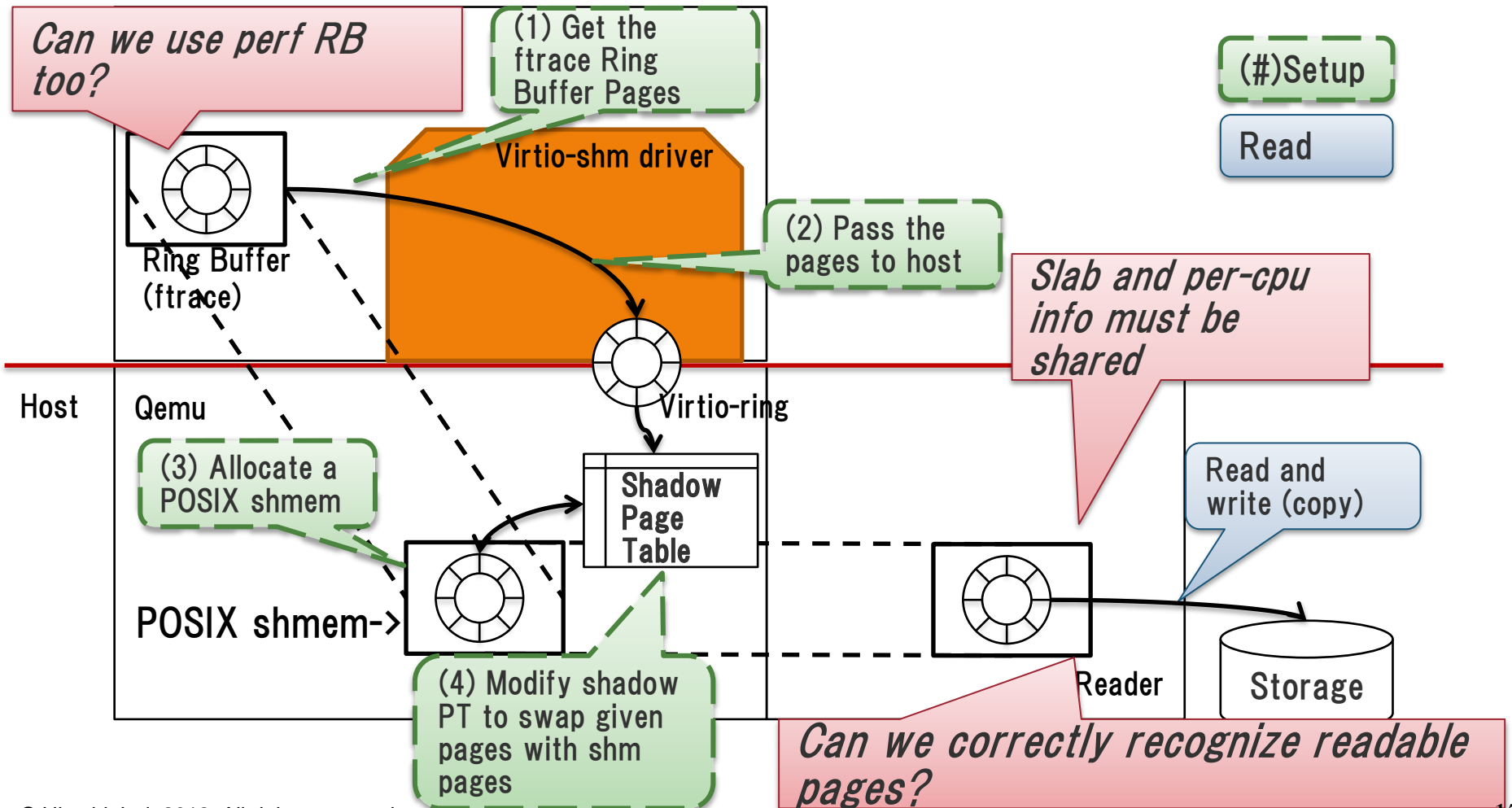
- Virtio device which provides APIs for assigning the guest pages to shared memory in host
  - Guest can assign any page to host's shmem
  - Qemu remaps original pages with the pages on shared memory
  - Similar to the ivshmem, but no big PCI address space required

