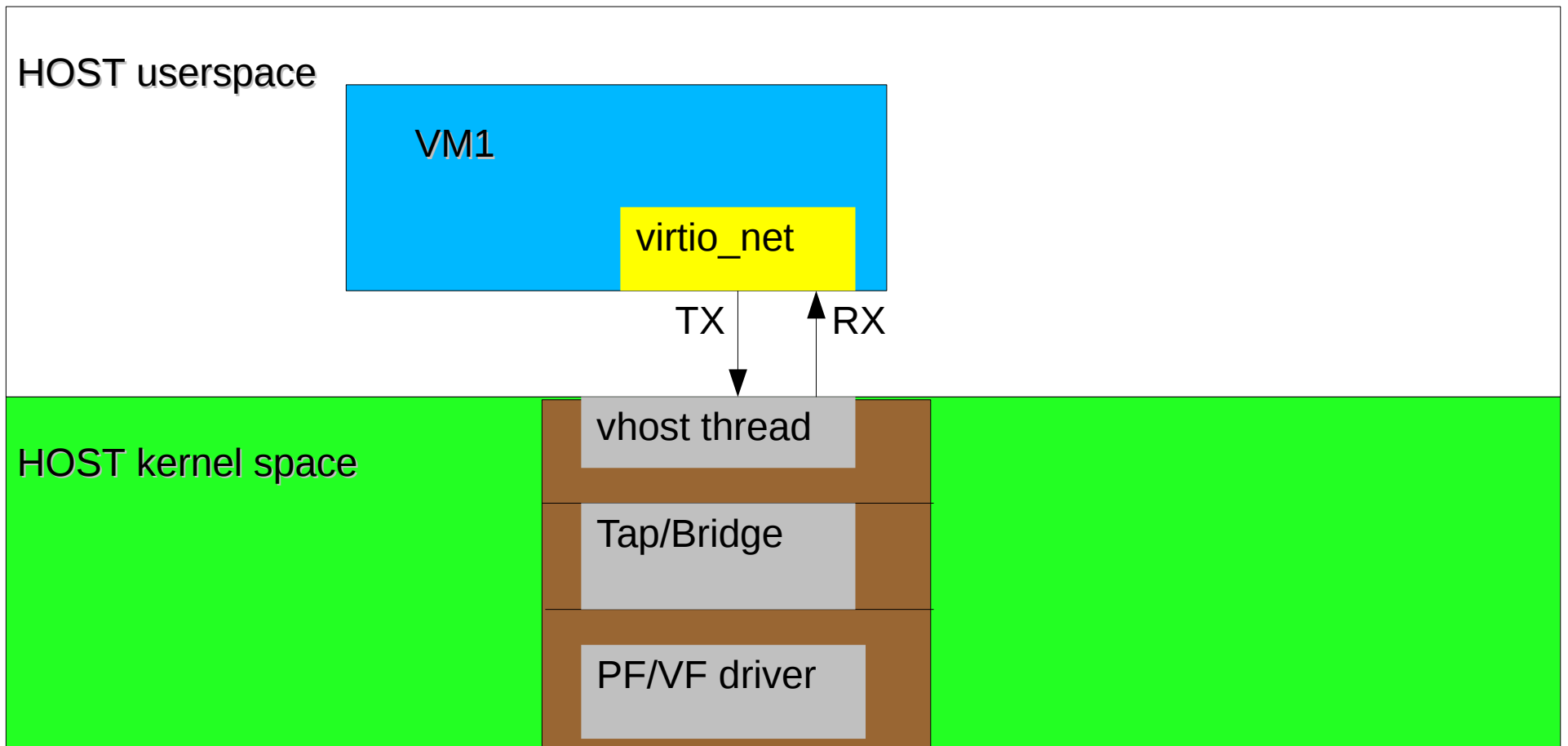


# KVM Performance – Vhost Scalability

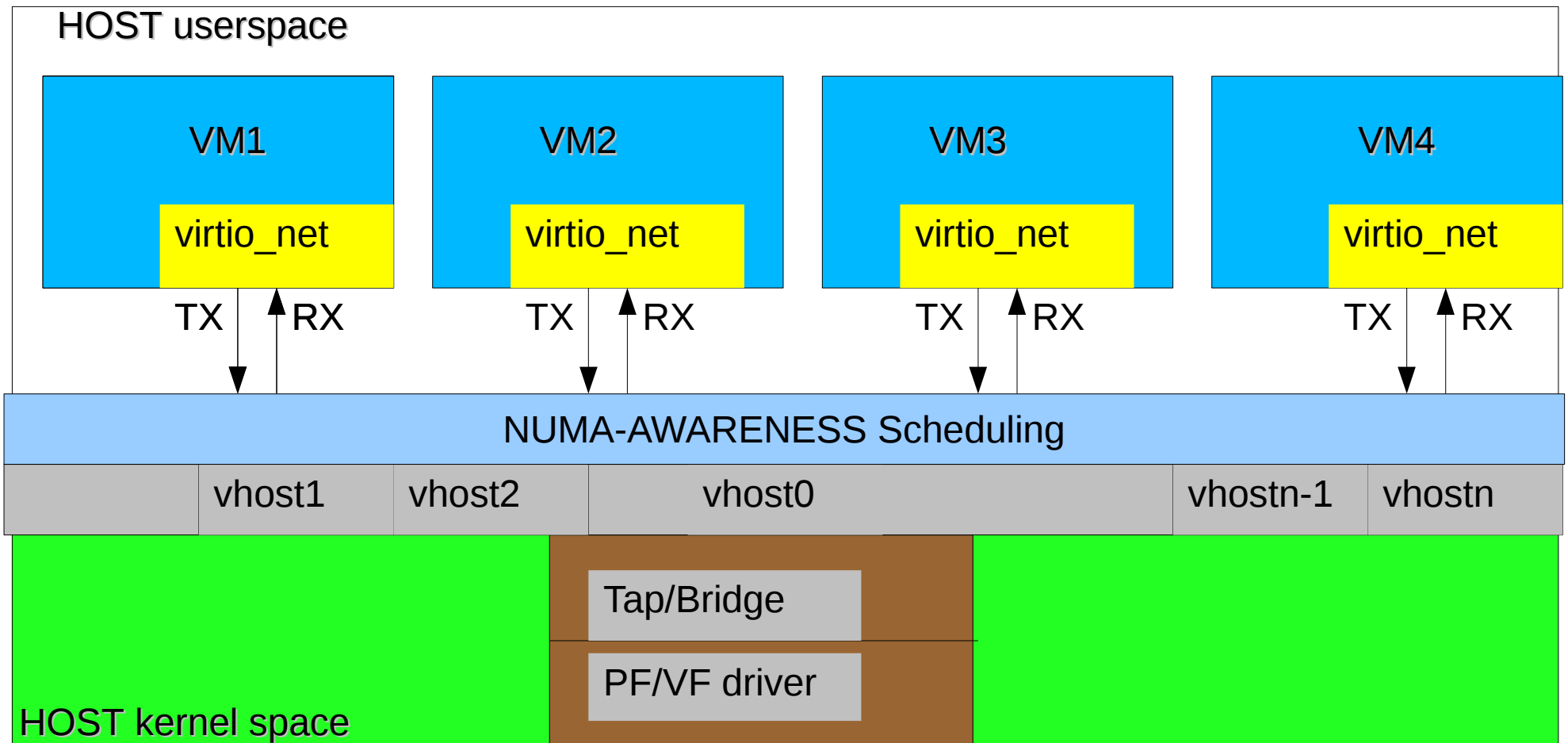
Shirley Ma, IBM  
John Fastabend, Intel

- Per virtio\_net device vhost thread



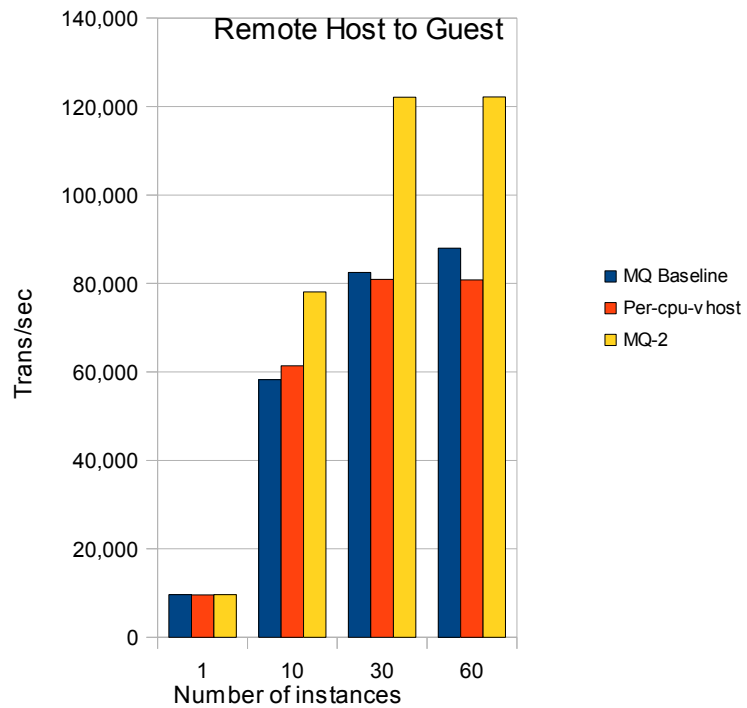
- The number of vhost threads depend on the number of virtio\_net device per VM.
- The number of vhost threads depend on the number of VMs.
- **Pros:**
  - Easy VM based cgroup control
- **Cons:**
  - Scheduling:
    - None numa-awareness
  - Performance
    - When increasing the number of virtio\_net devices per VM, or increasing the number of VMs, the performance does not scale
    - When the number of vhost threads are larger than the number of host cpus, there are lots of context switch overhead
    - vhost TX and RX are shared, so TX and RX work can't be processed simultaneously

