

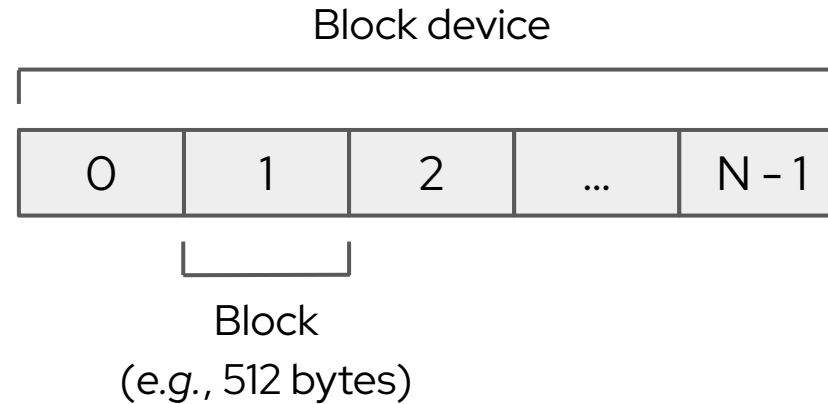
# libblkio

## Introducing the libblkio High-performance Block I/O API

Stefan Hajnoczi  
stefanha@redhat.com

Alberto Faria  
afaria@redhat.com

# Block devices



NVMe, SCSI, ATA, virtio-blk follow the block device model

**Read & Write** access data in units of blocks

**Flush** persists previously written data to permanent storage

**Discard** (TRIM) and **Write Zeroes** manage block allocation

# Where block device are used

## Databases

- MySQL & MariaDB

## File Systems

- Ceph Bluestore
- XFS, btrfs, etc

## Emulators & Hypervisors

- QEMU

## I/O Frameworks & Software-defined Storage

- SPDK

## Backup, Forensics, & Disk Imaging Tools

- casync
- qemu-img
- mkfs

# How libblkio came about

QEMU accumulated non-trivial block drivers

- io\_uring, NVMe userspace driver, ...

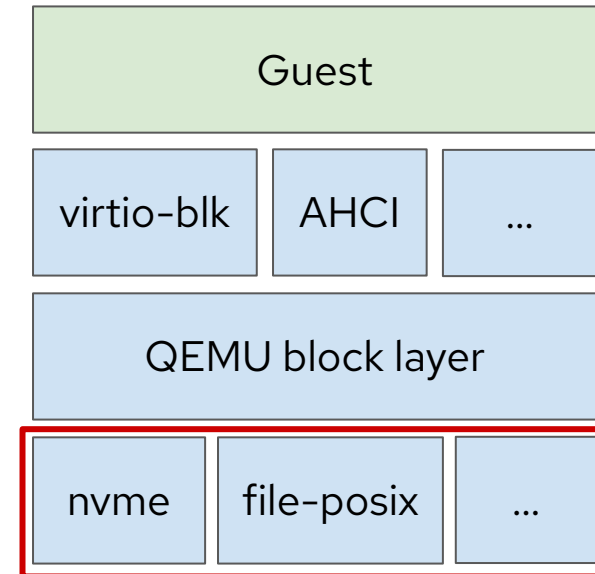
More drivers were needed:

- virtio-blk-vhost-vpda
- virtio-blk-vhost-user

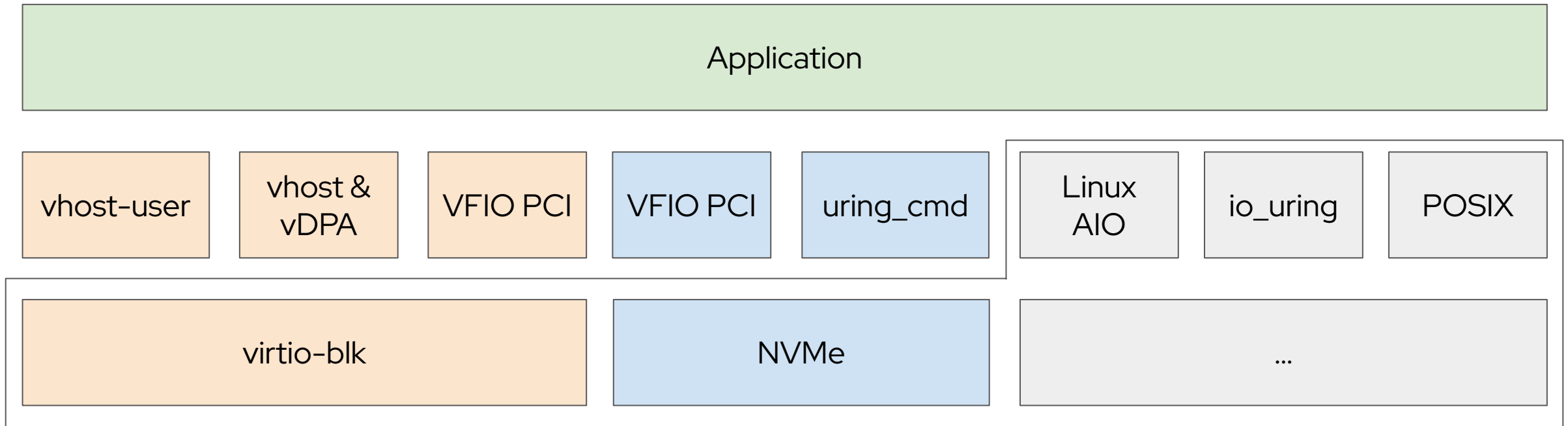
QEMU block drivers are only usable within QEMU