## ■ Virtio-serial with splice (enhancement)

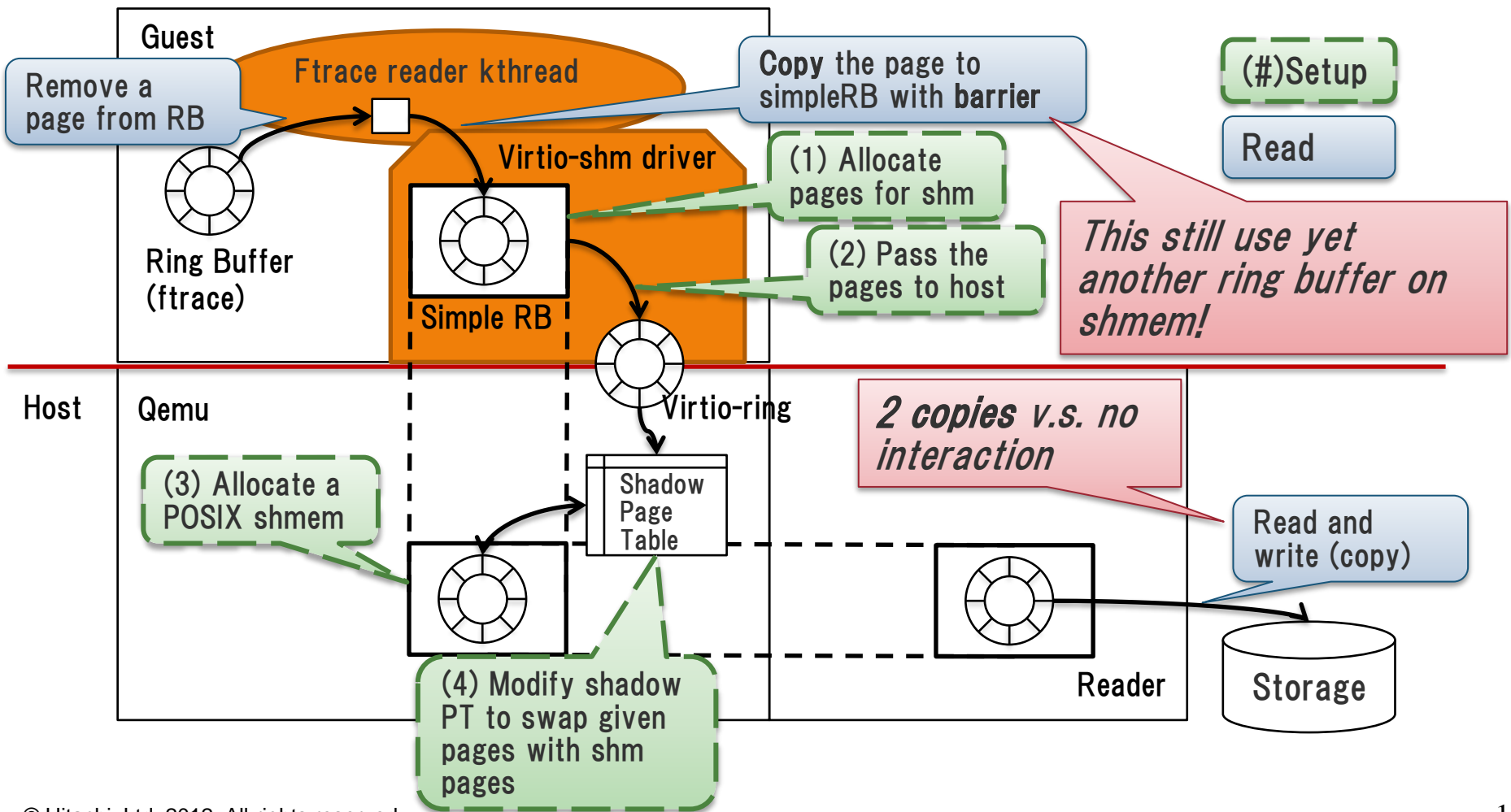
- Virtio device which provides chardev interface for guest
  - Guest can “splice” its data into the char device
  - Qemu copies data page to host-side pipe
  - Vhost can offload the copying process



