# Barrier Enabled IO Stack for Flash Storage

## Youjip Won

KAIST EE

# Motivation

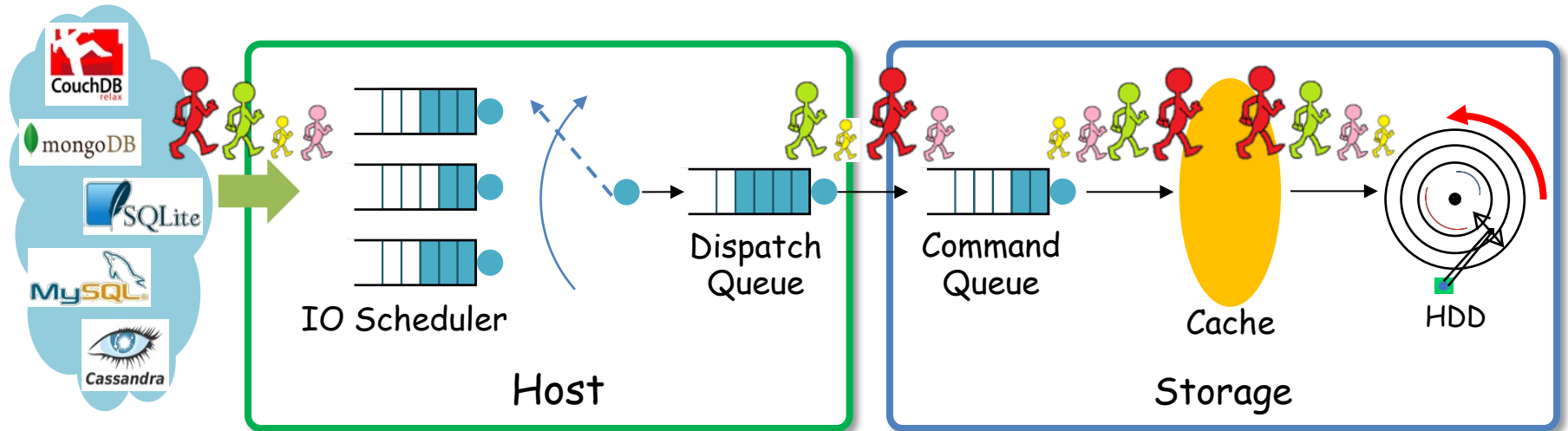# Modern IO Stack

## Modern IO stack is Orderless.

Issue ($I$)　　　　　　　　　　Dispatch ($D$)　Transfer ($X$)　Persist ($P$)



IO Scheduler

Dispatch Queue

Command Queue

Cache

HDD

Host

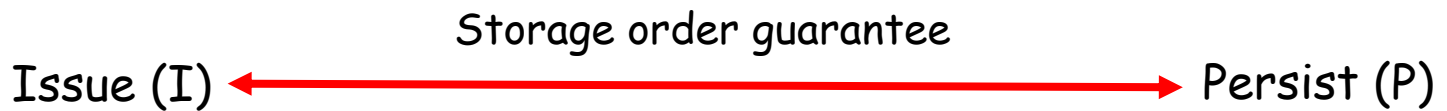Storage

$I \neq D$: IO Scheduling

$D \neq X$: Time out, retry, command priority

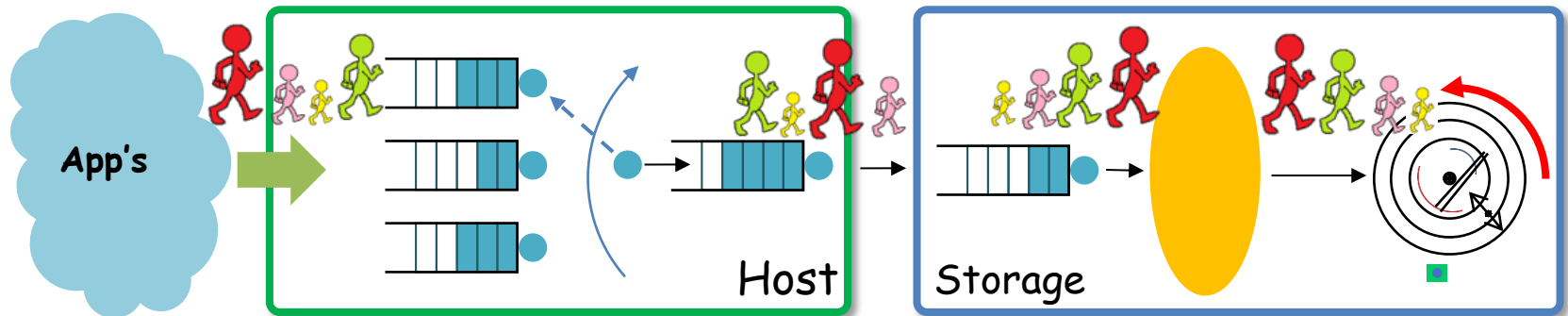$X \neq P$: Cache replacement, page table update algorithm of FTL

# Storage Order

Storage Order: The order in which the data blocks are made durable.
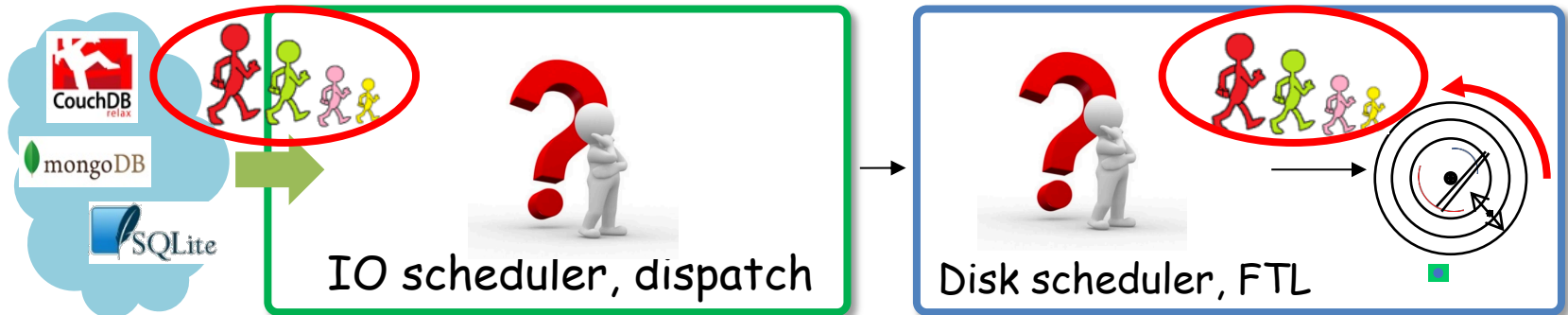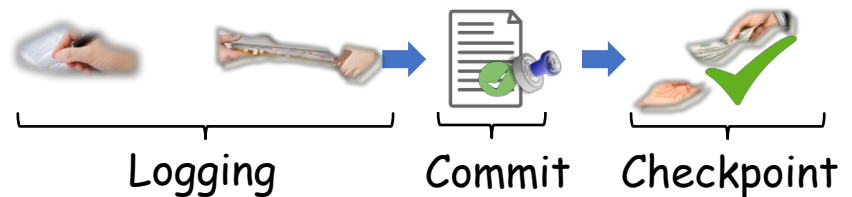
Guaranteeing the storage order