- Unable to reuse code in other programs

Decided to develop new drivers as a library instead



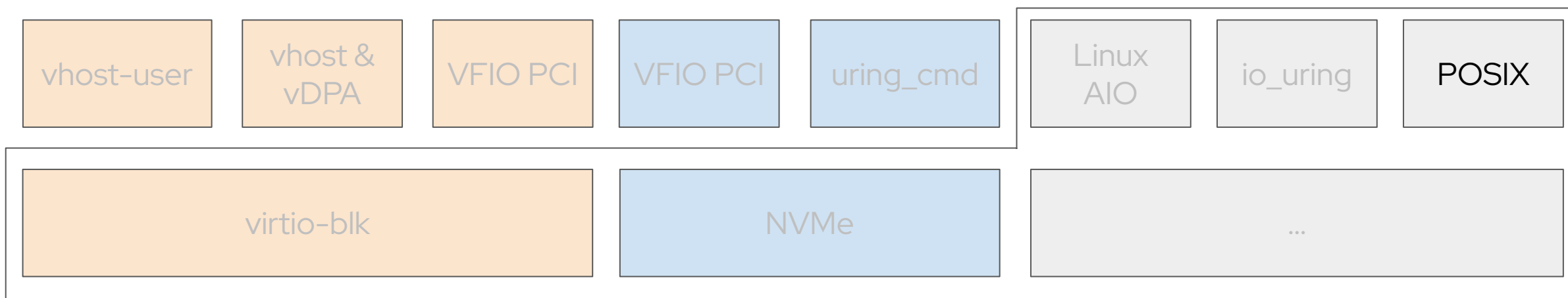
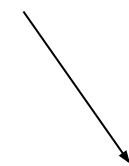
# Block I/O interfaces have proliferated!



How many can your application support?

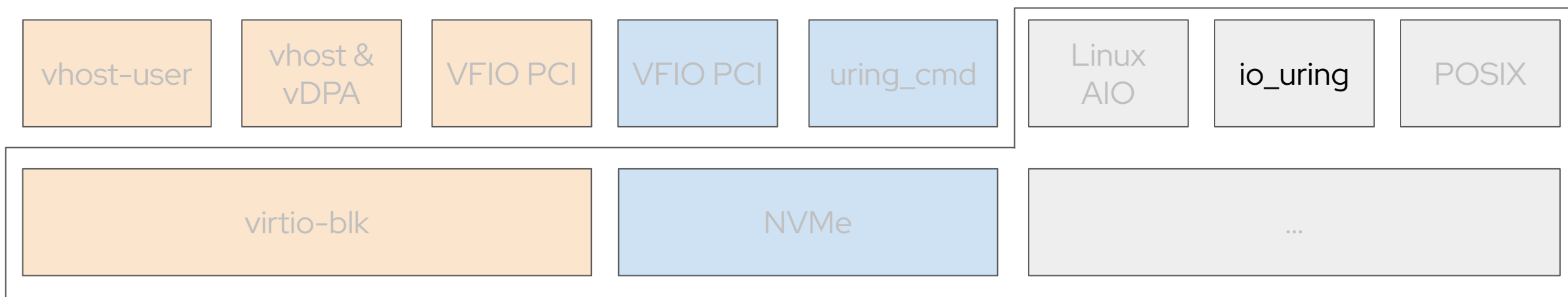
# Why use different interfaces?

**POSIX** is simple and device-agnostic



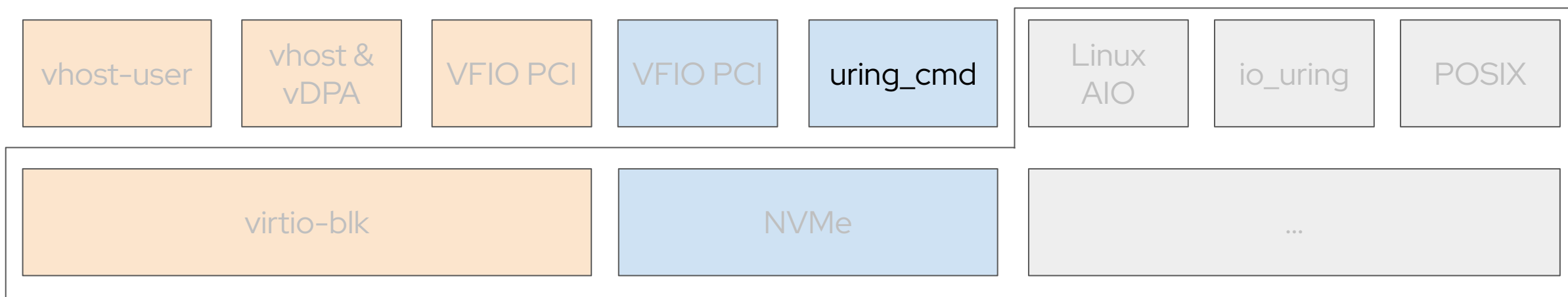
# Why use different interfaces?

**io\_uring** is asynchronous, less syscall overhead



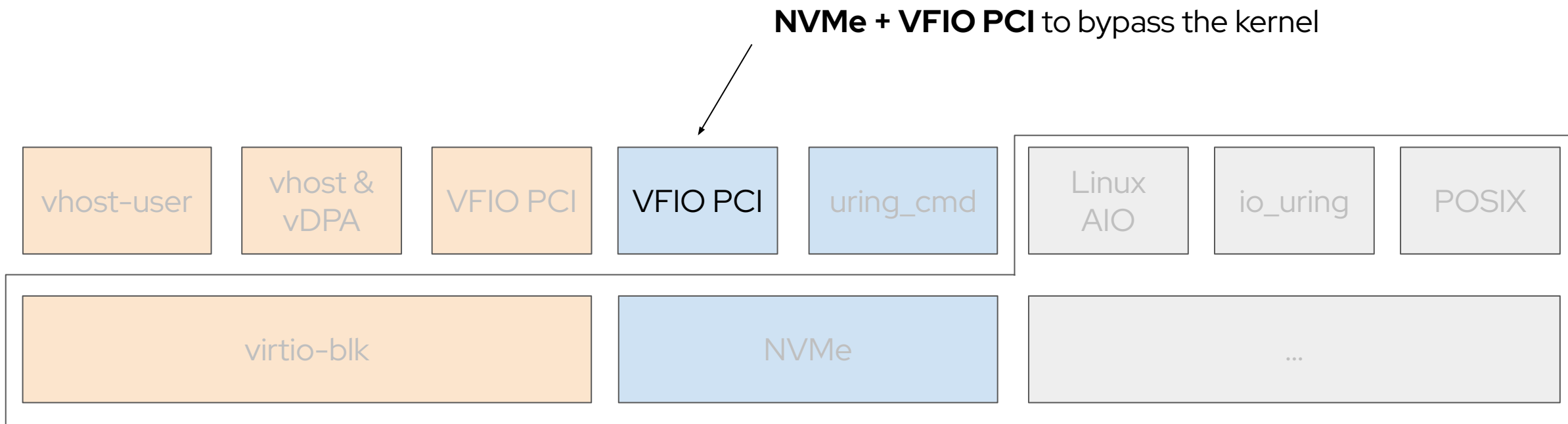
# Why use different interfaces?

