

# Empty Promise: Zero-copy Receive for vhost

Kalman Meth, Mike Rapoport, Joel Nider

{meth,joeln}@il.ibm.com

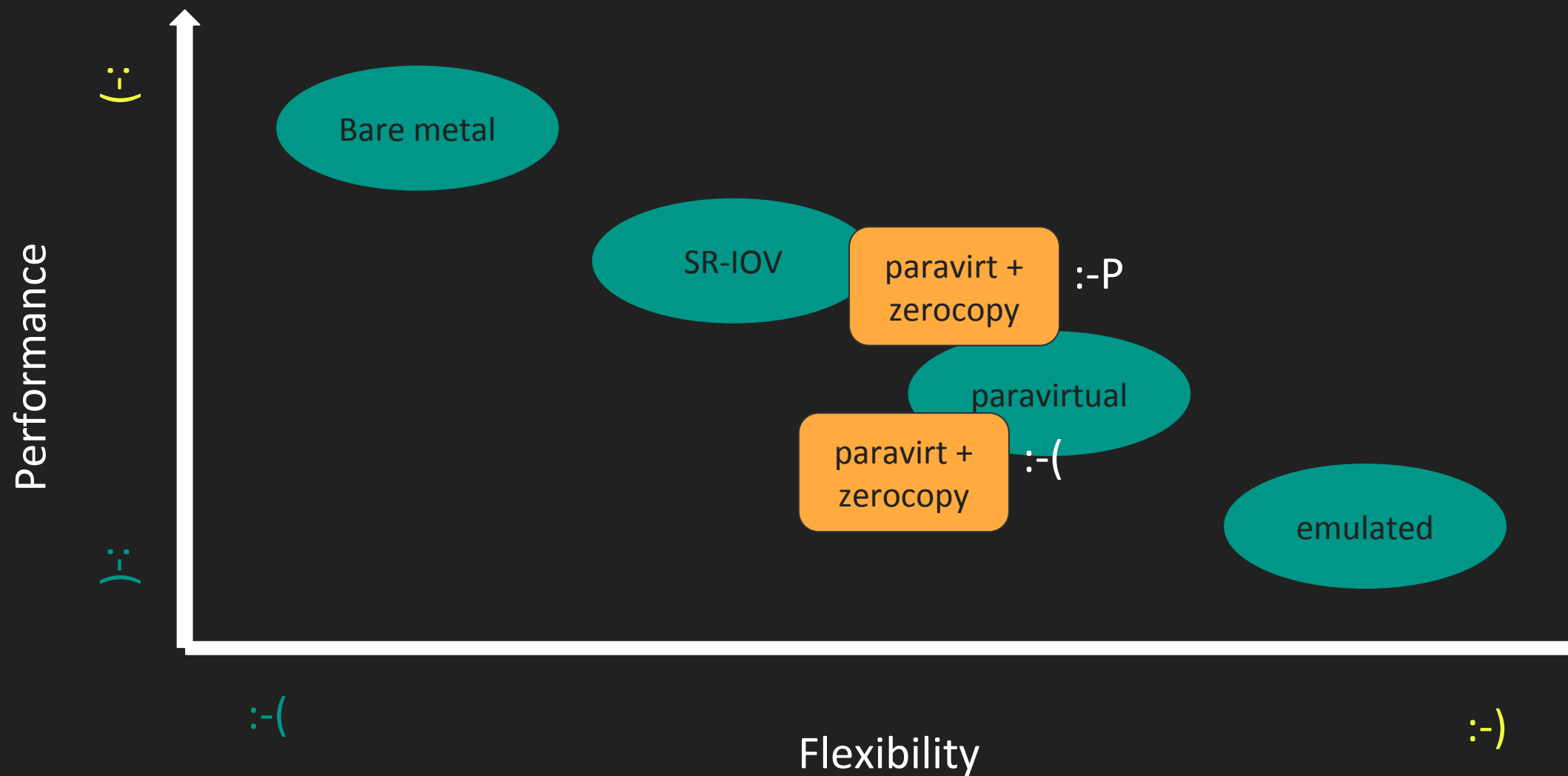
rppt@linux.vnet.ibm.com



This project has received funding from the *European Union's Horizon 2020 research and innovation programme* under grant agreements No 645402 and No 688386.