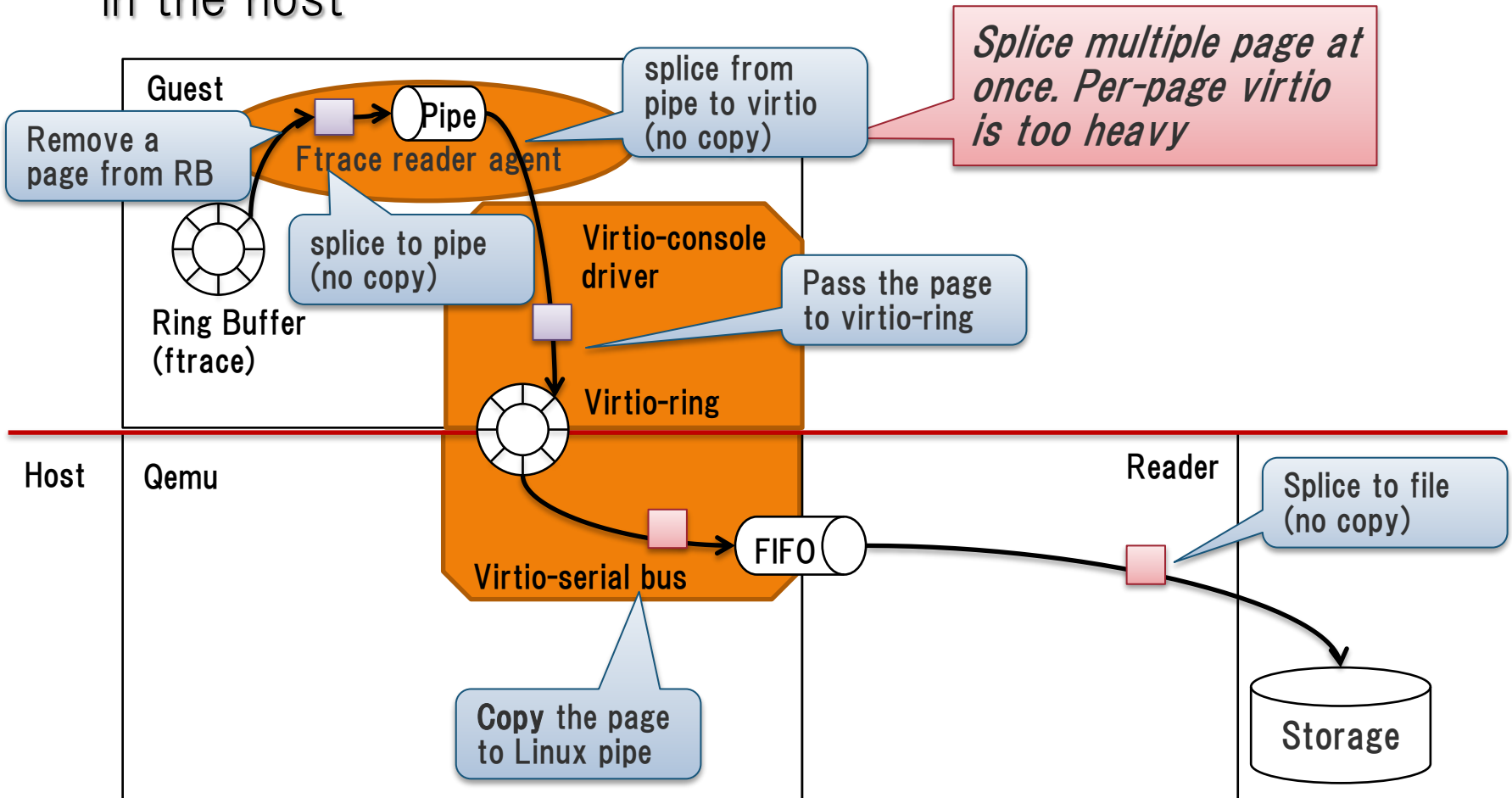
- Export Ftrace ring-buffer pages directly to host
  - How we can export per-cpu kmalloc object?



