

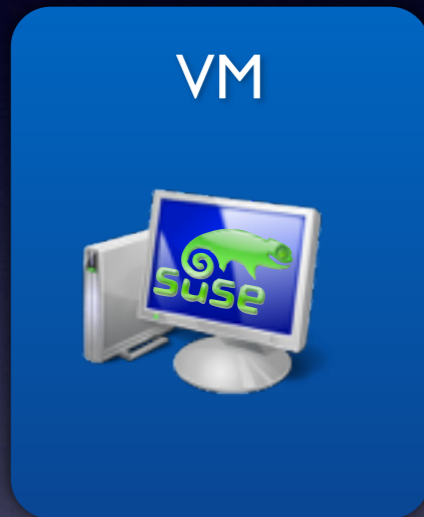
Migratable 40GBit/s Ethernet

About Me

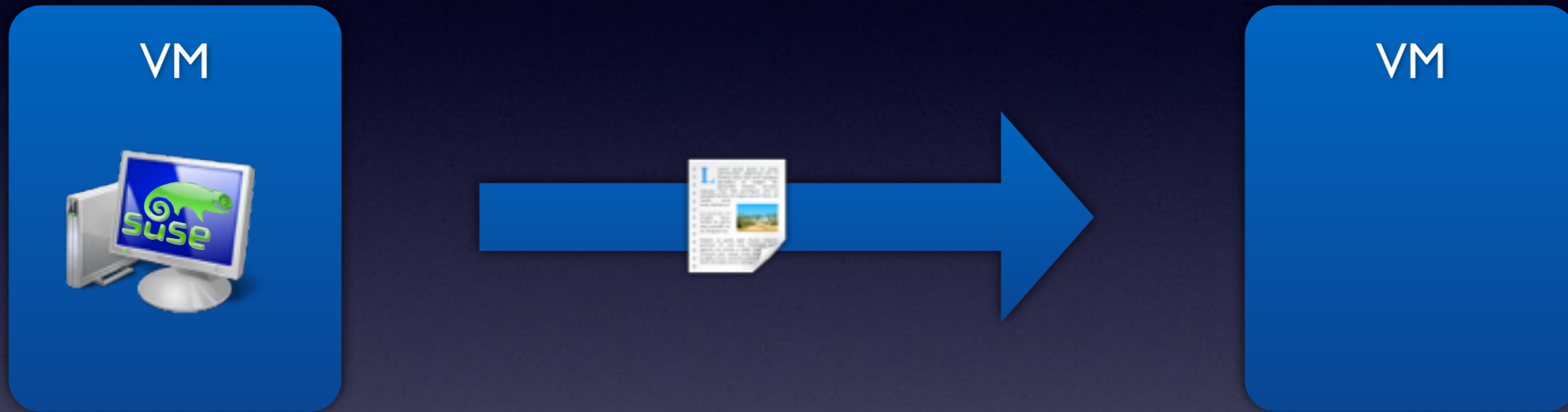
- Alexander Graf
- KVM and Qemu developer for SUSE
 - Server class PowerPC KVM port
 - S390x Qemu guest support
 - x86 Mac OS X in KVM
 - Nested SVM
 - ...

What is Live Migration

What is Live Migration



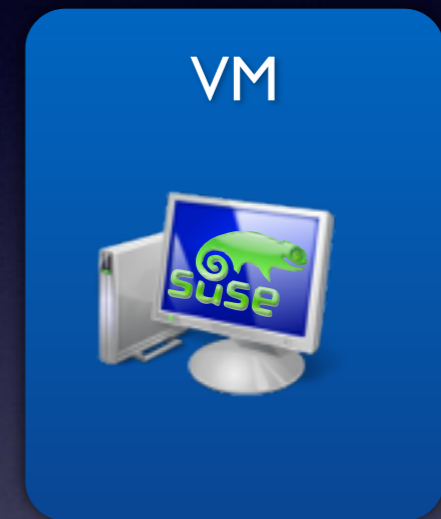
What is Live Migration



What is Live Migration



What is Live Migration

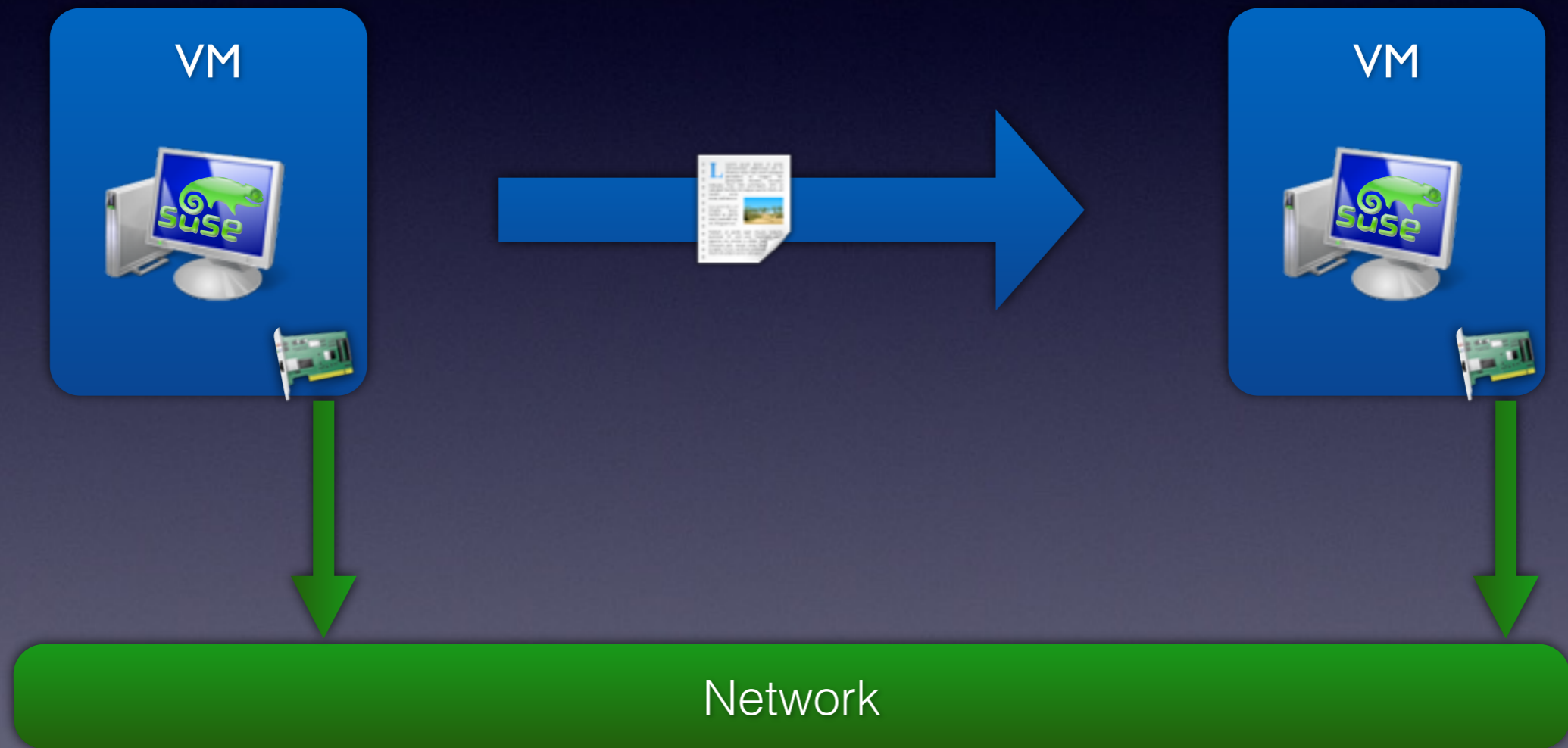


Network Live Migration

Network Live Migration

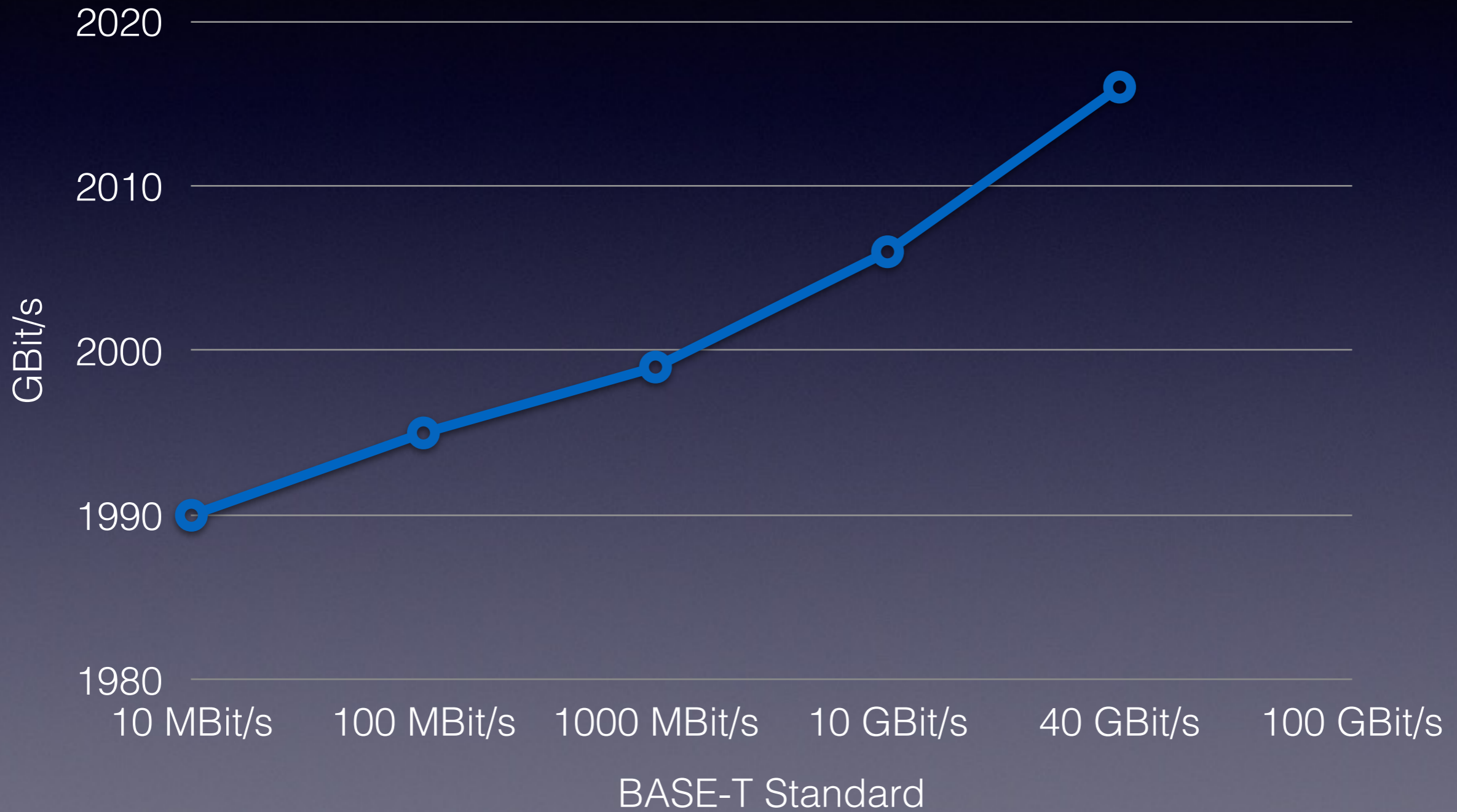


Network Live Migration

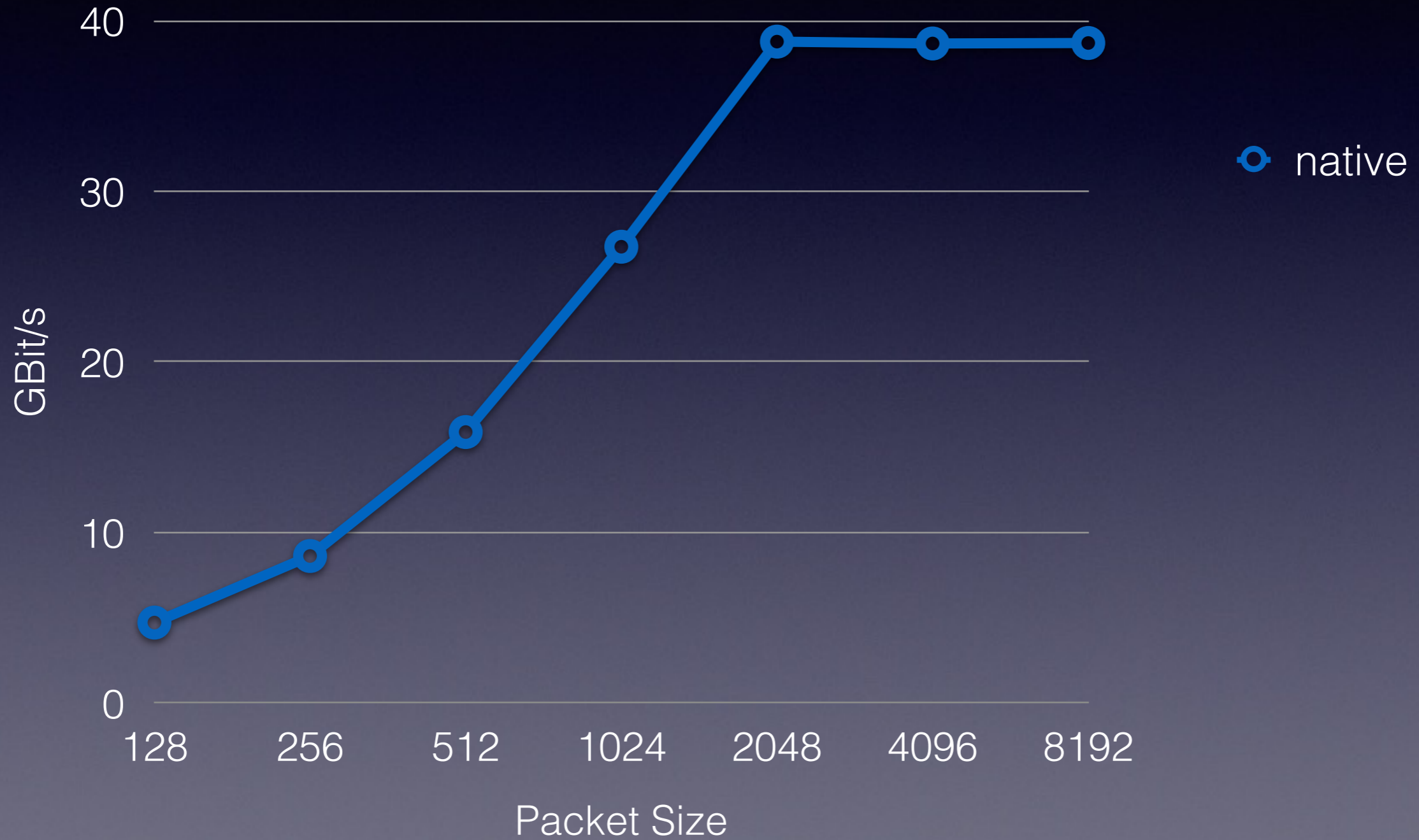


40 GBit/s

40 GBit/s

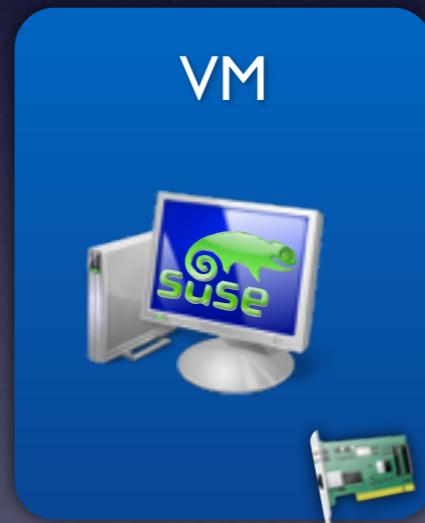


40 GBit/s

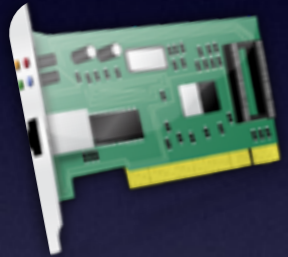


virtio

virtio

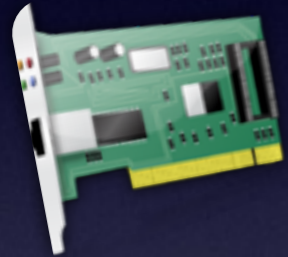


virtio



virtio-net

virtio



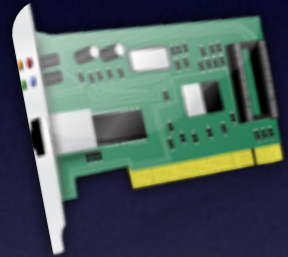
virtio-net

TX ring

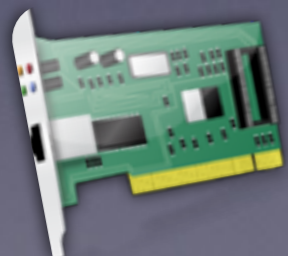
RX ring

Admin ring

virtio



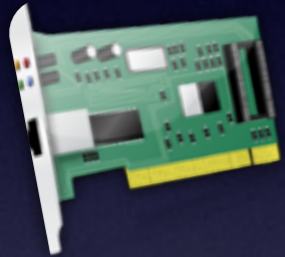
virtio-net



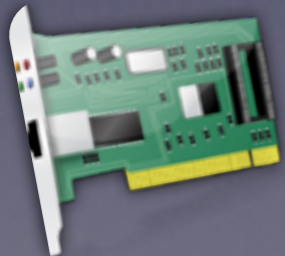
real NIC



virtio



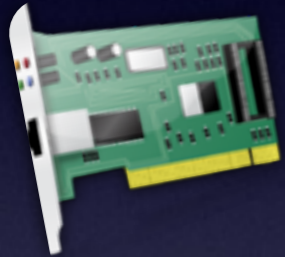
virtio-net



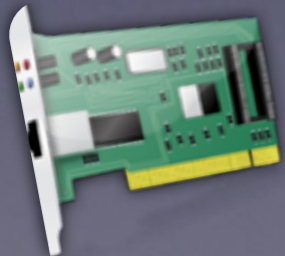
real NIC



virtio



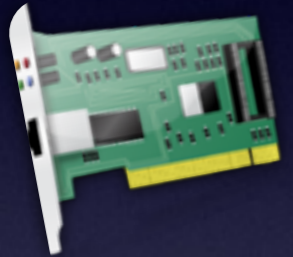
virtio-net



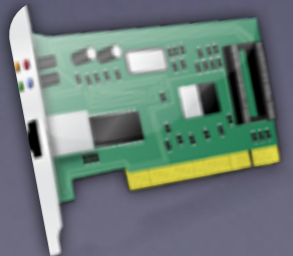
real NIC



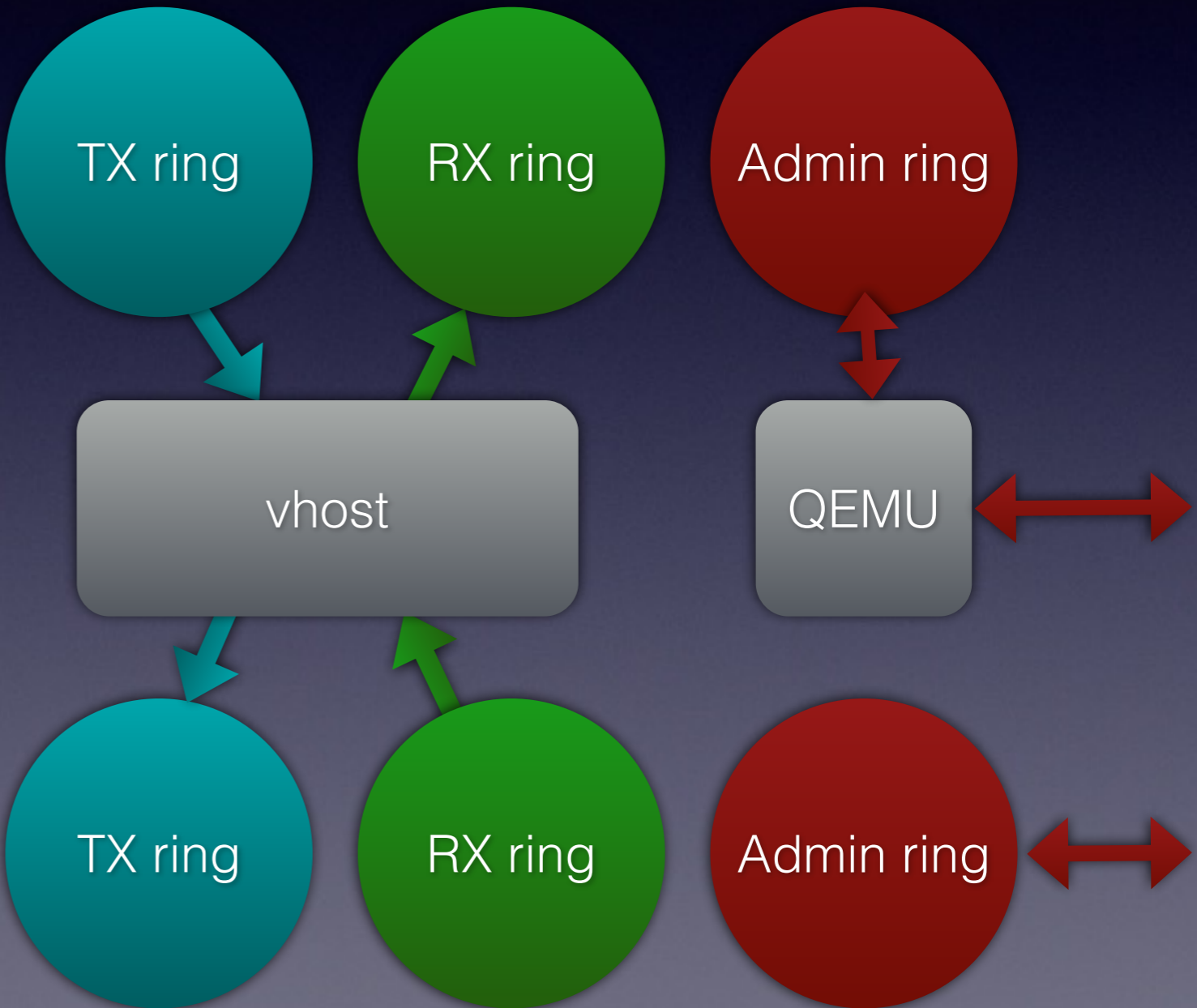
virtio



virtio-net

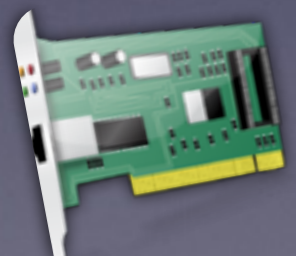
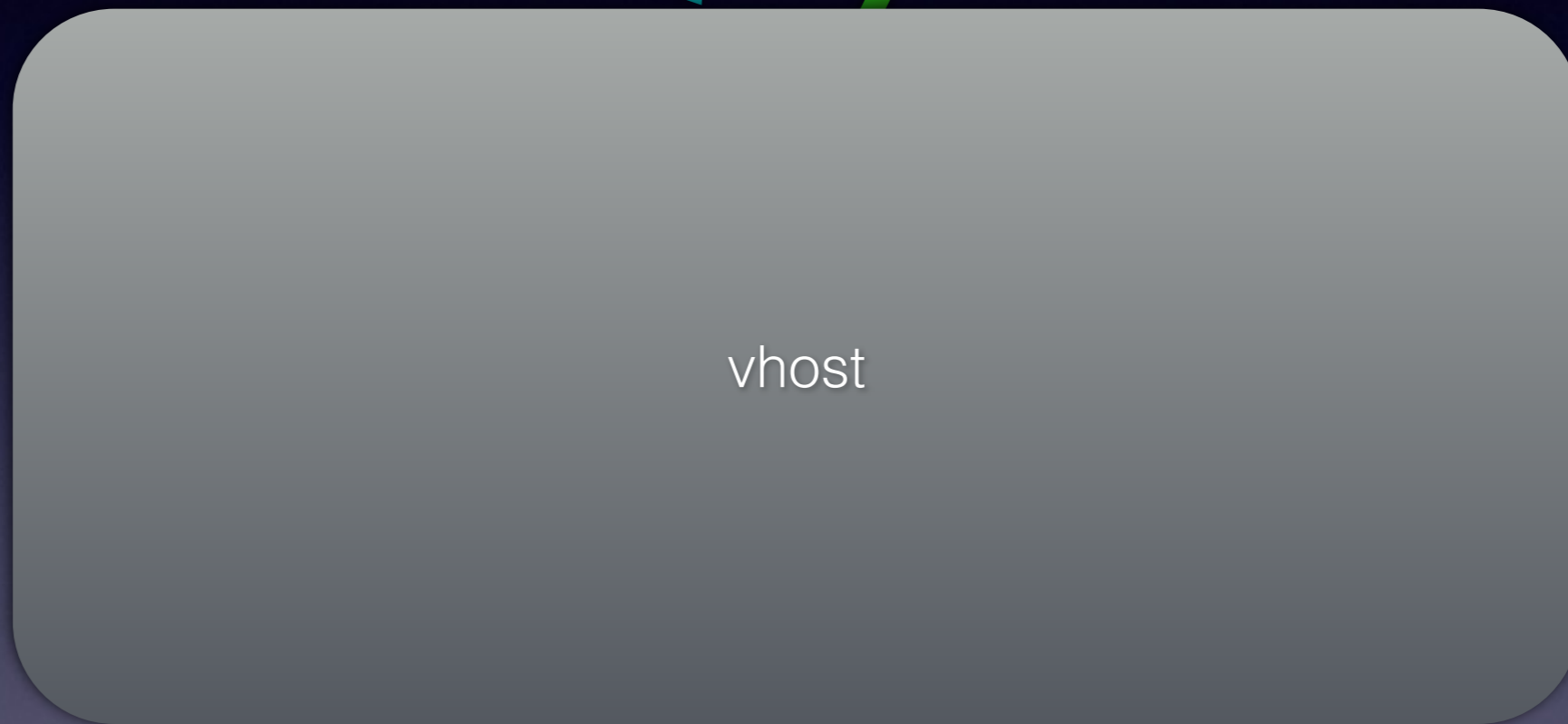


real NIC





virtio-net

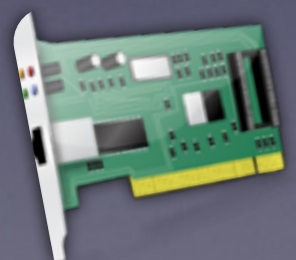
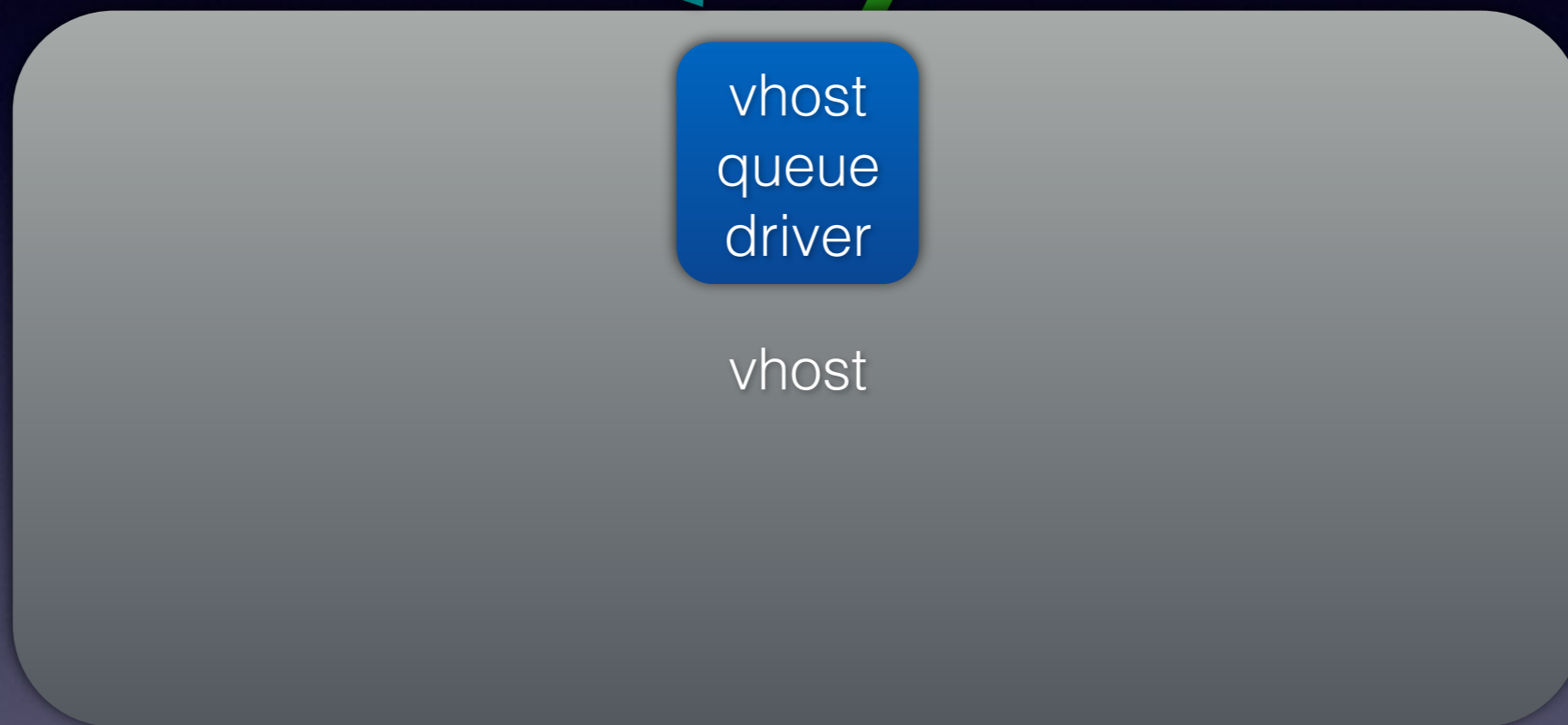
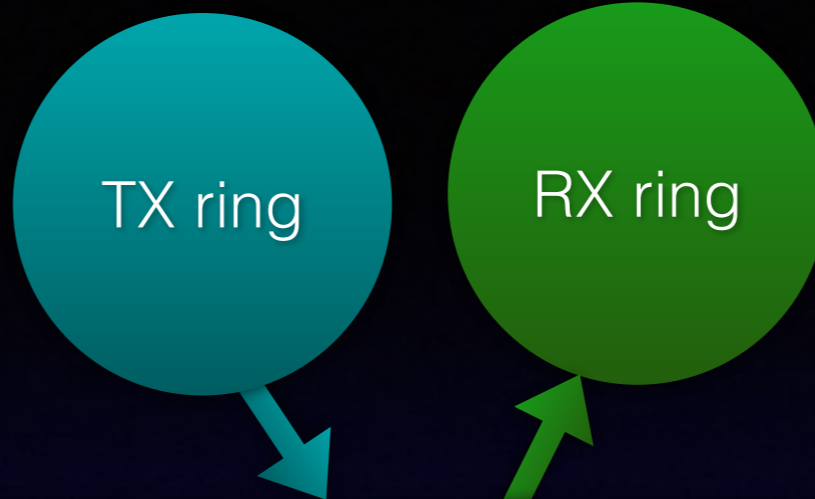