**NVMe + uring\_cmd** to bypass the VFS



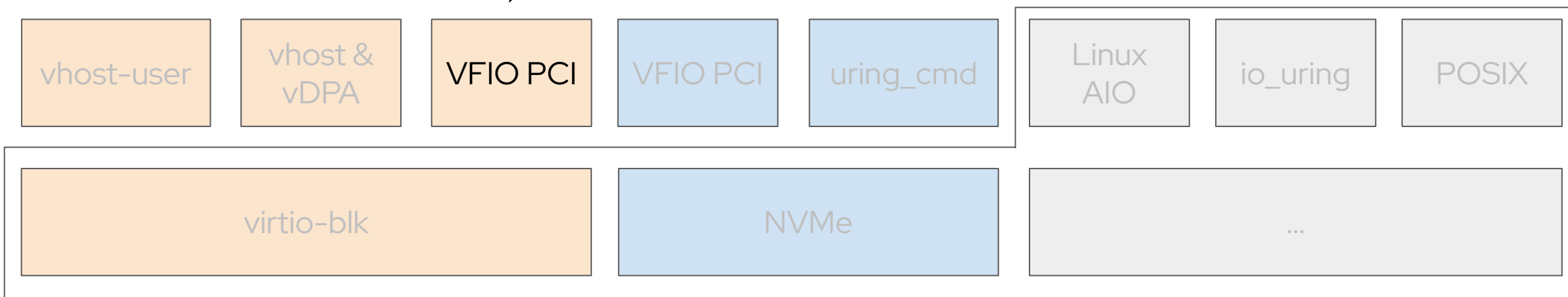


# Why use different interfaces?



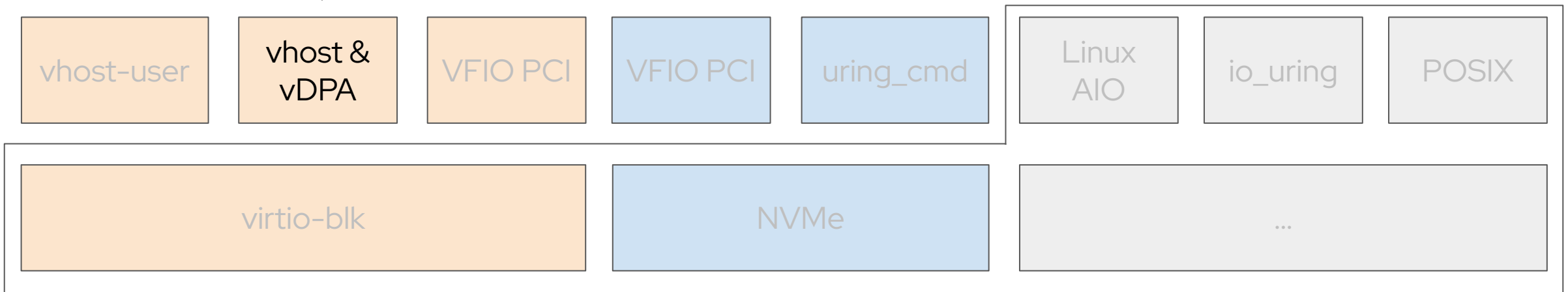
# Why use different interfaces?

**virtio-blk + VFIO PCI** to bypass  
the guest kernel's virtio-blk driver



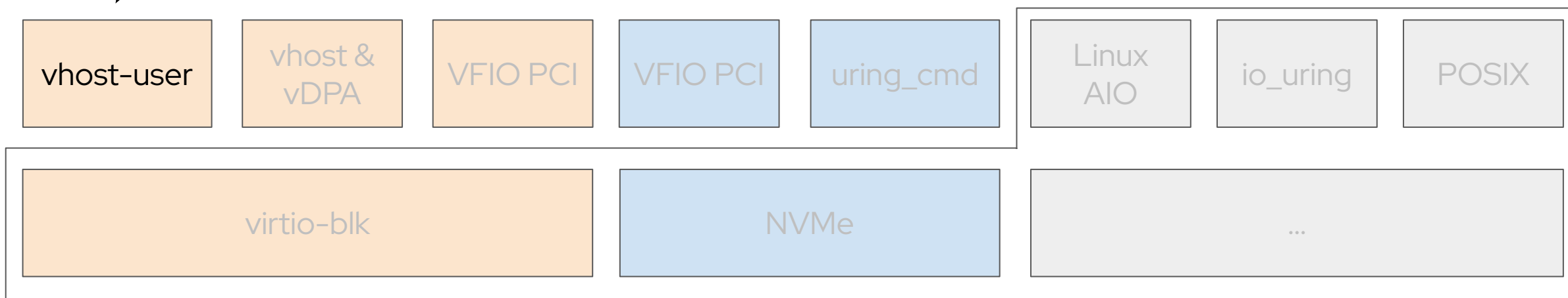
# Why use different interfaces?