- To generalize interface and guarantee memory coherency
  - At first, we have simple ring-buffer for memory coherency



- Virtio-serial opens chardev in the guest and FIFO(named pipe) in the host



# Comparison

	IVRing	Virtio-serial w/ splice	Virtio-shmem w/ simpleRB	Virtio-shmem
<b>Qemu-Reader</b>	Shmem	Pipe (1way)	Shmem	Shmem
<b>Guest-Qemu</b>	Shmem(PCI)	virtio	Shmem(EPT)	Shmem(EPT)
<b># of Copies</b>	1 (reader to file)	1 (virtio to pipe)	2 (in Guest, reader to file)	1 (reader to file)
<b>Guest-host interaction</b>	No	Once per I/O (16pages or more)	No	No
<b>Supported tracer</b>	SystemTap	Ftrace, user tools	Ftrace, and others?	Ftrace only
<b>No another RB</b>	IVRing	No	Simple RB	No
<b>Use virtio</b>	No	Yes	Hmm...	Hmm...
<b>Buffer Resize</b>	No	Support	Support	Support
<b>SMP scaling</b>	RB w/ lock	Per-cpu pipes	Per-cpu shmem	Per-cpu shmem
<b>VCPU hot-add</b>	Support (w/lock)	Add channels	Add shmem	Add shmem
<b>Live migration</b>	No	Possible	Possible	Possible
<b>Expectation</b>	Not scalable, not acceptable	Upstream Acceptable	Need Discussion	Need Discussion

# Comparison

	IVRing	Virtio-serial w/ splice	Vi sir	y/o
<b>Qemu-Reader</b>	Shmem	Pipe (1way)	Sh	
<b>Guest-Qemu</b>	Shmem(PCI)	virtio	Shmem(EPT)	Shmem(EPT)
<b># of Copies</b>	1 (reader to file)	1 (virtio to pipe)	2 (in Guest, reader to file)	1 (reader to file)
<b>Guest-host interaction</b>	No	Once per I/O (16pages or more)	No	No
<b>Supported tracer</b>	SystemTap	Ftrace, user tools	Ftrace, and others?	Ftrace only
<b>No another RB</b>	IVRing	No	Simple RB	No
<b>Use virtio</b>	No	Yes	Hmm...	Hmm...
<b>Buffer Resize</b>	No	Support	Support	Support
<b>SMP scaling</b>	RB w/ lock	Per-cpu pipes	Per-cpu shmem	Per-cpu shmem
<b>VCPU hot-add</b>	Support (w/lock)	Add channels	Add shmem	Add shmem
<b>Live migration</b>	No	Possible	Possible	Possible
<b>Expectation</b>	Not scalable, not acceptable	Upstream Acceptable	Need Discussion	Need Discussion

