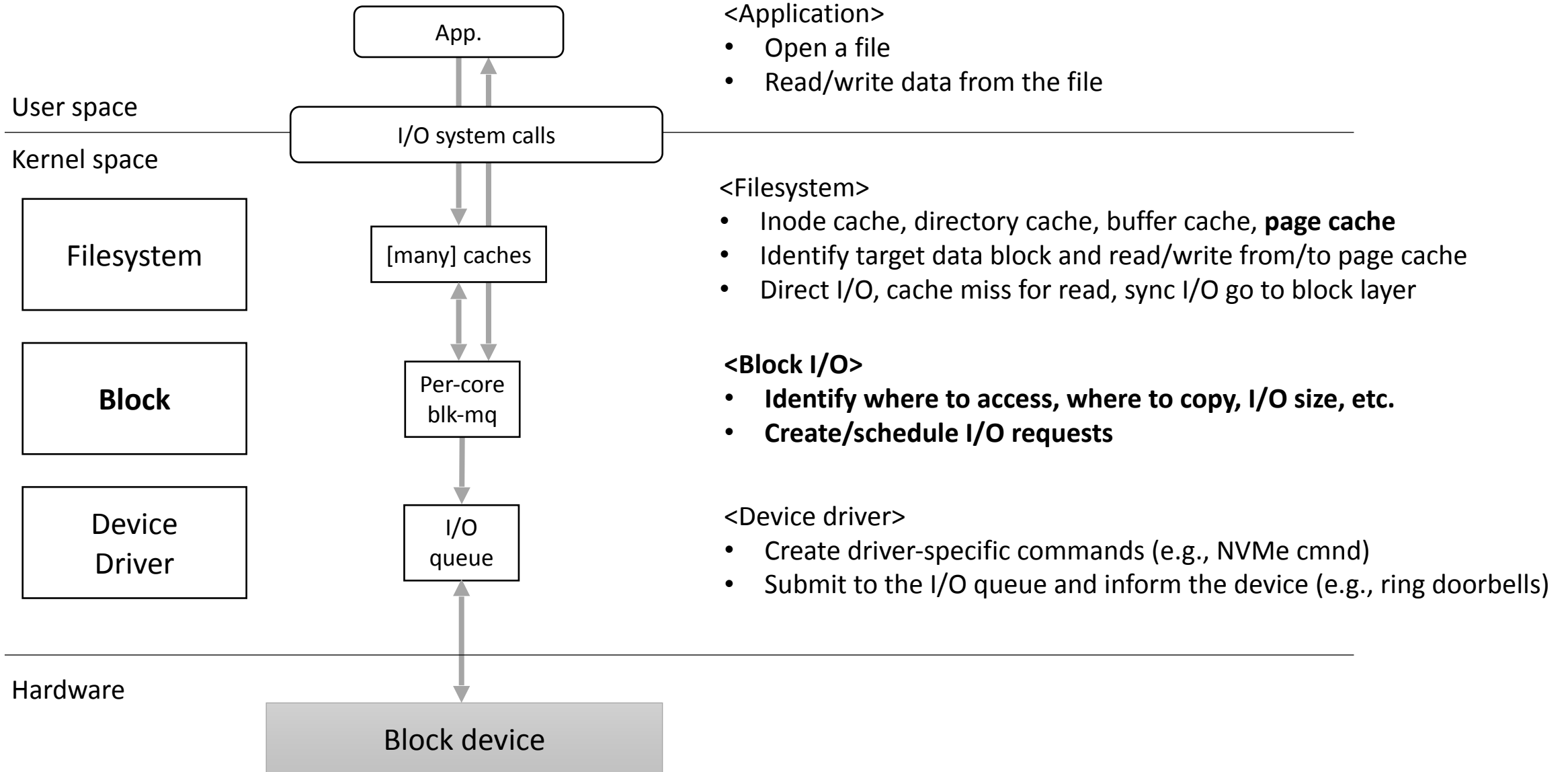# Linux kernel
# Block I/O Layer

# Overview: Accessing block devices

- **Block devices (e.g., HDD, SSD)**
  - Allow accessing fixed-size chunks of data
  - The fixed size chucks of data are called blocks
    - Block is the smallest logically addressable unit defined by filesystems (mostly 4KB)