**virtio-blk + vhost/vDPA** to access  
a slice of a physical virtio-blk device



# Why use different interfaces?

**virtio-blk + vhost-user** to communicate with a user-space storage server (e.g., SPDK, qemu-storage-daemon)



# Are block I/O interfaces similar?

## Same:

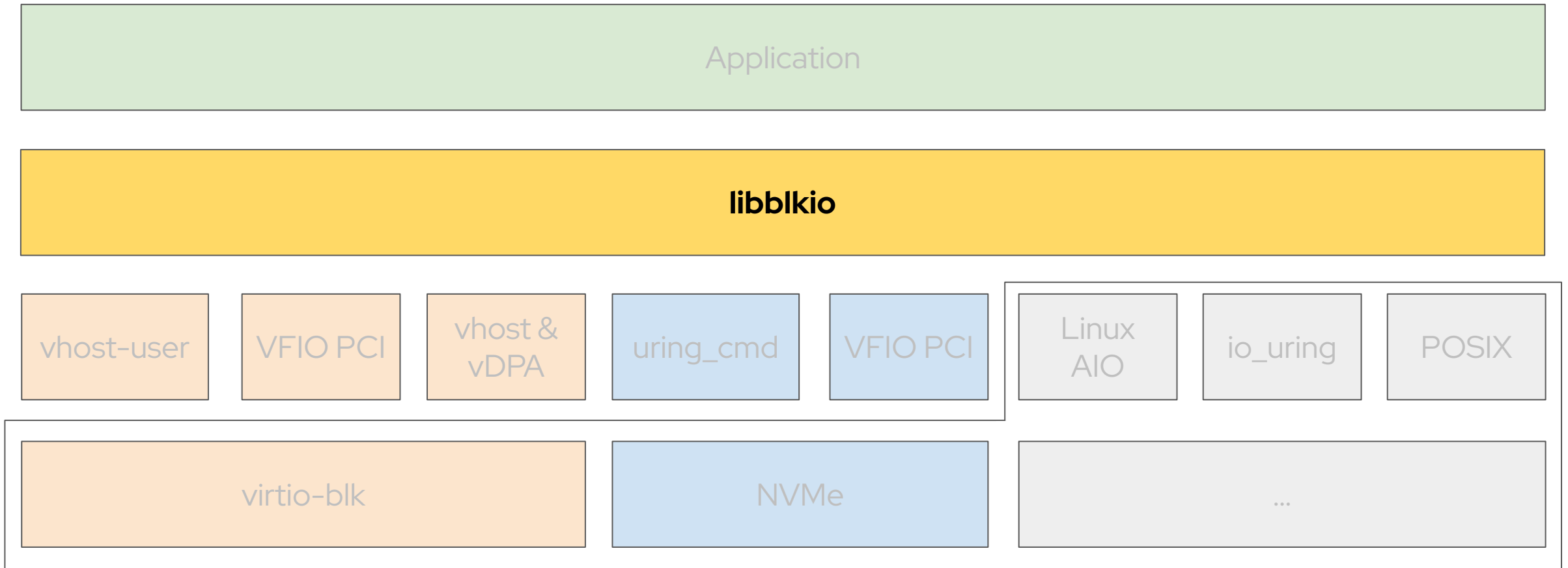
- ▶ Read, write, flush, discard, write zeroes
- ▶ Requests
- ▶ Queues

## Different:

- ▶ Synchronous or asynchronous
- ▶ Polling support
- ▶ I/O buffer memory constraints
- ▶ Queue memory layout
- ▶ Exact semantics and API details
- ▶ ...

How much effort is it to integrate a new interface into your application?

# libblkio provides a unified block I/O interface



## Where **libblkio** can be used

### Databases

- MySQL & MariaDB

### File Systems

- Ceph Bluestore
- XFS, btrfs, etc

### Emulators & Hypervisors

- QEMU

### I/O Frameworks & Software-defined Storage

- SPDK

### Backup, Forensics, & Disk Imaging Tools

- casync
- qemu-img
- mkfs

# libblkio

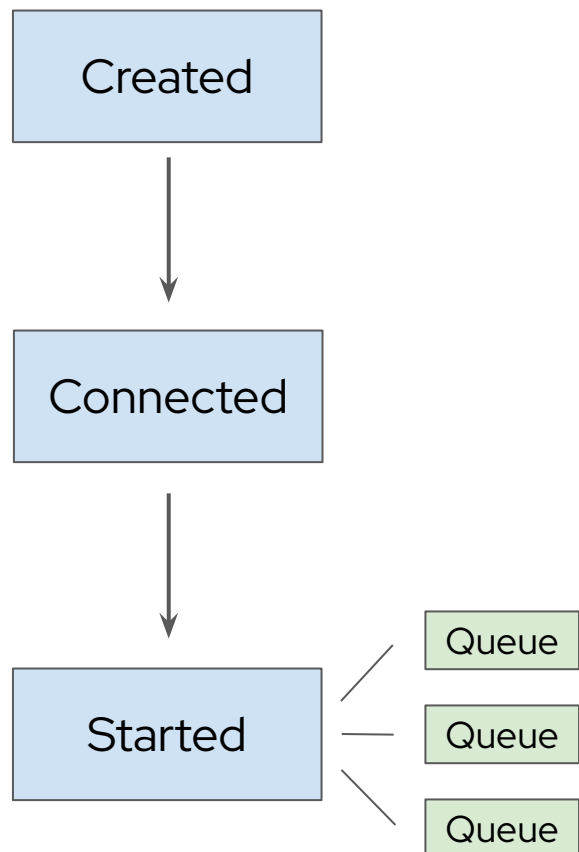
- ▶ C API (implemented in Rust)
- ▶ Provides several “**drivers**”
  - Each driver talks to a different underlying block I/O interface
- ▶ All drivers provide the **same API**
  - No code changes necessary to use a different driver



# libblkio

- ▶ Includes drivers for all environments
  - *"io\_uring"*, useful on **bare metal**, in **containers**, in **VMs**
  - *"virtio-blk-vfio-pci"*, useful to bypass guest kernel drivers in **VMs**
  - *"virtio-blk-vhost-vdpa"*, useful to split a physical virtio-blk device into many virtualized devices, e.g., for **containerized** environments
  - ...