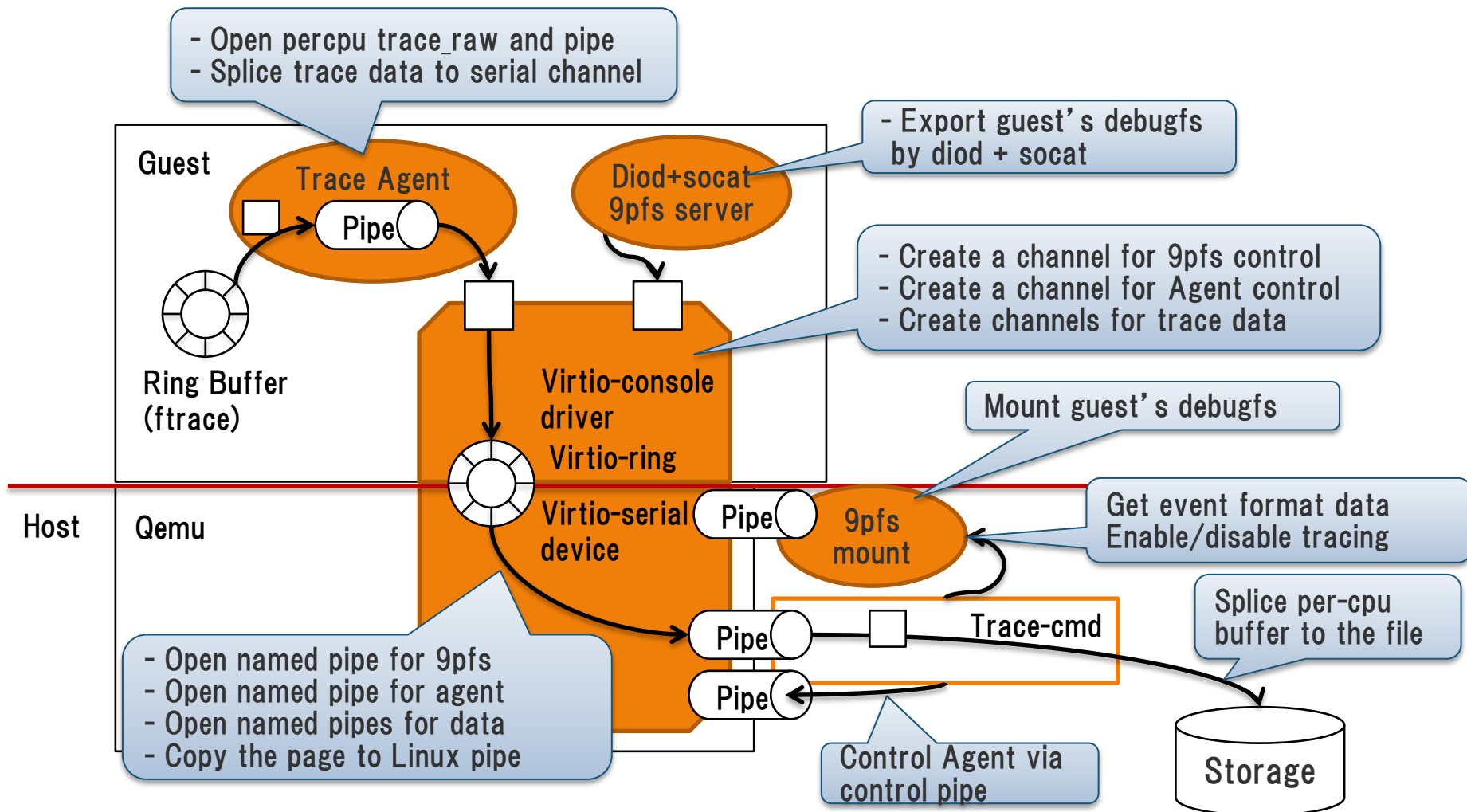
## Chosen method

- Simple implementation
- No claiming points
- Fit to virtio framework

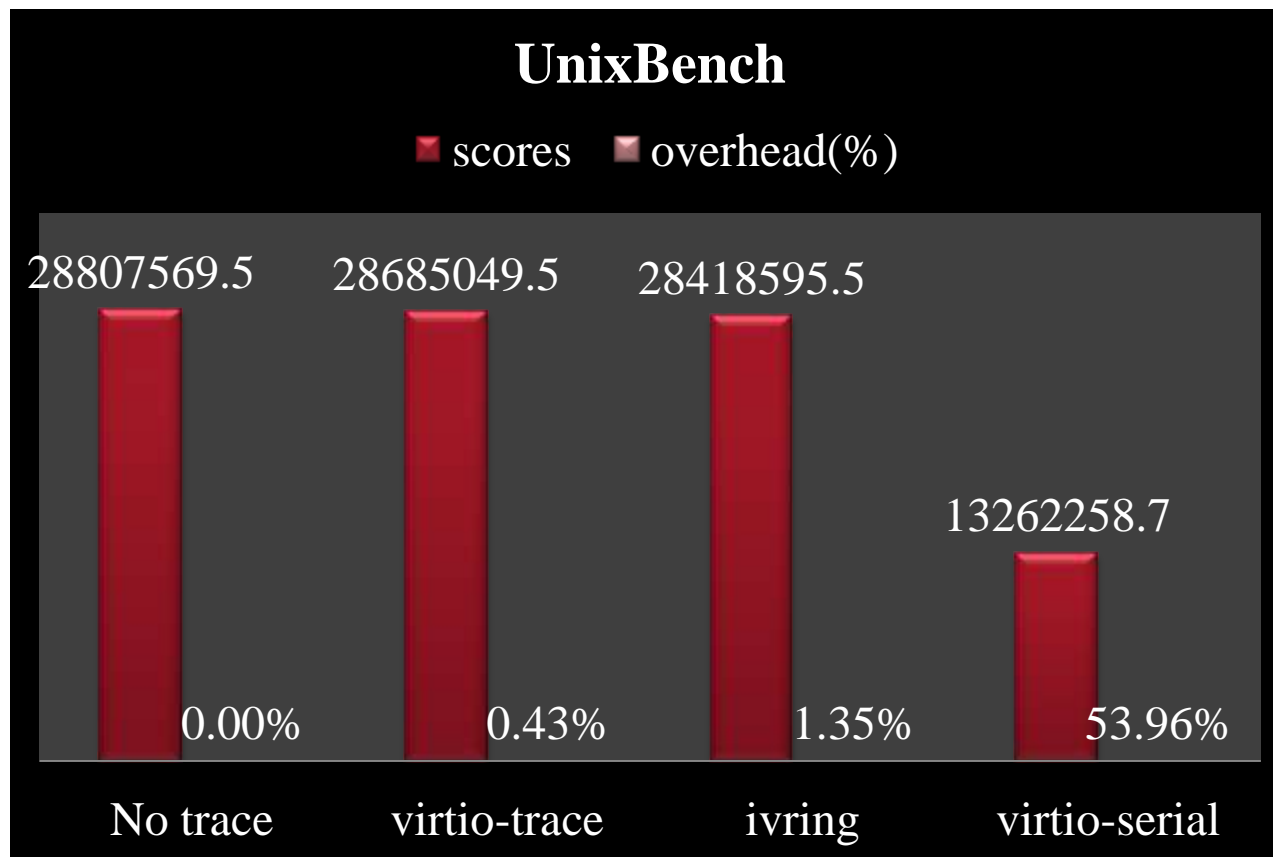


# Virtio-trace(serial w/ splice) Prototype

- Prototype consists of trace-agent, virtio-console driver, virtio-pipe device, and trace-cmd



- Compared with IVRing, virtio-trace is really fast?
  - Running UnixBench with tracing on guest VM



Host:  
Xeon x5660  
2.8Ghz/48GB

Guest:  
Single-Core VM  
1GB

Tracing:  
Scheduler and  
Syscall events

Compared with ivring, virtio-trace has lower overhead!

## ■ Maintainers accept it 😊

On Thu, 09 Aug 2012 21:30:29 +0900, Yoshihiro YUNOMAE <yoshihiro.yunomae.ez@hitachi.com> wrote:

> Hi All,  
>  
> The following patch set provides a low-overhead system for collecting kernel  
> tracing data of guests by a host in a virtualization environment.

Thankyou!

I've applied this, and it will head into linux-next in the next few days.

Cheers,  
Rusty.

## ■ Points

- NO “yet another ring buffer” in linux kernel
  - Use virtio ring for passing data
- More generic feature (not only for tracing)
  - Maybe useful for other use, like SPICE



## ■ So, what is the next step?

- Synchronize time-stamp between guest and host
  - Tsc-based trace\_clock is an option (for constant\_tsc machine)