Storage order guarantee

Issue (I) ←————————————————→ Persist (P)

$$(I = D) \wedge (D = X) \wedge (X = P)$$

Issue ($I$) ←——→ Dispatch ($D$) ←——→ Transfer ($X$) ←——→ Persist ($P$)
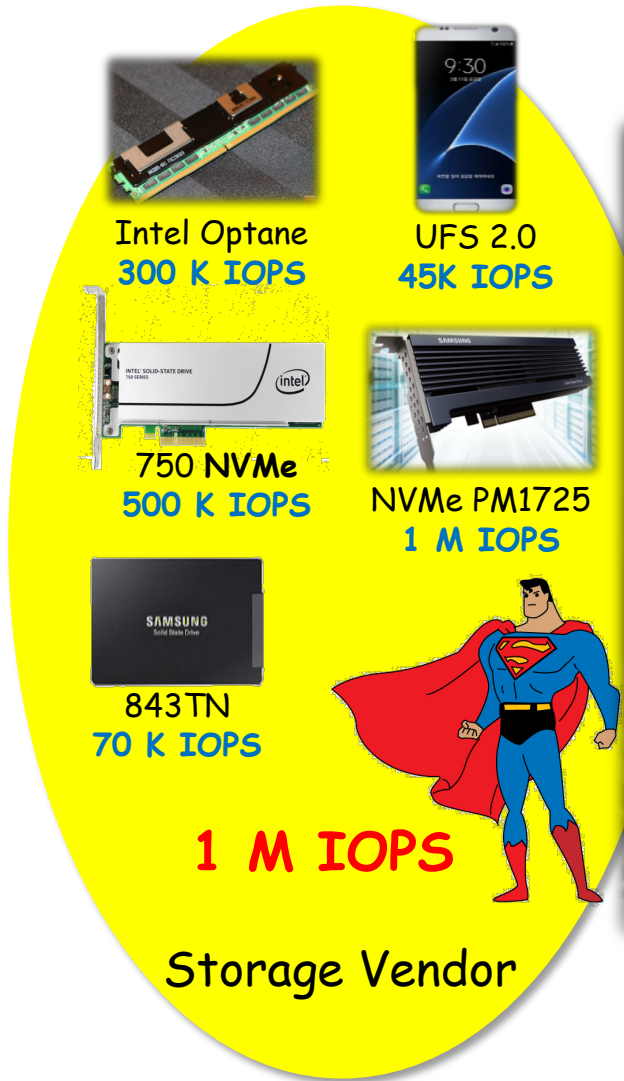


App's

Host

Storage

# Controlling the Storage Order

Applications need to control the storage order.

- Database logging
- Filesystem Journaling
- Soft-updates
- COW based filesystem



Logging     Commit     Checkpoint



IO scheduler, dispatch     Disk scheduler, FTL

# What's Happening Now....



Intel Optane
**300 K IOPS**

UFS 2.0
**45K IOPS**

750 **NVMe**
**500 K IOPS**

NVMe PM1725
**1 M IOPS**

843TN
**70 K IOPS**

**1 M IOPS**

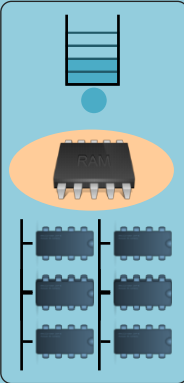## Storage Vendor

facebook
Google
CouchDB
SQLite
**274 OPS/s**
cassandra
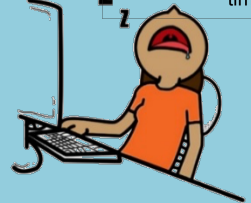**2,500 OPS/s**
MySQL
**1,132 OPS/s**

**I/O is Bottleneck!!!**
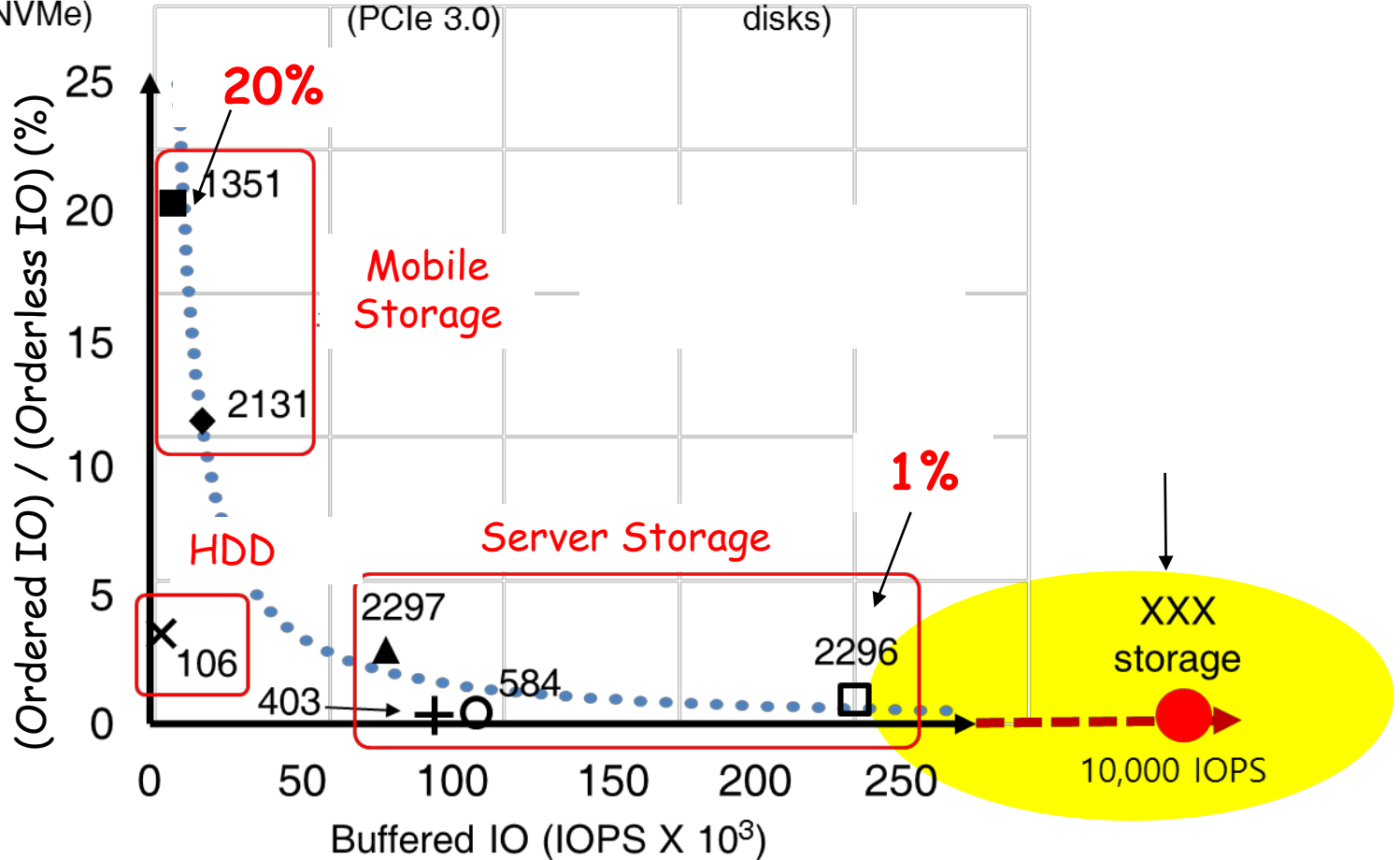
## Service Provider

## Storage

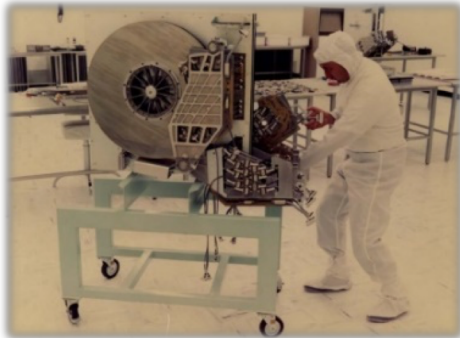# Overhead of storage order guarantee: write() + fdatasync()

**Legend:**
- ■ Galaxy S5 (eMMC 5.0)
- ◆ Galaxy S6 (UFS 2.0)
- ✚ Samsung 850 PRO (SATA3)
- ✕ HDD
- ○ Samsung 950 PRO (NVMe)
- ▲ OCZ RevoDrive 3 X2 (PCIe 3.0)
- □ OCZ RevoDrive 3 X2 (PCIe 3.0, RAID 0 of 4 disks)