real NIC





virtio-net

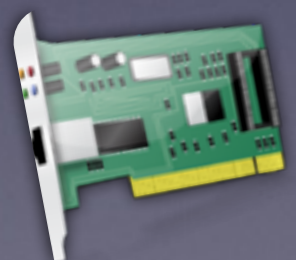
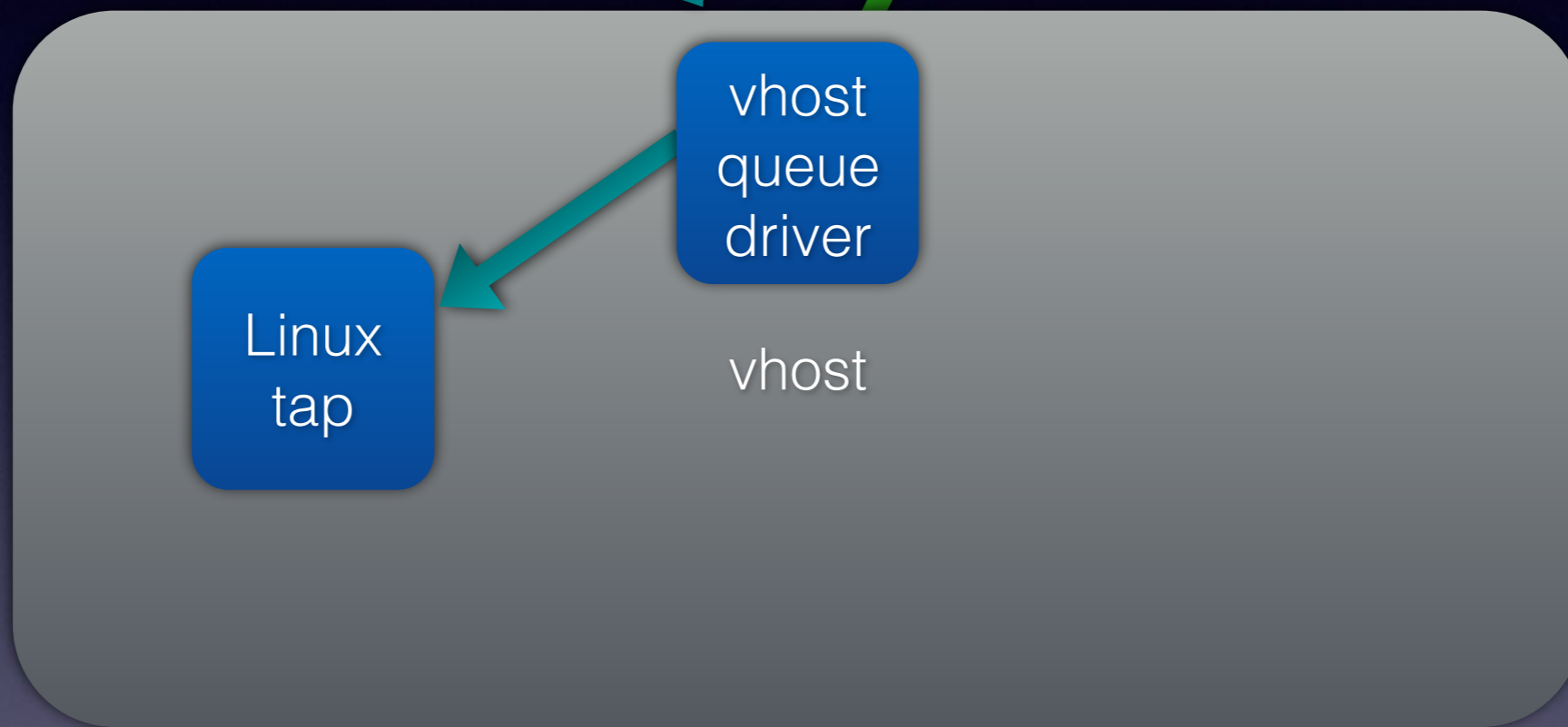
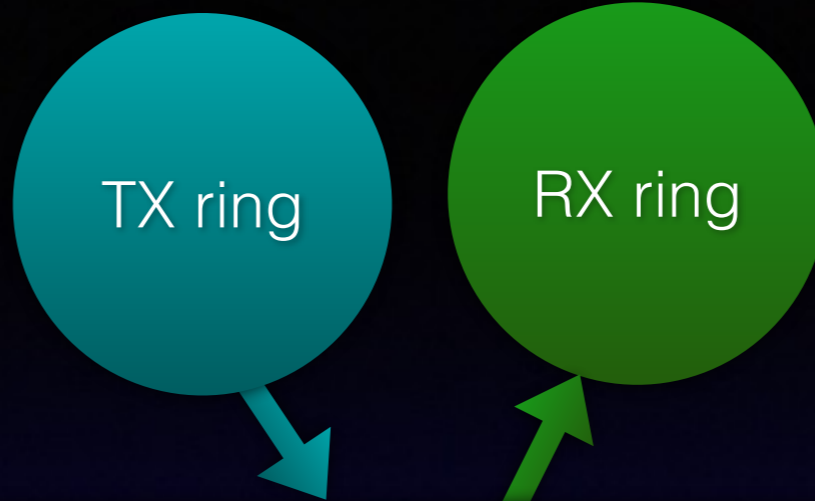


real NIC





virtio-net

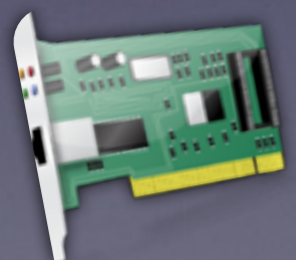
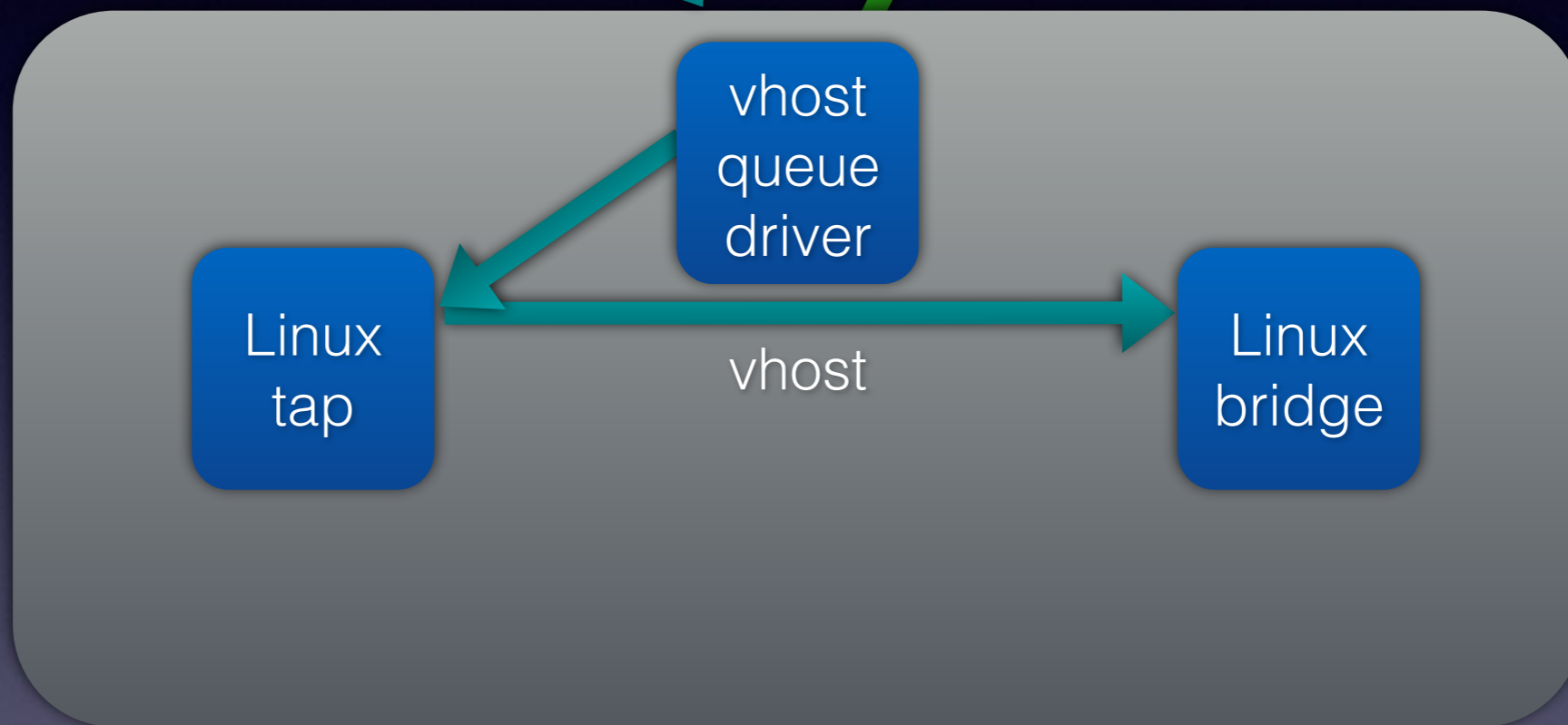
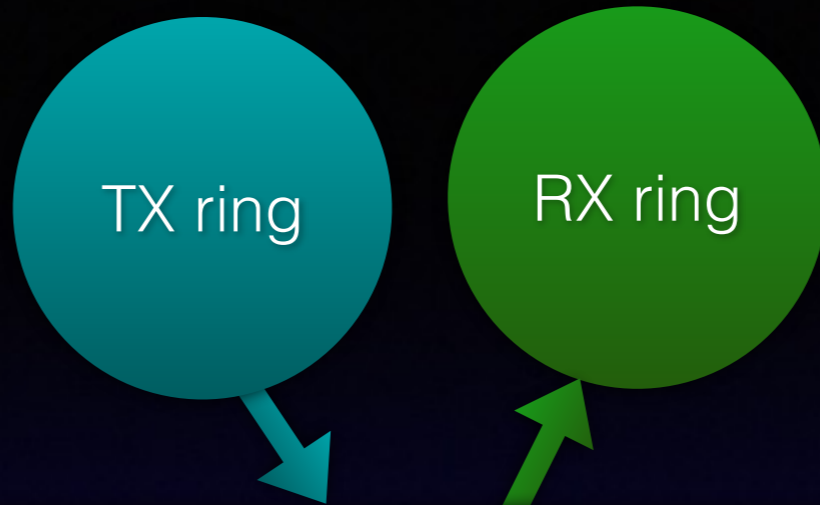


real NIC

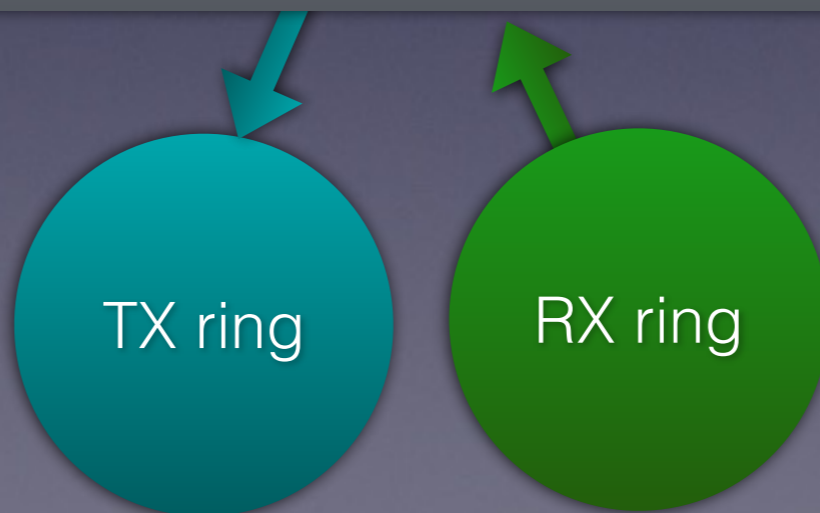




virtio-net

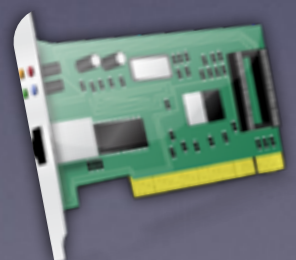
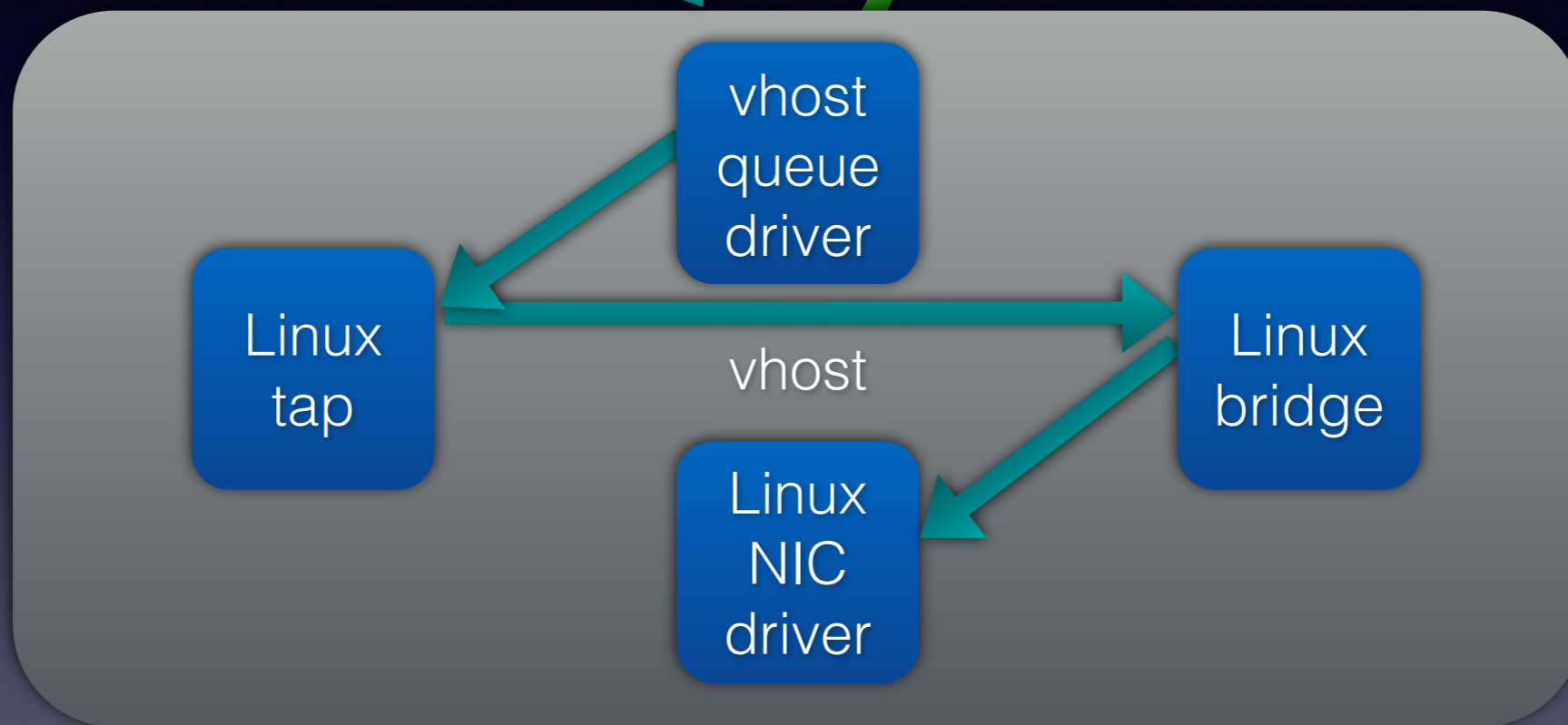
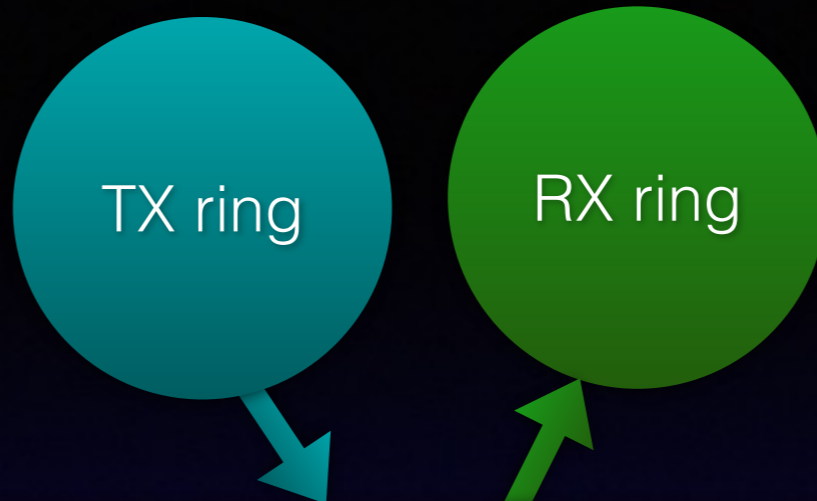


real NIC





virtio-net

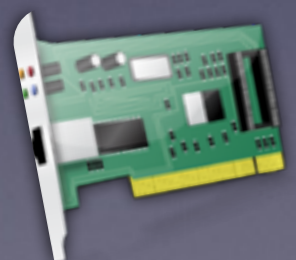
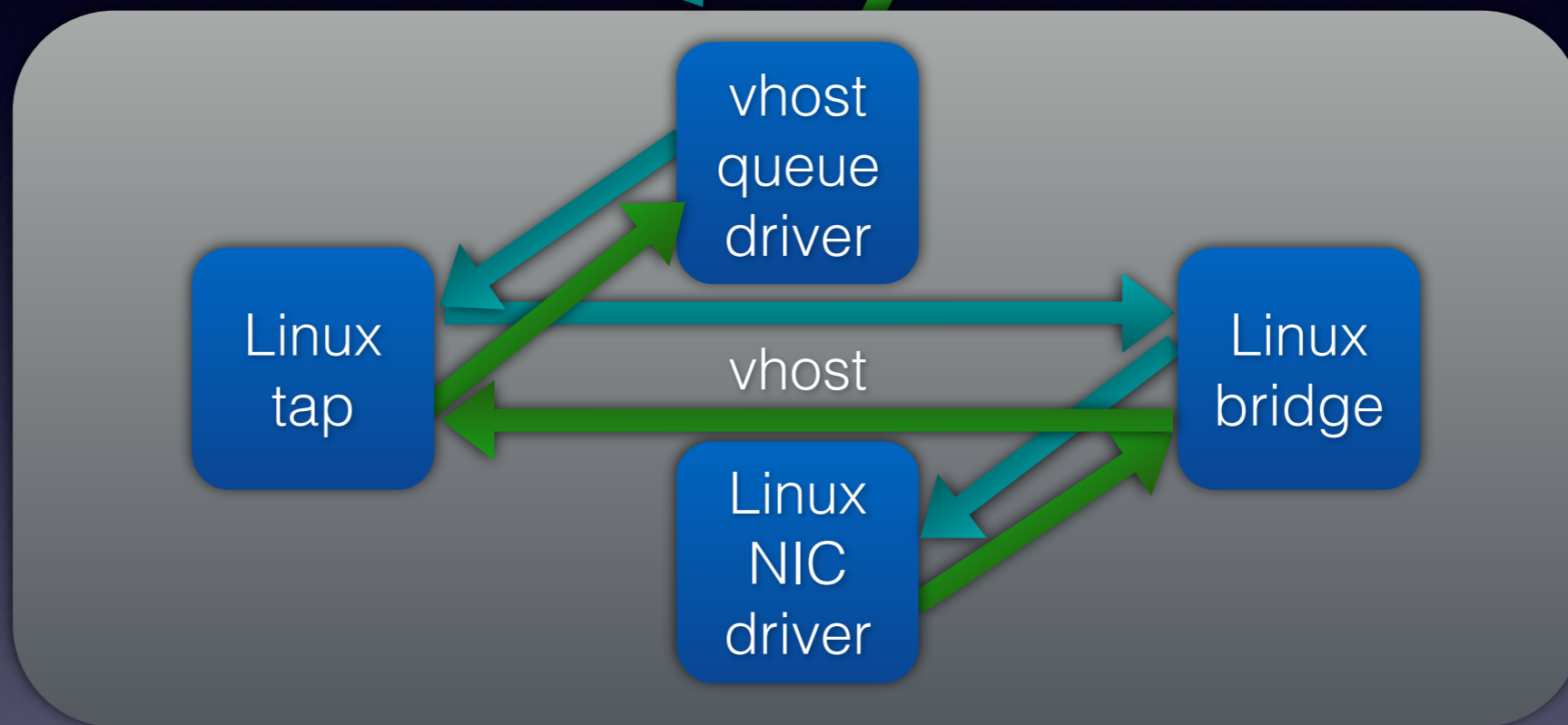
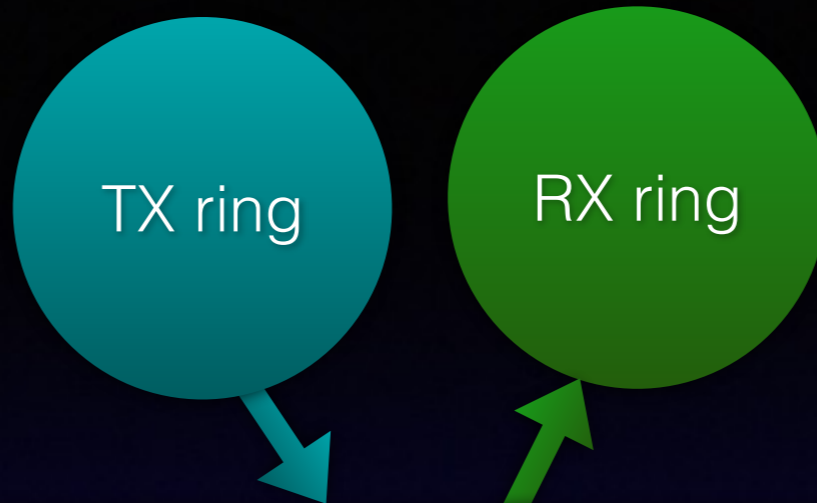