# Virtualization and IO



- No copy is better than copy
- Zerocopy TX without RX should feel lonely
- It was 8 years since the last attempt. Can we do better?

# More motivation



```
root@ubuntu-vm1: /home/perf
File Edit View Search Terminal Help
Samples: 99K of event 'cycles', Event count (approx.): 27191215670
Children      Self Command      Shared Object      Symbol
+ 100.00%     0.00% vhost-5119    [kernel.kallsyms]  [k] kthread
+ 100.00%     0.00% vhost-5119    [kernel.kallsyms]  [k] ret_from_fork
+ 100.00%     0.00% vhost-5119    [unknown]          [k] 0000000000000000
+ 99.91%      0.36% vhost-5119    [kernel.kallsyms]  [k] vhost_worker
+ 89.08%      0.57% vhost-5119    [kernel.kallsyms]  [k] handle_rx
+ 85.04%      0.01% vhost-5119    [kernel.kallsyms]  [k] handle_rx_net
+ 78.42%      0.05% vhost-5119    [kernel.kallsyms]  [k] macvtap_recvmmsg
+ 78.25%      0.60% vhost-5119    [kernel.kallsyms]  [k] macvtap_do_read
+ 71.53%      2.13% vhost-5119    [kernel.kallsyms]  [k] skb_copy_datagram_iter
+ 46.89%      46.89% vhost-5119    [kernel.kallsyms]  [k] copy_user_generic_strin
+ 12.86%      0.04% vhost-5119    [kernel.kallsyms]  [k] __softirqentry_text_sta
+ 12.80%      0.02% vhost-5119    [kernel.kallsyms]  [k] ret_from_intr
+ 12.77%      0.03% vhost-5119    [kernel.kallsyms]  [k] do_IRQ
+ 12.75%      0.04% vhost-5119    [kernel.kallsyms]  [k] net_rx_action
+ 12.69%      0.02% vhost-5119    [kernel.kallsyms]  [k] irq_exit
+ 12.68%      0.12% vhost-5119    [kernel.kallsyms]  [k] ixgbe_poll
+ 11.41%      1.98% vhost-5119    [kernel.kallsyms]  [k] ixgbe_clean_rx_irq
+ 8.64%       8.53% vhost-5119    [kernel.kallsyms]  [k] copy_page_to_iter
+ 6.80%       0.11% vhost-5119    [kernel.kallsyms]  [k] __kfree_skb
+ 5.67%       0.01% vhost-5119    [kernel.kallsyms]  [k] handle_tx_kick
For a higher level overview, try: perf report --sort comm,dso
```