**20%**

1351

Mobile Storage

2131

HDD

106

Server Storage

2297

1%

2296

XXX storage

403

584

10,000 IOPS

Y-axis: (Ordered IO) / (Orderless IO) (%)
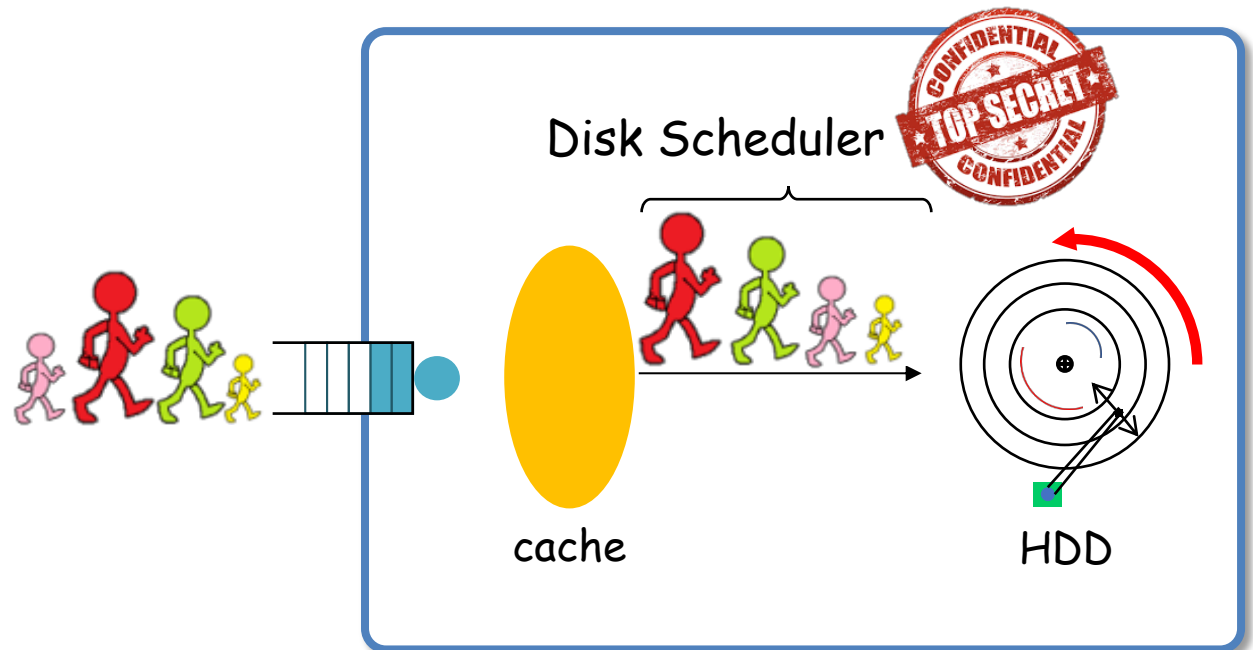
X-axis: Buffered IO (IOPS X $10^3$)

# Why has IO stack been orderless for the last 50 years?

## In HDD, host cannot control the persist order.

$$(I \times P) \equiv (I = D) \wedge (D = X) \wedge (X \times P)$$
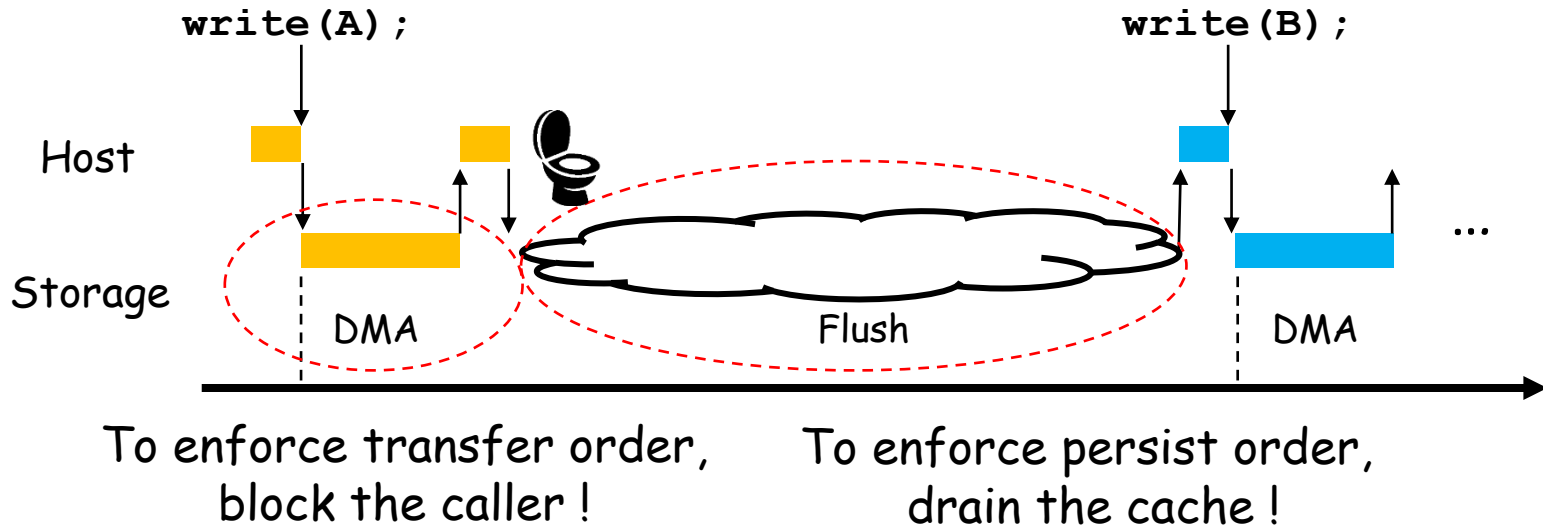
Disk Scheduler

250MB @ 1970's

cache

HDD

# Enforcing Storage Order in spite of Orderless IO Stack
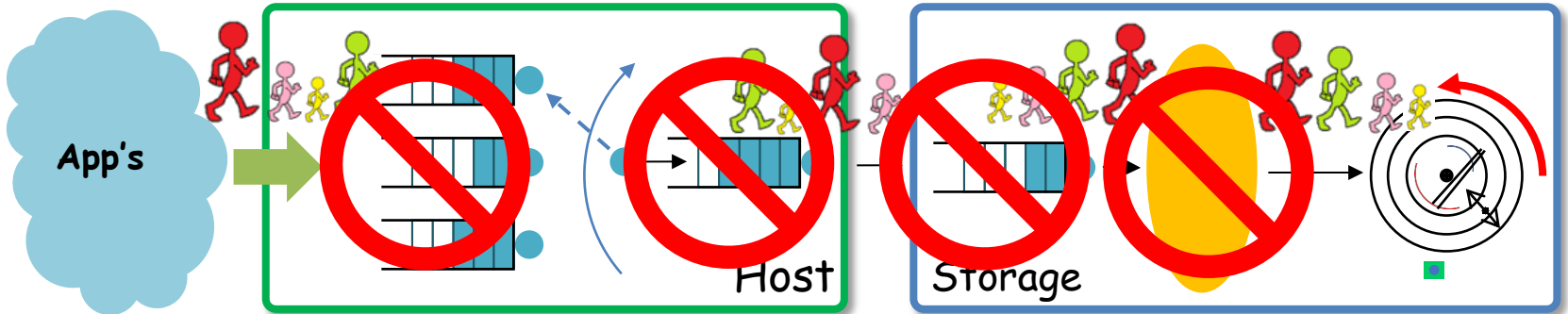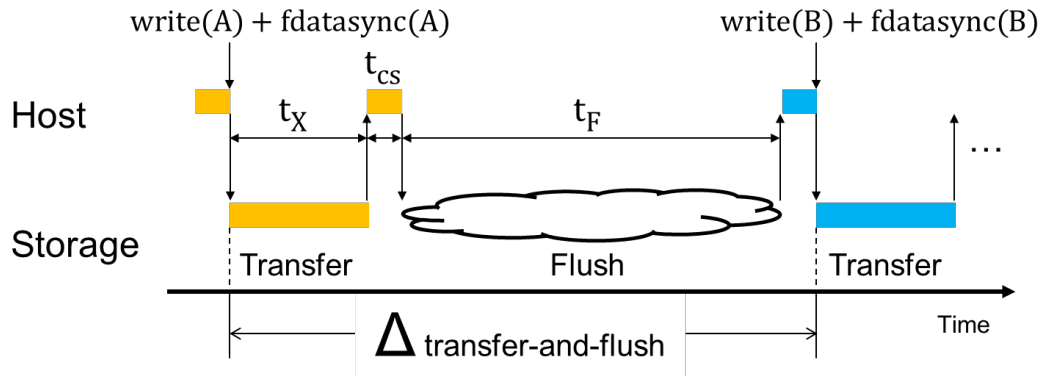
Interleave the write request with Transfer-and-Flush

```
write (A) ;          write (A) ;
write (B) ;    ──►   Transfer-and-flush;
                     write (B) ;
```