real NIC





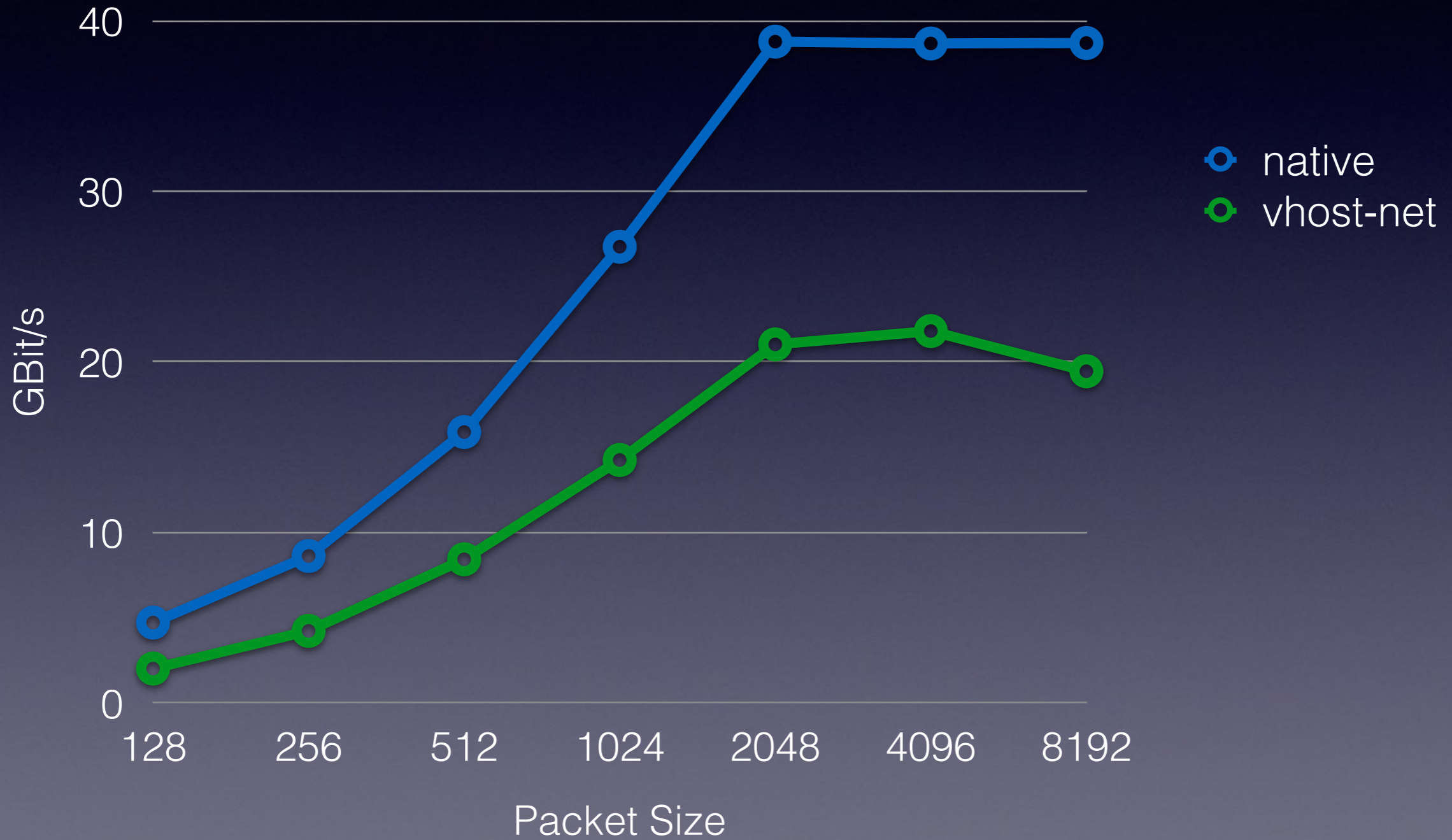
virtio-net



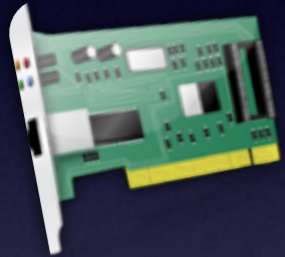
real NIC



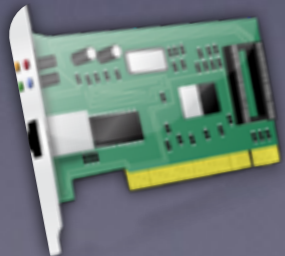
virtio



virtio



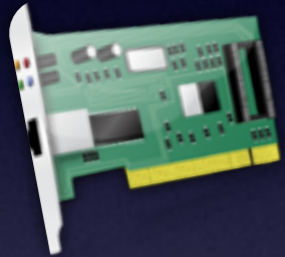
virtio-net



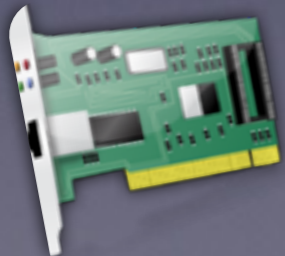
real NIC



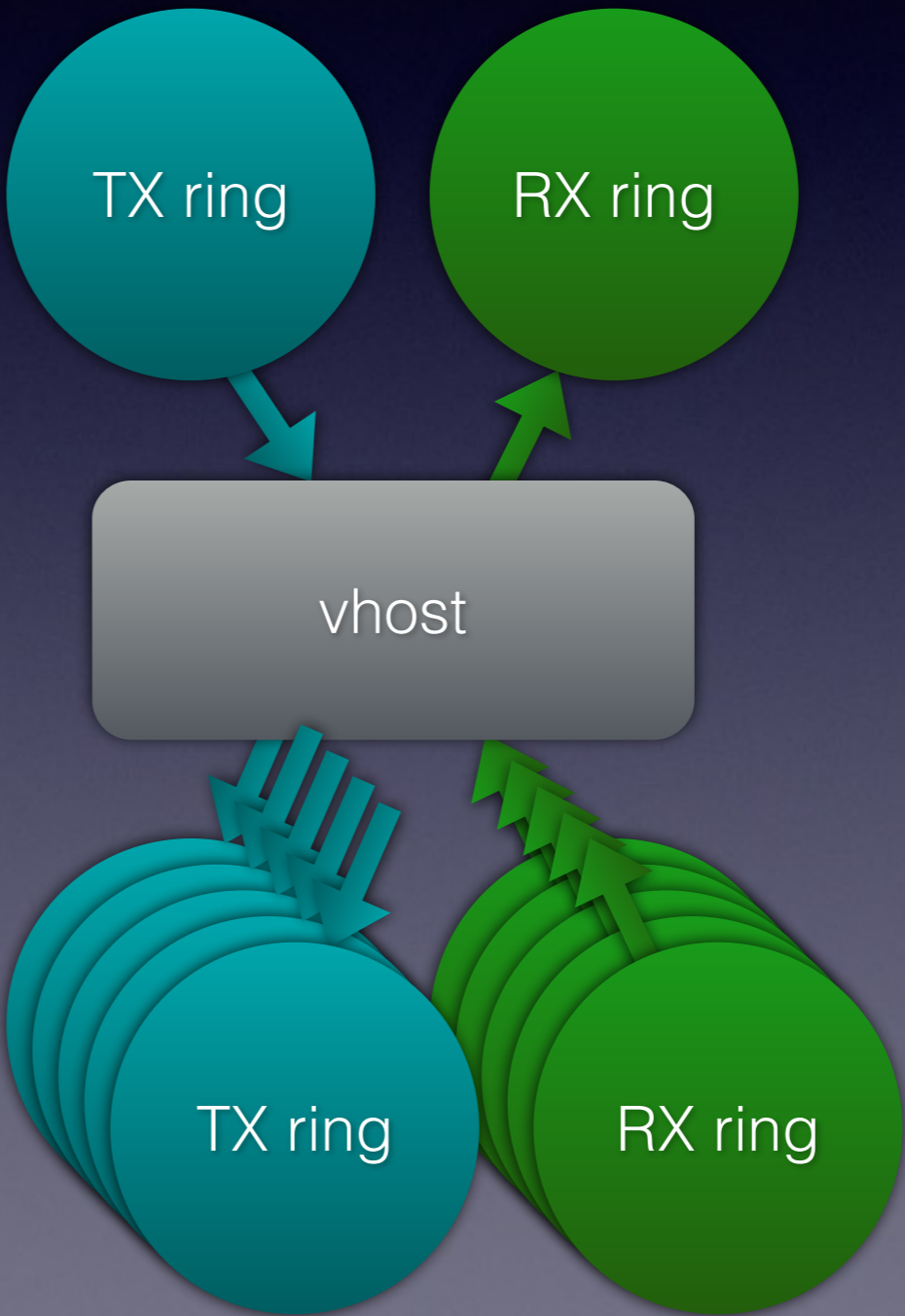
virtio



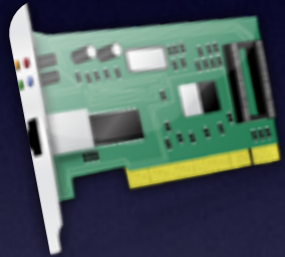
virtio-net



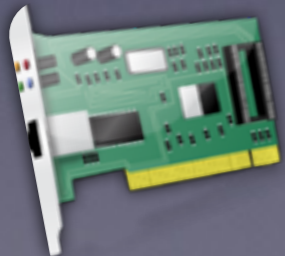
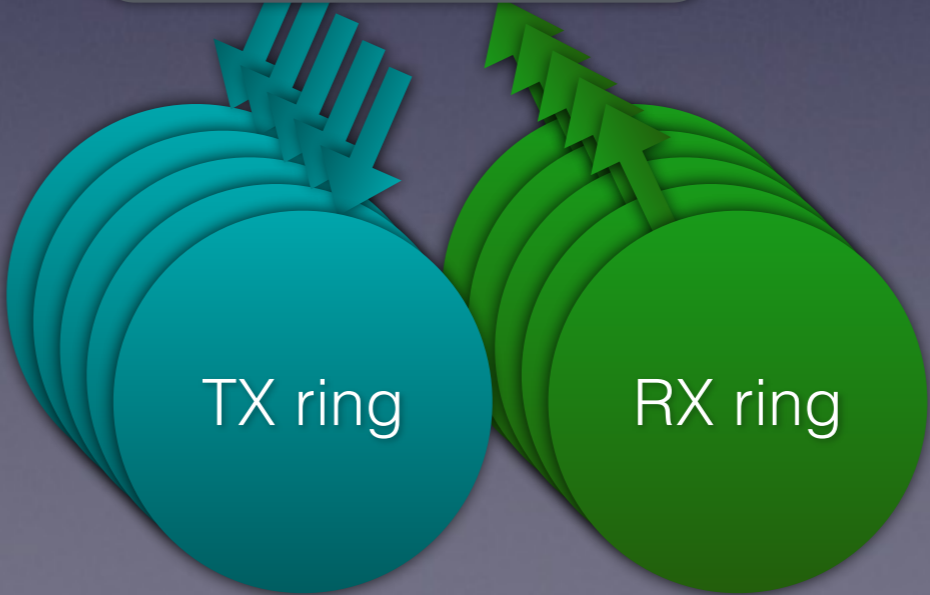
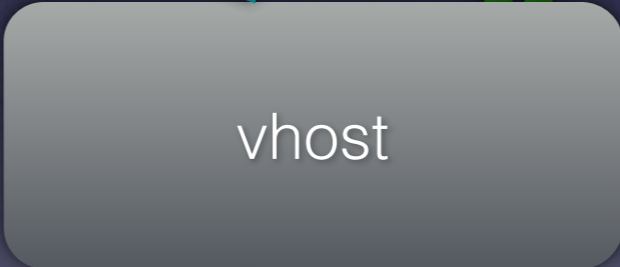
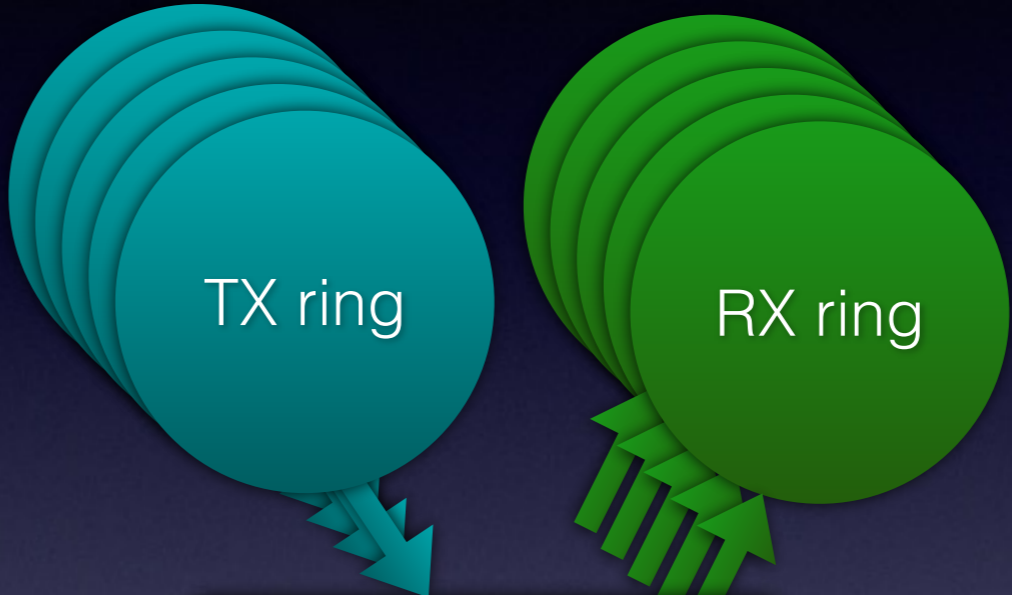
real NIC



virtio

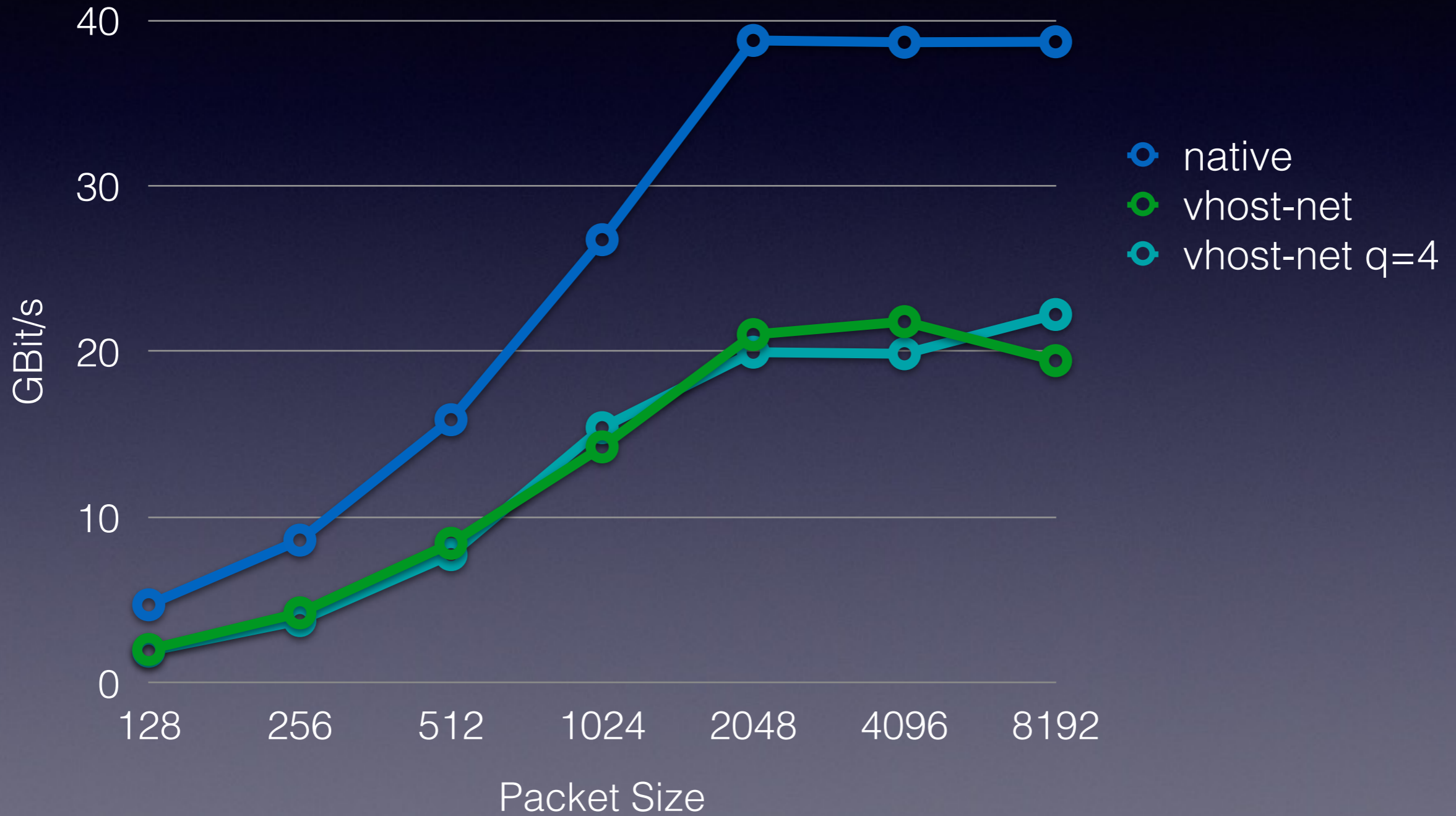


virtio-net



real NIC

virtio



virtio

- Bridge can be accelerated with flood=off learning=off
 - Requires VLAN
 - Shouldn't make a difference on direct connect

virtio

- Alternatives to bridge
 - Open vSwitch
 - rocker
 - macvtap
- Not benchmarked yet

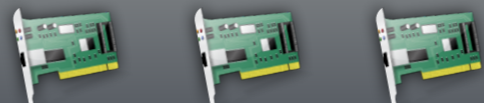
VFIO

VFIO

VM



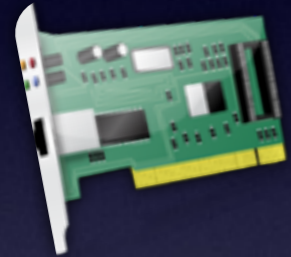
Host



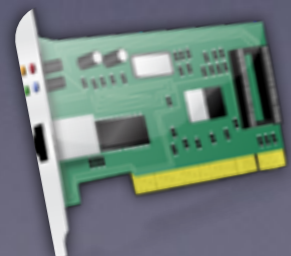
VFIO



VFIO



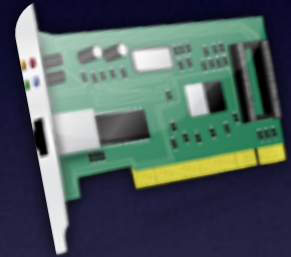
guest NIC



real NIC



VFIO



guest NIC



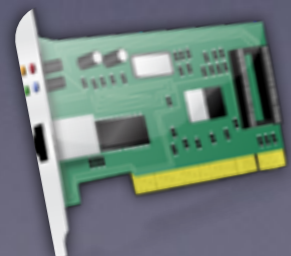
TX ring



RX ring

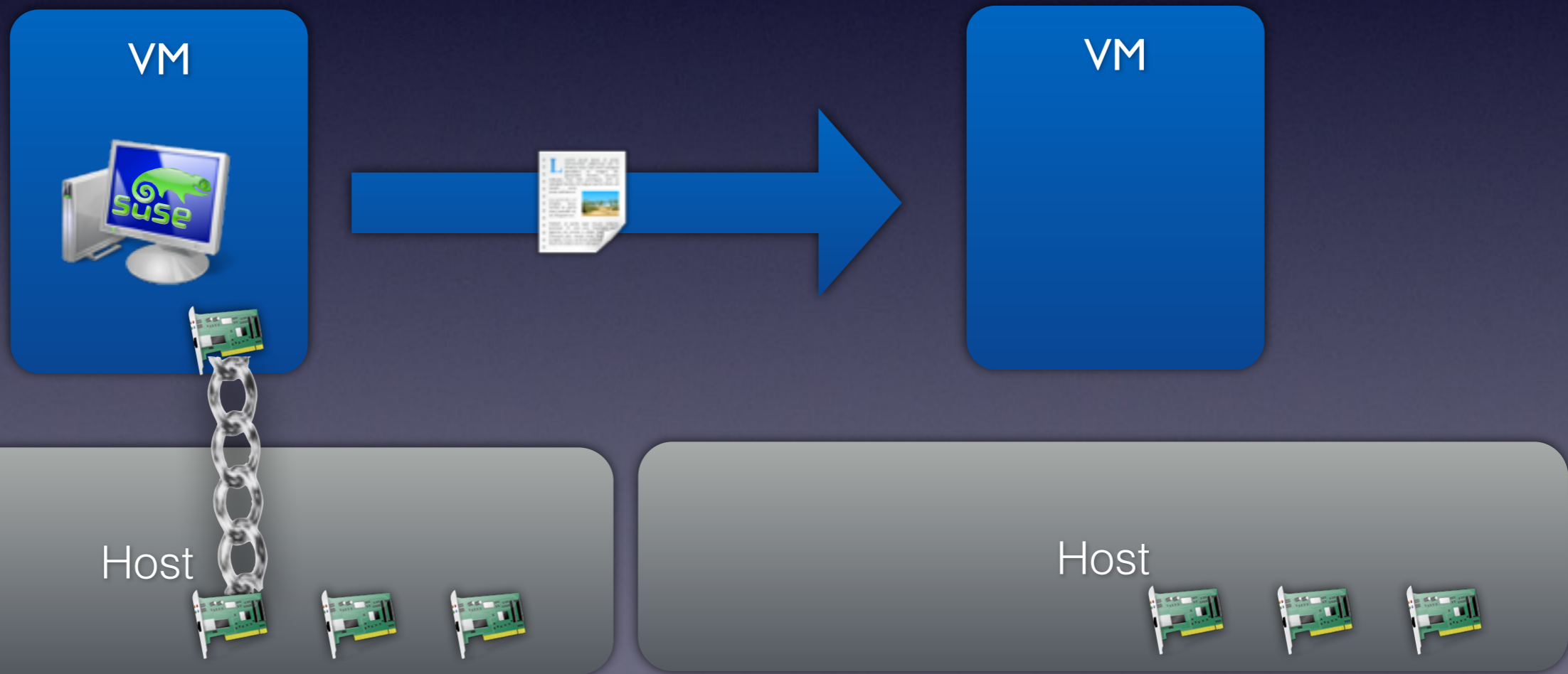


Admin ring

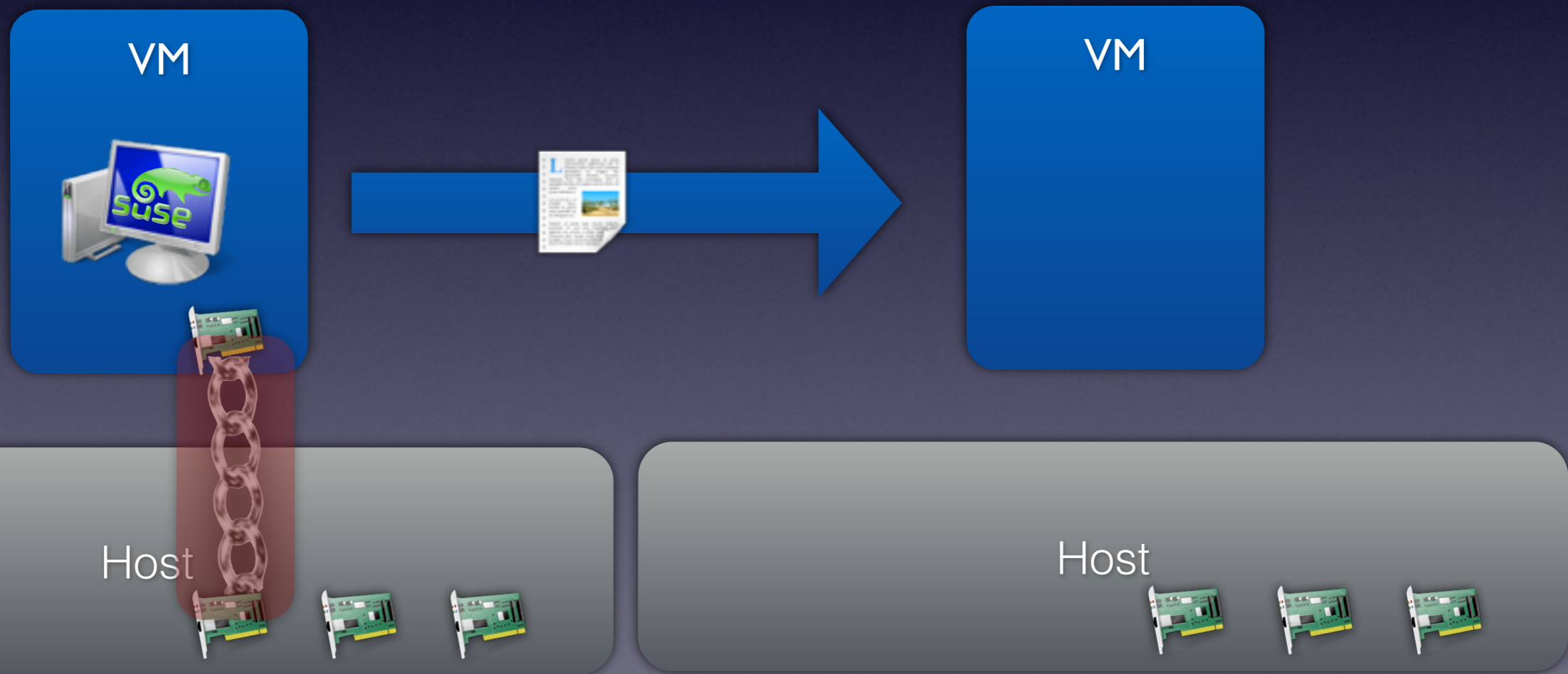


real NIC

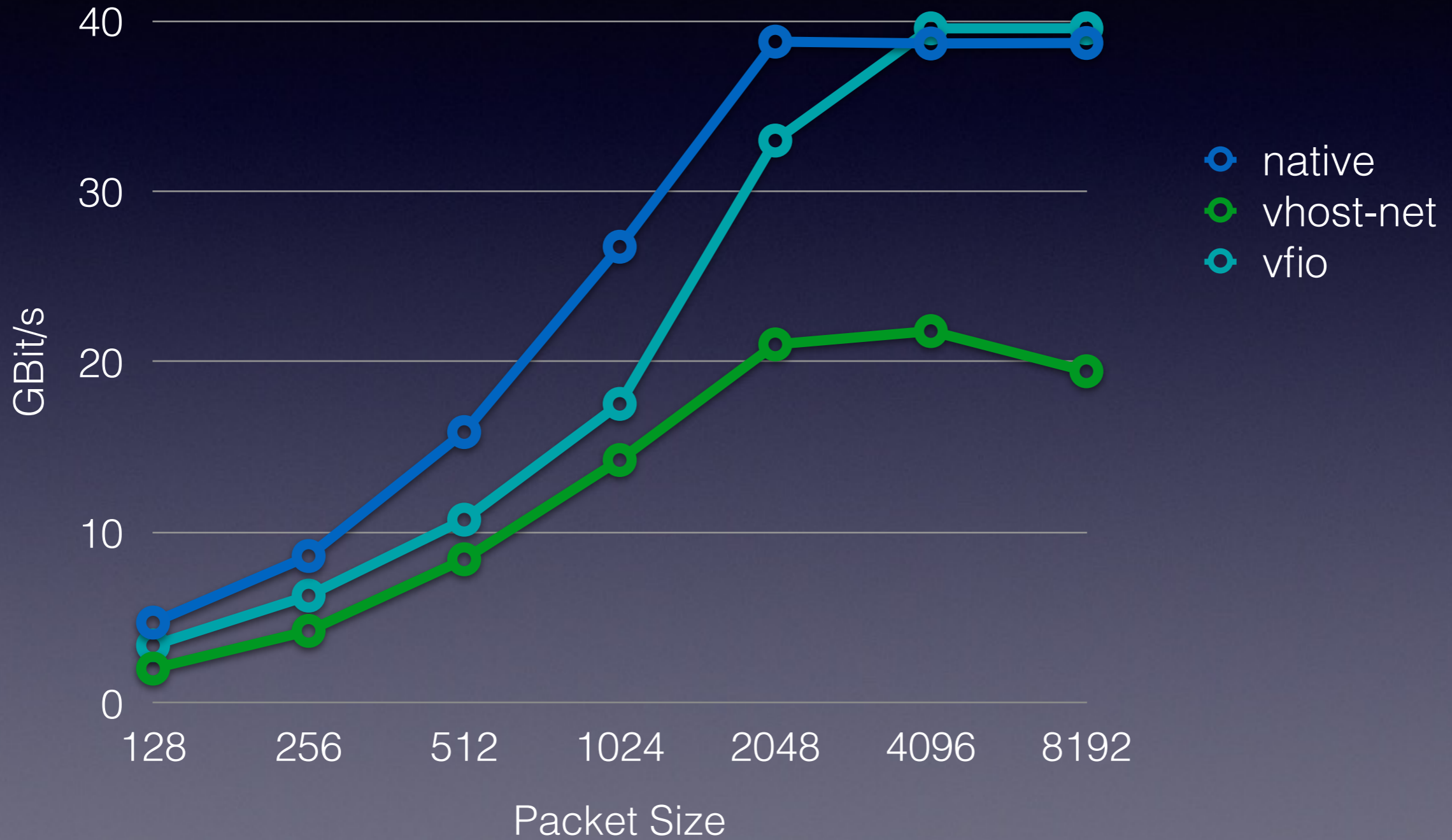
VFIO



VFIO



VFIO



VFIO

- Generic driver
- No chance for introspection

VFIO-i40evf

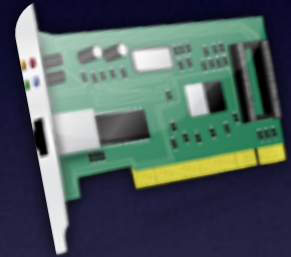
VFIO-i40evf



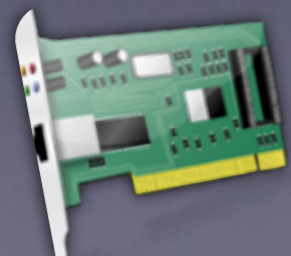
VFIO-i40evf



VFIO-i40evf



guest NIC



real NIC

VFIO-i40evf

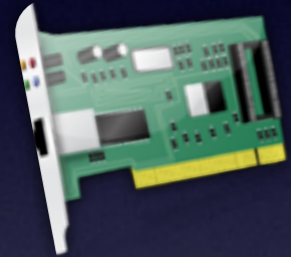


VFIO-i40evf



- Set MAC addresses
- Set RX/TX locations
- Set RX/TX properties
- ...