## Transmit

- Downstream routing is easy
- Memory is always at hand

## Receive

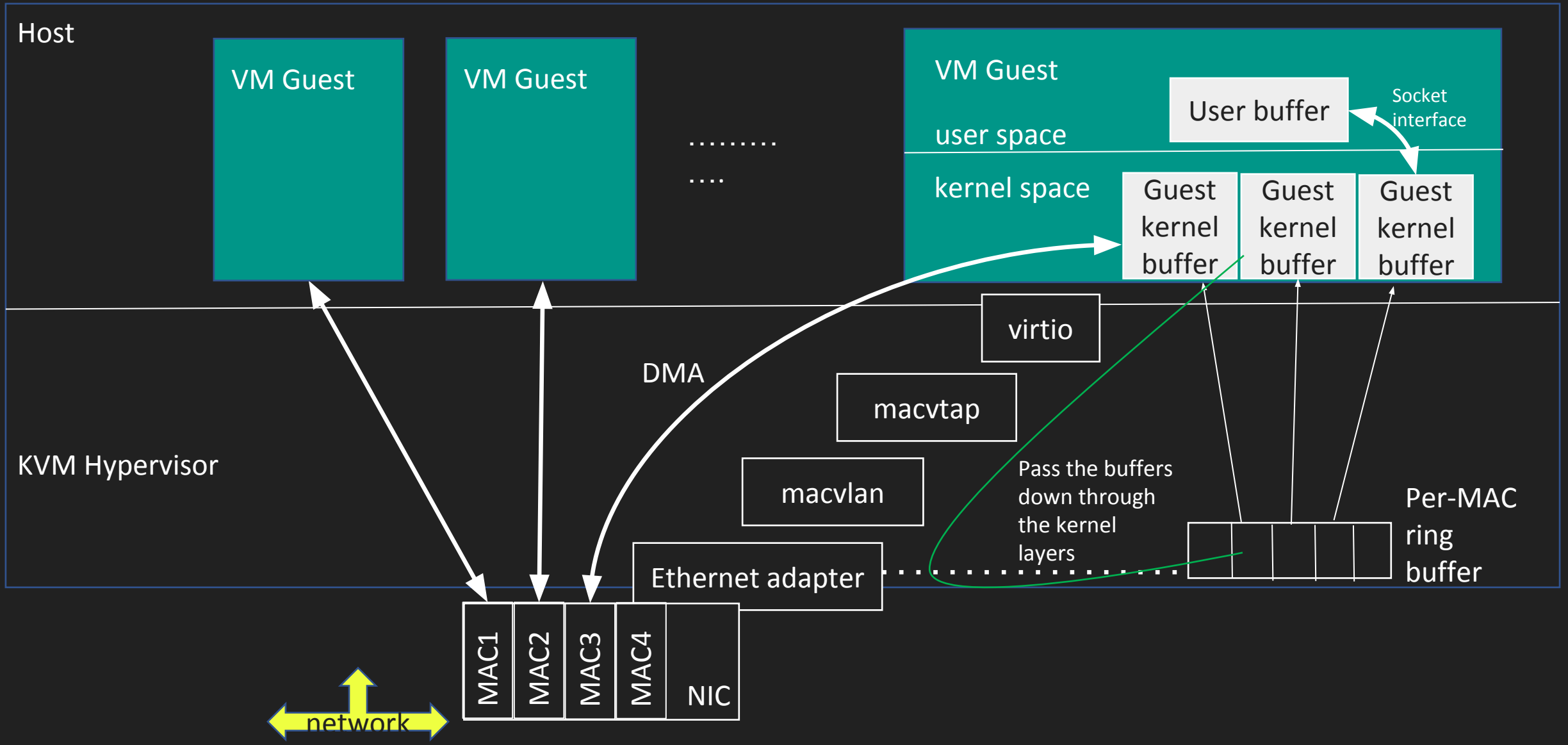
- Destination is not yet known
- Need memory for DMA
- Does not exist yet

# Assumptions



- Modern NICs are multiqueue
  - Dedicate queues to virtual NIC
- Guest allocates the buffers
  - Remapping DMA region to guest is more complex
- Tight coupling between physical and virtual NICs
  - Restrict zerocopy-RX to macvtap
- Minimal changes to guest

# Zero-Copy Rx Architecture







## macvtap

- `MSG_ZCOPY_RX_POST`
  - Control message from vhost-net to macvtap to propagate the buffers from guest to the lower levels
- `MSG_ZCOPY_RX`
  - Flag indicating that message contains preallocated buffers that should not be copied to userspace

## virtio-net

- `add_recvbuf_full_page()`
  - Ethernet adapter driver expects page size aligned buffers
  - Existing `add_recvbuf_*`() do not care since the data was always copied

- Isolate set of queues in physical NIC
- Create 1:1 correspondence between physical and virtual queues
- Clear RX descriptor ring
- Drop pre-allocated RX buffers in physical NIC driver

# Memory allocation



- virtio-net (guest)
  - Allocate buffers
    - DMA'able memory (`PAGE_SIZE` granularity and page aligned)
- vhost-net
  - Post buffers to macvtap
    - New control flag `MSG_ZCOPY_RX_POST` for `macvtap_recvmmsg()`
- macvtap
  - Allocate `skb`
  - Map `iovec` to `skb` (similar to `zerocopy_sg_from_iter`)
  - Pass the buffers to physical NIC
    - New method `ndo_post_rx_buffer()`
- Physical NIC driver adds new buffers to RX descriptor ring

- Physical NIC driver
  - DMA directly to the guest buffers
  - Setup `skb` structure
  - `netif_rx()` and friends
- `macvtap`
  - Queue `skb` as ready for the userspace
  - Inform `vhost-net` about the `virtio` descriptor associated with the `skb`.