- **Linux kernel has block I/O layer for accessing block devices**
  - Manage block devices
  - Create/schedule I/O requests
  - Interface with two layers
    - Upper layer: File System
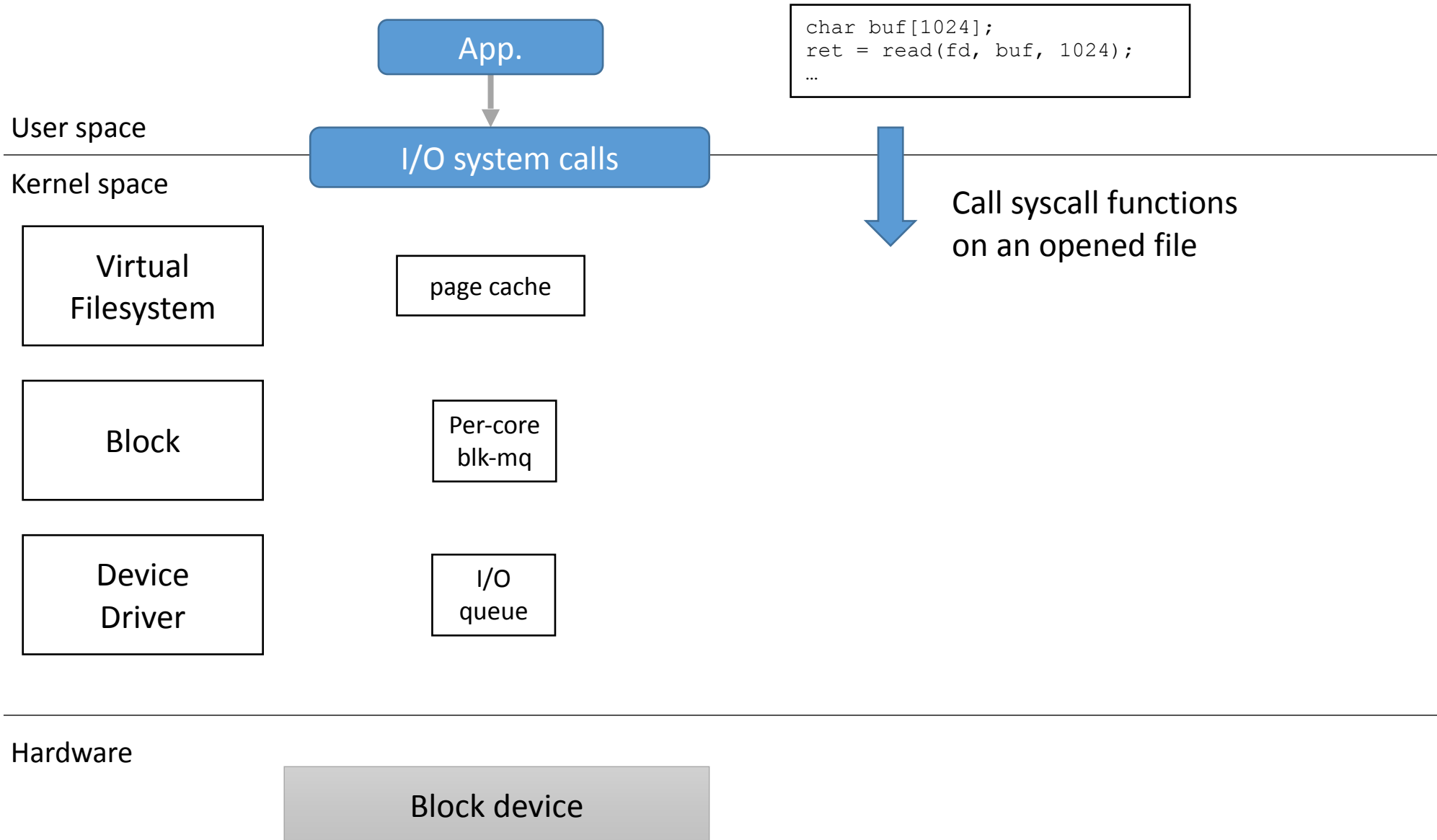    - Lower layer: Device Driver (such as NVMe)

# Overview: Accessing block devices

App.