VFIO-i40evf



guest NIC



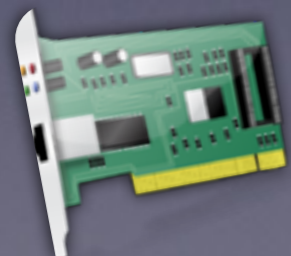
TX ring



RX ring

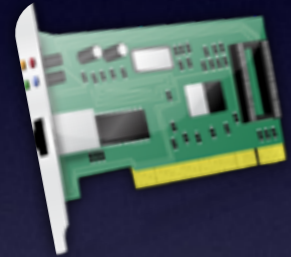


Admin ring

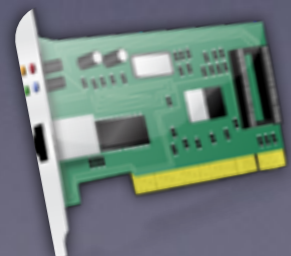


real NIC

VFIO-i40evf



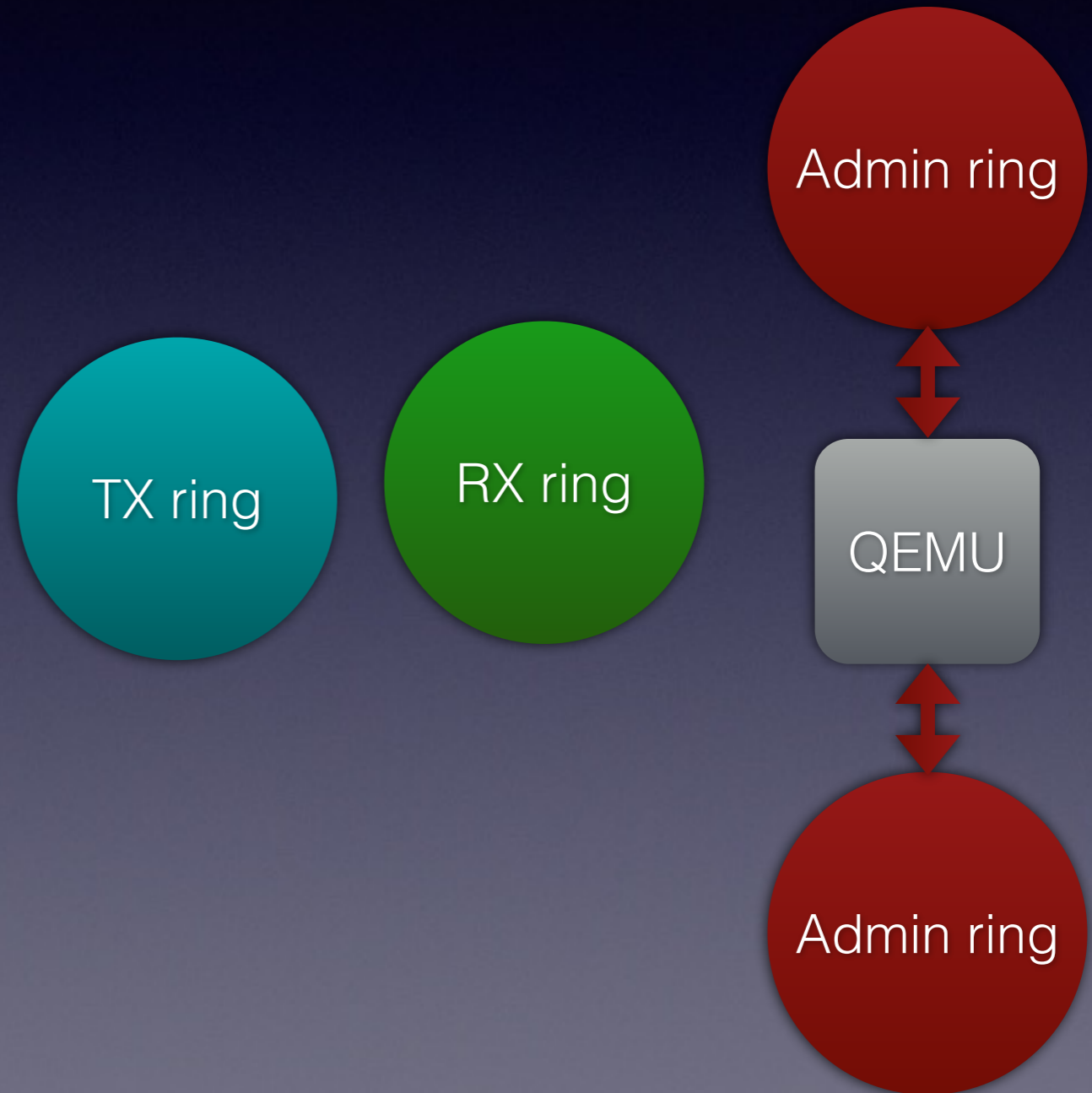
guest NIC



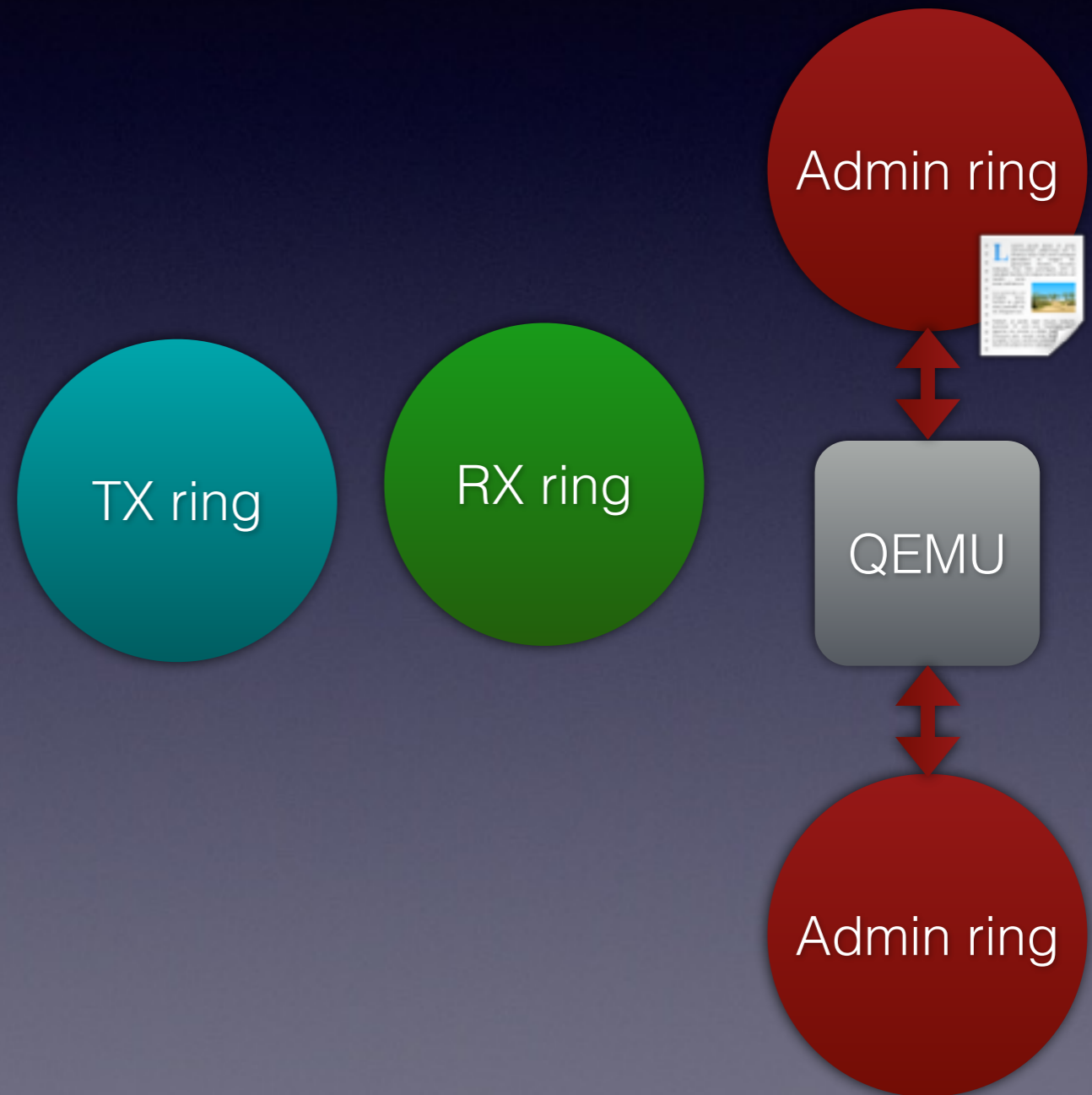
real NIC



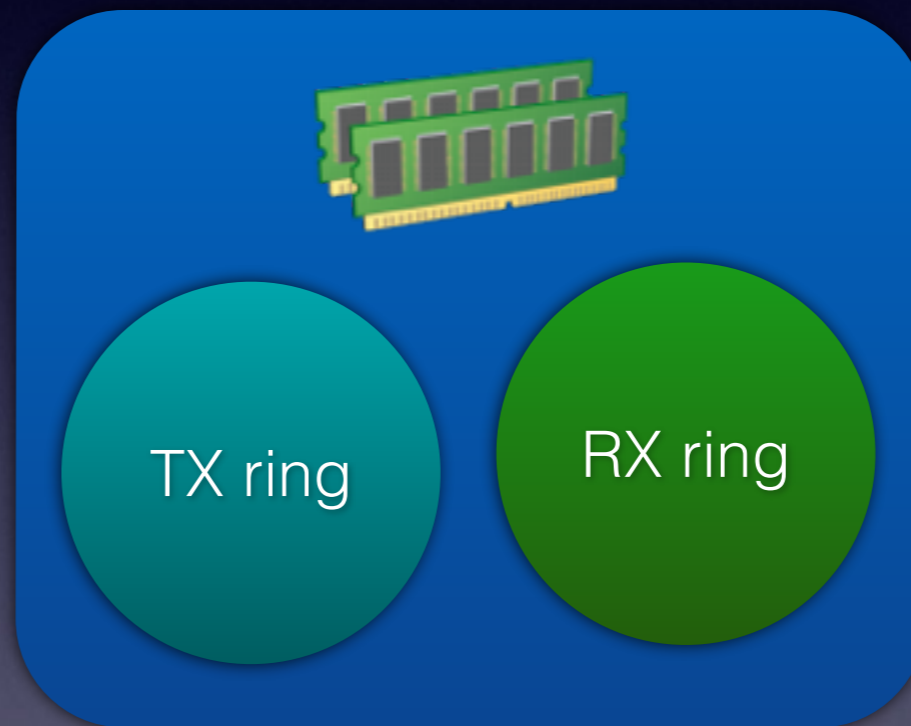
VFIO-i40evf



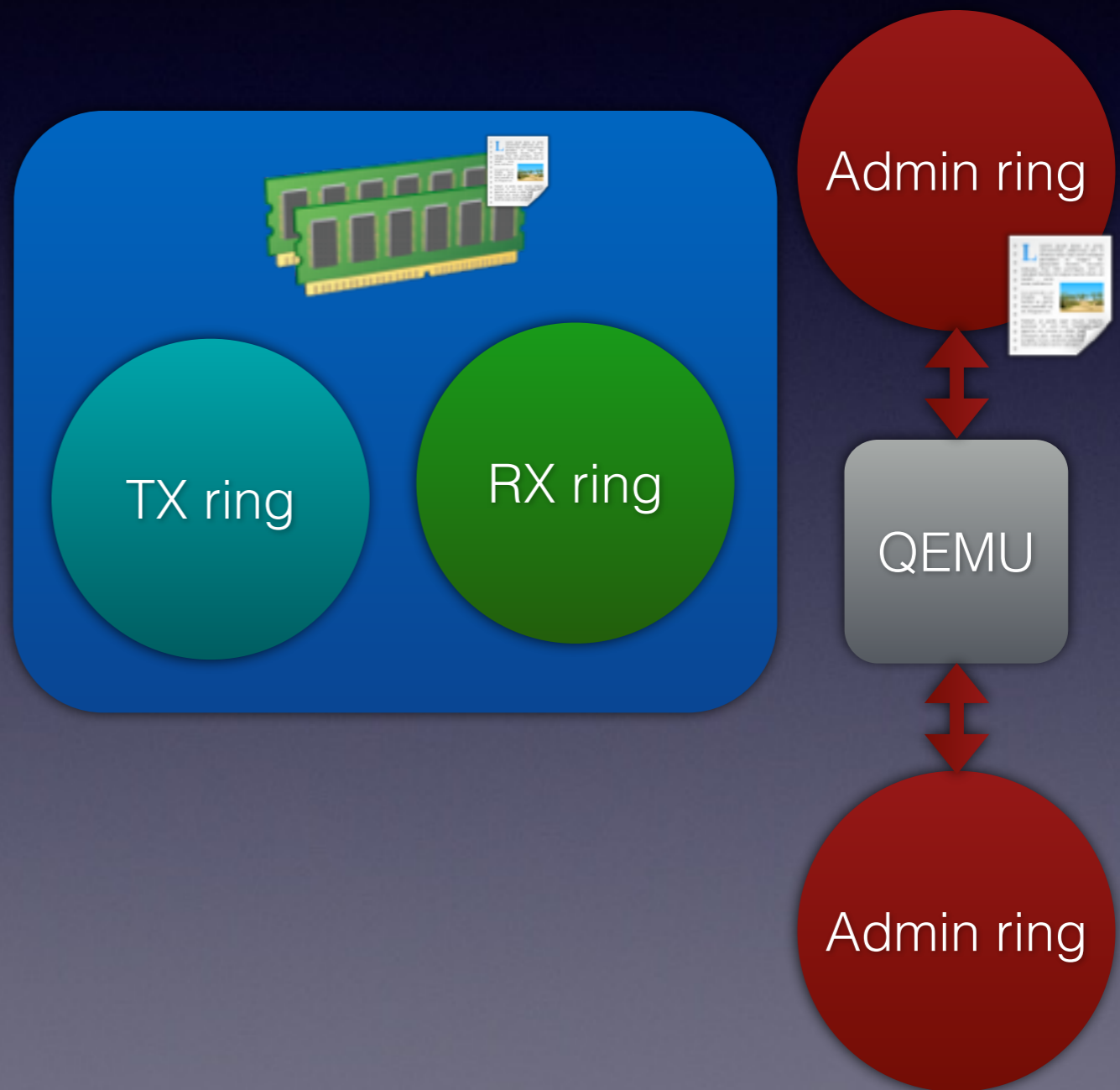
VFIO-i40evf



VFIO-i40evf



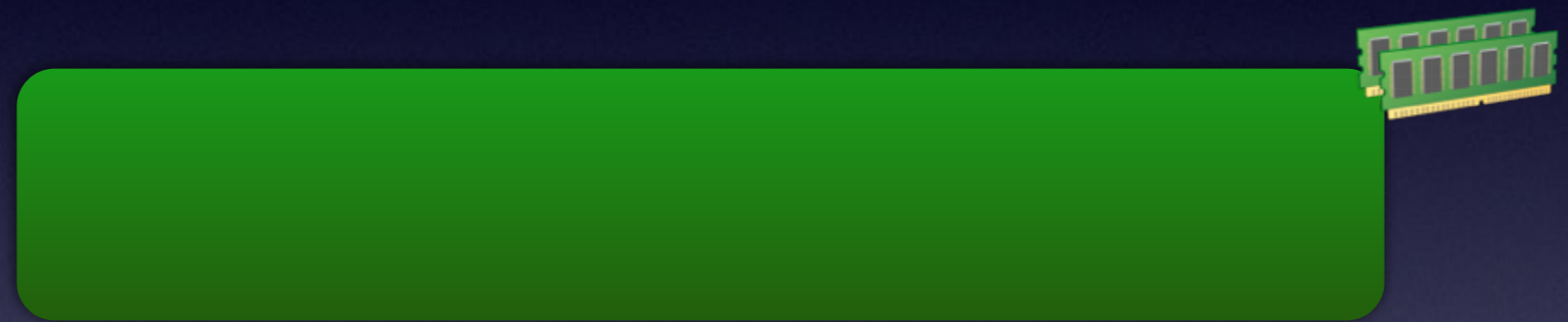
VFIO-i40evf



VFIO-i40evf



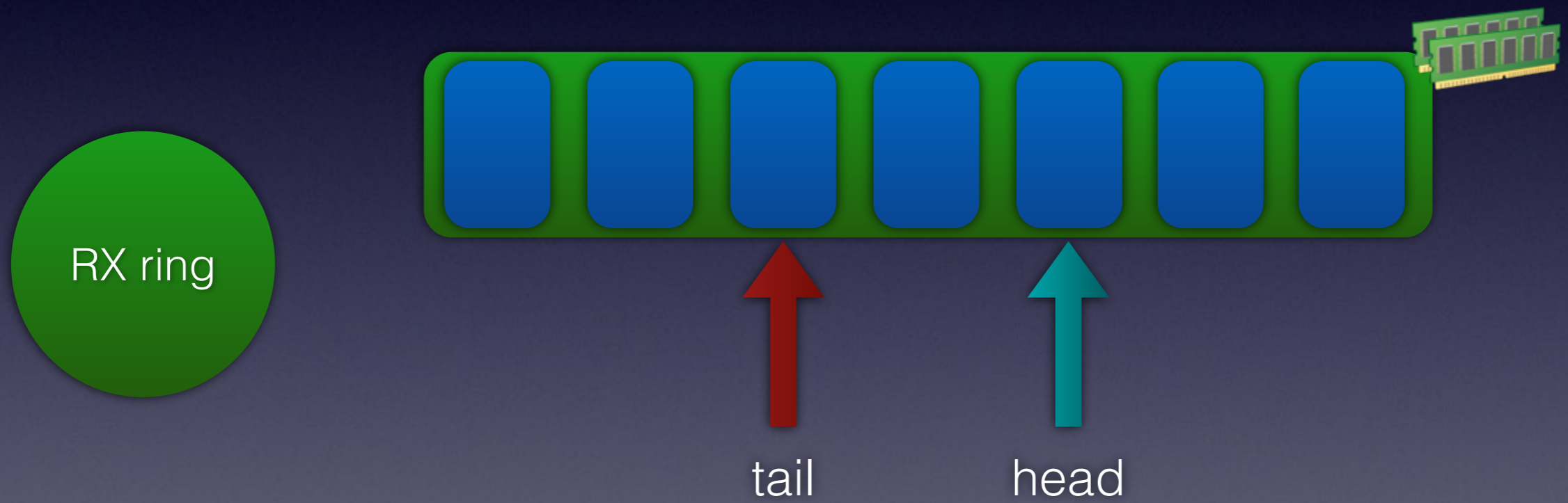
VFIO-i40evf



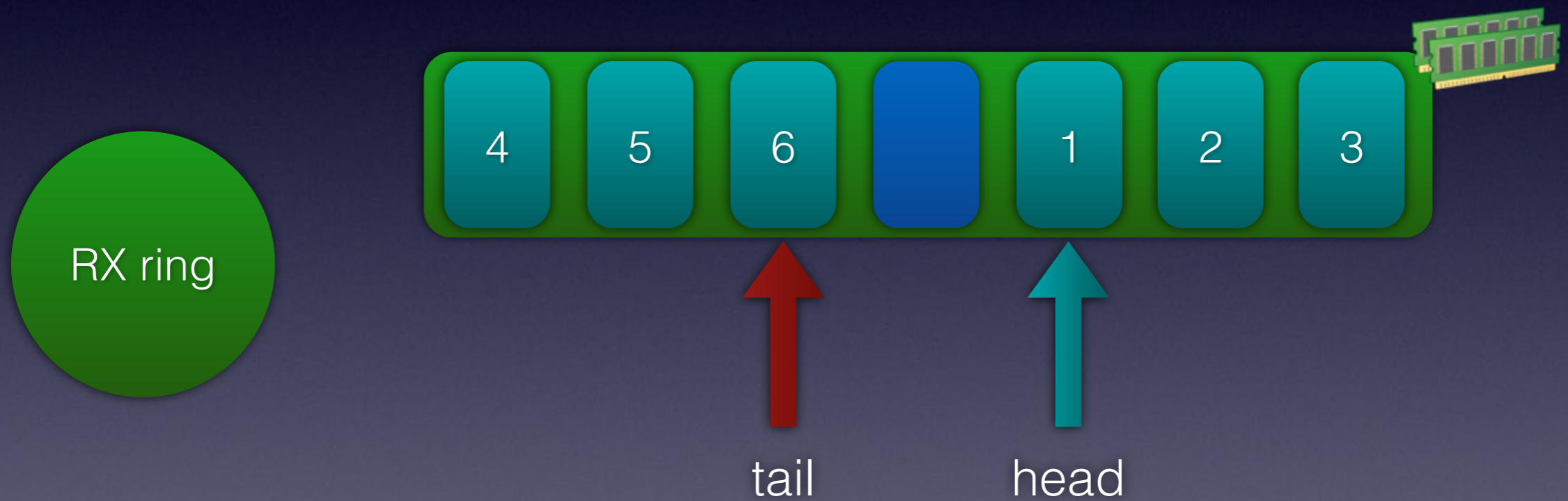
VFIO-i40evf



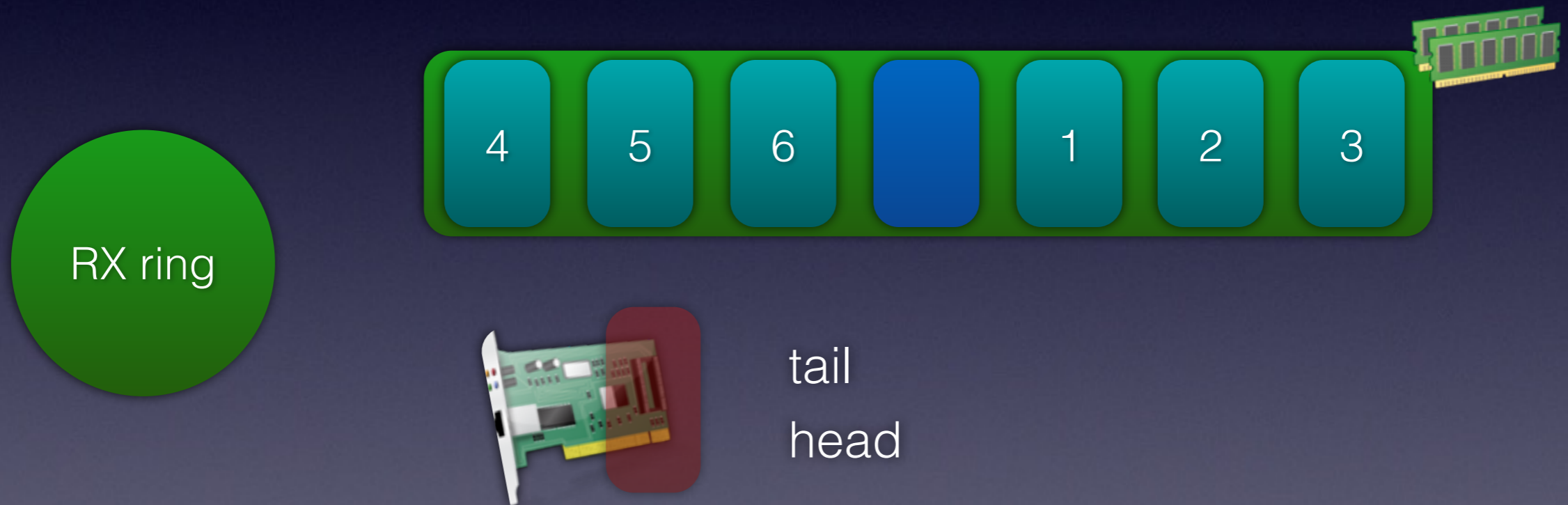