# Lifecycle

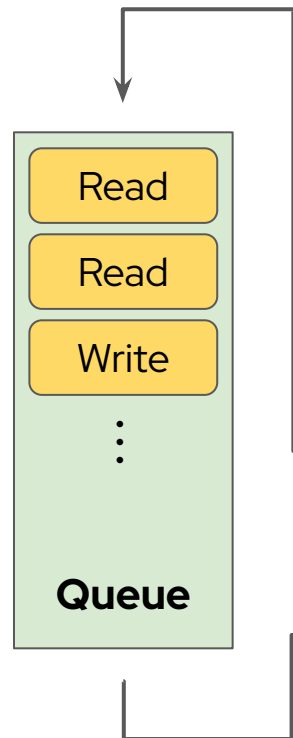


1. `blkio_create("io_uring")`

2. `blkio_set_str("path", "disk.img")`  
`blkio_connect()`

3. `blkio_set_int("num-queues", 3)`  
`blkio_start()`

# Queues



- ▶ Queues are independent of each other
  - Commonly assigned to a thread/core
- ▶ Workflow:
  1. Enqueue requests
  2. Submit requests and await completions
  3. Process completions

# Queues

- ▶ No request **ordering**
  - User must await completion before submitting next request to establish ordering
- ▶ Can **enqueue** any number of requests
  - But not all may be **in-flight** simultaneously
  - Drivers allow configuring this internal limit

# I/O modes

## Blocking I/O

1. Enqueue & submit requests
2. **Block** waiting for completions
3. Process completions

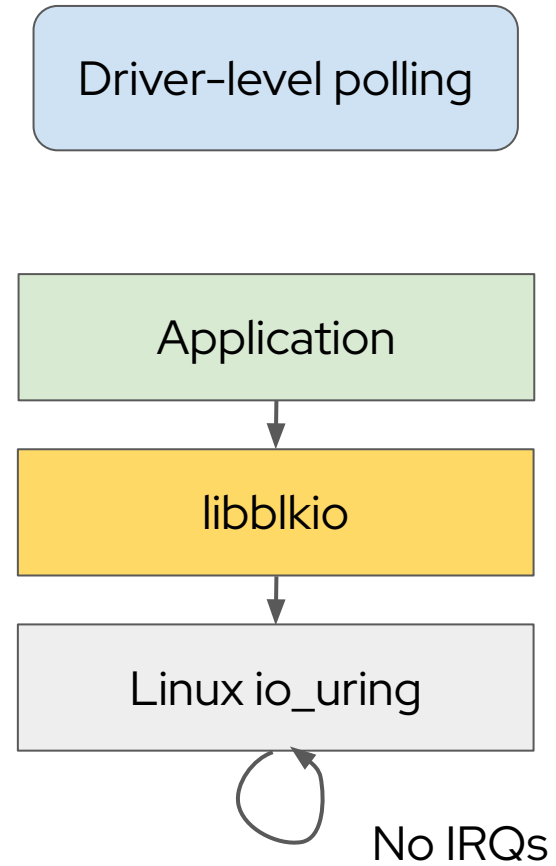
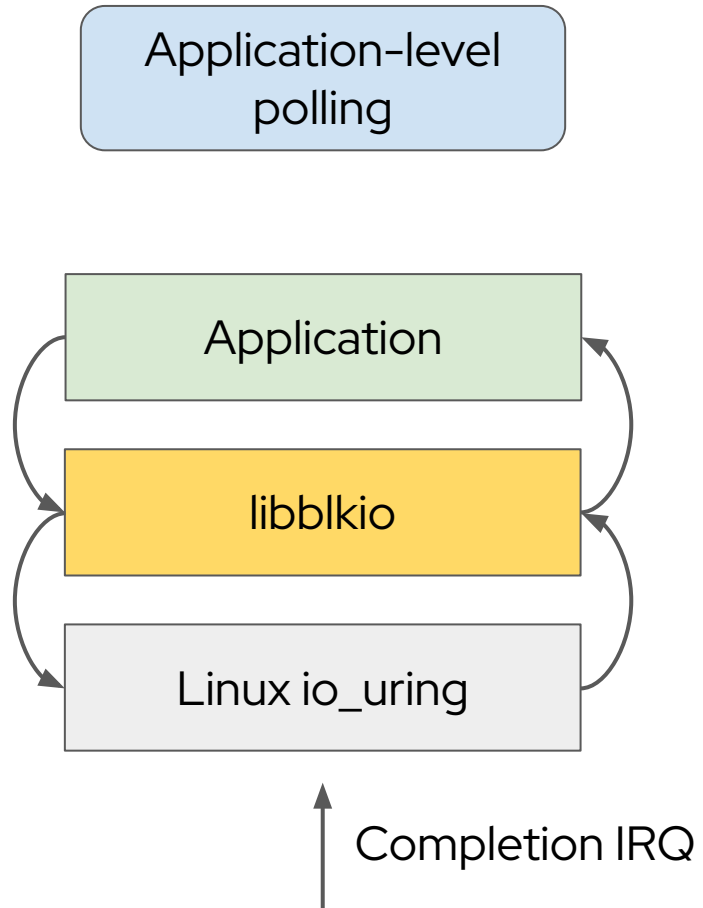
## Event-driven I/O

1. Enqueue & submit requests
2. Read/poll **eventfd**
3. Process completions

## Polled I/O

1. Enqueue & submit requests
2. **Loop** checking for completions
3. Process completions

# Polling modes

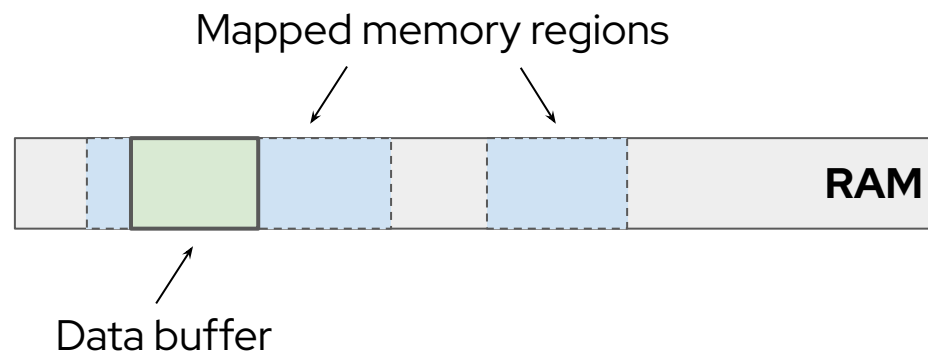


# Block limits and properties

- ▶ Devices/drivers may impose restrictions on requests
- ▶ Named **properties** expose this information
  - “max-transfer” – Maximum read/write request size
  - “buf-alignment” – Mandatory data buffer offset/size alignment
  - ...
- ▶ *Writable* properties are used for configuration
  - “path” – Device path for several drivers
  - “num-queues” – Number of queues to create
  - ...