- Per CPU vhost thread

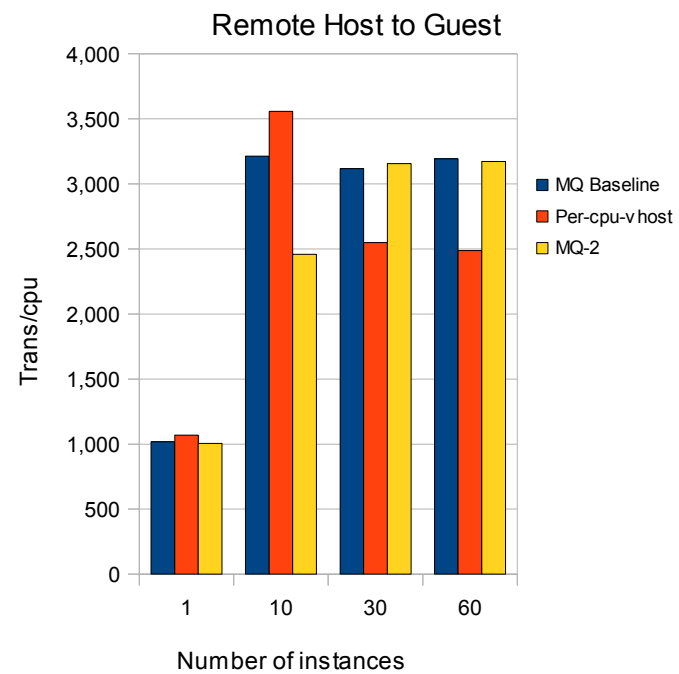


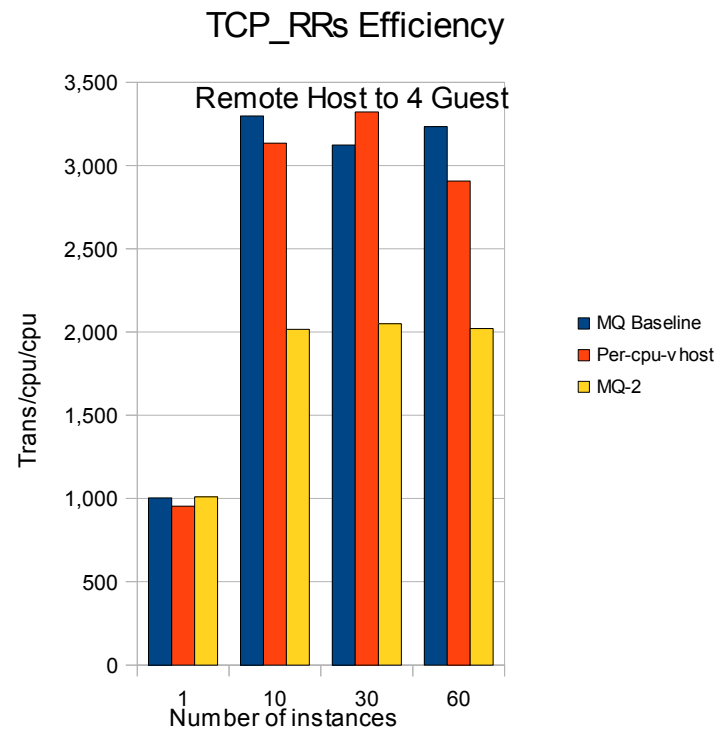
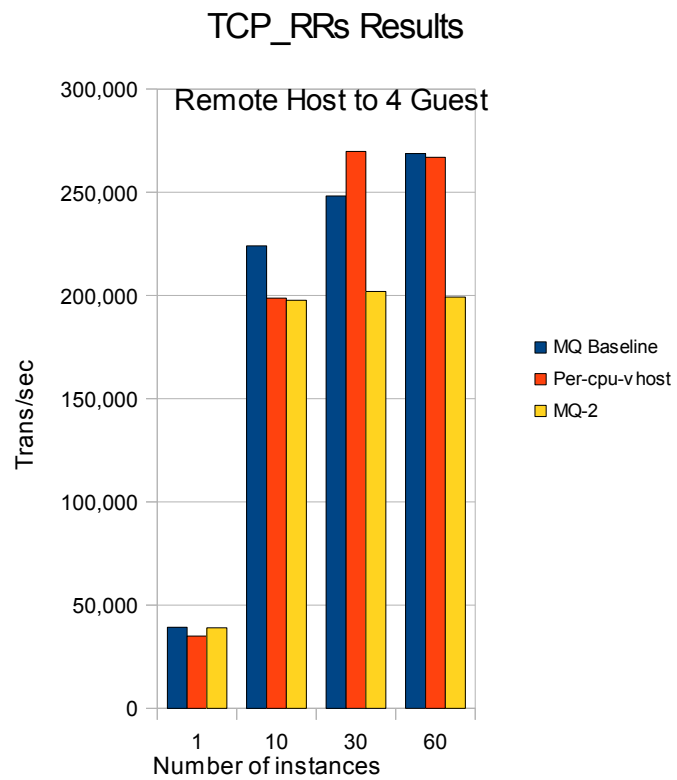
- Share vhost thread among VMs when number of VM/virtio\_net devices greater than the number of CPUs to avoid scheduling overhead
- Split vhost TX and RX work based on the workload
- NUMA-awareness scheduling
  - The vhost thread is picked up based on idelest allowed cpu in local numa node
- Cgroup control
  - The vhost thread is attached to the cgroup on the VM which the work comes from
  - When idle/need\_sched, the vhost thread is detached from the previous cgroup.