VFIO-i40evf



VFIO-i40evf



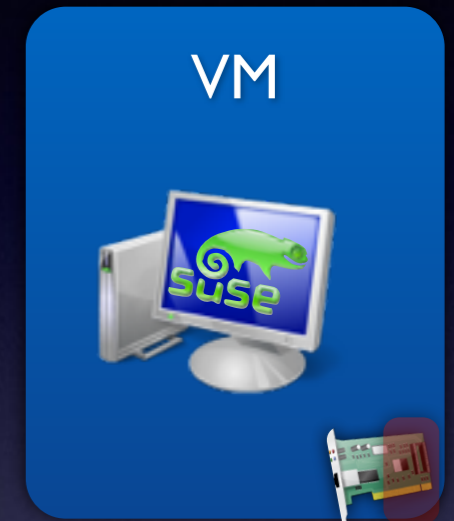
VFIO-i40evf



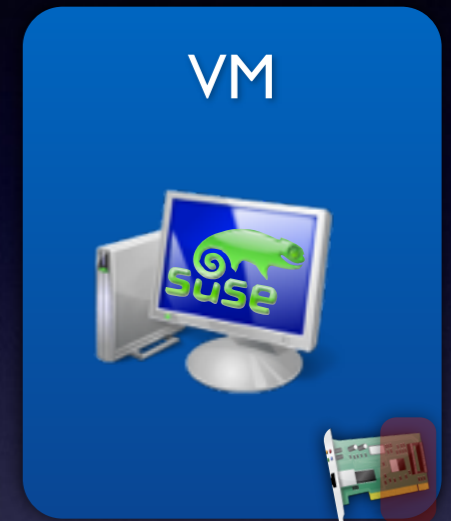
VFIO-i40evf



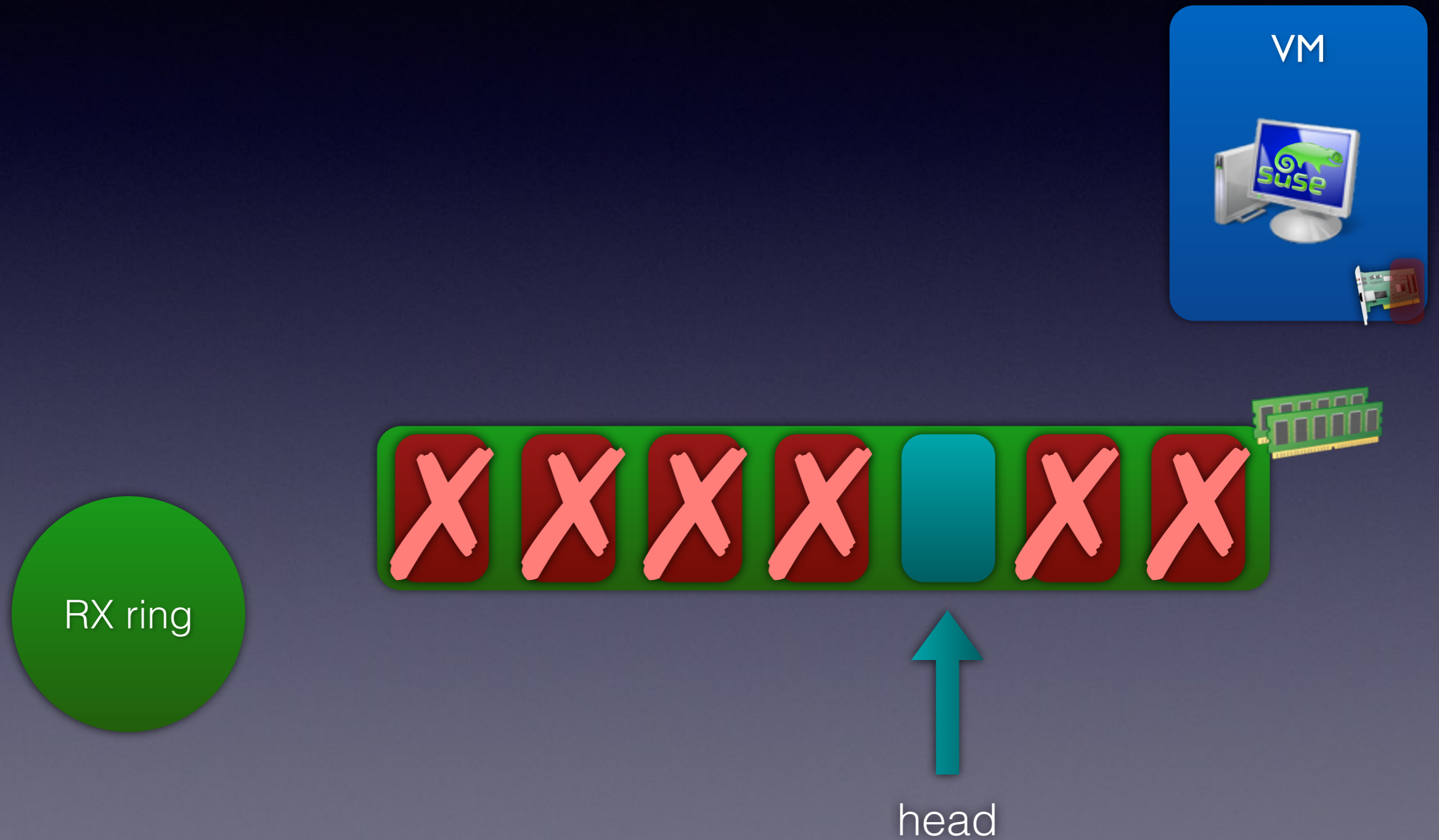
VFIO-i40evf



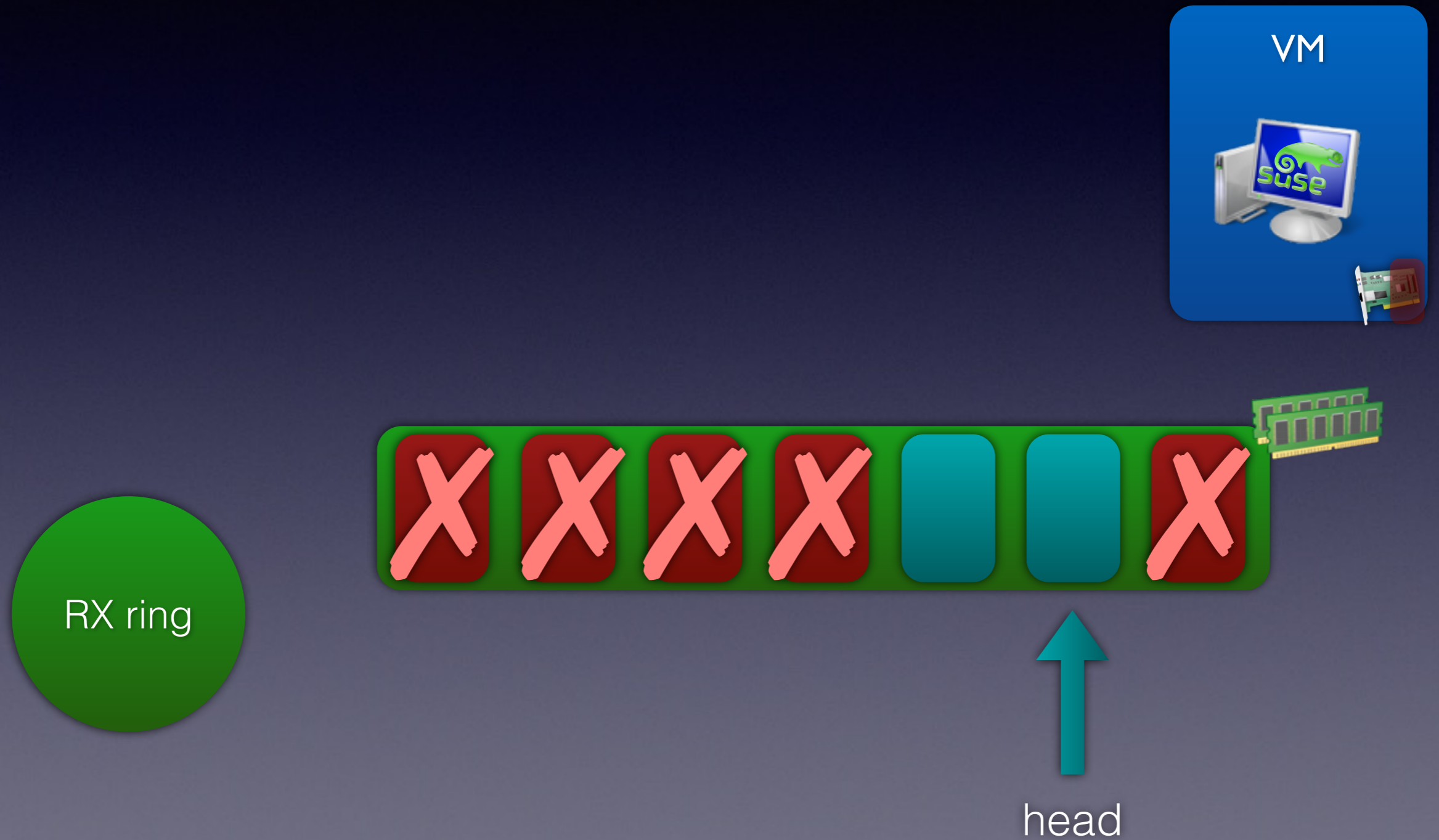
VFIO-i40evf



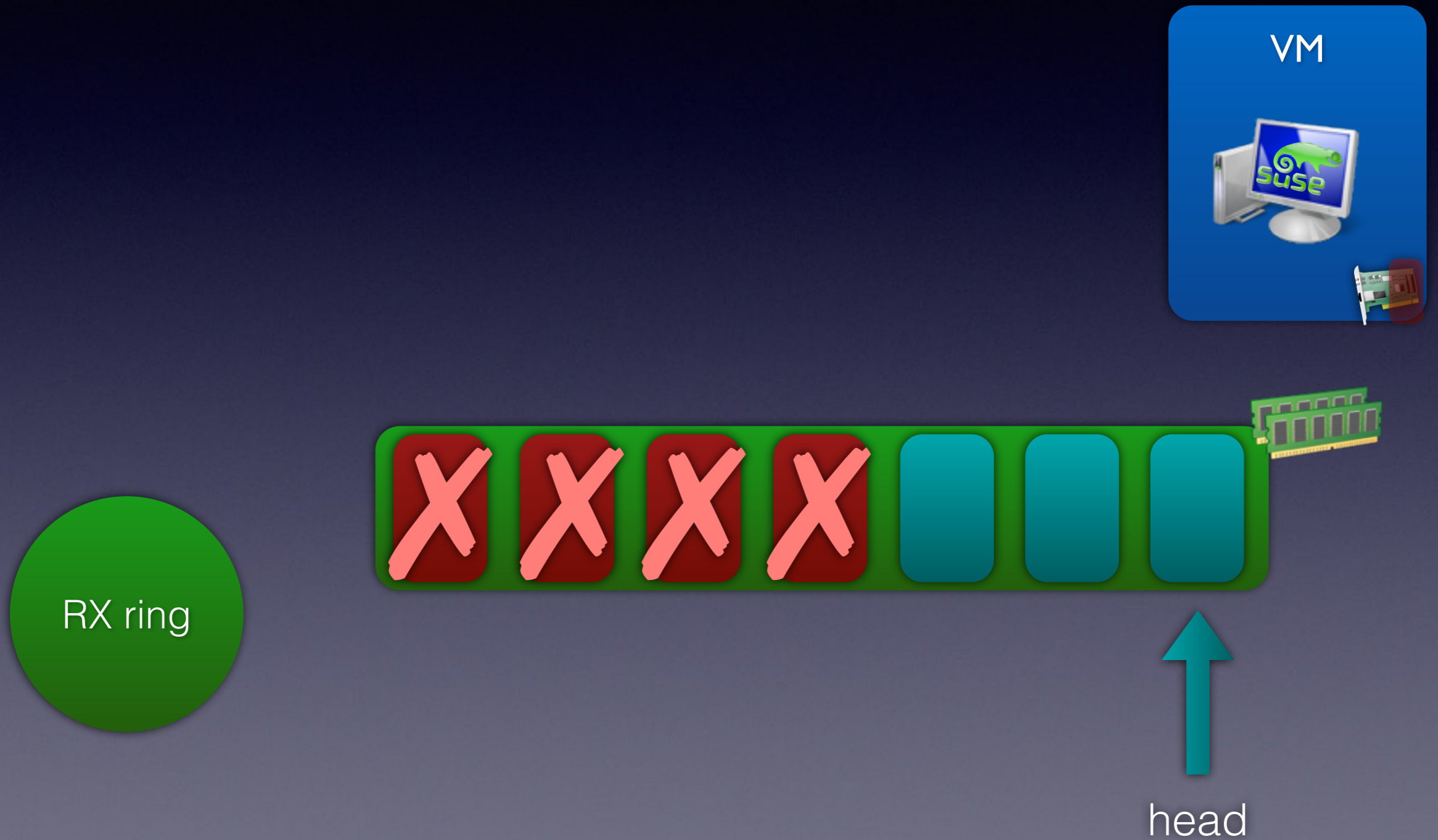
VFIO-i40evf



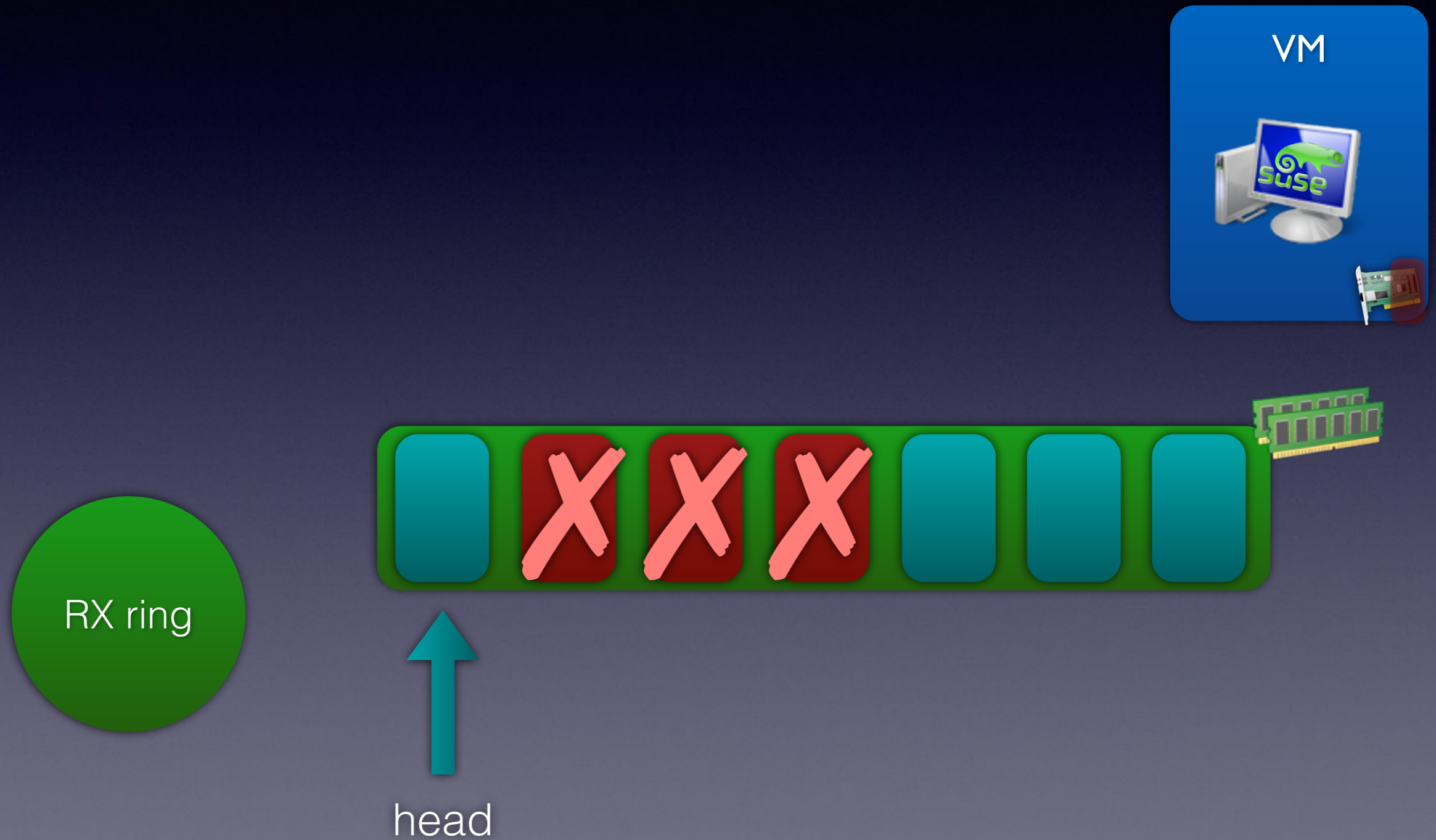
VFIO-i40evf



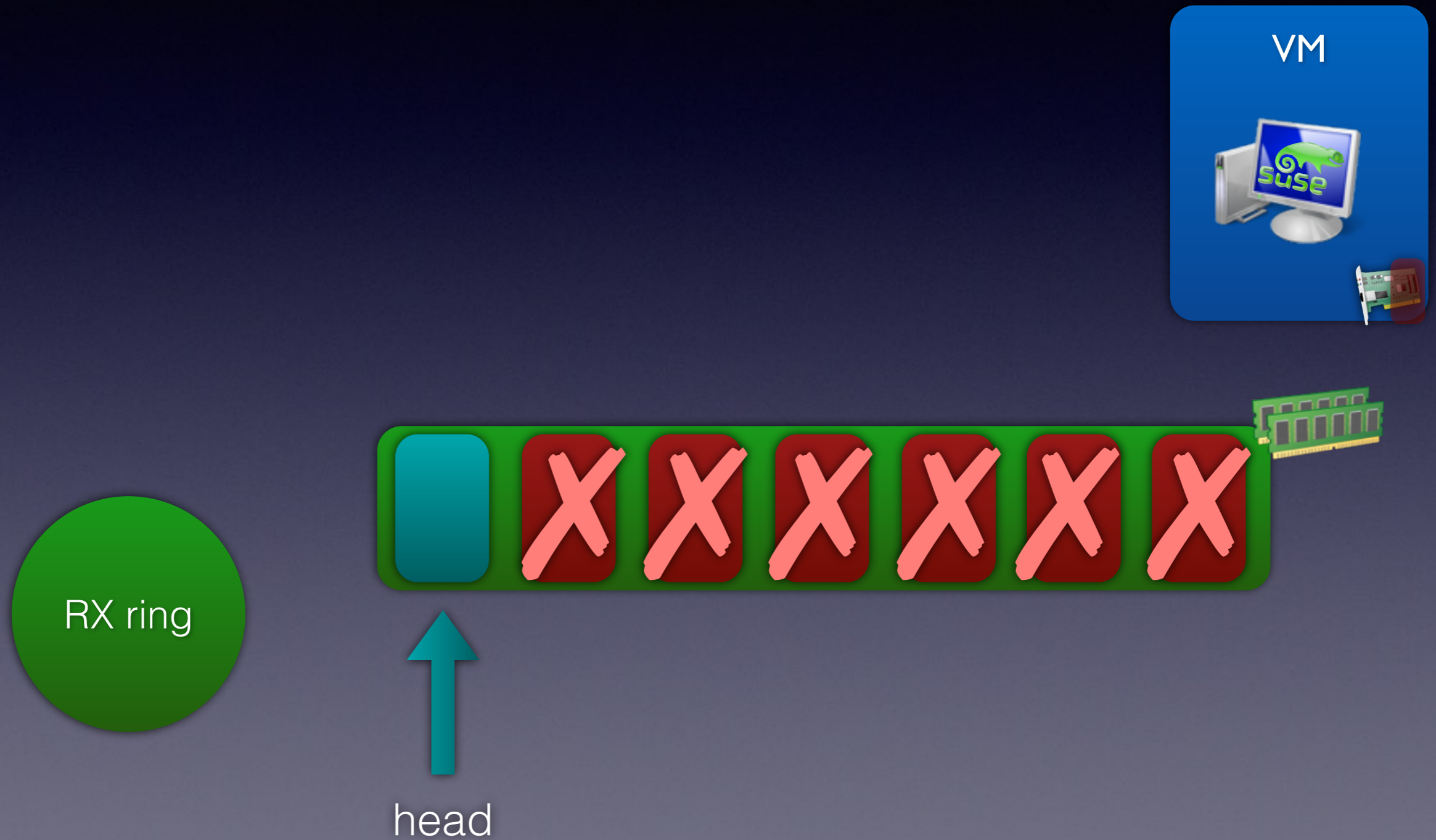
VFIO-i40evf



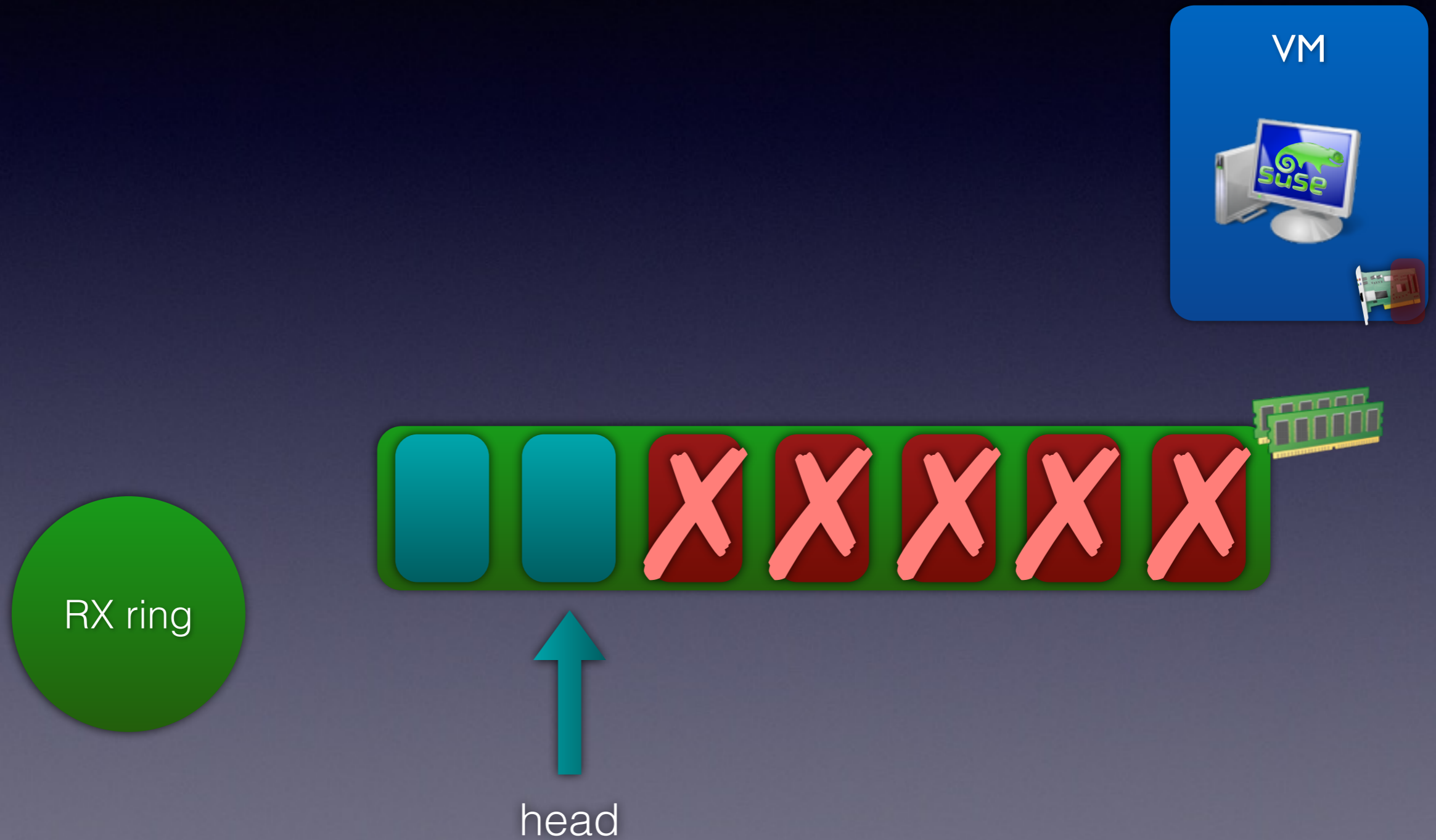
VFIO-i40evf



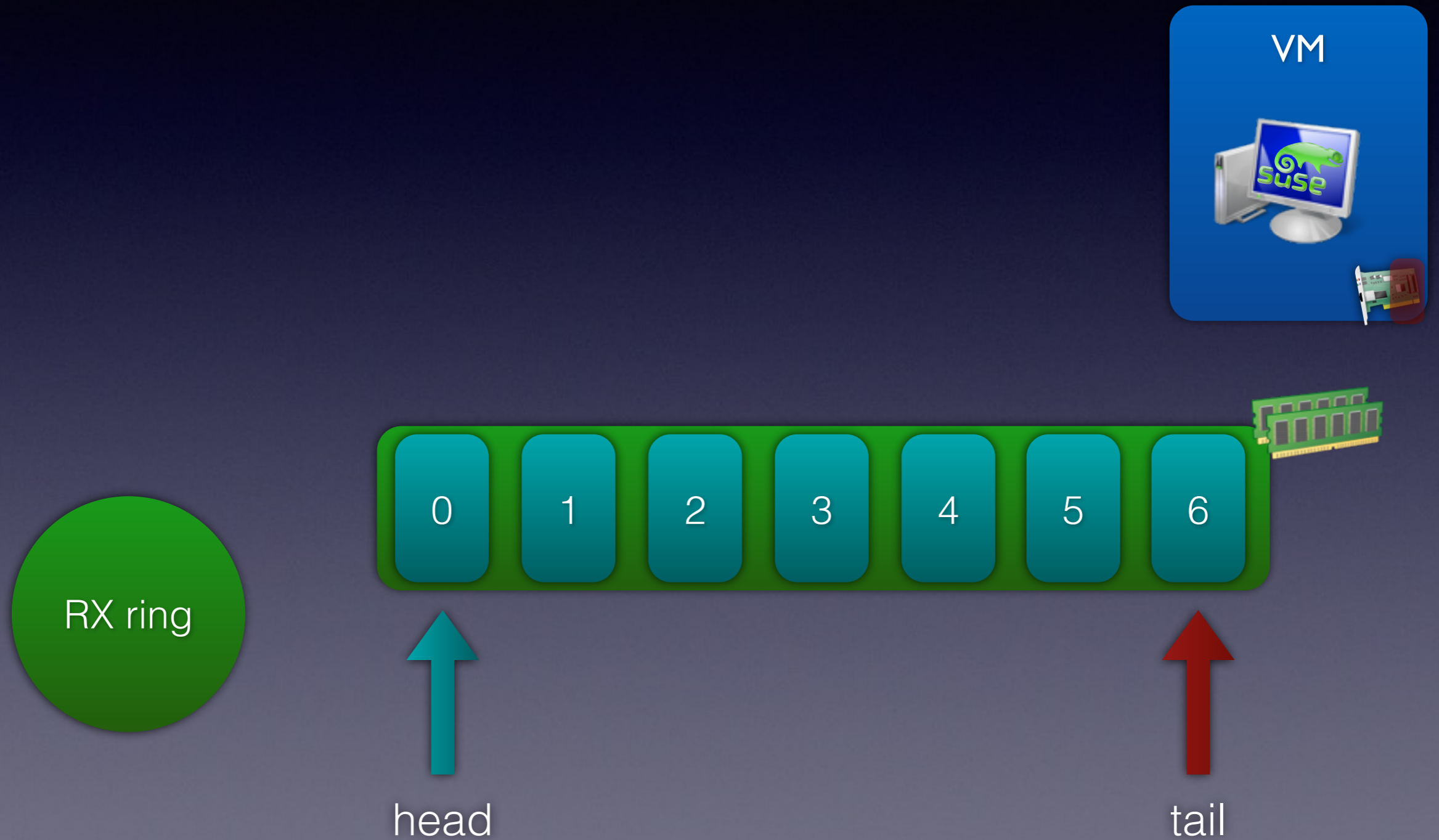
VFIO-i40evf



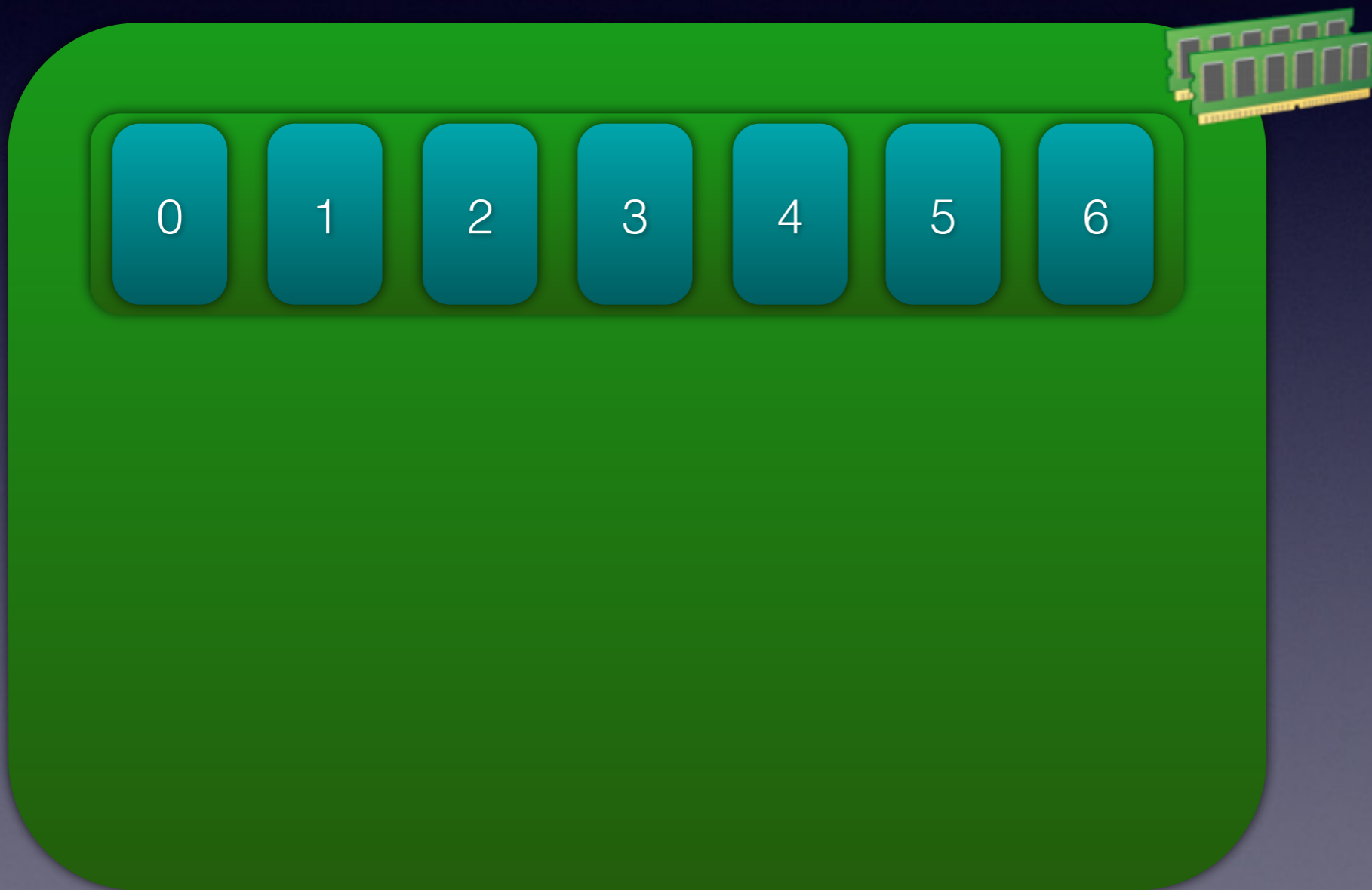
VFIO-i40evf



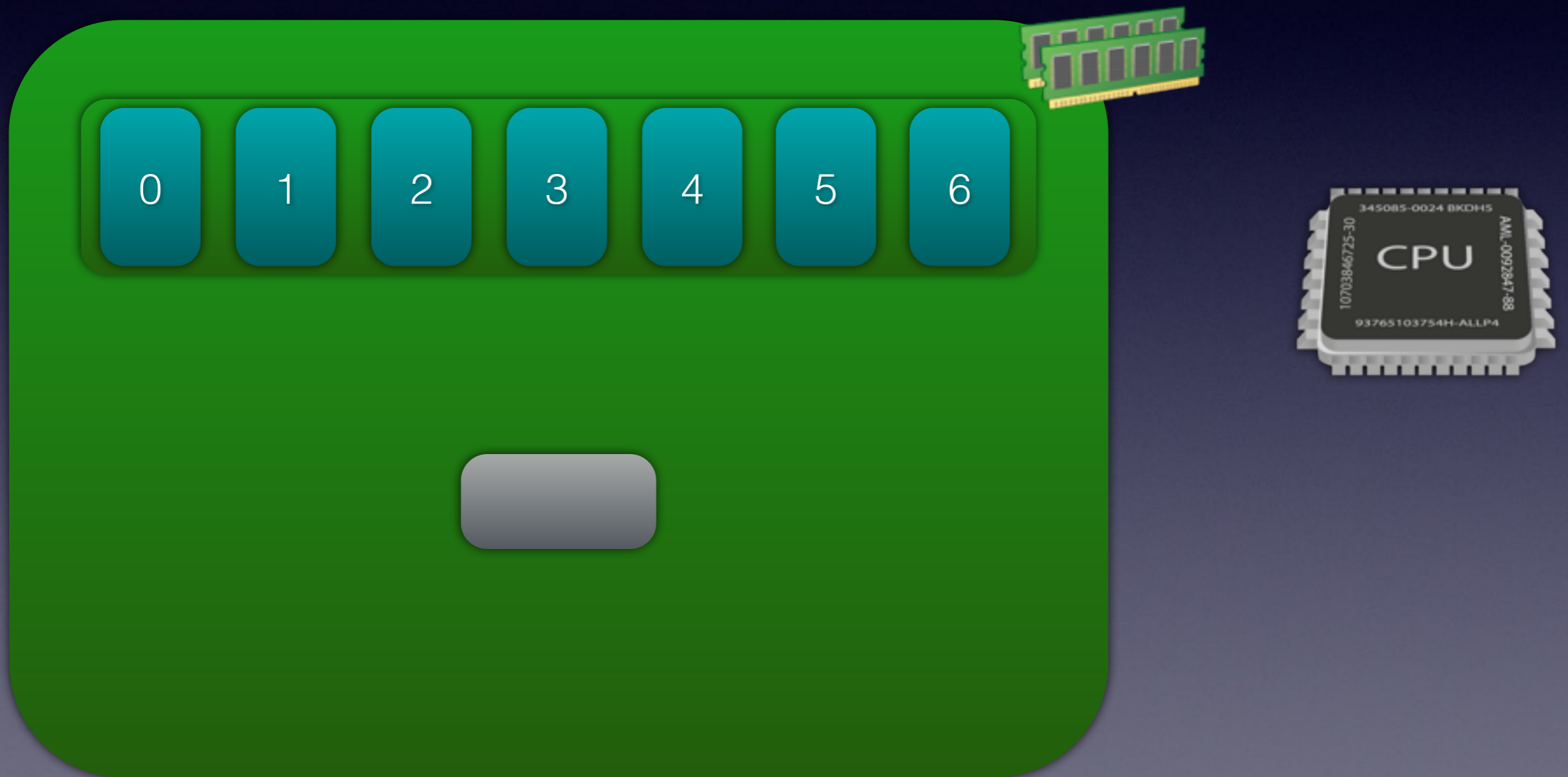
VFIO-i40evf



DMA