User space

I/O system calls

Kernel space

Filesystem

[many] caches

Block

Per-core blk-mq

Device Driver

I/O queue

Hardware

Block device

**\<Application\>**
- Open a file
- Read/write data from the file

**\<Filesystem\>**
- Inode cache, directory cache, buffer cache, **page cache**
- Identify target data block and read/write from/to page cache
- Direct I/O, cache miss for read, sync I/O go to block layer

**\<Block I/O\>**
- **Identify where to access, where to copy, I/O size, etc.**
- **Create/schedule I/O requests**

**\<Device driver\>**
- Create driver-specific commands (e.g., NVMe cmnd)
- Submit to the I/O queue and inform the device (e.g., ring doorbells)

# All these caches

- **inode cache and directory cache**
  - Enabling functionalities discussed in file systems
  - Faster access to information in inode and directory
    - Separated from "data" cache


- **Page cache**
  - Combines virtual memory and file data
  - Caches recently read data on persistent storage at the granularity of pages
  - Has a notion of "file"


- **Buffer cache**
  - Interfaces with block devices (hardware)
  - Caches recently read data blocks on persistent storage
  - Has no notion of "files"—just blocks

# Read I/O data path: Application

App.

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

User space

I/O system calls

Kernel space

Call syscall functions
on an opened file

Virtual
Filesystem

page cache

Block

Per-core
blk-mq

Device
Driver

I/O
queue

Hardware

Block device

# Read I/O data path: Virtual Filesystem

App.

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

[Input]
• fd: file descriptor
• buf: user buffer
• buffer length: 1024B

User space

I/O system calls

Kernel space

**Virtual Filesystem**

page cache

kiocb

(1) Create **kiocb** that includes:
• file instance looked-up by fd
• I/O vector to copy 1024B read-data to buf
• pos: where to read/write in the file

Block

Per-core blk-mq

Device Driver

I/O queue

Hardware

Block device

# Read I/O data path: Virtual Filesystem

App.