### TCP\_RRs Results

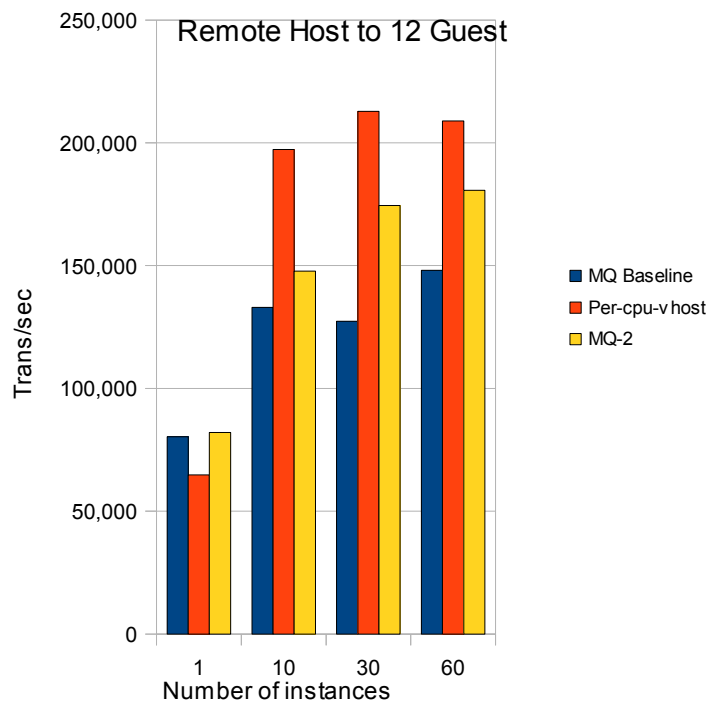


### TCP\_RRs Efficiency

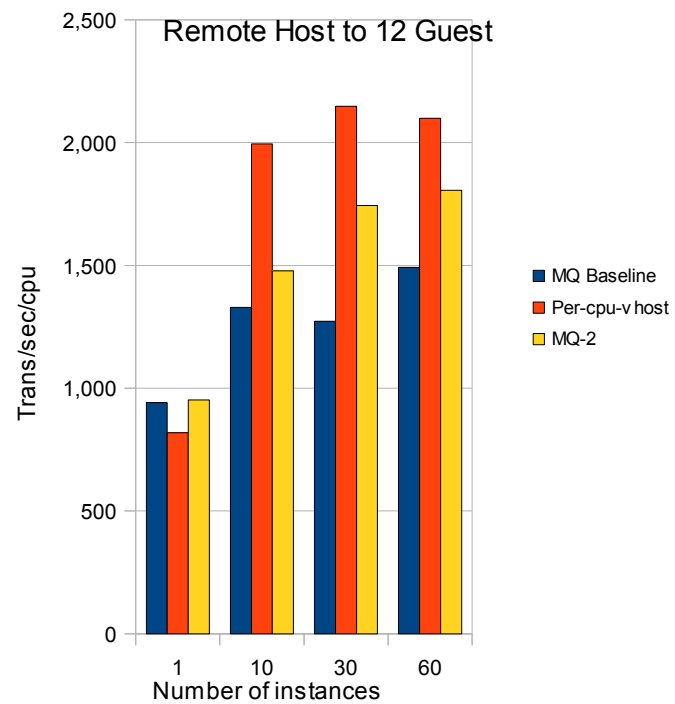




### TCP\_RRs Results

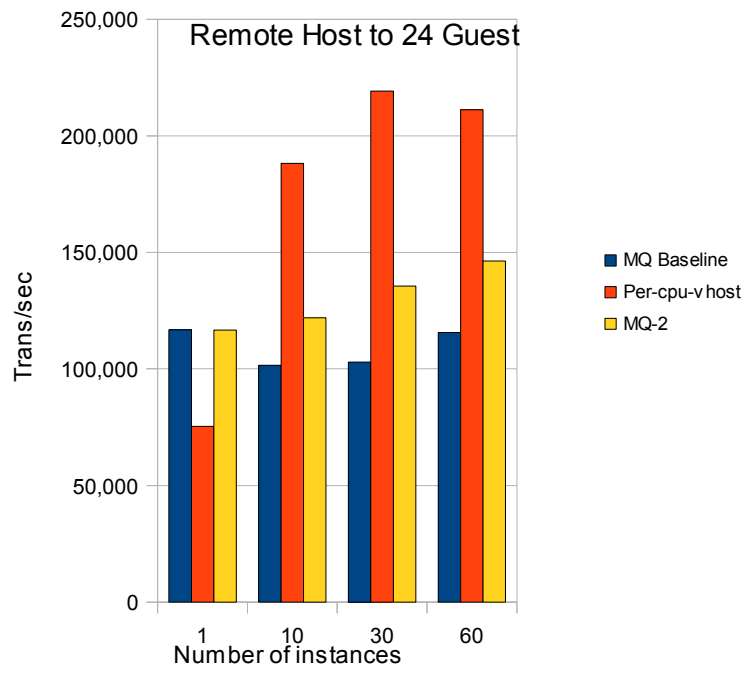


### TCP\_RRs Efficiency

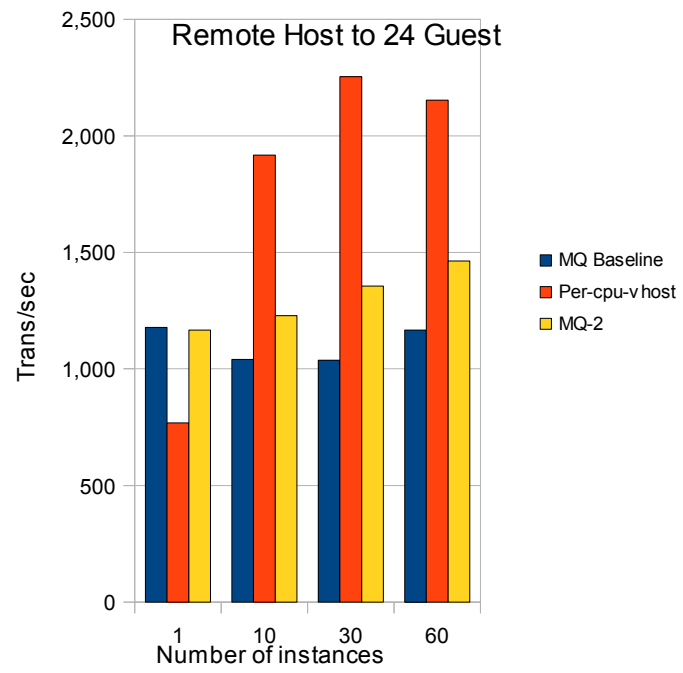