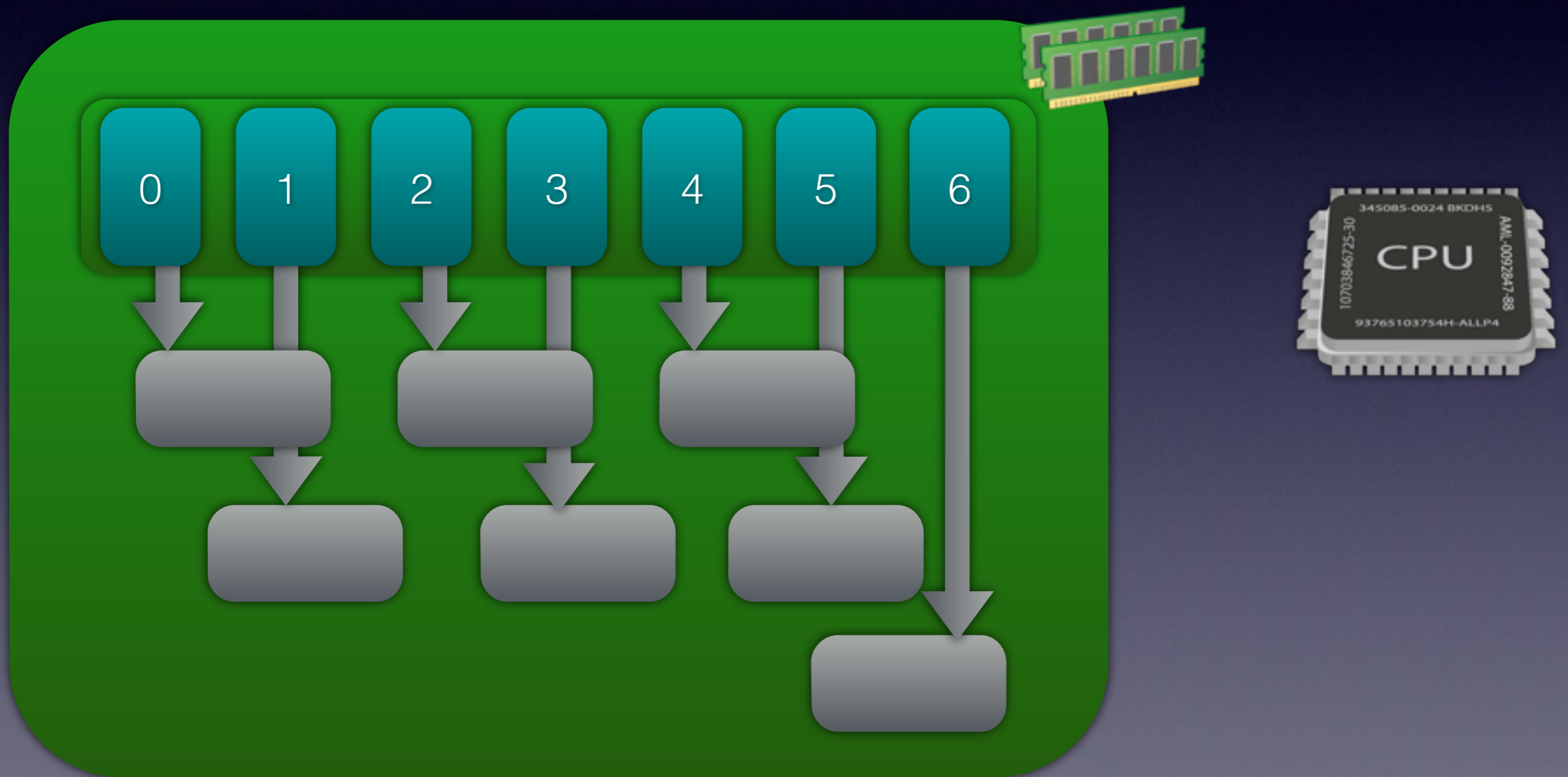
DMA



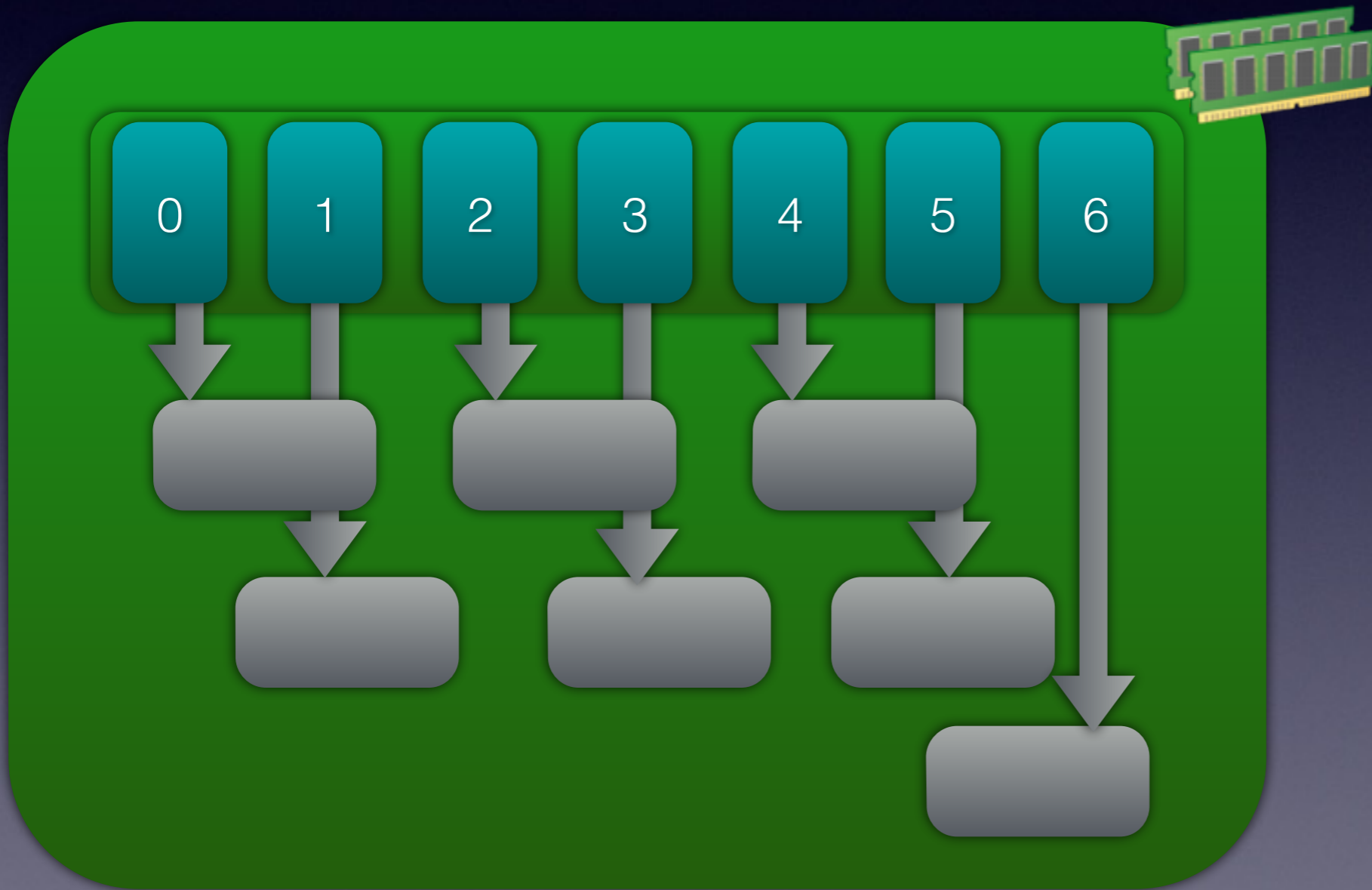
DMA



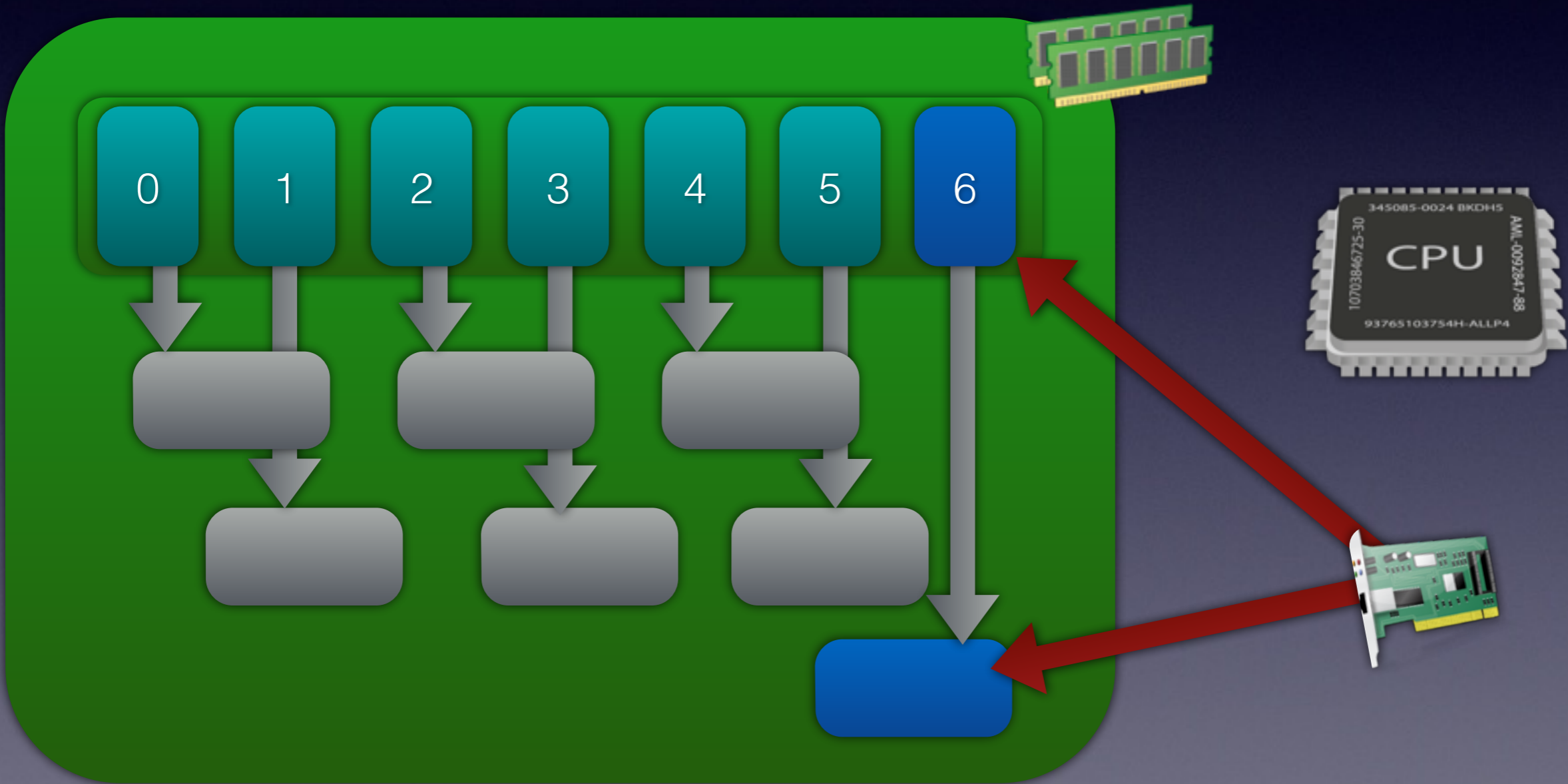
DMA



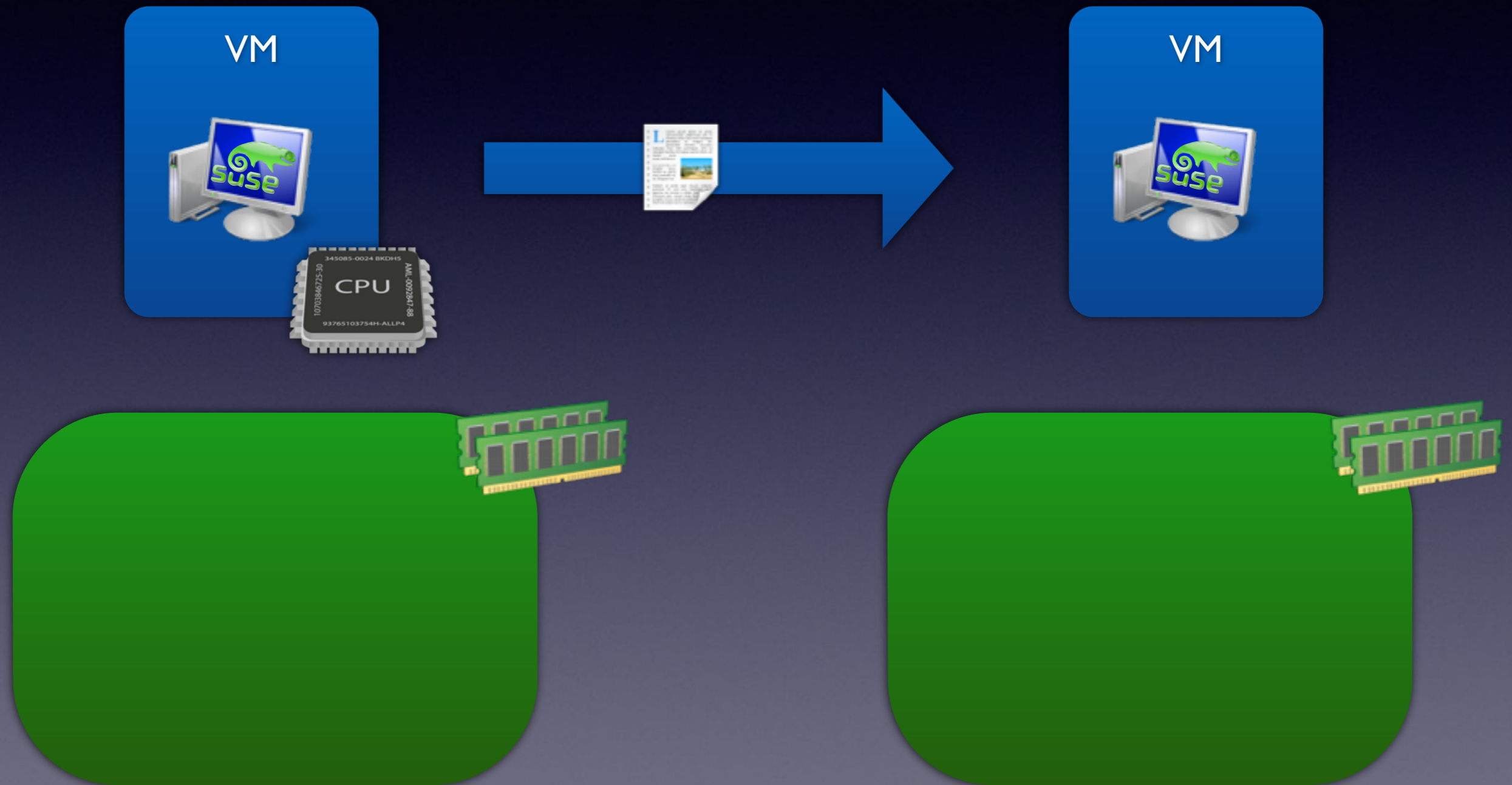
DMA



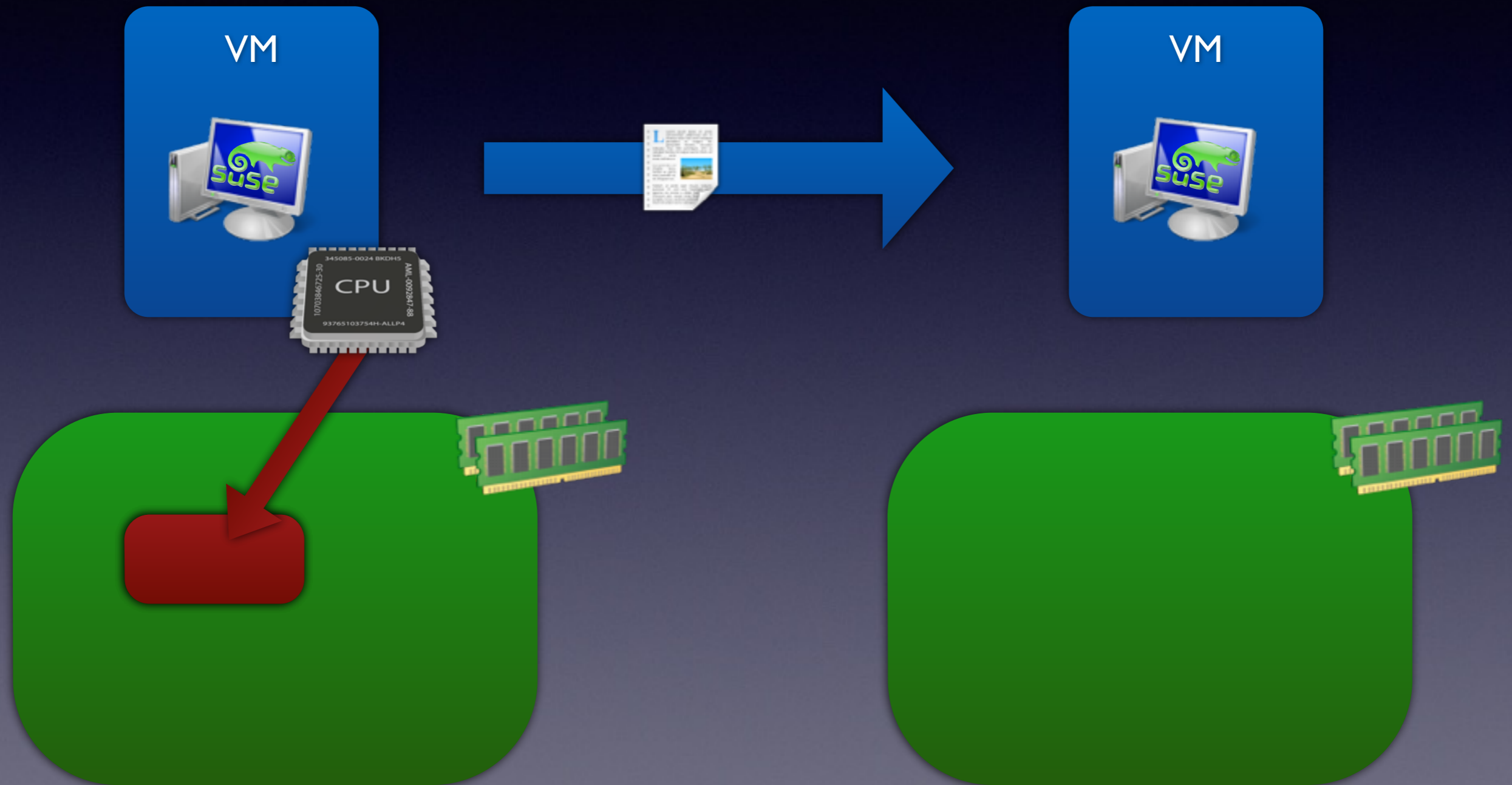
DMA



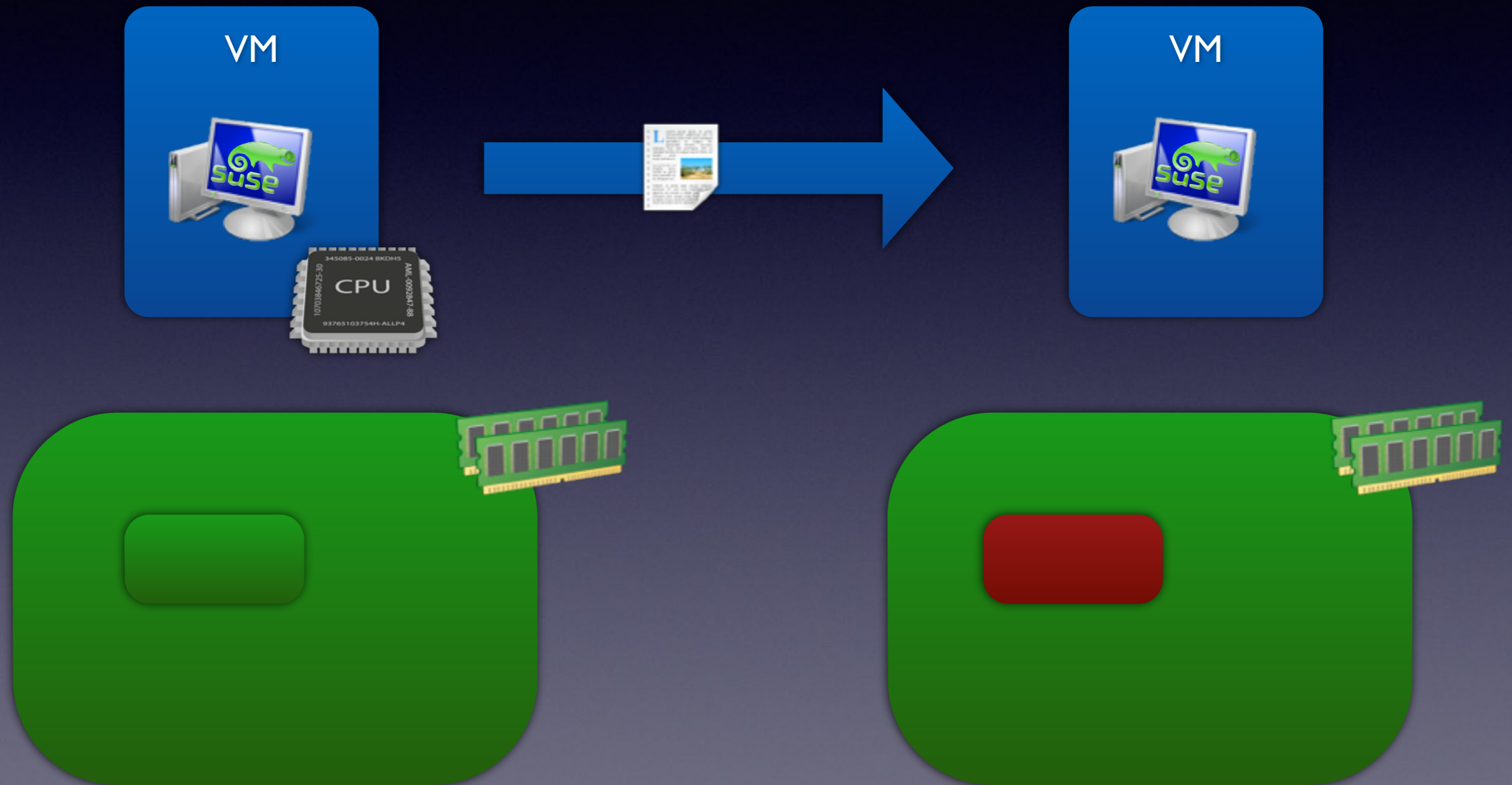
VFIO-i40evf



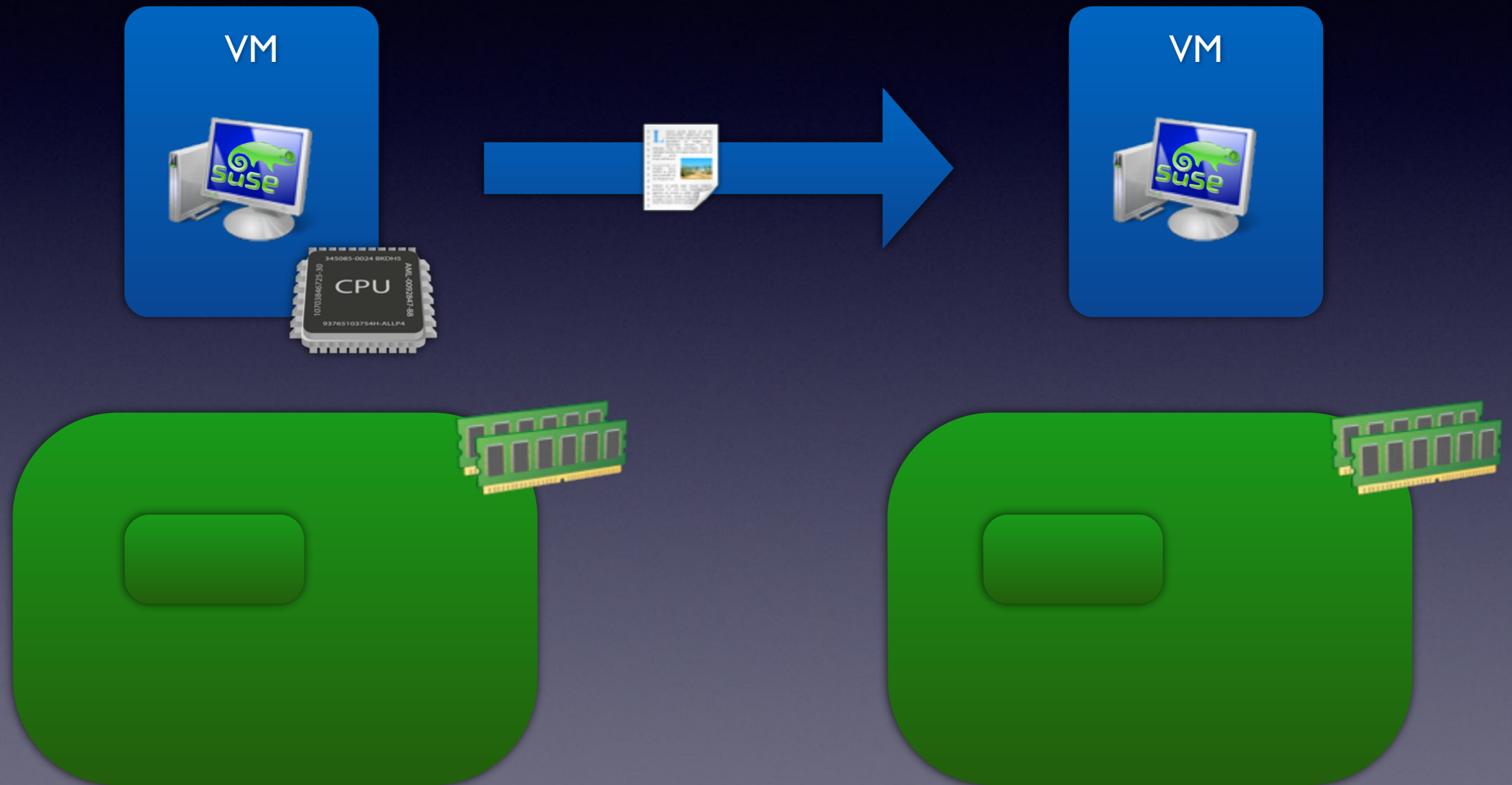
VFIO-i40evf



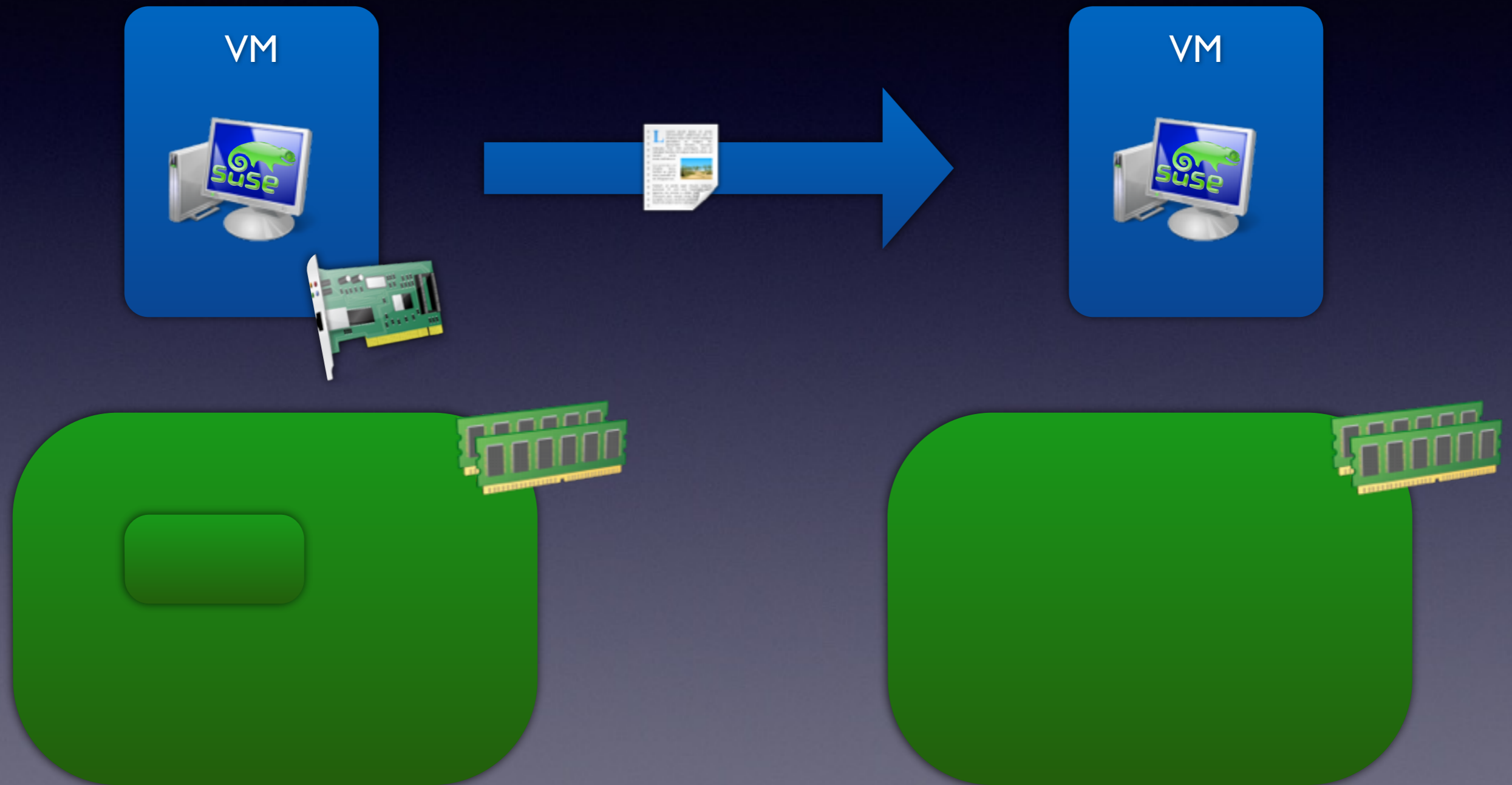
VFIO-i40evf



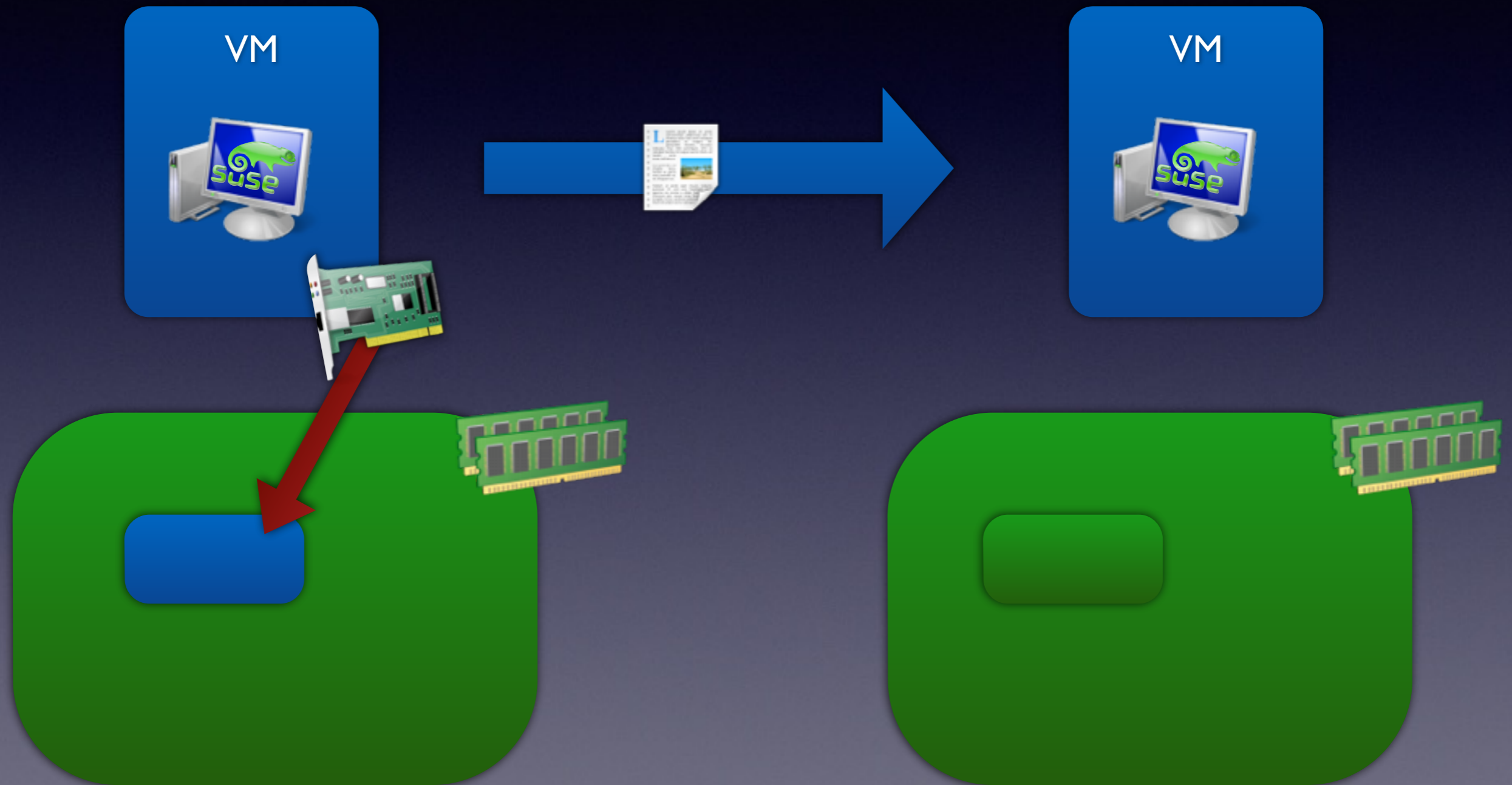
VFIO-i40evf



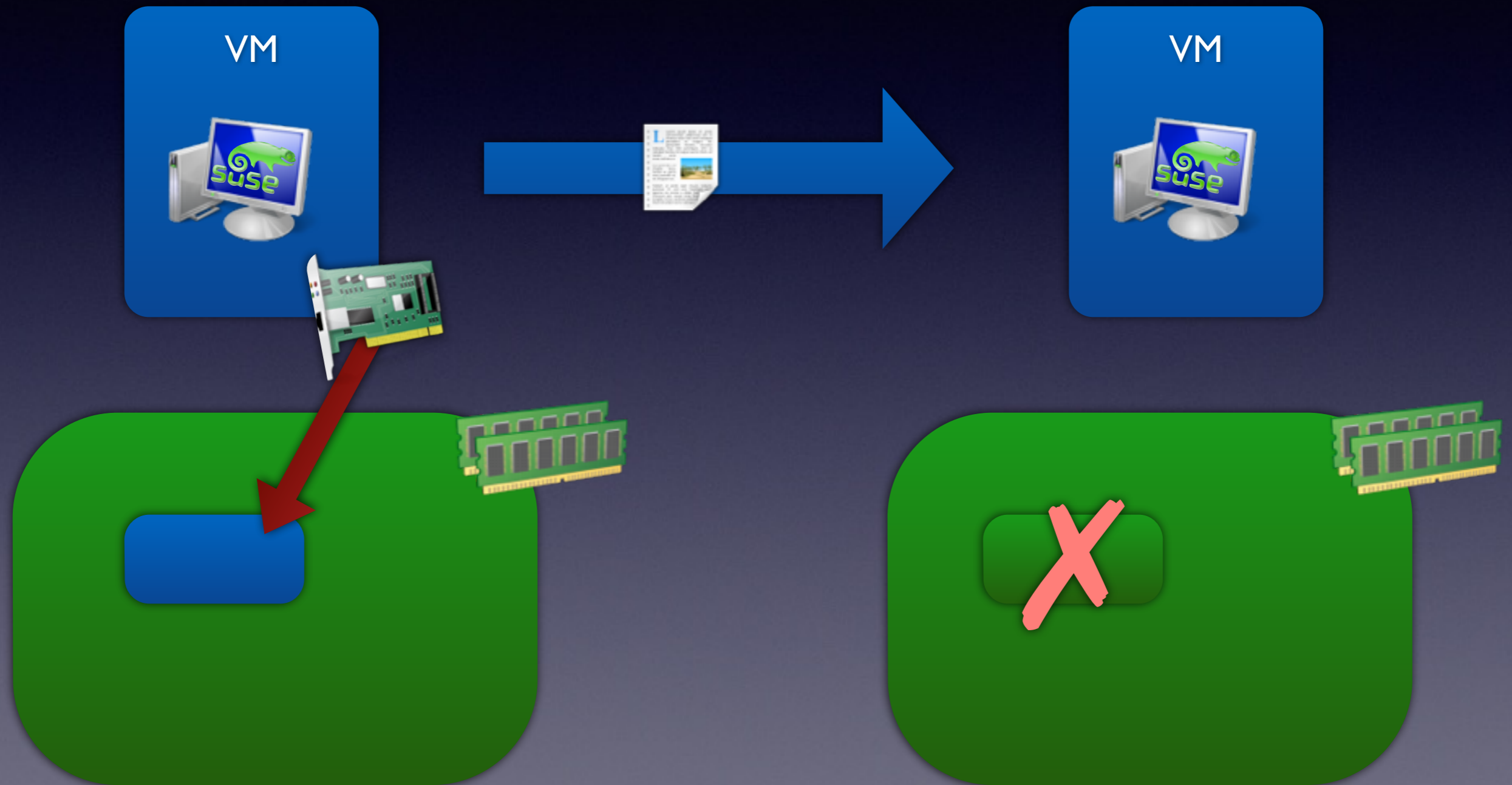
VFIO-i40evf



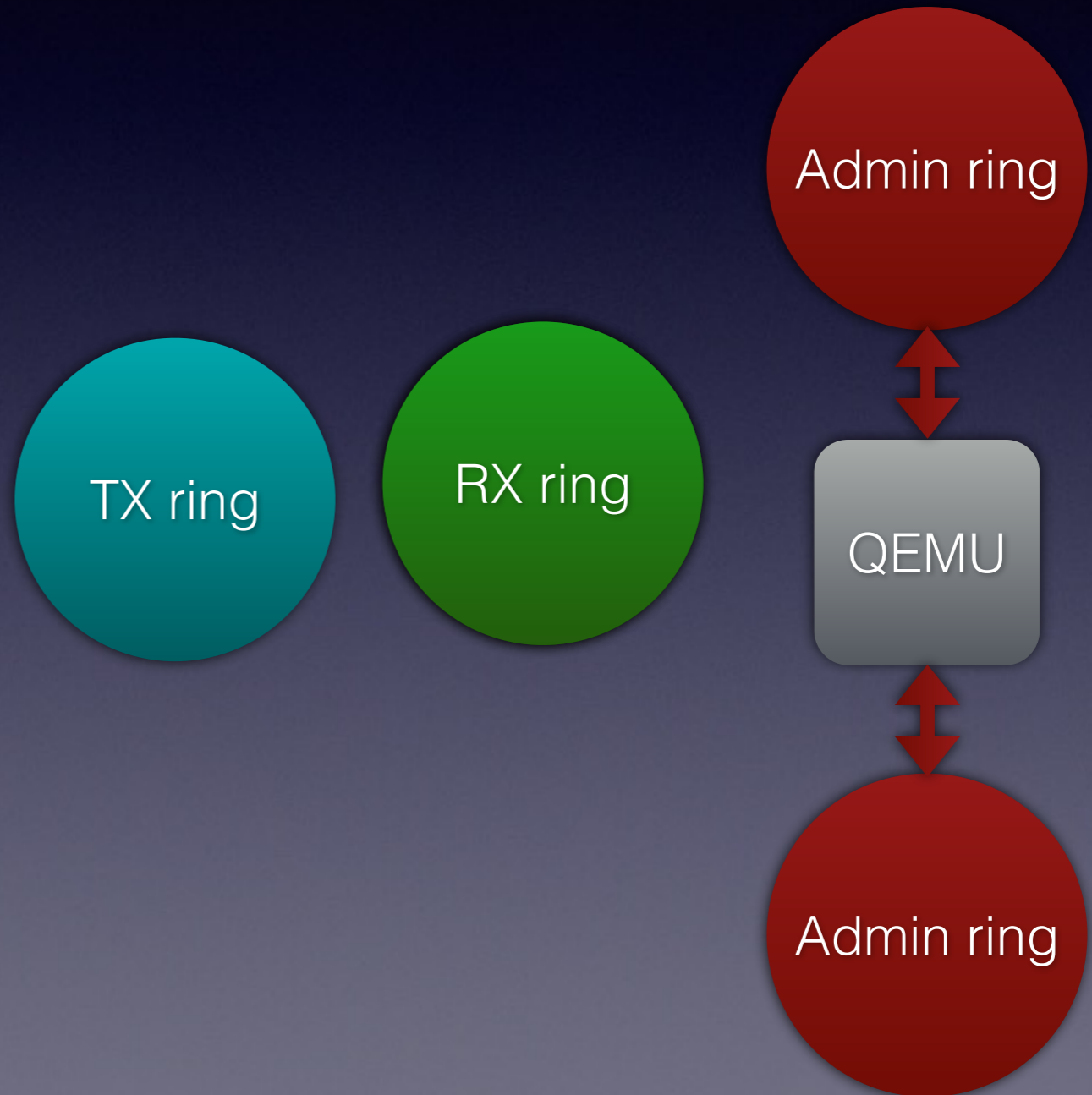
VFIO-i40evf



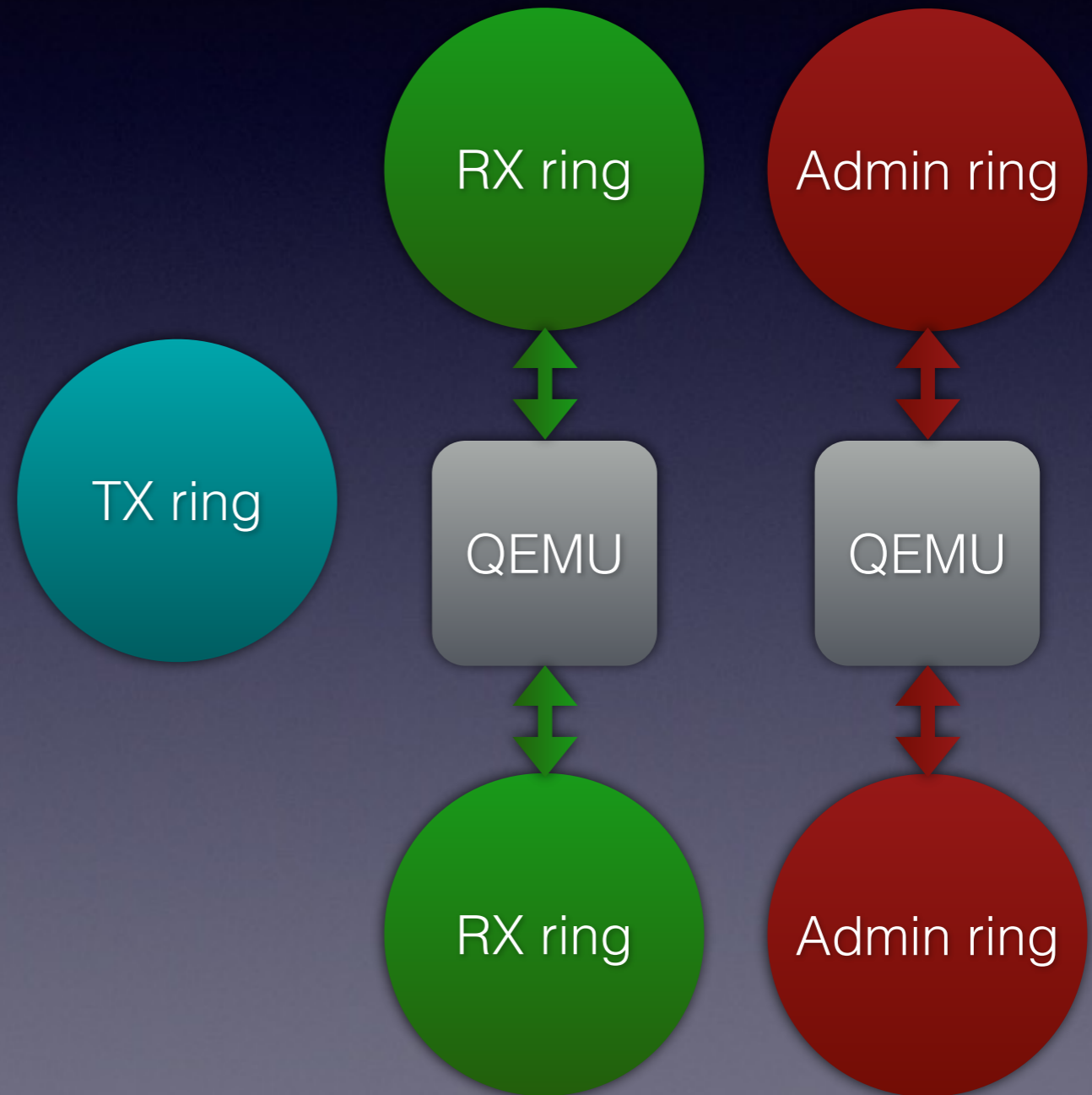
VFIO-i40evf



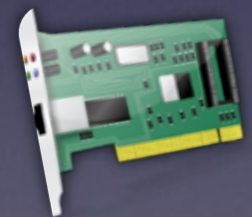