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

**[Input]**
- fd: file descriptor
- buf: user buffer
- buffer length: 1024B

User space

I/O system calls

Kernel space

**Virtual Filesystem**

page cache

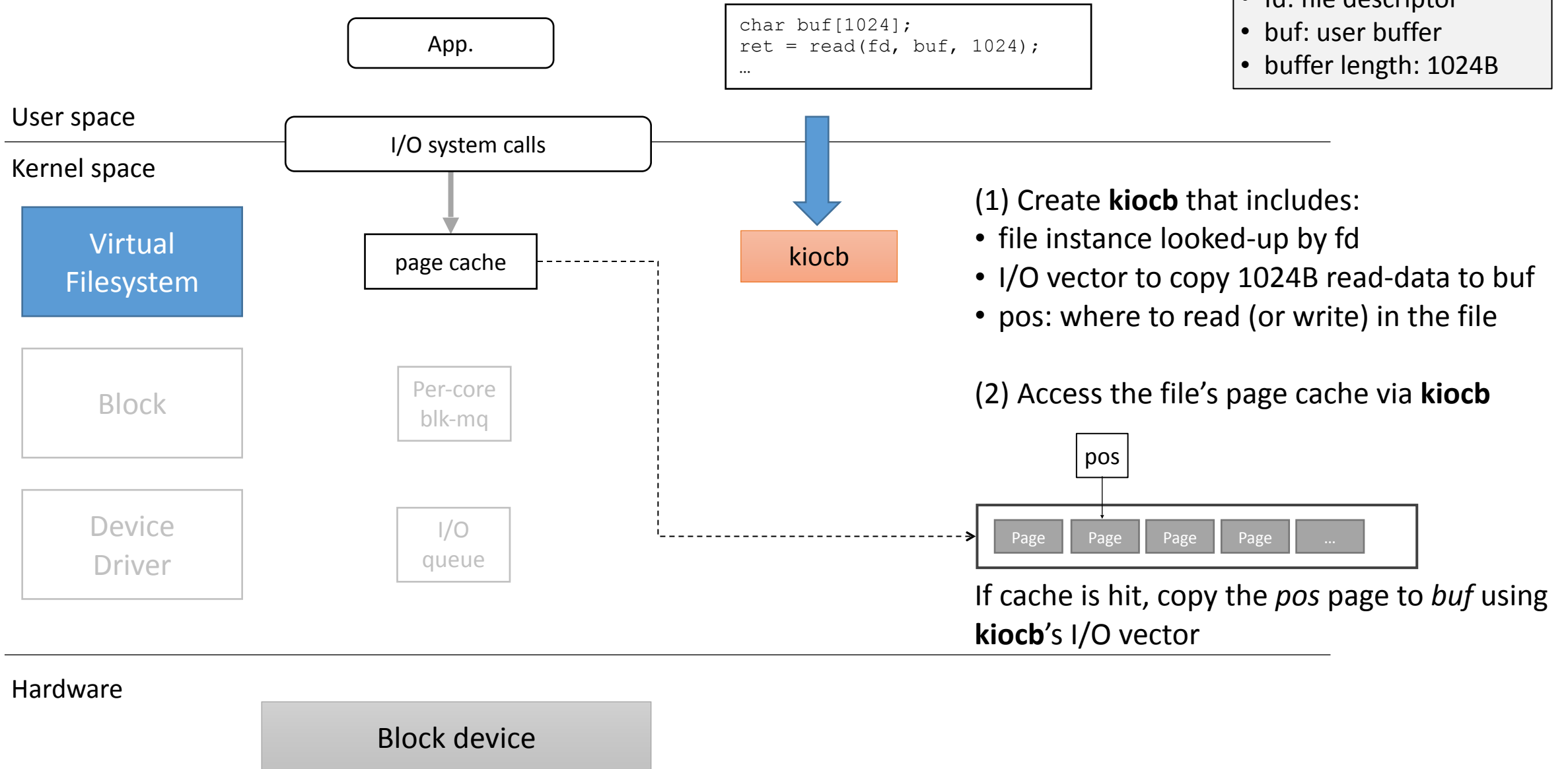kiocb

Block

Per-core blk-mq

Device Driver

I/O queue

(1) Create **kiocb** that includes:
- file instance looked-up by fd
- I/O vector to copy 1024B read-data to buf
- pos: where to read (or write) in the file

(2) Access the file's page cache via **kiocb**

pos

| Page | Page | Page | Page | … |

If cache is hit, copy the *pos* page to *buf* using **kiocb**'s I/O vector

Hardware

Block device

# Read I/O data path: Block (init. bio)

**App.**

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

**[Input]**
- file & pos
- Target page address

User space

I/O system calls

Kernel space

Virtual Filesystem

page cache

**kiocb**

Create **bio** that includes:
- Data block of the device to access
- I/O vector to copy the data block

**Block**

Per-core blk-mq

**bio**

[page cache]

I/O vector | 0 | 1 | … |

| Page | Page | Page | Page | … |

Target page address

Device Driver

I/O queue

(1) Find the data block to read using filesystem info (inode, pos, etc.)

Hardware

(2) The data block will be DMAed to the pages that bio's I/O vector points to

**Block device** data

(For direct I/O, bio's I/O vector would point to the user buffer => next slide)

# Read I/O data path: Block (init. bio)

[buf]

| Page | Page | Page | Page | ... |
|------|------|------|------|-----|

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

App.

User space

Kernel space

I/O system calls

Virtual Filesystem

page cache

kiocb

Block

Per-core blk-mq

bio

I/O vector

| 0 | 1 | ... |
|---|---|-----|

Device Driver

I/O queue

Target page address

**Direct I/O: page cache is not involved**
- In general, page cache has a single cache policy (such as LRU)
- Applications like DB might want to implement their own cache policy with direct I/O

Hardware

Block device  data

# Read I/O data path: Block (init. request)

App.