# Memory regions

- ▶ Some drivers require **pre-registering** memory for data buffers
  - e.g., for establishing IOMMU mappings
- ▶ libblkio provides a unified **memory region** abstraction for this
  - Can be mapped/unmapped dynamically
  - Drivers may require data buffers to belong to mapped regions





# Memory regions

- ▶ There may be a limit on **concurrently** mapped regions
  - Mapping/unmapping might not be cheap
  - Ideally would map once and use many times
- ▶ Drivers may impose further restrictions on data buffers
  - Memory alignment, file descriptor-backed memory, ...
- ▶ Utilities for **allocating** suitable memory are provided

# Feature summary

- ▶ Unified, multi-queue block I/O API
  - ▶ Blocking I/O, event-driven I/O, polled I/O
    - Fits your app's I/O model
  - ▶ Properties
  - ▶ Memory regions
- ▶ Drivers are modular
    - Low integration effort
    - Can be contributed to <https://gitlab.com/libblkio/libblkio>

# Case study & evaluation

## Case study: libblkio in QEMU

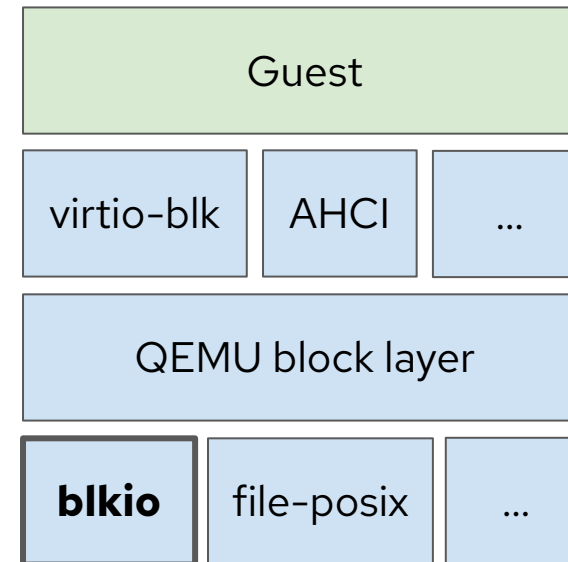
- ▶ QEMU (<https://www.qemu.org/>) is an emulator with a full block layer
- ▶ New QEMU block driver adds:
  - `-blockdev io_uring,filename=test.img,...`
  - `-blockdev virtio-blk-vhost-user,path=vhost-user-blk.sock,...`
  - `-blockdev virtio-blk-vhost-vdpa,path=/dev/vhost-vdpa-0,...`
- ▶ ~700 source lines of code (SLOC)
  - Applications typically need less code
- ▶ Expected in QEMU 7.2 release

# QEMU block drivers

Guests (VMs) submit I/O requests to *emulated storage controllers*.

The block layer hands requests to drivers.

Block driver integrates libblkio with QEMU.



# I/O buffer memory in QEMU

Guest  
RAM

Guest RAM is long-lived with occasional hotplug.

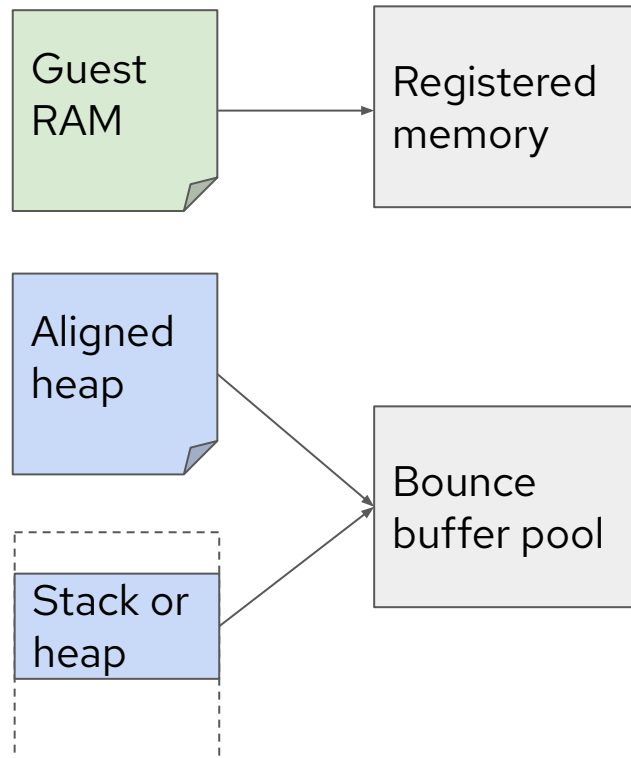
Aligned  
heap

Block drivers like crypto or qcow2 allocate internal I/O buffers.

Stack or  
heap

Small I/O from within QEMU may be arbitrary stack or heap memory.