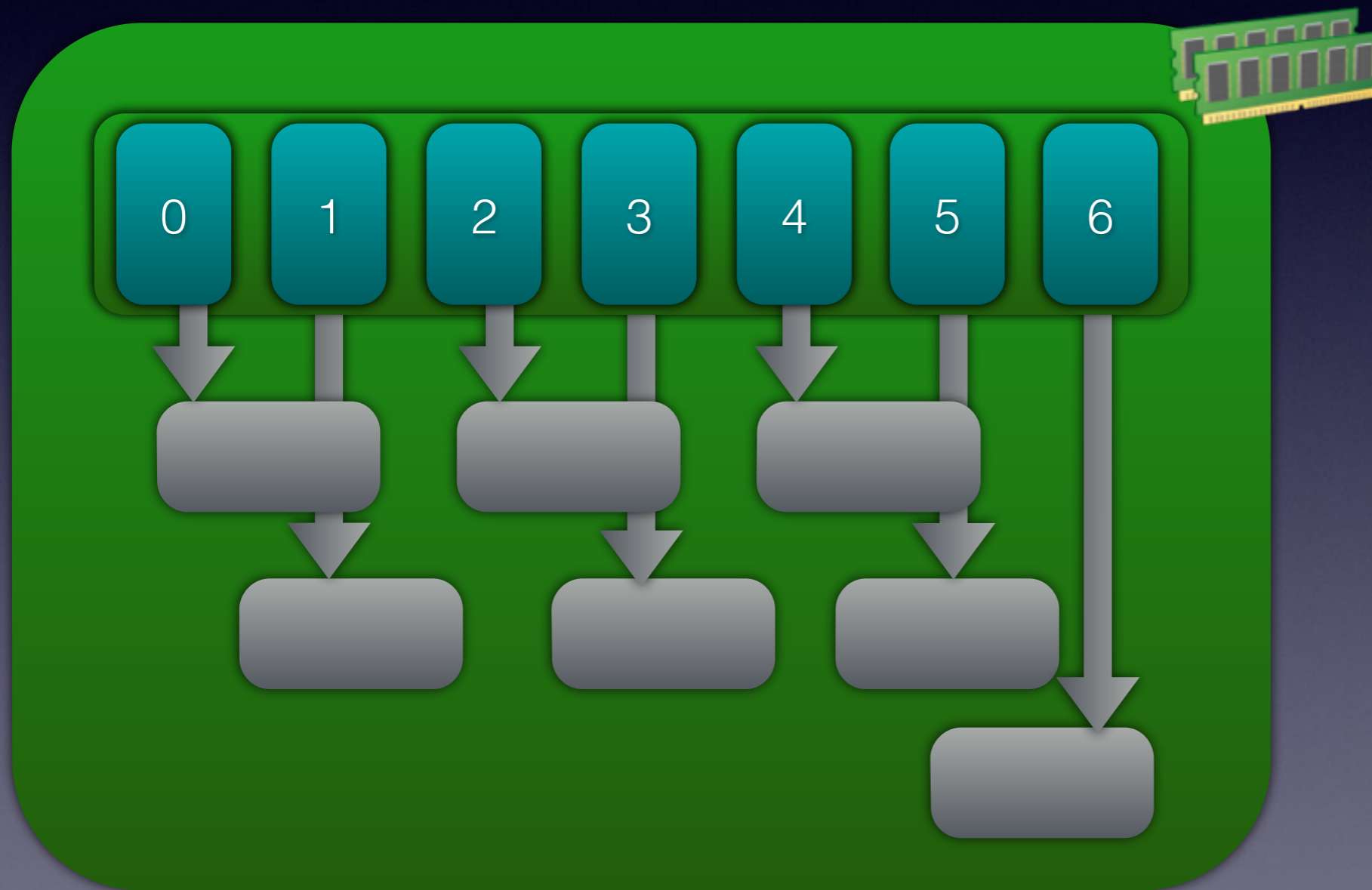
VFIO-i40evf



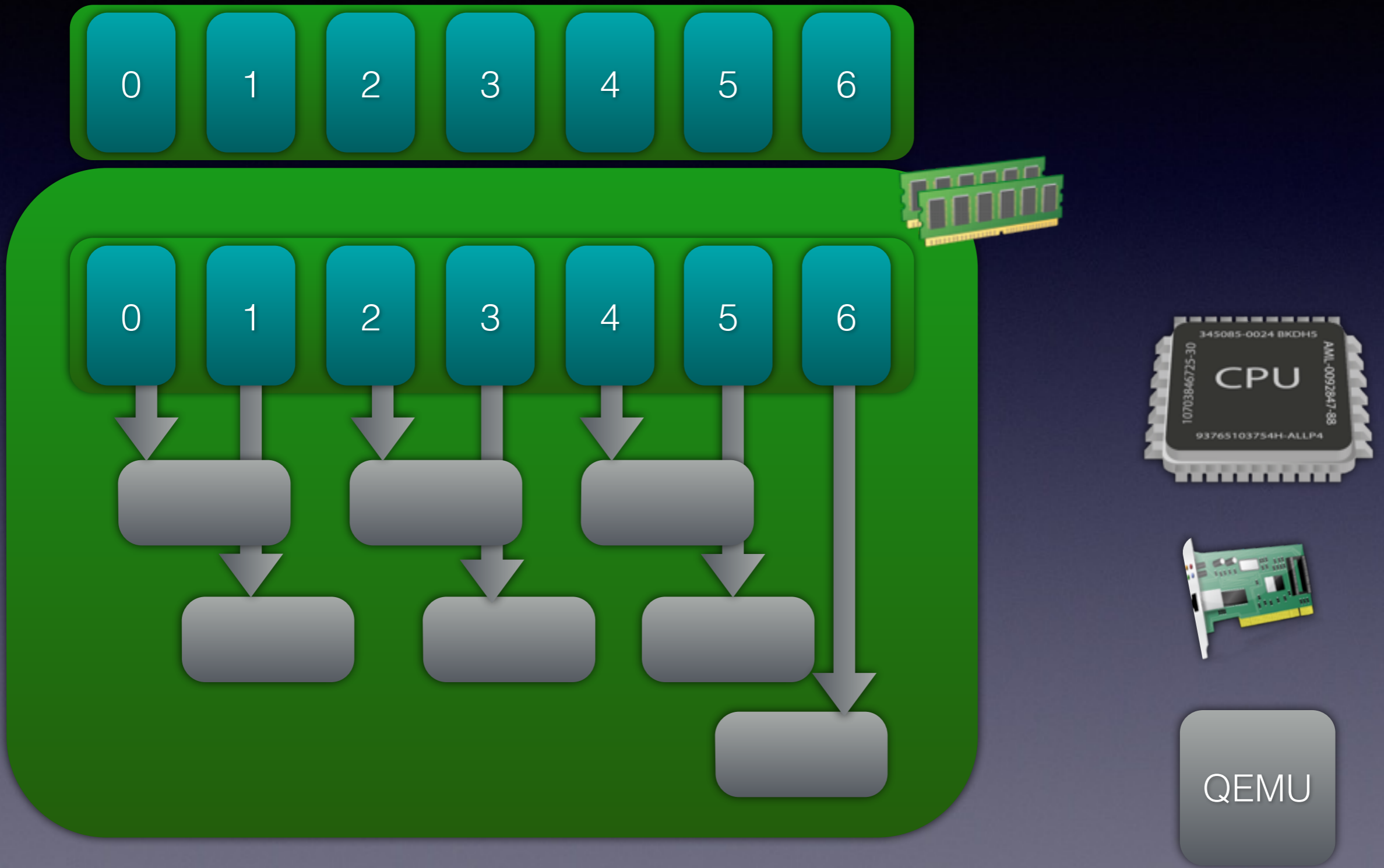
VFIO-i40evf



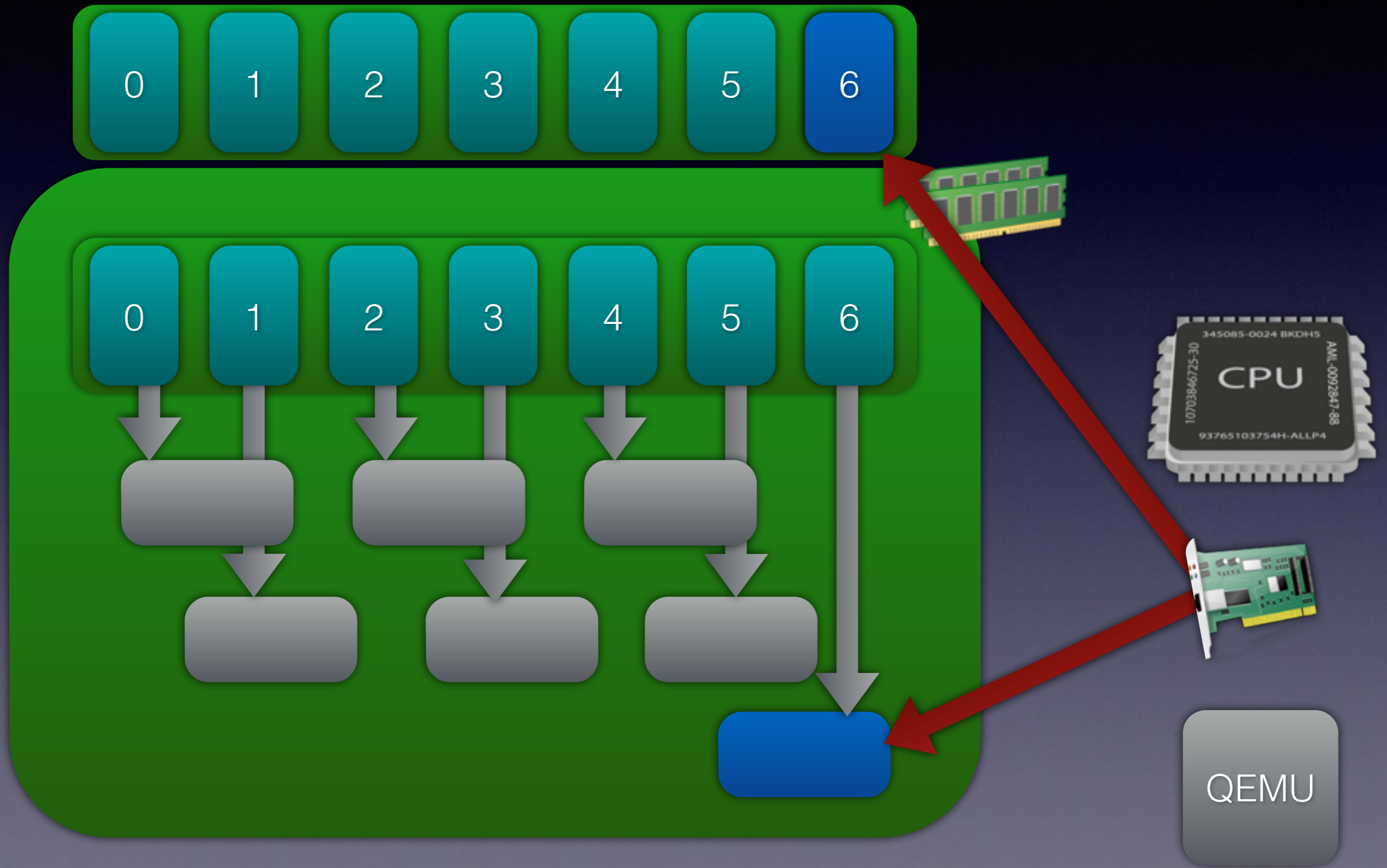
VFIO-i40evf



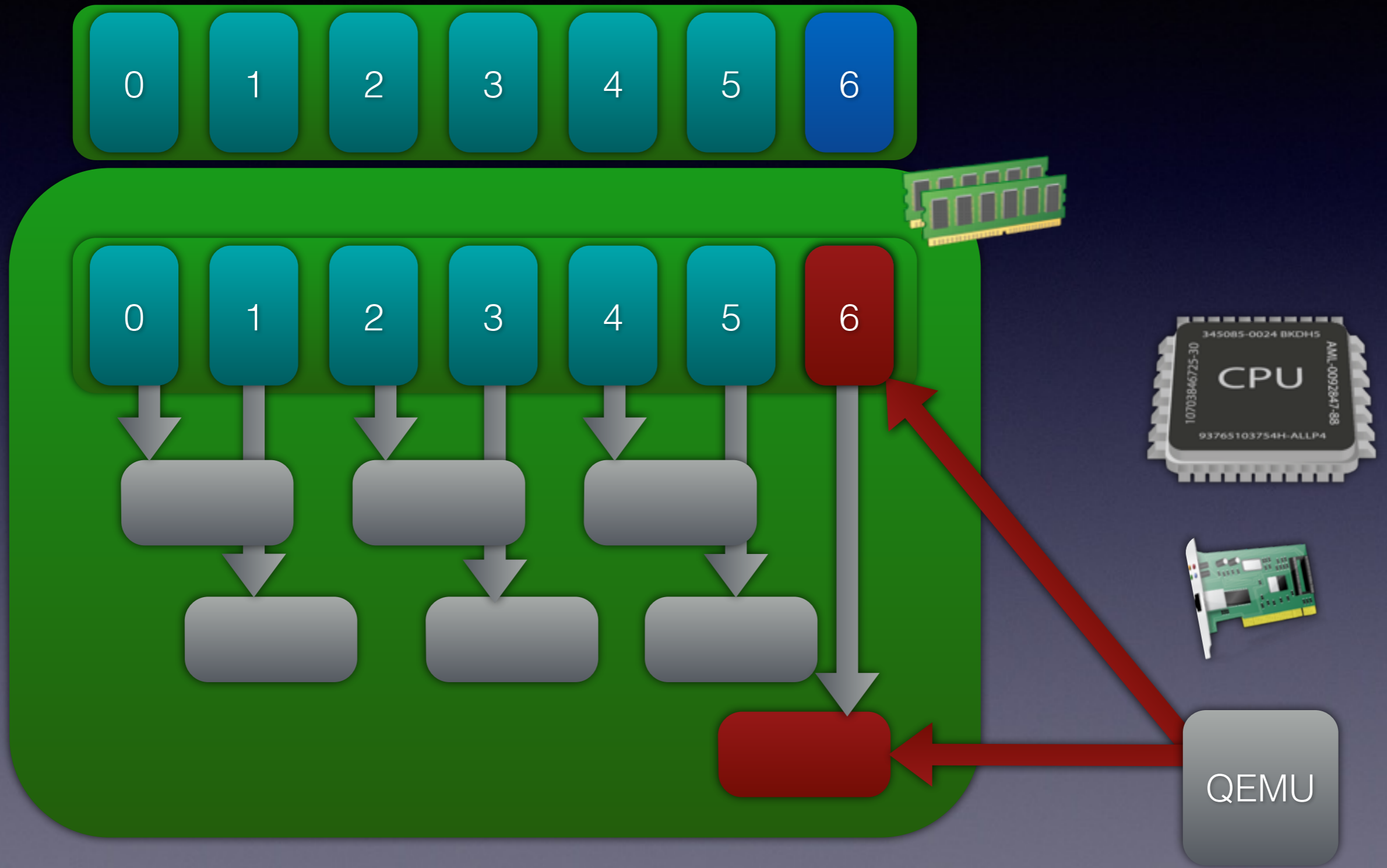
VFIO-i40evf



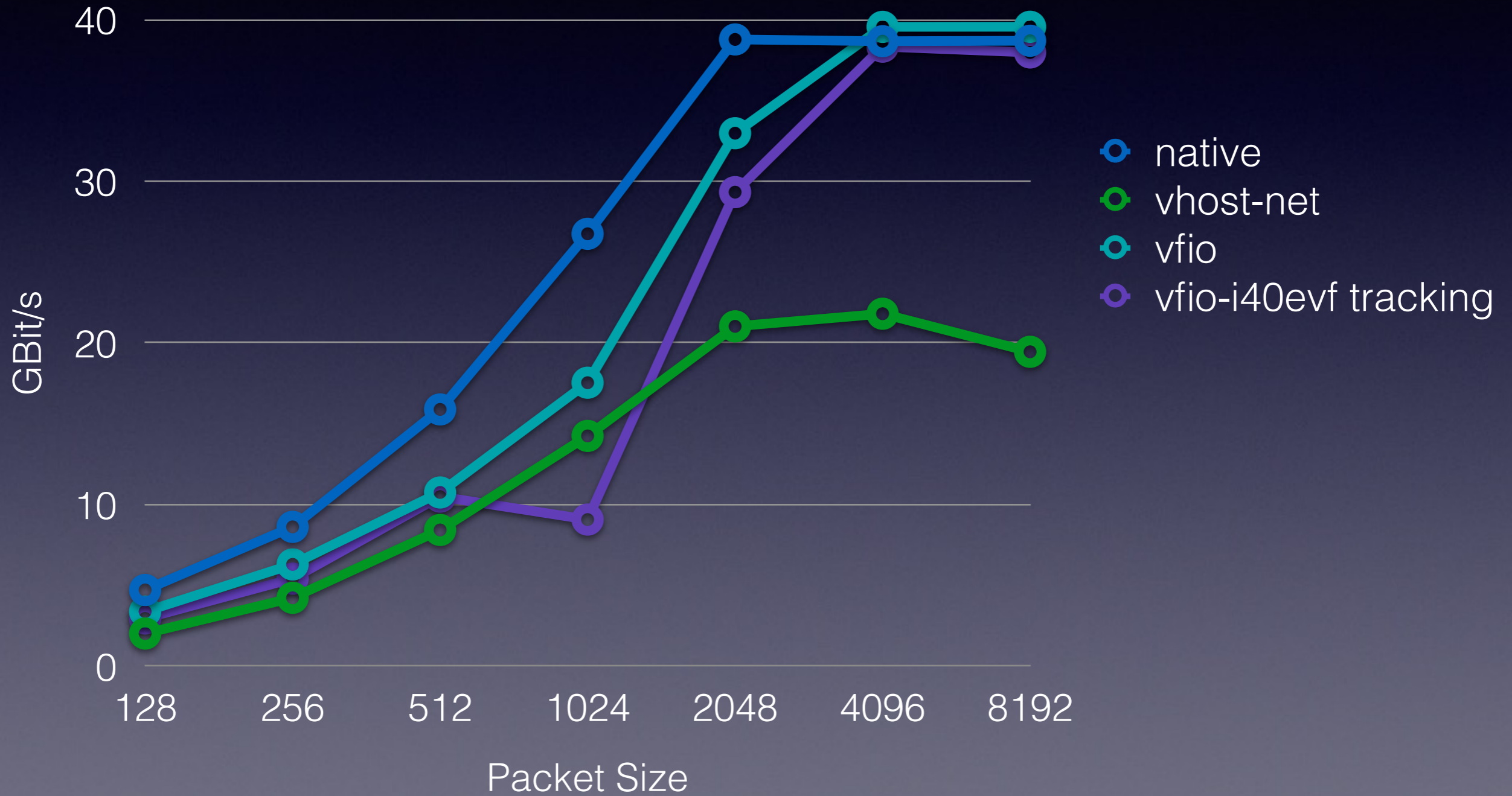
VFIO-i40evf



VFIO-i40evf



VFIO-i40evf

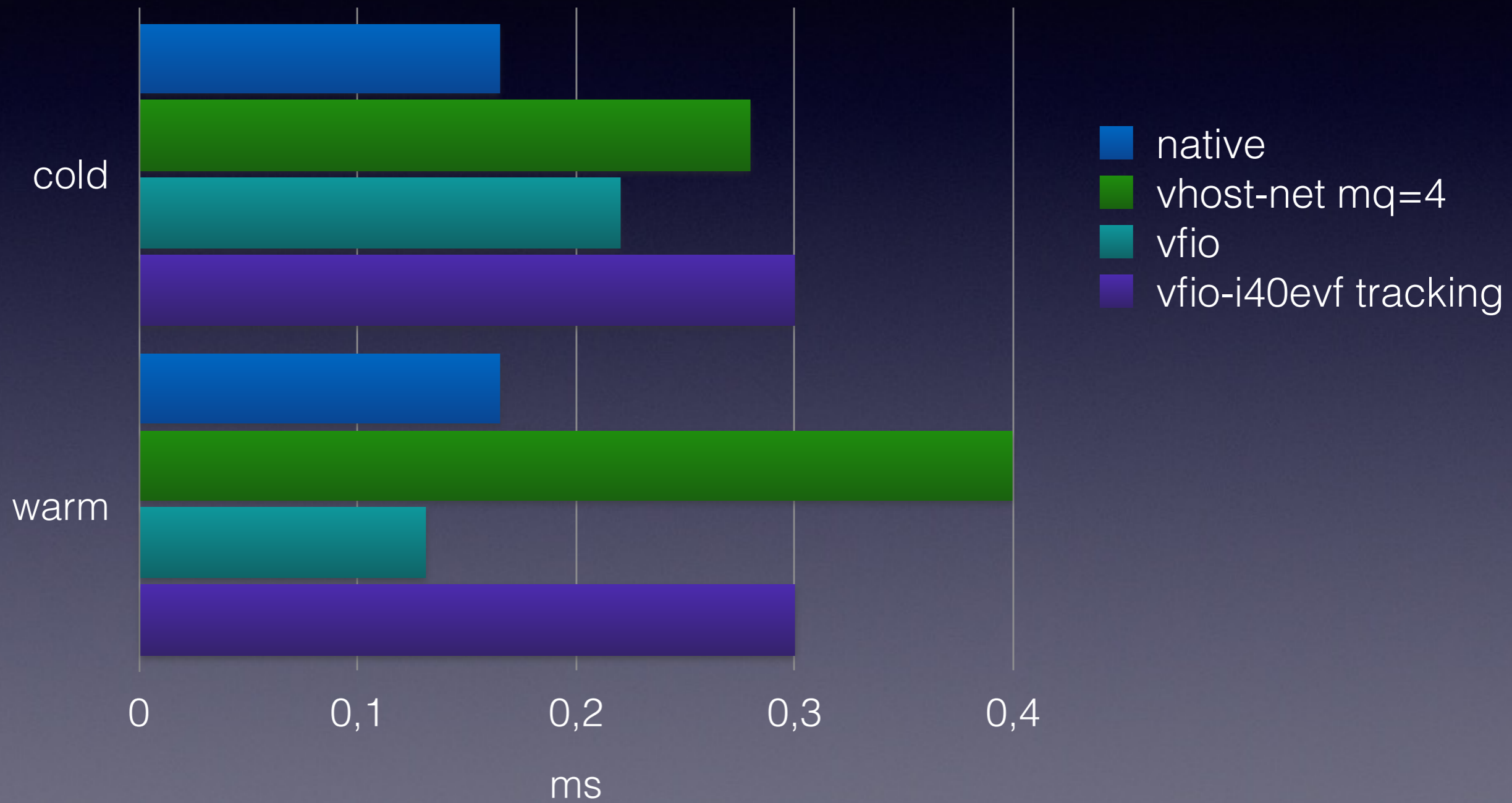


VFIO-i40evf

- Live Migration works
- No memory overcommit
- Performance identical to VFIO in normal operation
- Throughput close to VFIO during migration

Latency

Latency

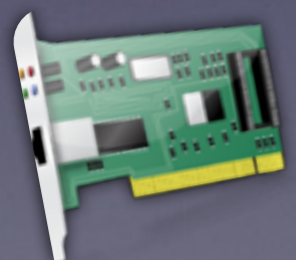
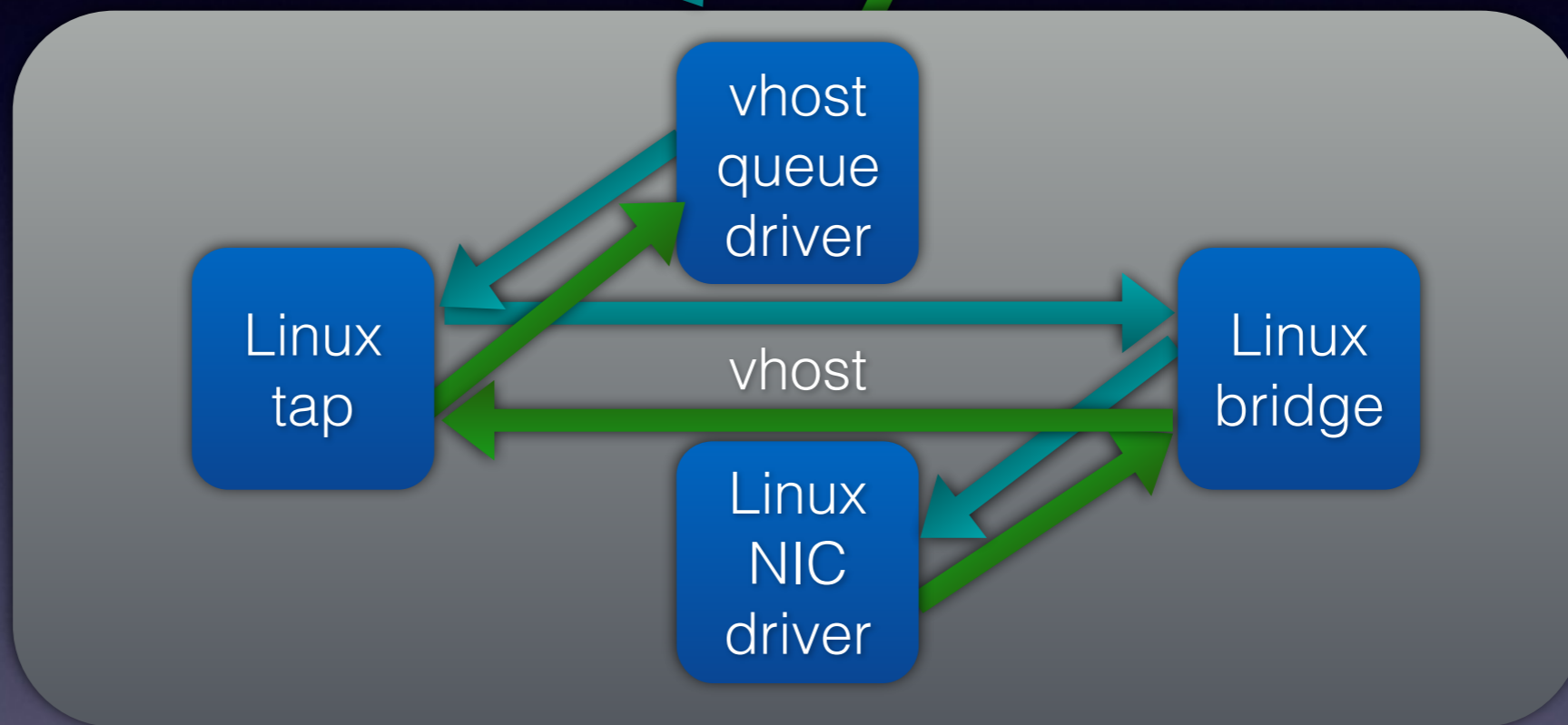
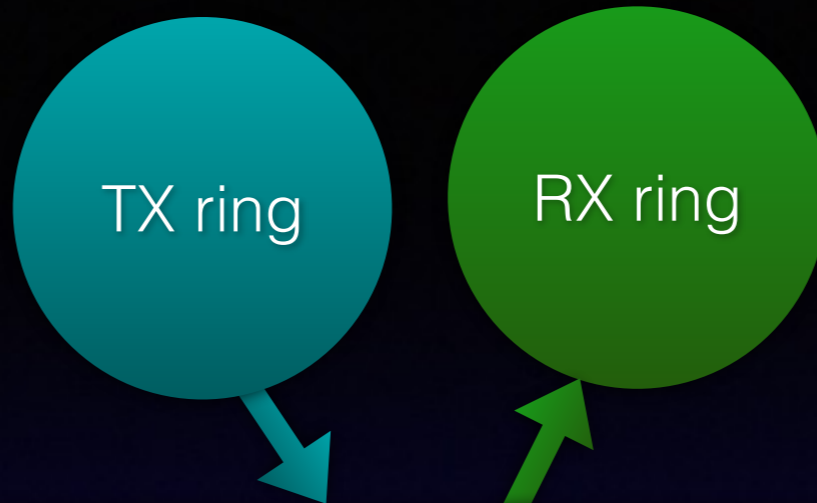