**write(A);**                              **write(B);**

Host

Storage

DMA                    Flush                    DMA

To enforce transfer order,          To enforce persist order,
    block the caller !                  drain the cache !
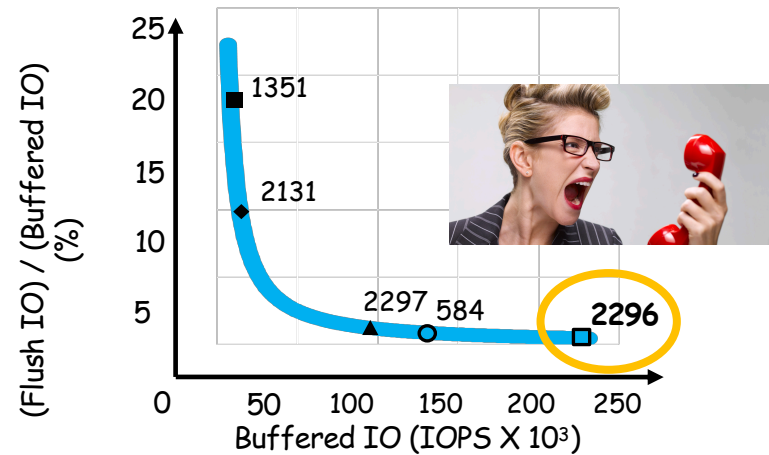
# Transfer-and-Flush

# Overhead of Transfer-and-Flush



**NVMe PM1725**
**120K IOPS**

Ordering Guarantee
**< 2%**

**NVMe PM1725**
**2K IOPS**



Storage Performance (IOPS)

- X25-M **35 K IOPS** (2009)
- 830 PRO **80 K IOPS** (2012)
- 850 PRO **100 K IOPS** (2014)
- Intel 750 **440 K IOPS** (2015)
- PM1725 **1 M IOPS**

(Flush IO) / (Buffered IO) (%) vs Buffered IO (IOPS X 10³)

1351, 2131, 2297, 584, **2296**

# Developing Barrier-enabled IO Stack

In the era of HDD
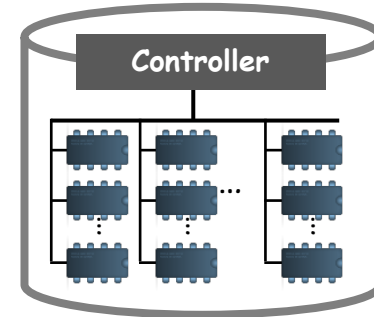(circa 1970)

In the era of SSD
(circa 2000)



Controller

Seek and rotational delay.

➡ The host cannot control persist order.

➡ the IO stack becomes orderless.
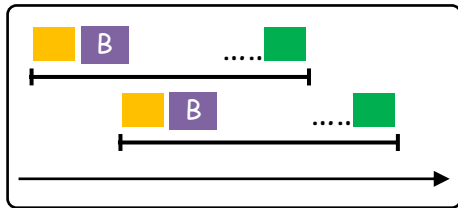
➡ **use transfer-and-flush** to control the storage order

~~Seek and rotational delay~~

➡ The host may control persist order.

➡ The IO stack may become order-preserving.

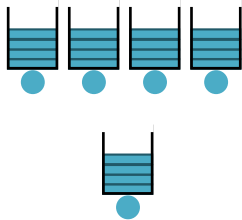➡ Control the storage order **without Transfer-and-Flush**

It is a time to re-think the way to control the storage order.

# Barrier-enabled IO Stack
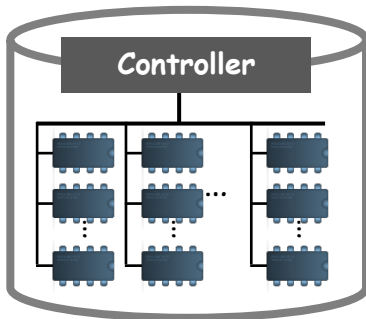


**BarrierFS**

- Dual-Mode Journaling
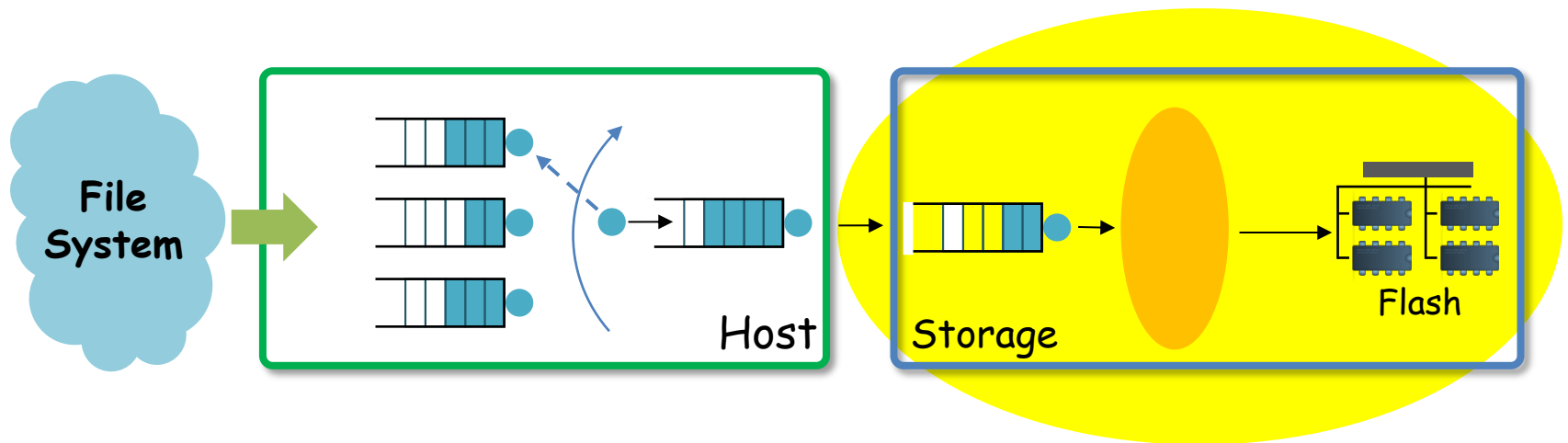- fbarrier() / fdatabarrier()

**Order-preserving Block Device Layer**

- Order-preserving dispatch
- Epoch-based IO scheduling

**Barrier-enabled Storage**

Controller

- Barrier write command

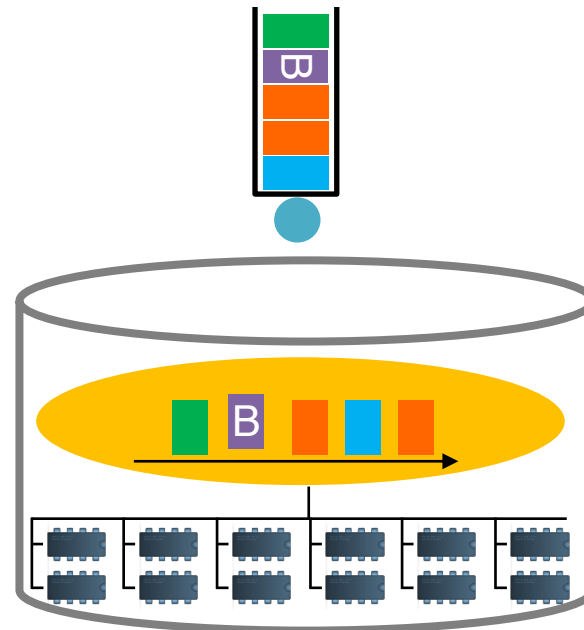# Barrier-enabled Storage



File System → Host → Storage → Flash
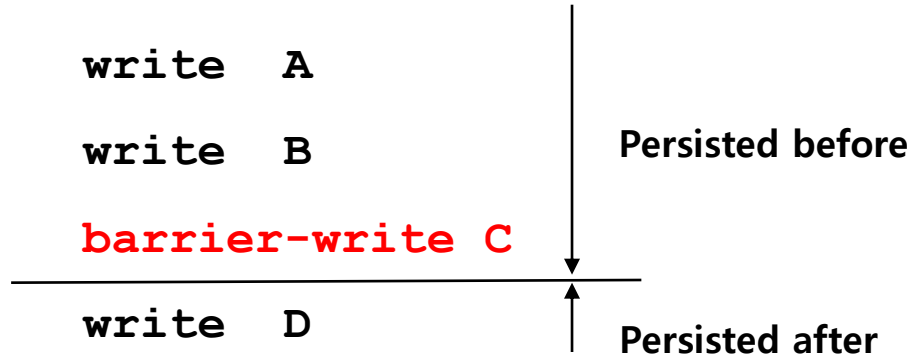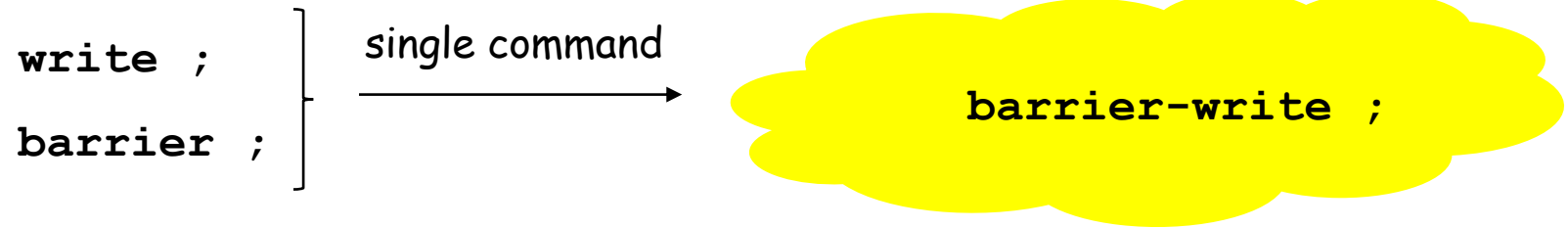
# To Control the Persist Order, X = P

barrier command (2005, eMMC)

```
write (A) ;
write (B) ;
write (C) ;
barrier;
write (D) ;
```

