

Zero-copy Receive for vhost

Kalman Meth, Mike Rapoport, Joel Nider

{meth,joeln}@il.ibm.com

rppt@linux.vnet.ibm.com

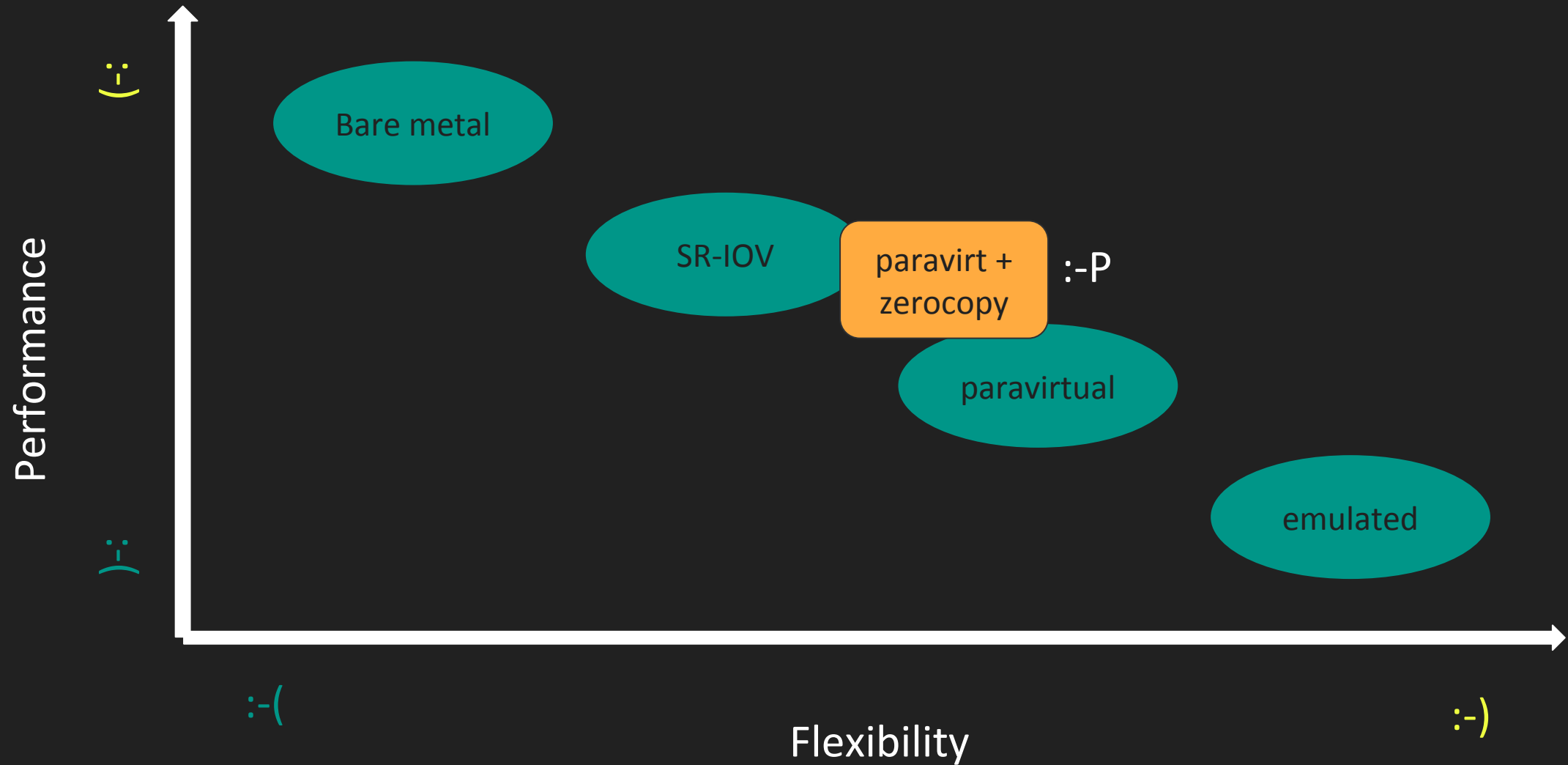


MIKELANGELO

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



Virtualization and IO



- No copy is better than copy
- Zerocopy TX without RX should feel lonely
- It was 7 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
```

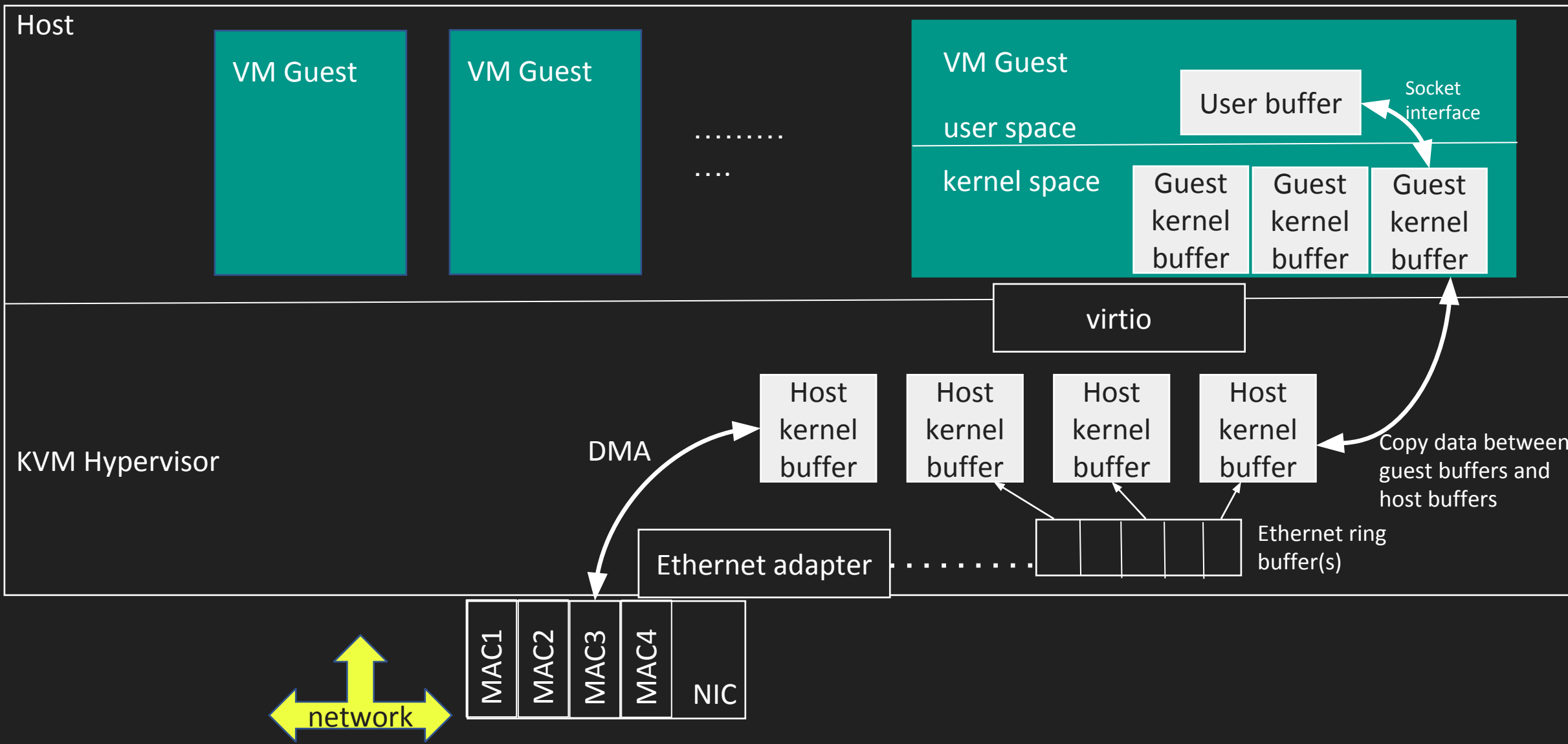
Transmit

- Downstream routing is easy
- Memory is always at hand
- ~~Still has problems~~
 - ~~head-of-line blocking~~

Receive

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

Virtio RX path

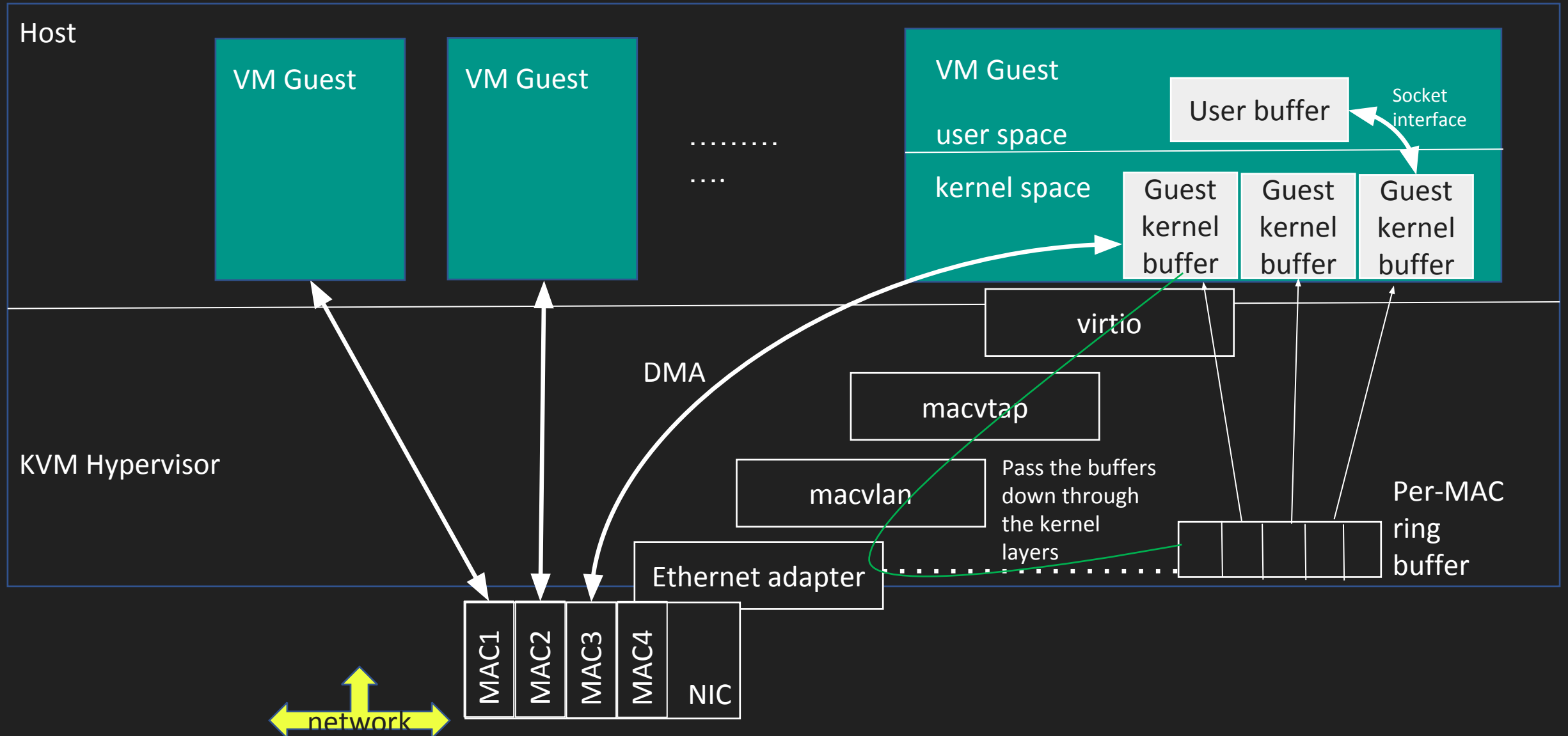


Assumptions



- 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

Zero-Copy Rx Architecture



- 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_recvmsg()`
- 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

netdev

- `->ndo_set_zerocopy_rx(struct net_device *pdev, struct net_device *vdev)`
 - Pass `vdev` down the stack to the ethernet adapter to bind physical and virtual queues.
 - Similar to `->ndo_dfwd_add_station()`, maybe just add flags there...
- `->ndo_post_rx_buffer(struct net_device *dev, struct sk_buff *skb)`
 - Passes a single (page aligned) buffer to the ethernet adapter
 - `skb` contains pointer to the upper level device and `ubuf_info`

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

- No unified mechanism for RX memory allocation
 - Each driver wraps `alloc_page()` differently
 - Page pool is a savior?
- What to do when running out of buffers
 - Drop?
 - Fallback to copy?
- virtio does not know how to deal with fragment offsets
 - `skb_copy_datagram_iter()` takes care of `frag->page_offset`
- Allocating page for Ethernet frame is wasteful
 - Even worse for architectures with `PAGE_SIZE > 4K`

Implementation status



- Progress is slower than expected
- Initial implementation for ixgbe, plans to add other NICs later on
- Zerocopy packets reach the guest :-)
- Not stable enough to run benchmarks :-)

Thank you!