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

User space

I/O system calls

Kernel space

Virtual Filesystem

page cache

kiocb

Block

Per-core blk-mq

bio

request

Device Driver

I/O queue

Hardware
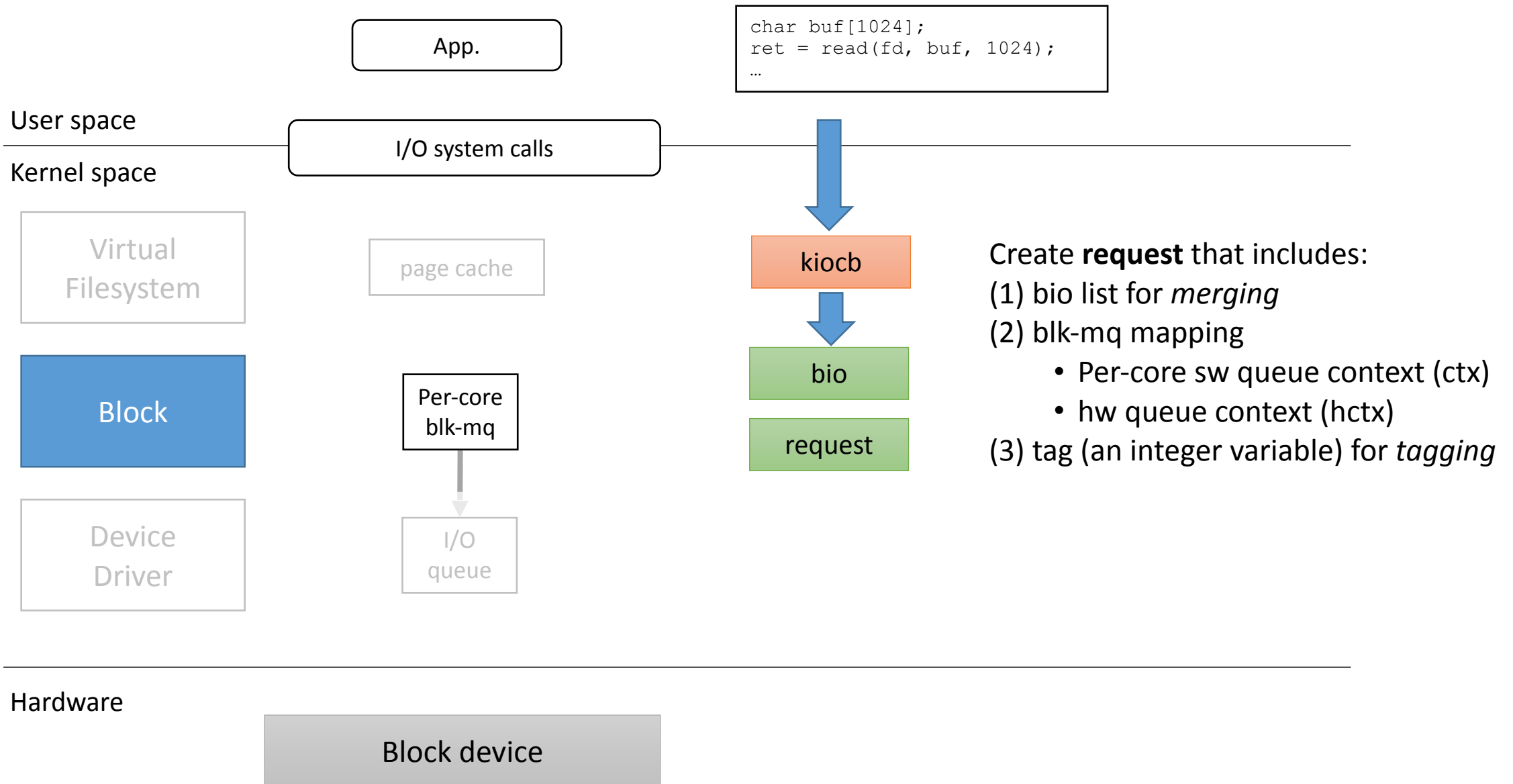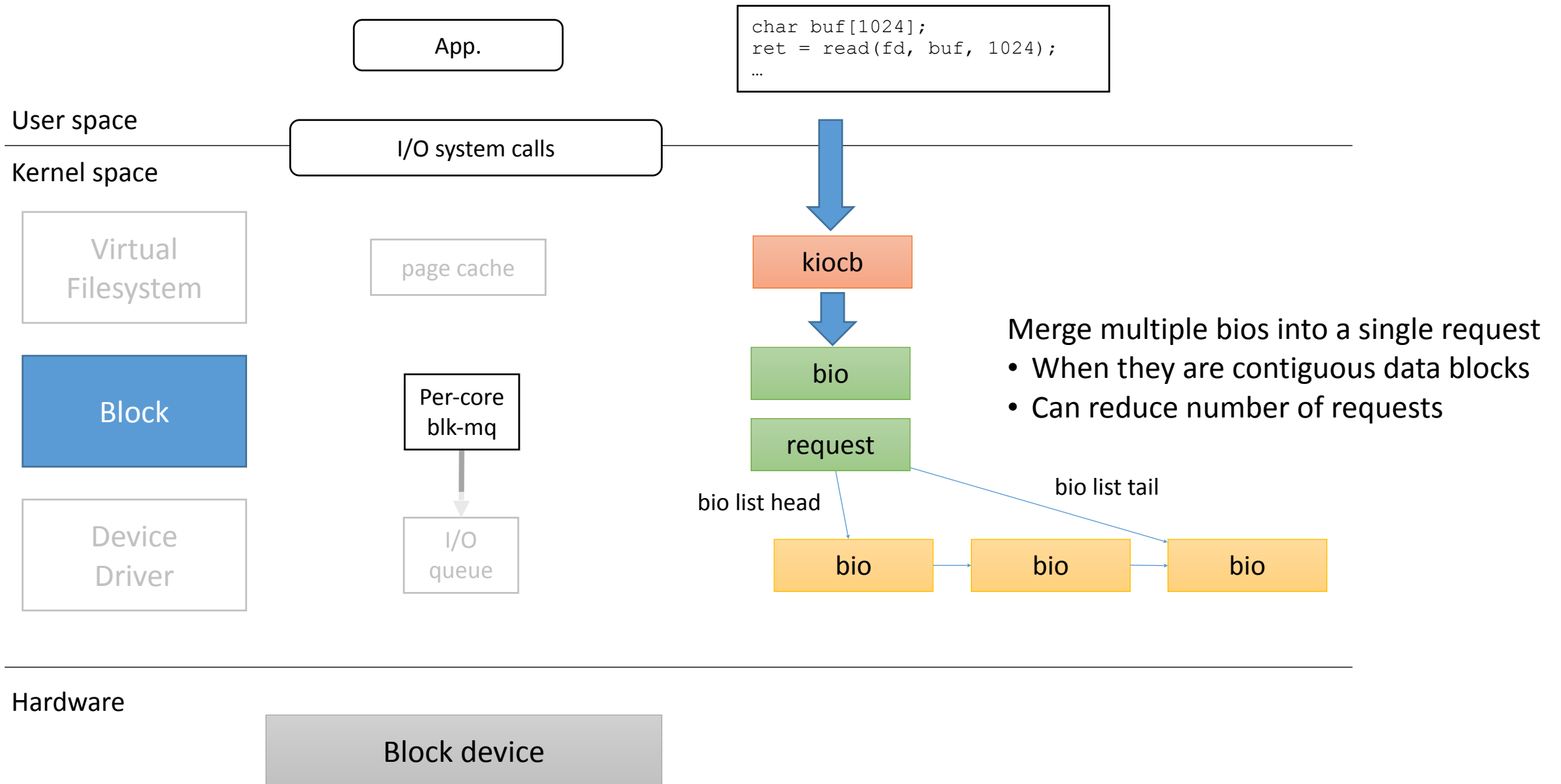
Block device