# Mapping QEMU I/O buffers



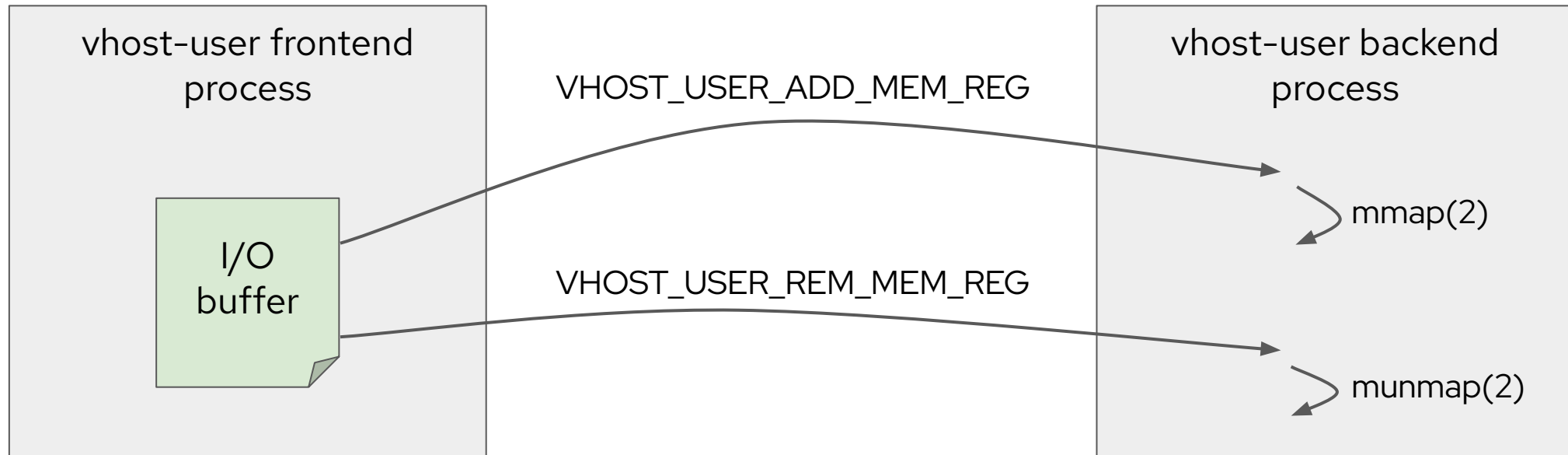
Guest RAM is permanently mapped to libblkio.

Bounce buffer pool is permanently mapped to libblkio

- Incurs copy overhead
- Alternative: temporary mappings?
- Alternative: intercept heap buffer allocation?

Existing applications need to make similar decisions about mapping I/O buffer memory.

## vhost-user map/unmap operation



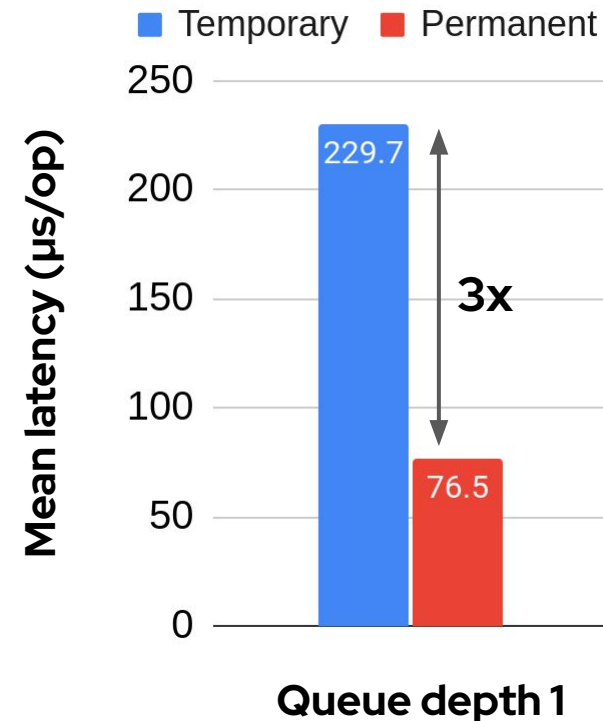
Mapping involves file descriptor passing over an AF\_UNIX socket.



# Measuring vhost-user map/unmap overhead

- ▶ Permanent mapping **vs** temporary mapping
  - 1 vCPU, 512-byte random reads
  - Mean latency (**lower is better**)
- ▶ qemu-storage-daemon vhost-user-blk export
- ▶ Out-of-the-box config without guest & vhost-user-blk polling

The I/O buffer mapping strategy is important!



# Real-world performance: QEMU with libblkio

How does libblkio's io\_uring driver compare against QEMU's io\_uring implementation?

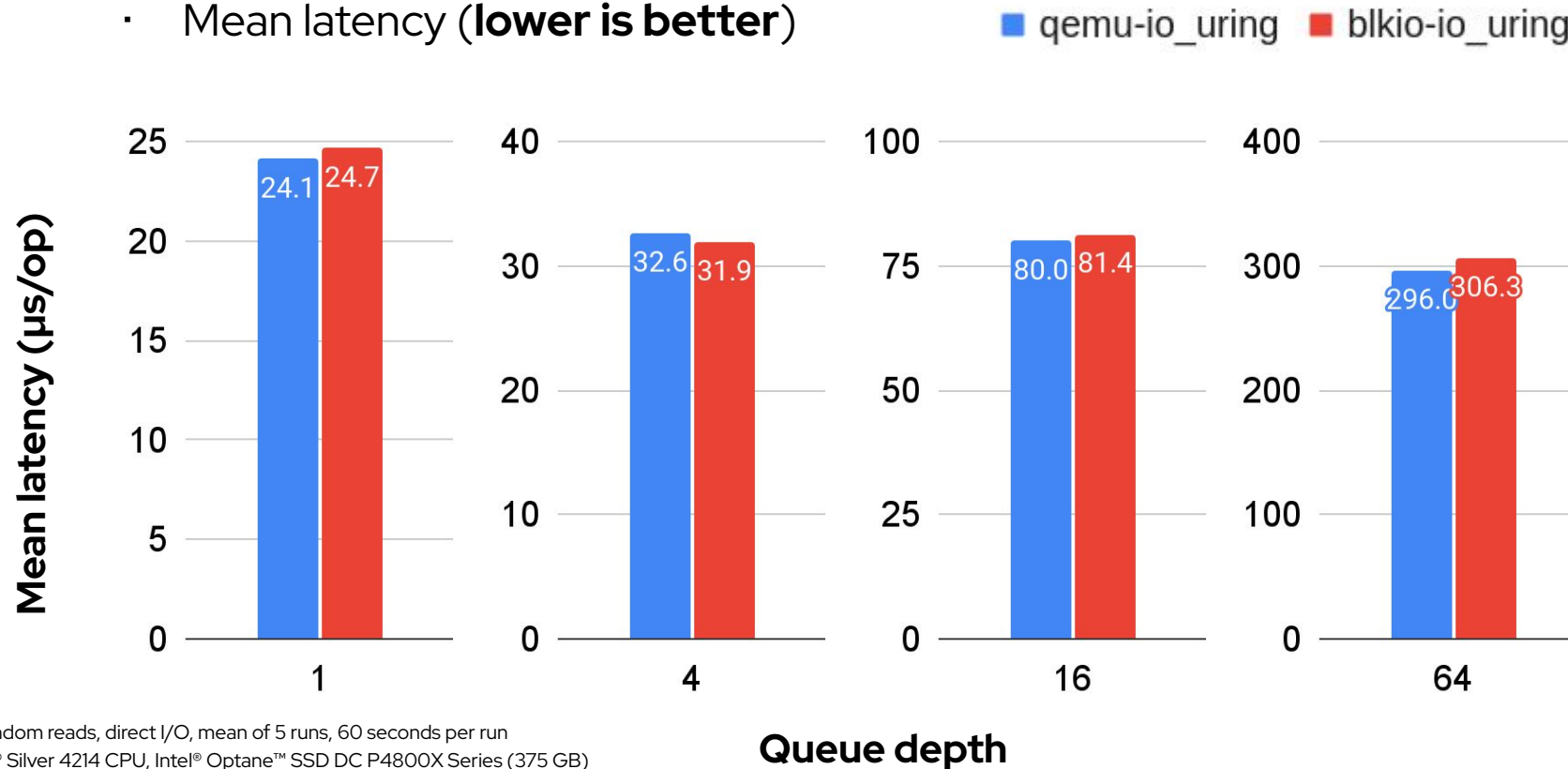
**QEMU** ▶ `-blockdev host_device,filename=/dev/nvme0n1,aio=io_uring,...`