# Barrier Write

```
write ;
barrier ;
```

single command →

**barrier-write ;**

```
write   A
write   B
barrier-write C
write   D
```

Persisted before

Persisted after

With Barrier Write command,

host can control the persist order <span style="color:red">without flush</span>.

$$(I \xcancel{\to} P) \equiv (I \xcancel{\to} D) \wedge (D \xcancel{\to} X) \wedge (X \xcancel{\to} P)$$
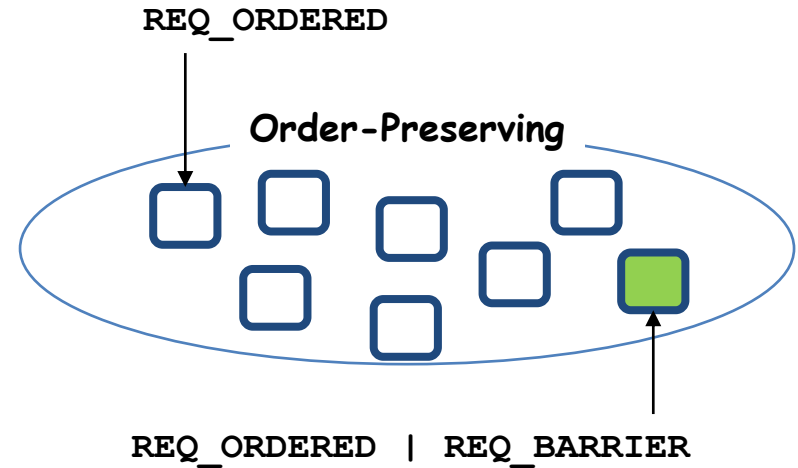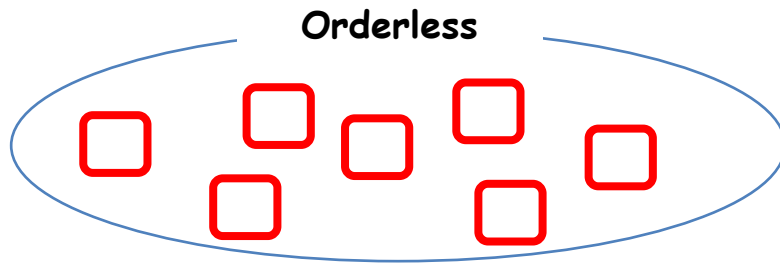
# Order-preserving Block Device Layer

# Order Preserving Block Device Layer

- ✓ New request types

- ✓ Order Preserving Dispatch

- ✓ Epoch Based IO scheduling

# Request Types



**Orderless**

**Order-Preserving**

REQ_ORDERED

REQ_ORDERED | REQ_BARRIER

{A,B,E} → {G,H}

H G F E D C B A → block layer ? → E F H D A C B E → Cache

# Order Preserving Dispatch Module (for D = X)

➢ Ensure that the barrier request is serviced in-order.

Set the command priority of 'barrier' type request  to ORDERED.

```
write A
barrier-write B //set the command priority to 'ORDERED'
write C
```

Storage

# SCSI Command Priority
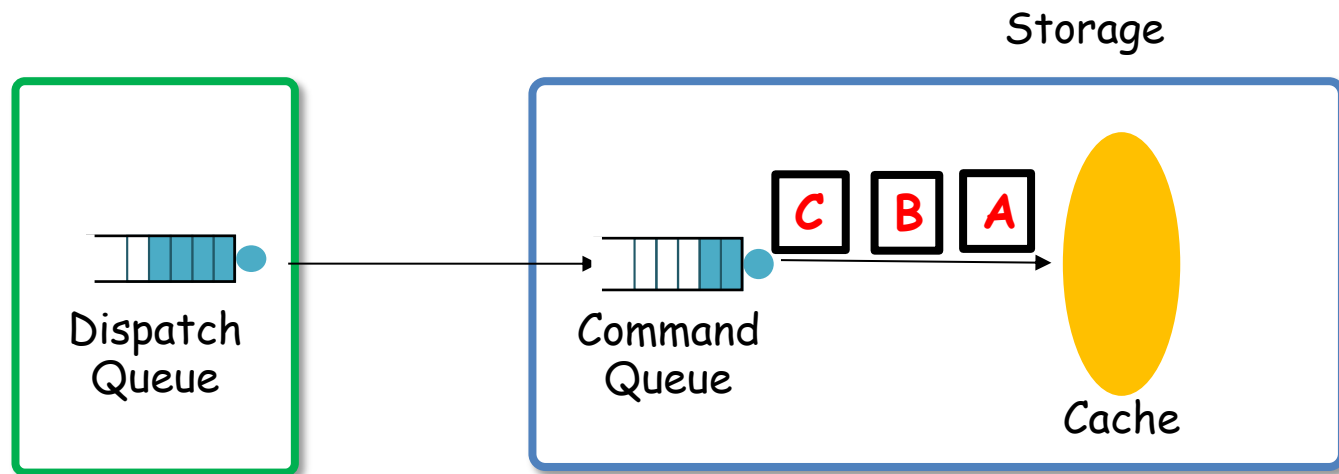
✓ **Head of the Queue**

Dispatch Queue → (HoQ) → Command Queue

✓ **Ordered (Barely being used)**

Dispatch Queue → (Ordered) → Command Queue

✓ **Simple (Default)**

anywhere

Dispatch Queue → (Simple) → Command Queue

# Order Preserving Dispatch



Legacy Dispatch

```
write(A);  write(B);
```

Host

Storage

DMA        DMA        ...

Caller blocks.

DMA transfer overhead

Order Preserving Dispatch

```
write(A);  // "ordered"
    write(B);  //"simple"
```

Host

Storage

DMA        DMA        ...

Caller does not block.  👍

No DMA transfer overhead  👍

With Order Preserving Dispatch, host can control the transfer order without DMA transfer.