# What should a "request" contain?

- **Where is the request going?**
  - Which device
  - But devices have multiple "queues"
    - Which queue
    - Identified by a "hardware context"

- **Where should the request response be directed?**
  - Which CPU core and which application
  - Identified by a "software context"

- **A request identifier**
  - tag

# Read I/O data path: Block (init. request)

App.

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

User space

I/O system calls

Kernel space

Virtual Filesystem

page cache

Block

Per-core blk-mq

Device Driver

I/O queue

kiocb

bio

request

Create **request** that includes:
(1) bio list for *merging*
(2) blk-mq mapping
  - Per-core sw queue context (ctx)
  - hw queue context (hctx)
(3) tag (an integer variable) for *tagging*

Hardware

Block device

# Read I/O data path: Block (merging)

App.

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

User space
_____
Kernel space

I/O system calls

Virtual Filesystem

page cache

kiocb

Block

Per-core blk-mq

bio

request

Merge multiple bios into a single request
- When they are contiguous data blocks
- Can reduce number of requests

Device Driver

I/O queue

bio list head                          bio list tail

bio  →  bio  →  bio

Hardware

Block device

# Read I/O data path: Block (blk-mq mapping case 1)

#I/O queues = #cores
(default configuration)

App.

User space
Kernel space

I/O system calls

Block

ctx[0]

hctx[0,0]

request

ctx[1]

hctx[1,0]

ctx[2]

hctx[2,0]

Device
Driver

I/O
queue[0]

I/O
queue[1]

I/O
queue[2]

Hardware

Block device

# Read I/O data path: Block (blk-mq mapping case 2)

#I/O queues = #cores x 2

App.

User space

Kernel space

I/O system calls

Current Linux kernel supports three hctx types:
- read/write/poll

Block

ctx[0]

ctx[1]

ctx[2]

hctx[0,0]  hctx[0,1]

hctx[1,0]  hctx[1,1]

hctx[2,0]  hctx[2,1]

In this example:
- hctx[*,0] for read I/O
- hctx[*,1] for write I/O

request

Device Driver

I/O queue[0]   I/O queue[3]

I/O queue[1]   I/O queue[4]

I/O queue[2]   I/O queue[5]

Hardware

Block device

# Read I/O data path: Block (scheduling)

App.

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

User space

I/O system calls

Kernel space

Virtual Filesystem

page cache

Block

Per-core blk-mq

I/O queue

Device Driver

kiocb

bio

request

When request initialization is done, perform I/O scheduling at per-req granularity
- Linux supports Noop/deadline/BFQ/Kyber (Noop skips I/O scheduling step)
- Existing schedulers were mostly designed for HDD (slow devices)
- For NVMe SSD (fast devices), I/O scheduling may add extra latency, so Noop is chosen in many cases

Hardware

Block device

# Read I/O data path: Device Driver (NVMe)

App.

```
char buf[1024];
ret = read(fd, buf, 1024);
…
```

User space

I/O system calls

Kernel space

Virtual Filesystem

page cache

kiocb

Block

Per-core blk-mq

bio

request

Create **cmnd** that includes:
- request instance
- NVMe-specific data structures (command header, etc.)

Device Driver

I/O queue

cmnd

Perform NVMe operations (next slide)
- As I/O queues, NVMe defines submission queue (SQ) and completion queue (CQ) per core

Hardware

Block device

# Read I/O data path: Device Driver (NVMe)

# Read I/O data path: Response

App.

User space
Kernel space

I/O system calls

(2) For sync I/O,
• Return to app

(3) For async I/O,
• Add a completion event in the ring buffer (kioctx)
• Signal the app to wake-up

Virtual Filesystem

page cache

kiocb - - - - - - - - → kioctx

Block

Per-core blk-mq

bio

request

Device Driver

I/O queue

(1) From the completion entry,
• Extract the tag
• Find the corresponding request instance from the tag
• Call block-layer completion callbacks

Hardware

Block device

# Write I/O path (difference from Read)

App.

```
char buf[1024];
ret = write(fd, buf, 1024);
…
```

User space

Kernel space

Virtual Filesystem

I/O system calls

page cache

kiocb

(1) Write the user data in the page cache

(2) For async writes,
- Return to app immediately
- Write to device later

Block

Per-core blk-mq

Device Driver

I/O queue

pos

| Page | Page | Page | Page | … |

Hardware

Block device

(3) For sync writes,
- Write to device immediately through block layer