Other ideas

Hardware aware vhost



virtio-net

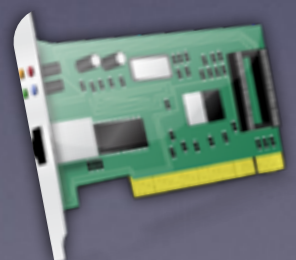
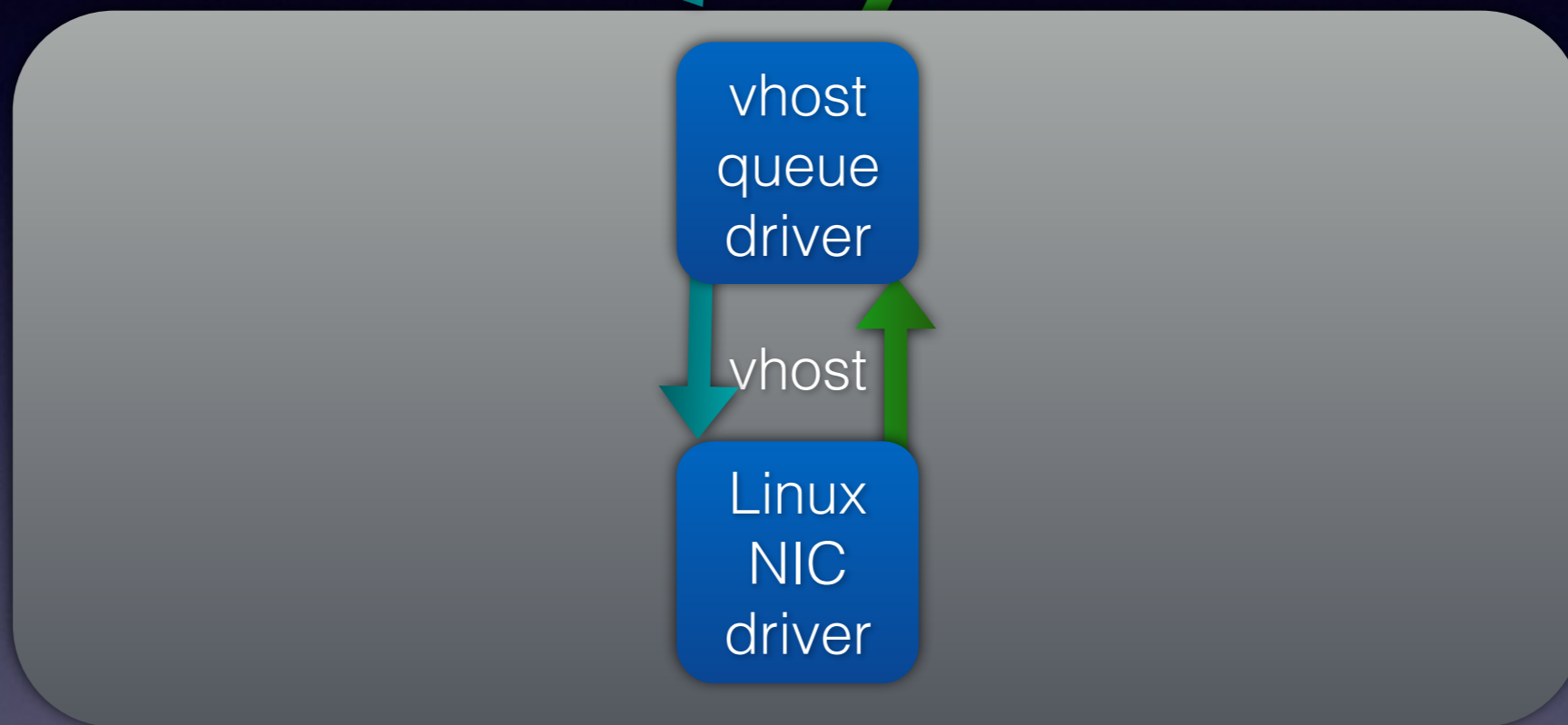
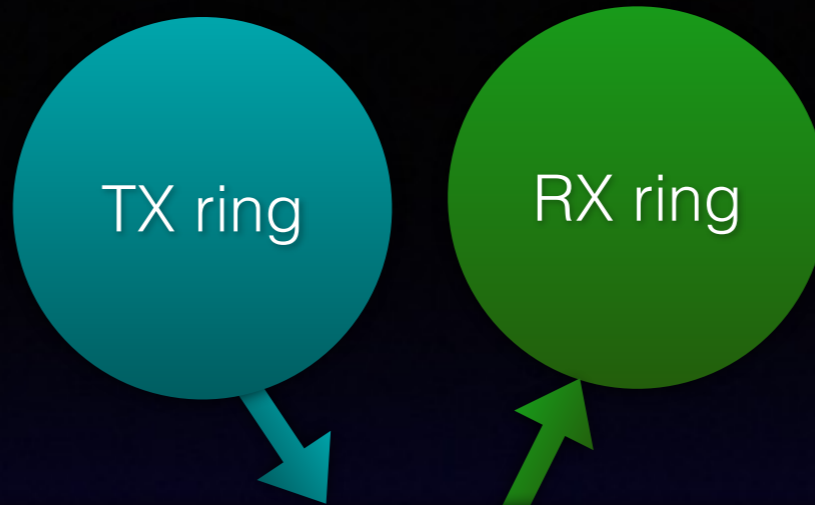


real NIC





virtio-net



real NIC

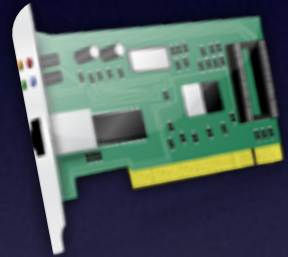


Hardware aware vhost

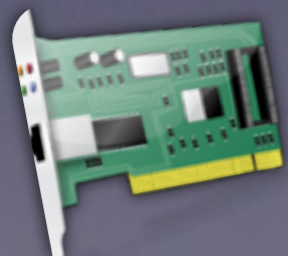
- Vhost talks directly to NIC driver
- NIC driver allocates queues for Vhost
- No guest exposure to real NIC

Virtio Queue Format Extension

Virtio Queue Format Extension



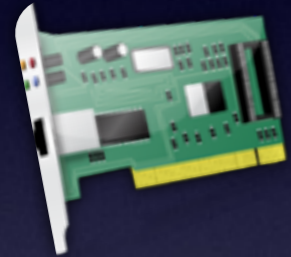
virtio-net



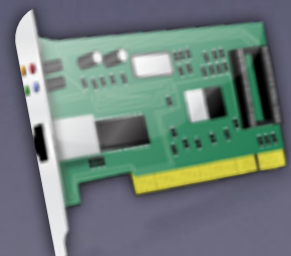
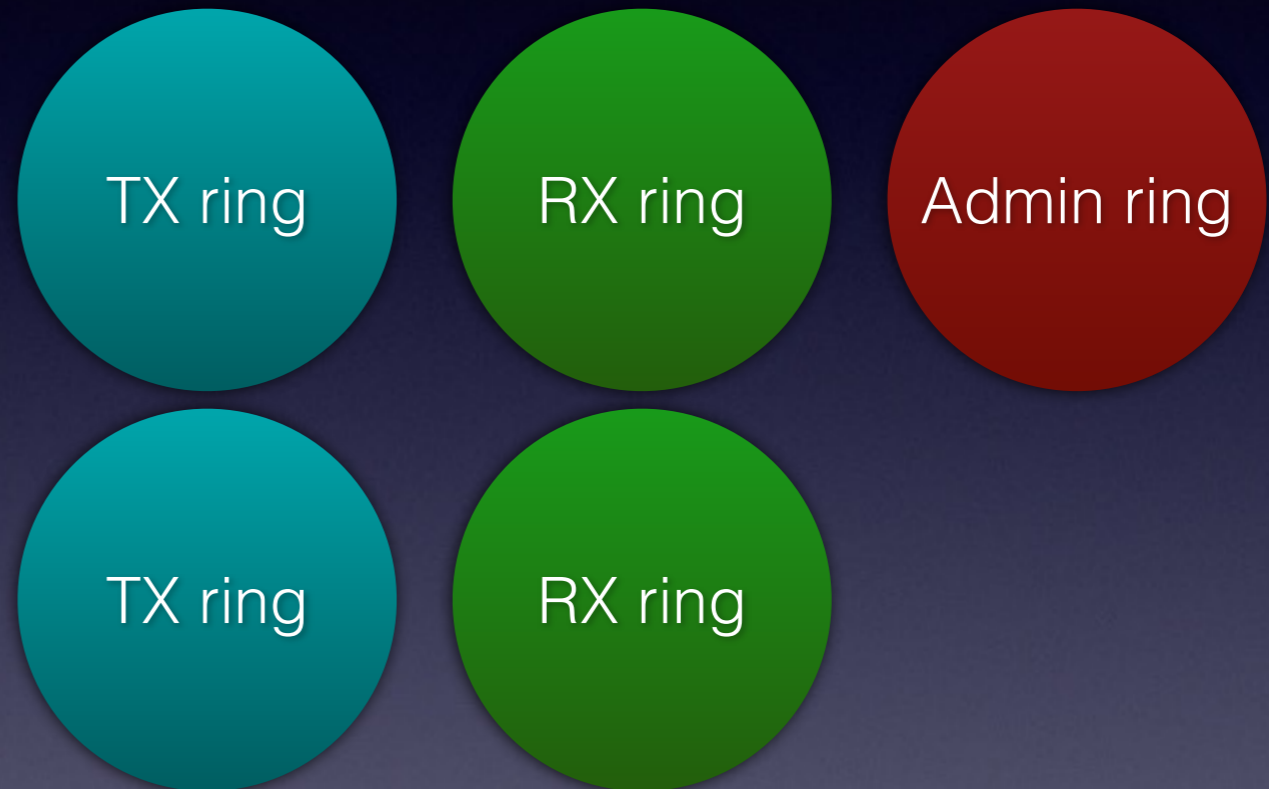
real NIC



Virtio Queue Format Extension



virtio-net



real NIC

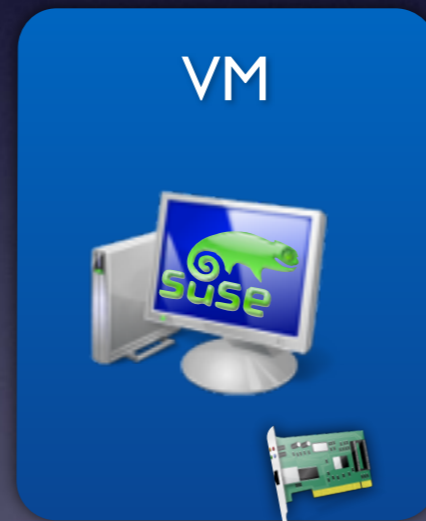


Virtio Queue Format Extension

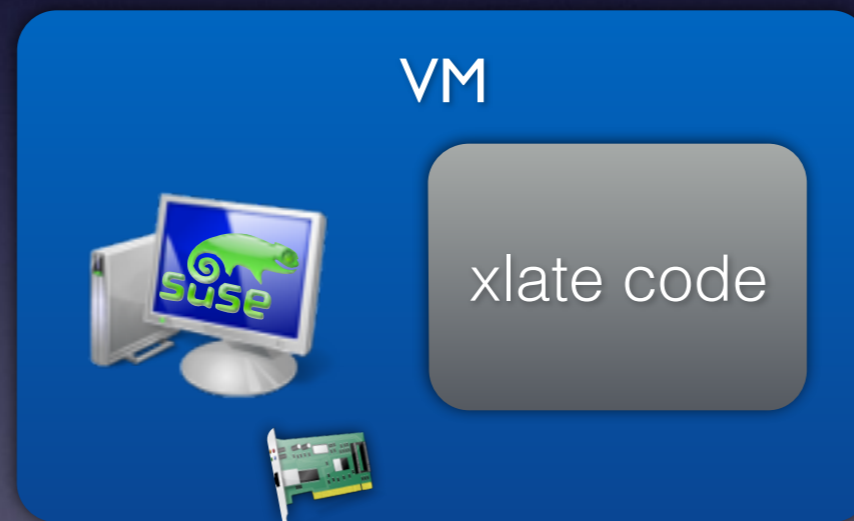
- Guest sees virtio device with additional BAR
- New BAR provides direct access to HW queues

Guest injected driver

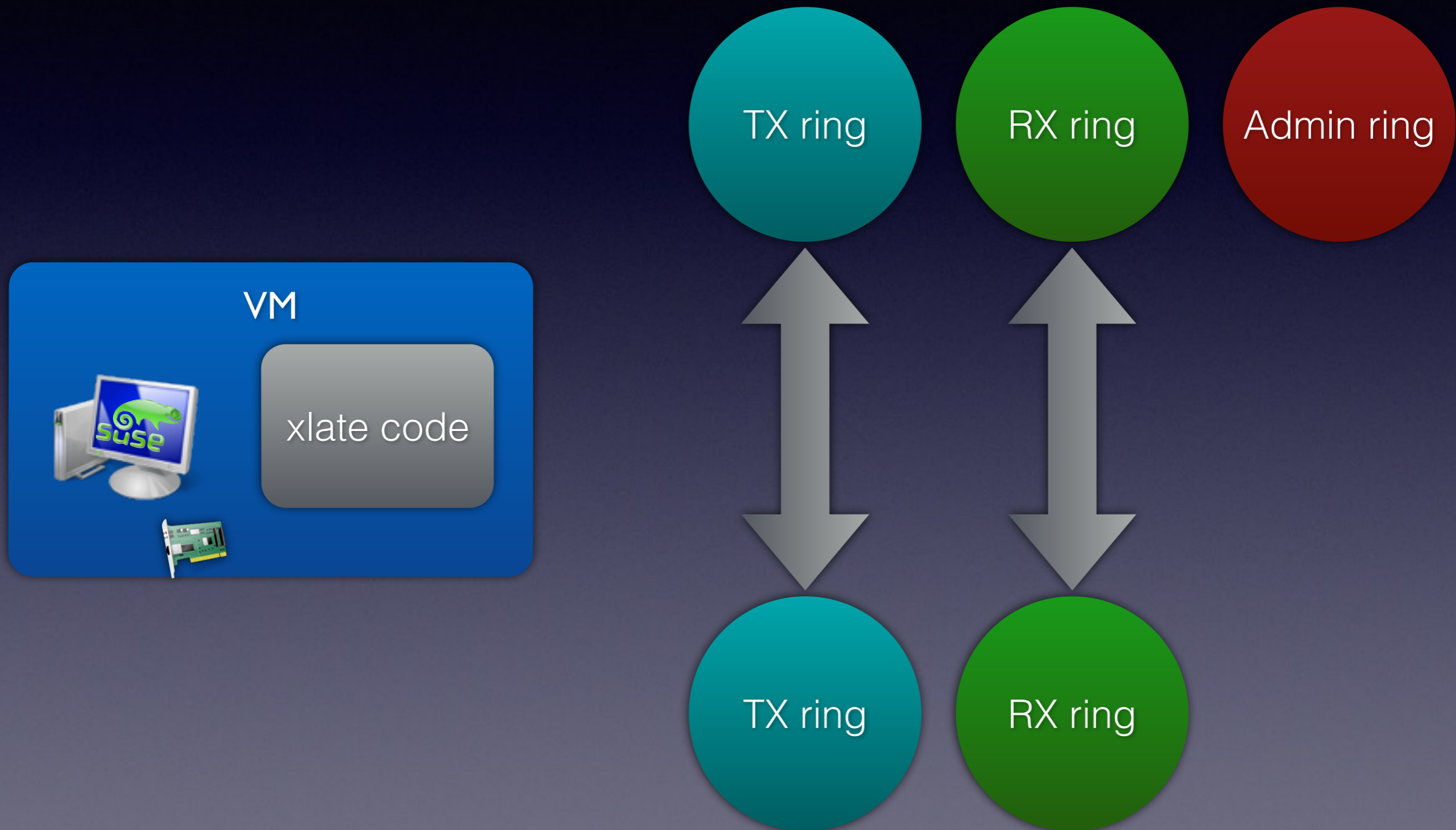
Guest injected driver



Guest injected driver



Guest injected driver

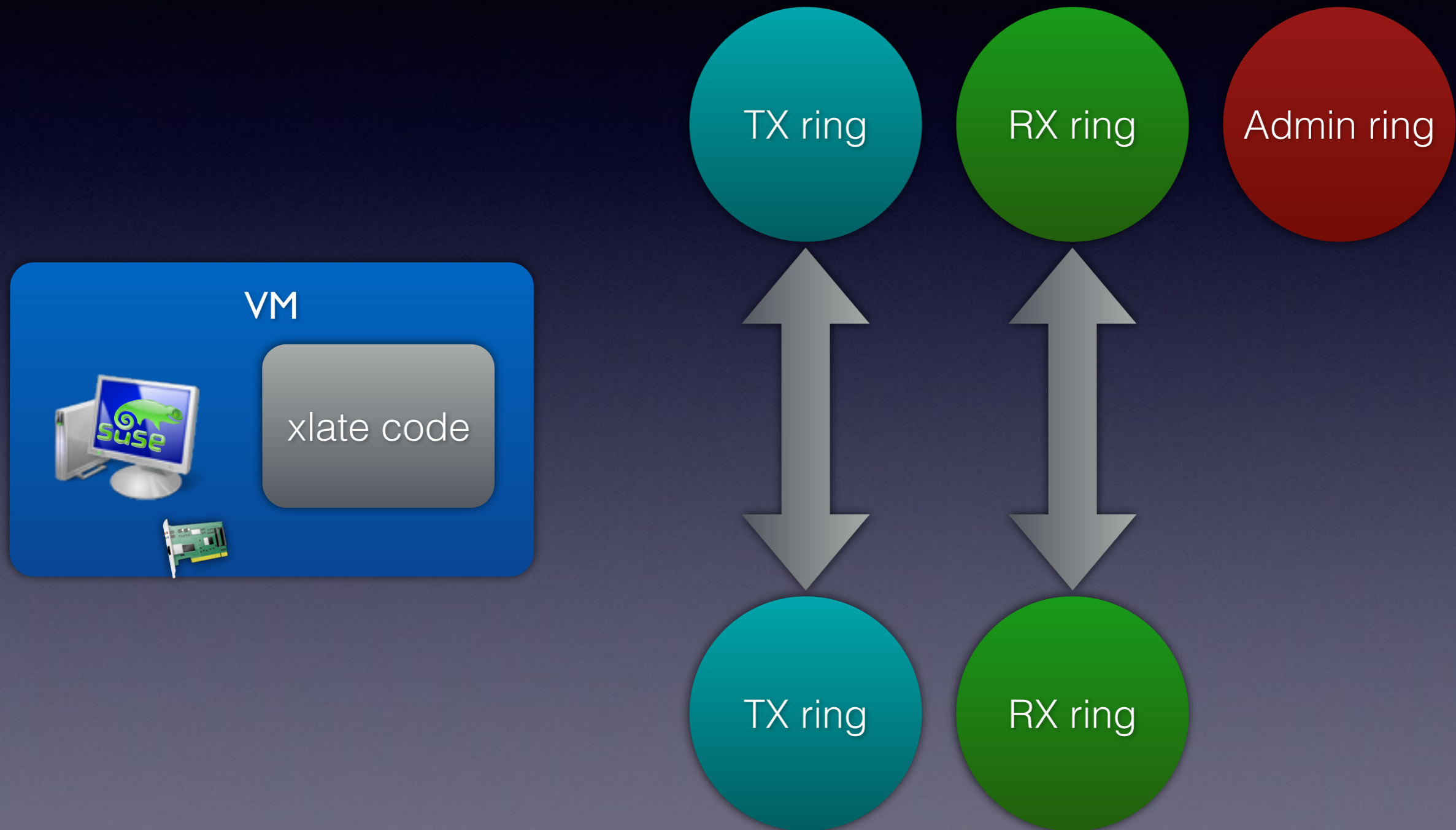


Guest injected driver

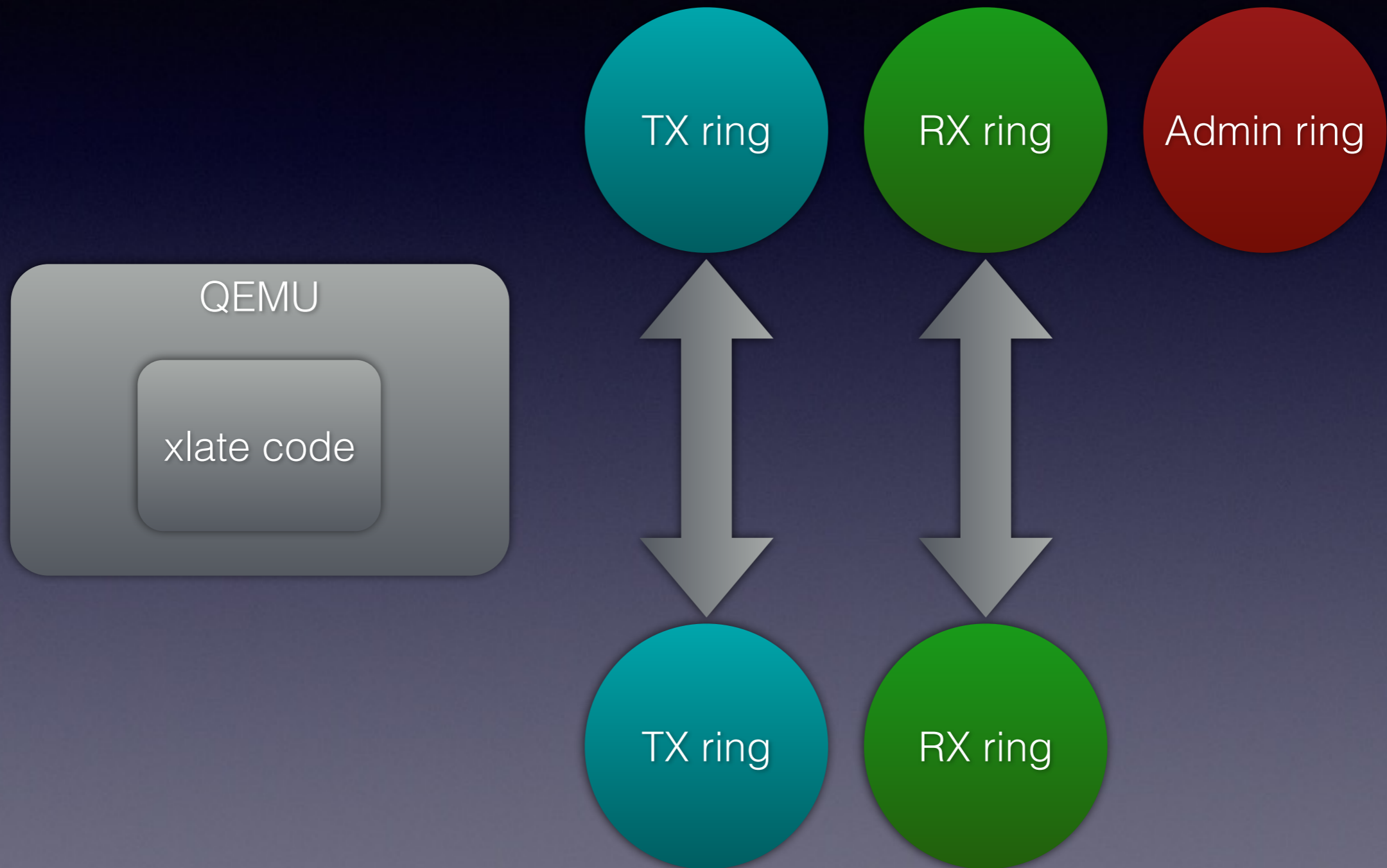
- Guest sees extended virtio device
- Code to convert hw queues to virtio queues comes from hypervisor
- Benefits:
 - Less code in privileged context
 - Faster IRQ path

QEMU based virtio
queue conversion

QEMU based virtio queue conversion



QEMU based virtio queue conversion



QEMU based virtio queue conversion

- QEMU allocates real hardware
- QEMU has specific HW/virtio conversion driver
- Guest only sees virtio
- Probably bad for latency

Conclusions

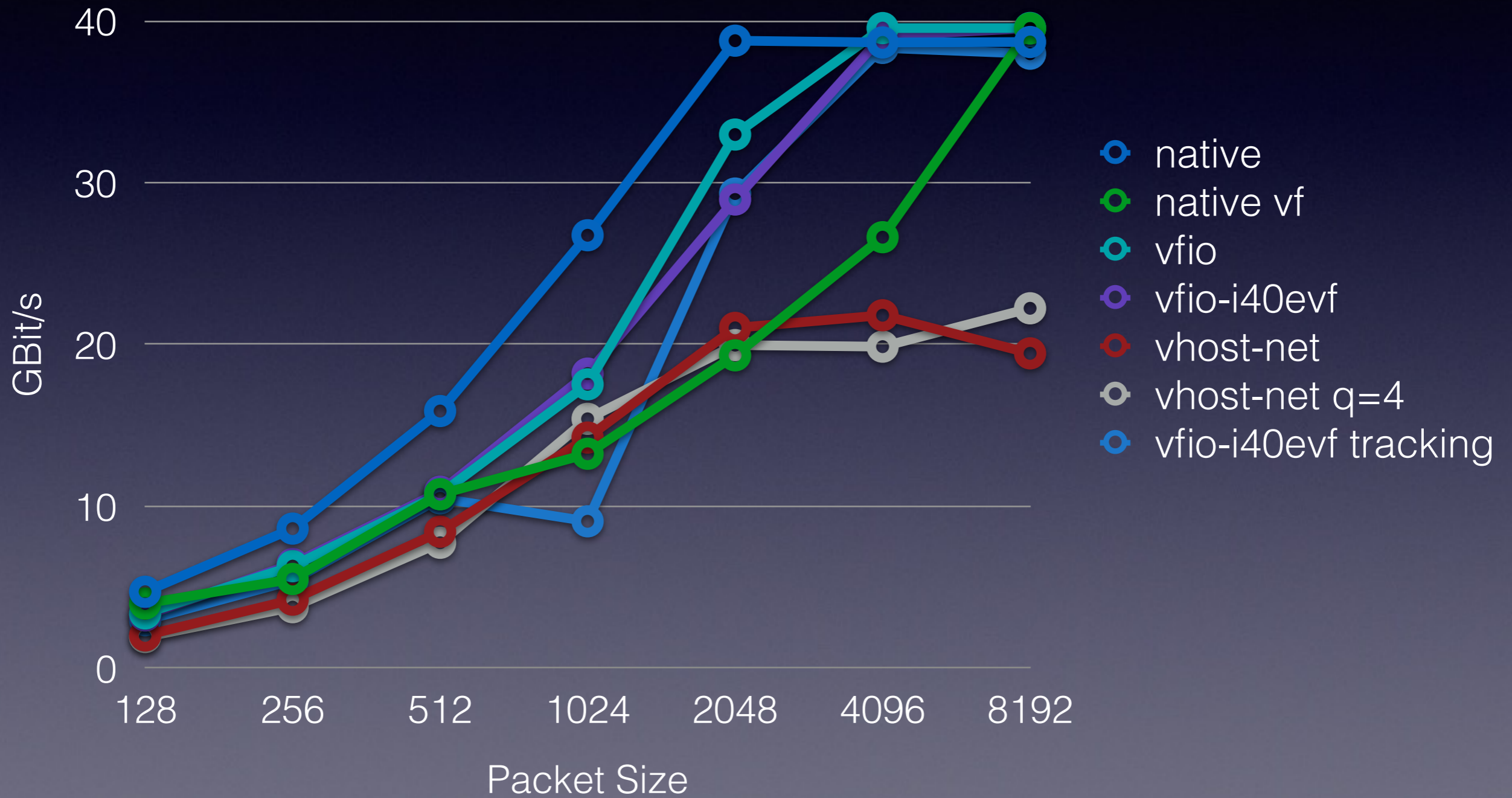
- Virtio is fast enough for most cases
- Room for improvement exists
- Live migration of real hardware works

Thank You

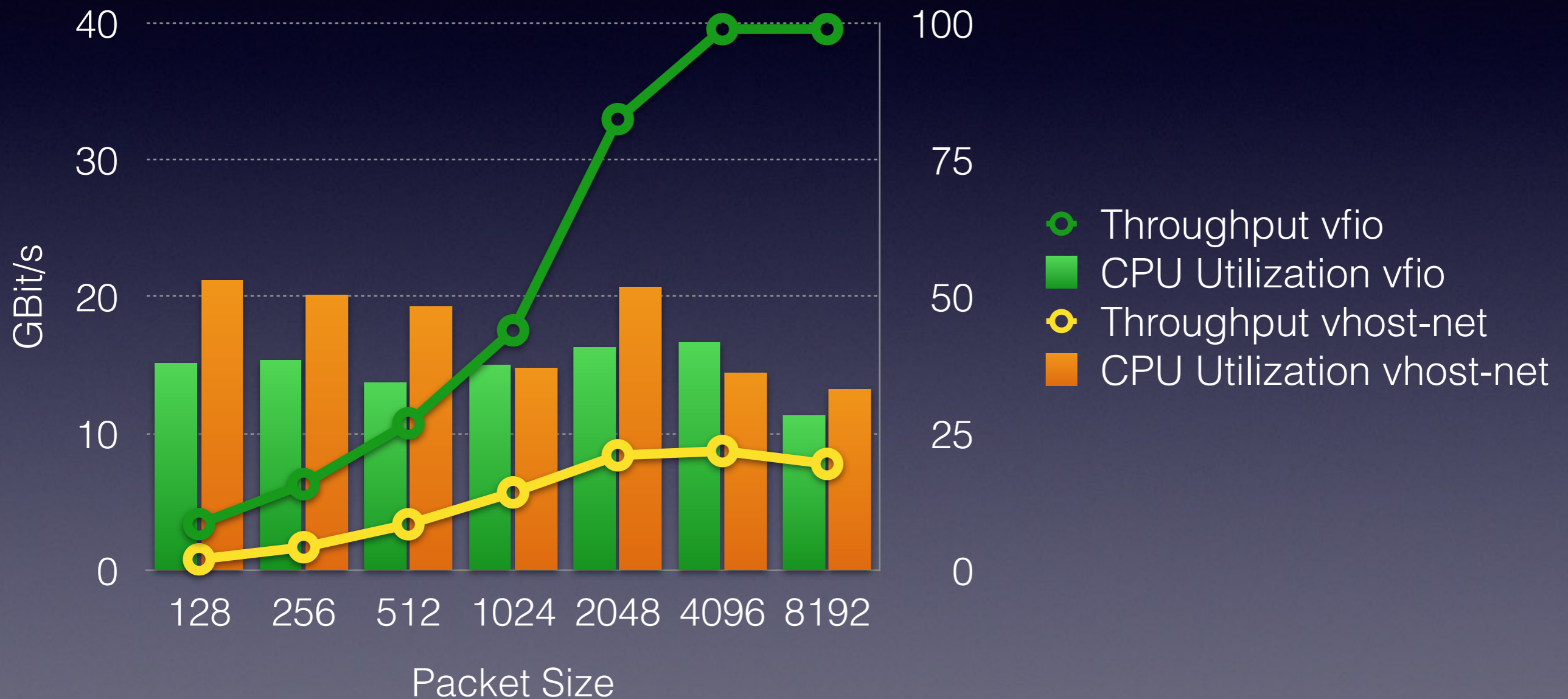
[git://github.com/agraf/qemu.git](https://github.com/agraf/qemu.git) vfio-i40vf

Performance Backup

Throughput



Throughput vs CPU Utilization



Throughput vs Latency