$$(I \quad P) \equiv (I \quad D) \wedge (D \quad X) \wedge (X = P)$$

# Epoch Based IO scheduler (for I = D)

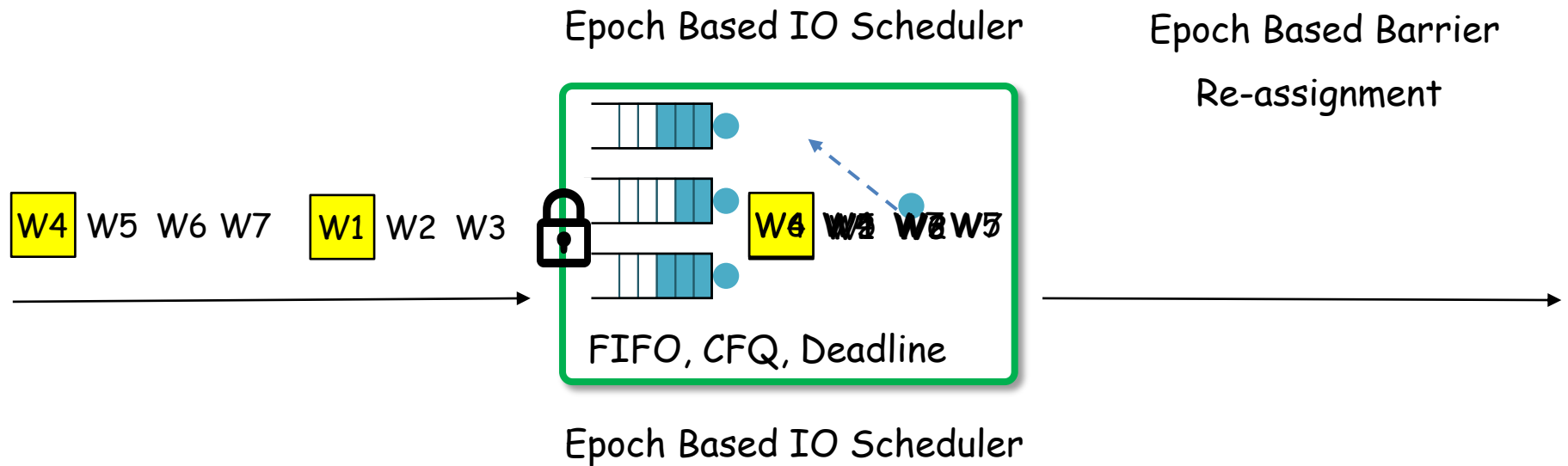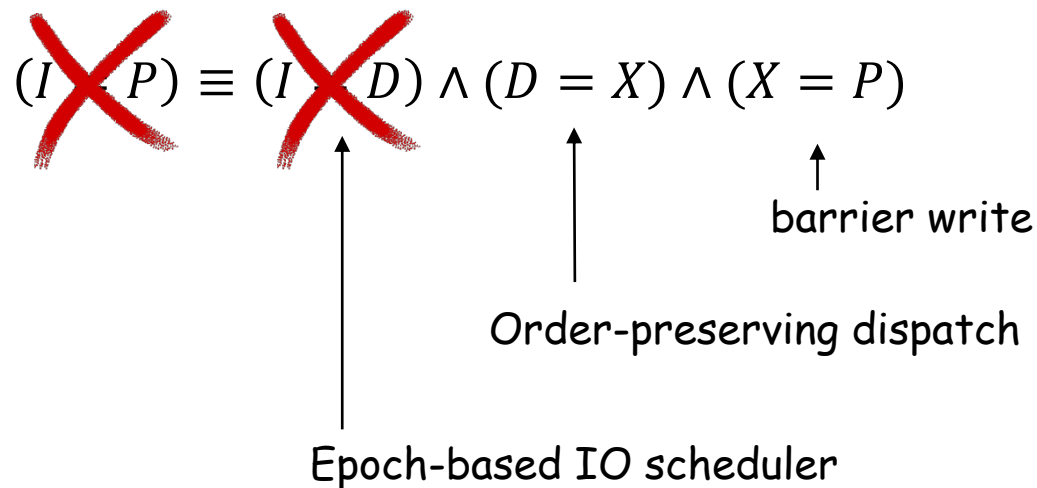➢ Ensure that the OP requests between the barriers can be freely scheduled.

➢ Ensure that the OP requests does not cross barrier boundary.

➢ Ensure that orderless requests can be freely scheduled independent with barrier.

Epoch Based IO Scheduler

Epoch Based Barrier Re-assignment

W4 W5 W6 W7    W1 W2 W3

W4  W4 W6W7W5

FIFO, CFQ, Deadline

Epoch Based IO Scheduler

With Epoch Based IO Scheduling, host can control the dispatch order with existing IO scheduler.

$$(I \cancel{\quad} P) \equiv (I \cancel{\quad} D) \wedge (D = X) \wedge (X = P)$$

barrier write

Order-preserving dispatch

Epoch-based IO scheduler

# Order Preserving Block Device Layer

**Control Storage Order without Transfer-and-Flush !**



File system

IO Scheduler

Dispatch Queue

Host

Command Queue

Cache

Flash

Storage

# Enforcing the Storage Order

## Legacy Block Layer (With Transfer-and-Flush)

`write(A);`                                                              `write(B);`

Host

Storage    DMA                    Flush                         DMA     ...

## Order Preserving Block Layer

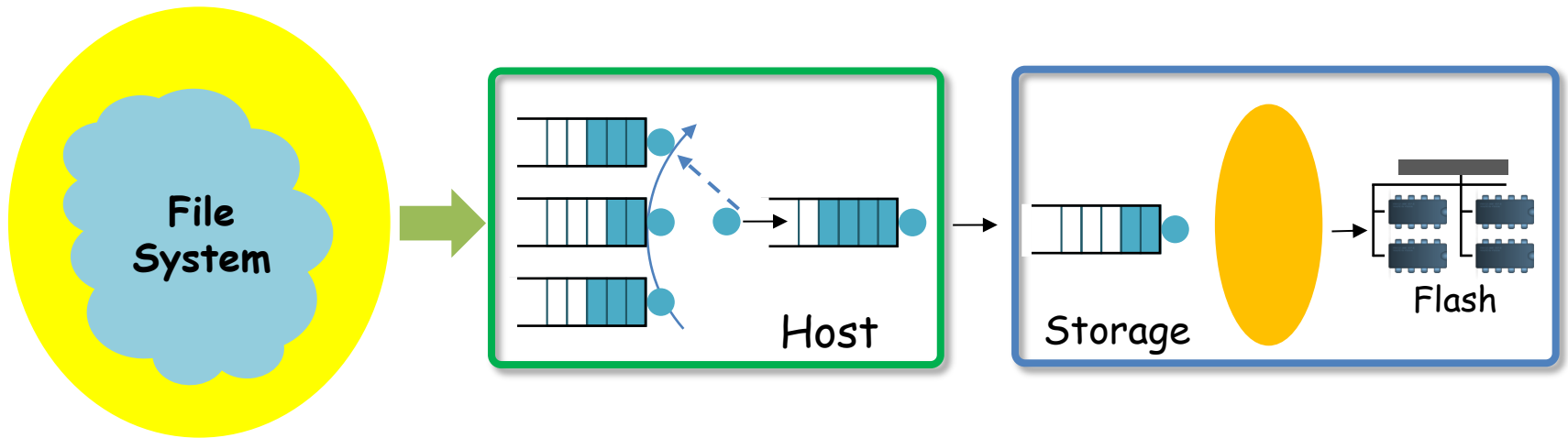`bwrite(A);`    `bwrite(B);`

Host

Storage

No Flush !

No DMA !

No Context Switch !

# Barrier-enabled Filesystem

# New primitives for ordering guarantee

| | Durability guarantee | Ordering guarantee |
|---|---|---|
| Journaling | ✓ fsync()<br>➢ Dirty pages<br>➢ journal transaction<br>➢ Durable | ✓ fbarrier()<br>➢ Dirty pages<br>➢ Journal transaction<br>➢ ~~durable~~ |
| No journaling | ✓ fdatasync()<br>➢ Dirty pages<br>➢ durable | ✓ fdatabarrier()<br>➢ Dirty pages<br>➢ ~~durable~~ |