VS

**libblkio** ▶ `-blockdev io_uring,filename=/dev/nvme0n1,...`

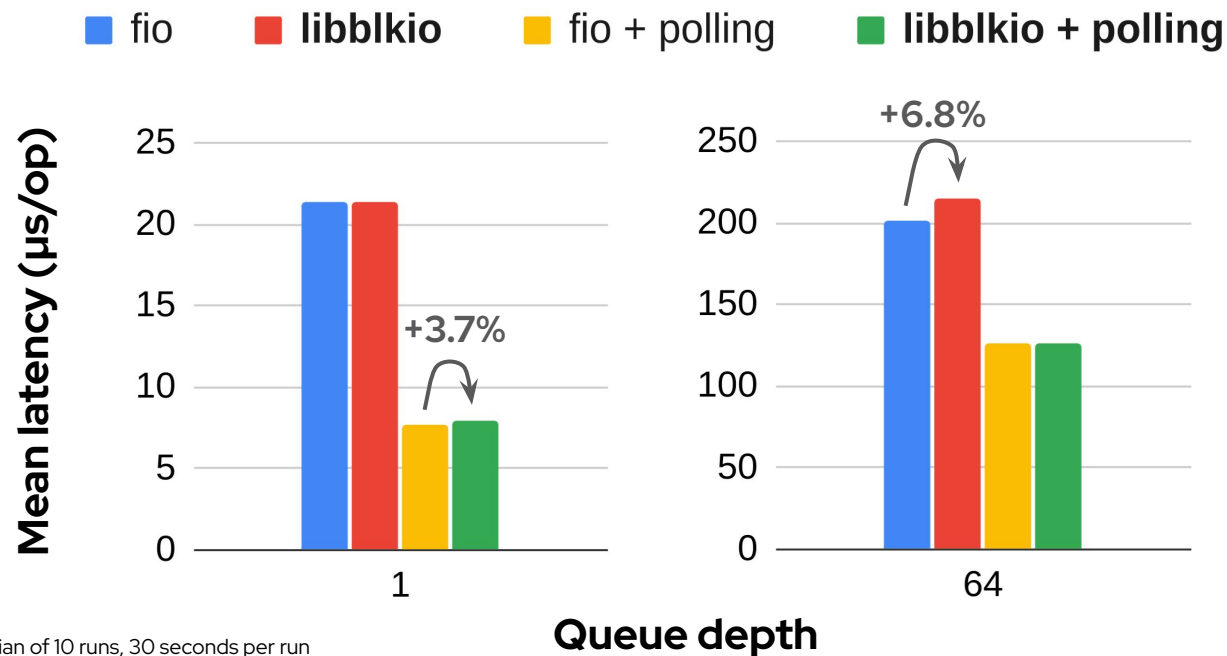
# QEMU io\_uring randread

- ▶ Native io\_uring **vs** libblkio io\_uring QEMU block drivers
  - 1 vCPU, 512-byte random reads
  - Mean latency (**lower is better**)



# fio micro-benchmarks

- ▶ fio's io\_uring engine **vs** libblkio fio engine using "io\_uring" driver
  - 1 core, 512-byte random reads, w/o and w/ driver-level polling (*i.e.*, poll queues)
  - Mean latency (**lower is better**)



# Future work

## Future work: Stable blkio Rust crate

Native Rust API should be idiomatic and safe.

Current API exposes MemoryRegion and raw iovec pointers

- ▶ Better abstractions needed for safety

How to manage request lifetime across do\_io() loop?

- ▶ Caller is trusted to keep the iovecs alive for the duration of the request

Hanna Reitz has been experimenting with the blkio crate in Rust applications.

# Future work: Network storage

NVMe over TCP, NBD, iSCSI, etc could be added.

Data path APIs are asynchronous but control path APIs are synchronous.

May need async control path API to avoid hangs.

## Future work: Queue passthrough

Emulators like QEMU present a storage device like virtio-blk to the guest.

They emulate the storage controller and invoke corresponding libblkio APIs.

When the underlying device is the same type as the guest device, passing through the queues bypasses the emulator and libblkio for better performance.

Currently in development for virtio-blk devices by Stefano Garzarella.



# Thank you

Red Hat is the world's leading provider of enterprise open source software solutions. Award-winning support, training, and consulting services make Red Hat a trusted adviser to the Fortune 500.



[linkedin.com/company/red-hat](https://www.linkedin.com/company/red-hat)



[youtube.com/user/RedHatVideos](https://www.youtube.com/user/RedHatVideos)



[facebook.com/redhatinc](https://www.facebook.com/redhatinc)



[twitter.com/RedHat](https://twitter.com/RedHat)

# Extended fio micro-benchmark results