### TCP\_RRs Results



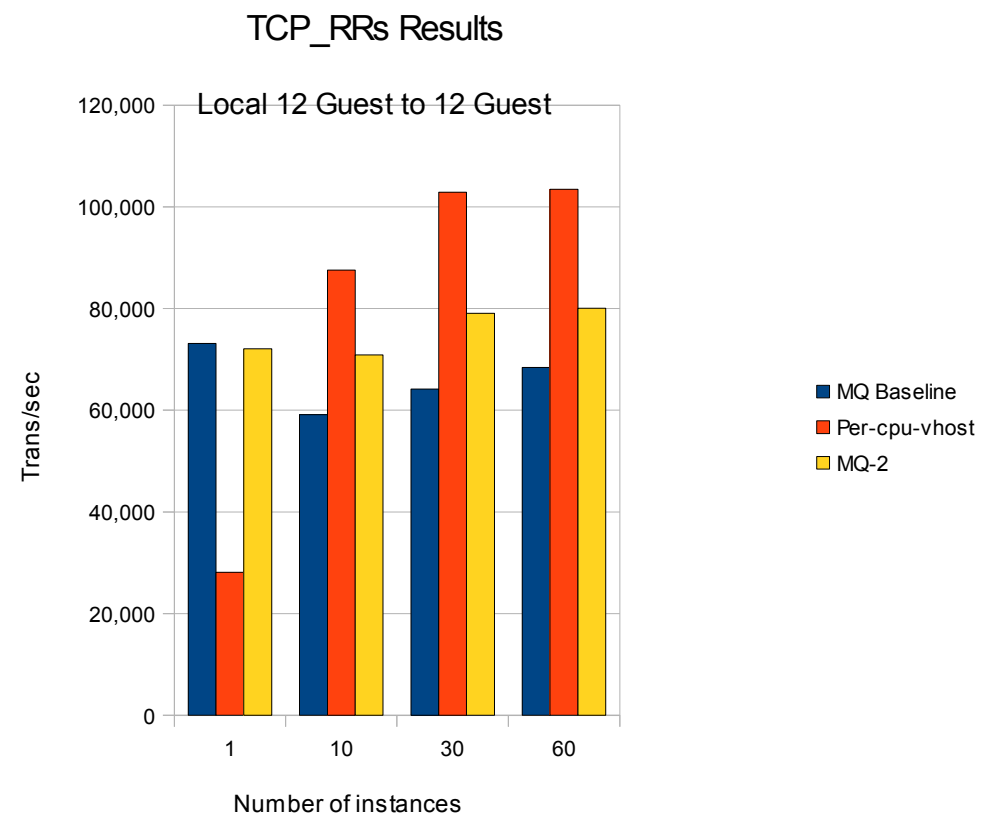
### TCP\_RRs Efficiency



- Per CPU Vhost thread shows win with many Vms

- Crossover occurs at about

- Similar trends for
  - TCP\_STREAM
  - UDP\_STREAM, UDP\_
  - Local (east-west) traffic





Thanks to Tom Lendacky for perf data!

Questions/Comments?

<http://github.com/jrfastab/>

- IBM x3650-M2
  - Intel E5530 2.4 Ghz Nehalem processors
  - dual socket with 4 cores/socket
  - Hyperthreading disabled
  - NIC X540-SR1 10GbE (ixgbe)