# fsync() in EXT4

{Dirty Pages ($D$), Journal Logs ($JD$)} → {Journal Commit ($JC$)}

- Two Flushes
- Three DMA Transfers
- A number of Context switches



fsync () start

fsync () end

Filesystem

JBD

Storage

DMA  D  DMA  JD  Flush  DMA  JC  FUA

# fsync() in BarrierFS

- write Dirty pages 'D' with order-preserving write

- write Journal Logs 'JD' with barrier write

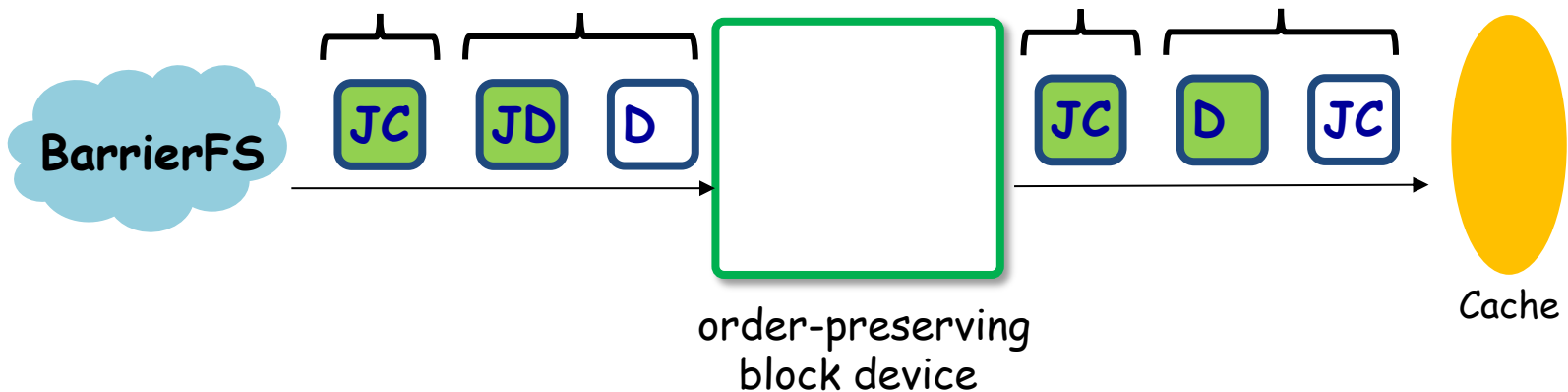- write Journal Commit Block 'JC' with barrier write
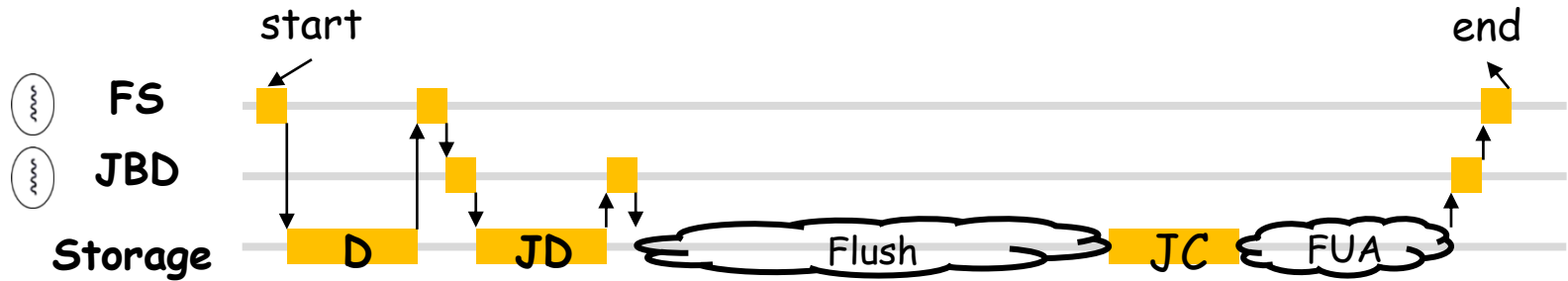
- flush

☐   order-preserving write (`REQ_ORDERED`)

🟩   barrier write (`REQ _ORDERED|REQ_BARRIER`)

## {D,JD} → {JC}
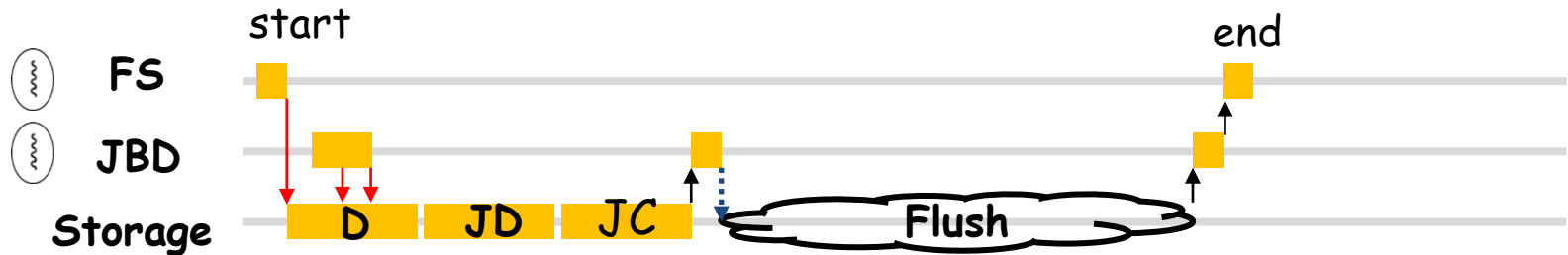


order-preserving block device

Cache

# Efficient fsync() implementation

✓ **fsync() in EXT4**
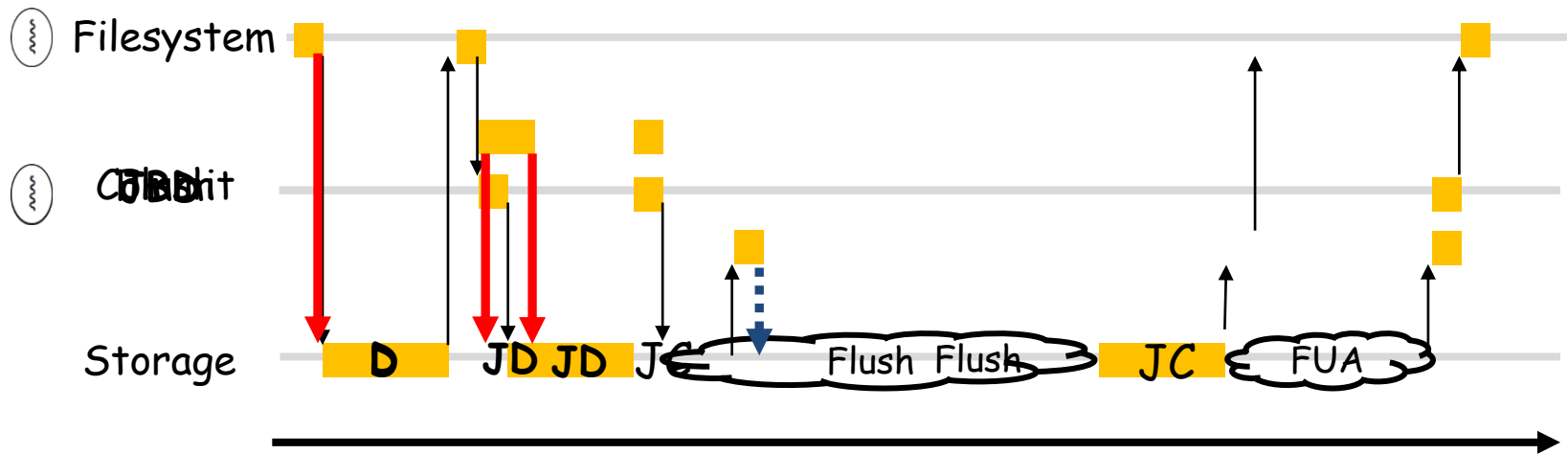


✓ **fsync() in BarrierFS**
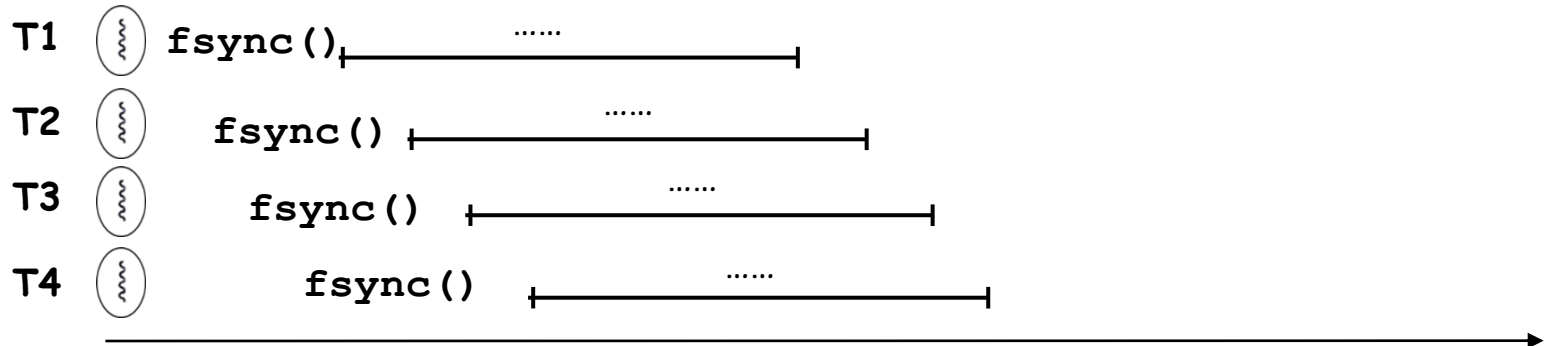
# Dual Mode Journaling

- Journal Commit

  - Dispatch 'write JD' and 'write JC'        → Control plane

  - Make JD and JC durable                    → Data Plane

- Dual Mode Journaling

  - separate the control plane activity and the data plane activity.