  - Or agent gives the guest's trace\_clock offset (but how?)
    - In generic, vmexit pattern matching can give us a hint.
- Consolidate 9pfs server to trace-agent
  - For simplicity and ease of use (reducing steps of setup)
  - We can prepare setup script for guest tracing
- Fix some issues on Qemu's chardev(serial backend)
  - Hotplug issue
  - Guest blocking issue
- Live migration support

- IVRing is not acceptable for upstream
- Improved virtio-serial to support splice is accepted
  - Build the prototype and measure the performance
- Still under development...
  - Make setup easier
  - Time synchronizing
  - Clarify “chardev” issues on Qemu
  - Fix Qemu chardev for CPU-hotplug and non-blocking

**TO BE  
CONTINUED** 

- See, IVTrace slide @ LinuxCon Japan 2012
  - Low-Overhead Ring-Buffer of Kernel Tracing  
[http://events.linuxfoundation.org/images/stories/pdf/lcjp2012\\_yunomae.pdf](http://events.linuxfoundation.org/images/stories/pdf/lcjp2012_yunomae.pdf)
  - Tracing Across Host OS and Guest OS  
[http://events.linuxfoundation.org/images/stories/pdf/lcjp2012\\_nagai.pdf](http://events.linuxfoundation.org/images/stories/pdf/lcjp2012_nagai.pdf)

- Linux is a registered trademark of Linus Torvalds.
- UNIX is a registered trademark of The Open Group.
- All other trademarks and copyrights are the property of their respective owners.