# Packet receive (cont)

- vhost-net
  - `handle_rx_zero_copy()`:
    - Update virtqueue
    - Kick macvtap with `->recvmsg(MSG_ZCOPY_RX)`
- macvtap (again)
  - `macvtap_do_read_zero_copy()`
    - `skb_array_consume`
    - Cleanup

# Implementation status



- Initial implementation
  - still sub-optimal
- Stable enough to benchmark
- Source:
  - <https://github.com/mikelangelo-project/linux-zecorx>
  - <https://github.com/mikelangelo-project/qemu-zecorx>

# Test setup



- 2x IBM System x3550 M4 Server
  - Intel® Xeon® Processor E5-2660
    - 8 cores, 16 threads
  - 56G RAM
  - Intel 82599ES 10-Gigabit Network Connection (ixgbe)
- Back to back connection for host NICs
- VM with 4 vCPUs, 2G RAM
- Linux v4.8
- netperf 2.6.0



```

rppt@marva: ~/data/perf
rapopor... x rppt@m... x rppt@m... x rapopor... x rapopor... x rapopor... x
# Event 'cycles:pp'
#
# Baseline      Delta      Shared Object      Symbol
# .....      .....      .....
#
  49.86%    -46.79%    [kernel.vmlinux]   [k] copy_user_generic_string
  12.35%
  5.57%     -5.14%    [kernel.vmlinux]   [k] copy_page_to_iter
  5.21%     +0.01%    [vhost]             [k] skb_release_data
  3.40%
  2.77%
  1.60%     +1.52%    [vhost]             [k] vhost_get_vq_desc
  1.60%
  1.30%     -1.03%    [kvm]               [k] skb_copy_datagram_iter
  1.25%     -1.09%    [kernel.vmlinux]   [k] free_pcpages_bulk
  1.20%     -0.15%    [kernel.vmlinux]   [k] vhost_signal
  0.88%
  0.82%     -0.35%    [kernel.vmlinux]   [k] free_hot_cold_page
  0.68%     -0.55%    [kvm]               [k] __apic_accept_irq
  0.68%     -0.68%    [vhost_net]        [k] skb_set_owner_w
  0.51%     +0.32%    [kernel.vmlinux]   [k] _raw_spin_lock
  0.46%     +0.77%    [vhost]             [k] macvtap_do_read
                        [k] __free_page_frag
                        [k] kvm_irq_delivery_to_apic_fast
                        [k] handle_rx
                        [k] __slab_free
                        [k] translate_desc
:

```

```
traffic-gen$ netperf -H guest
```

	<b>v4.8</b>	<b>v4.8 + ZCRX</b>
Throughput (Mbit/sec)	9255.22	4396.34
System Utilization (%)	10.9095	12.2941
CPU usage - vhost (%)	74.676	95.979
CPU usage - qemu (%)	99.6685	100.83

- RX interrupt, and vhost are on the same CPU
- Retry with forced CPU affinity

```
vm-host$ taskset -cp 4 $(pgrep vhost)
vm-host$ echo 5 > /proc/irq/111/smp_affinity_list
traffic-gen$ netperf -H guest
```

	<b>v4.8</b>	<b>v4.8 + ZCRX</b>
Throughput (Mbit/sec)	9255.22	8492.17
System Utilization (%)	10.9095	12.121
CPU usage - vhost (%)	74.676	86.1722
CPU usage - qemu (%)	99.6685	100.103



# What happened



- New bottleneck is in DMA mapping
- Latency has grown
  - Measured with `rdtsc()` in `vhost-net::handle_rx()`
    - Copy:  $\approx 2$  tscs/byte
    - No copy:  $\approx 3.5$  tscs/byte
- Page recycling in `ixgbe` replaced with `ndo_post_buffers()`
  - Sequential instead of parallel
  - Move frequent `dma_map*()` / `dma_unmap()`

# Can we do better?



- Maybe
- tl;dr
  - Try to re-parallel DMA mapping and RX processing
  - Better batching for memory allocation and dma\_map/unmap
  - Major changes to virtio ring
  - AF\_XDP based virtio backend in user space

Thank you!