  - Separate thread to each

    - Commit Thread (Control Plane)
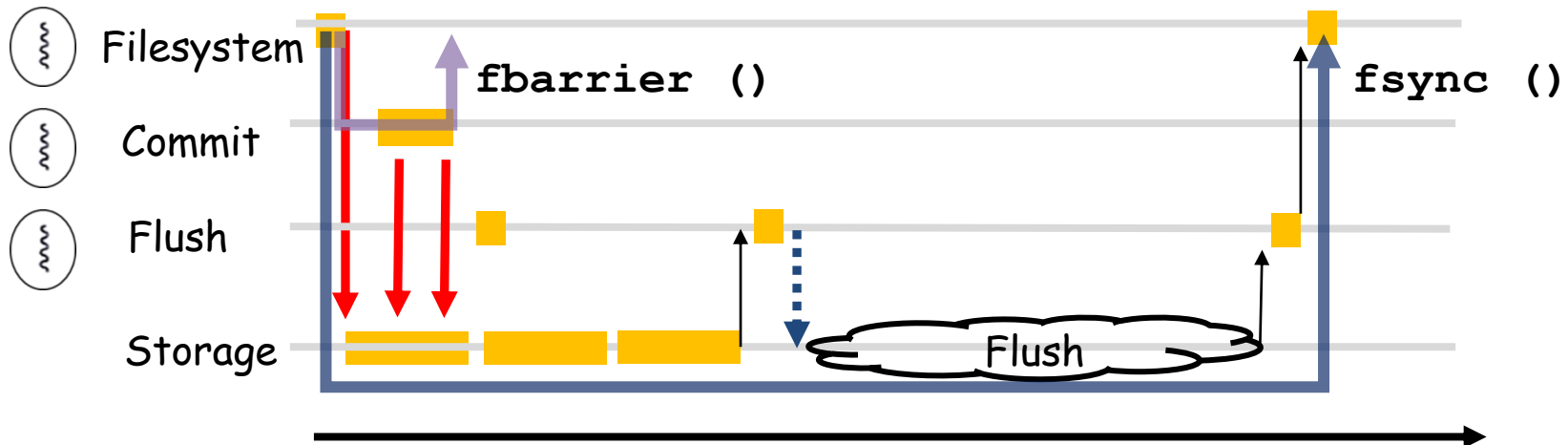
    - Flush Thread (Data Plane)

# Implications of Dual Thread Journaling

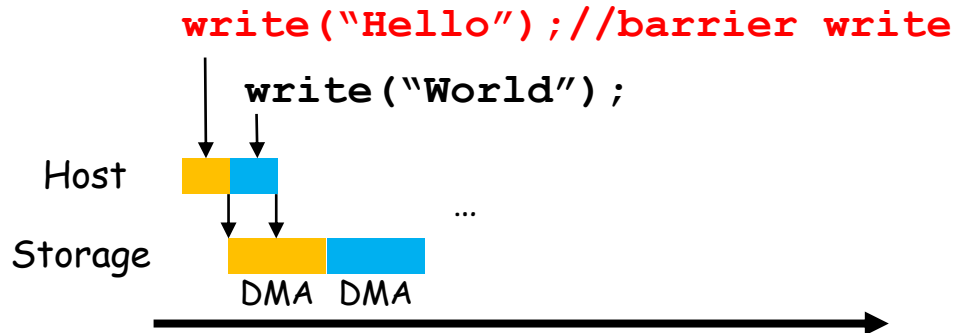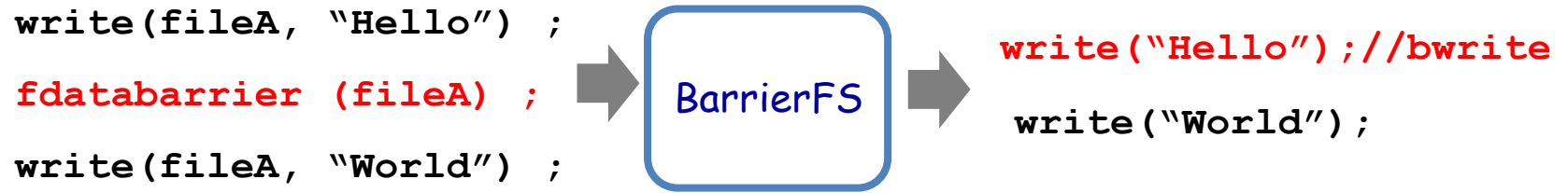✓ Journaling becomes concurrent activity.

T1 fsync() ...... 

T2 fsync() ......

T3 fsync() ......

T4 fsync() ......

✓ Efficient Separation of Ordering Guarantee and Durability Guarantee

Filesystem **fbarrier ()** **fsync ()**

Commit

Flush

Flush

Storage

# fdatabarrier()

- write Dirty pages 'D' with order-preserving write

```
write(fileA, "Hello") ;
fdatabarrier (fileA) ;
write(fileA, "World") ;
```
→ BarrierFS →
```
write("Hello");//bwrite
write("World");
```

**write("Hello");//barrier write**
write("World");

Host

Storage

...

DMA  DMA

DMA transfer overhead  NO
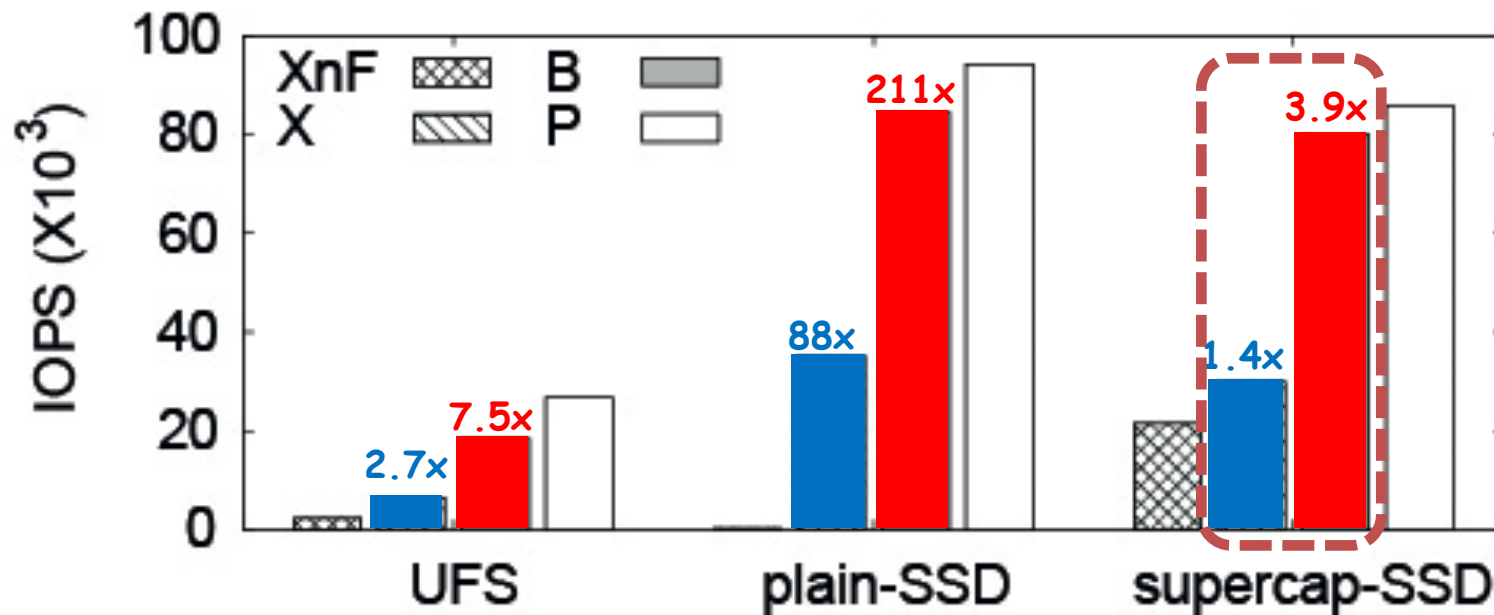
Flush overhead  NO

Context switch  NO

# Experiments

- Platforms: PC server (Linux 3.16), Smartphone (Galaxy S6 Linux 3.10)

- Flash Storages:

    - Mobile-SSD(UFS2.0, 2ch), Plain-SSD (SM 850, 8ch), Supercap-SSD (SM843, 8ch)

- Workload

    1. Micro benchmark: Mobibench, FxMark (Microbenchmark)

    2. Macro Benchmark: Mobibench(SQLite), filebench(varmail), sysbench(MySQL)

- IO stack

    1. Durability guarantee: EXT-DR(fsync()), BFS-DR(fsync())

    2. Ordering guarantee: EXT4-OD (fdatasync(), NO-barrier), BFS-OD (fdatabarrier())

# Benefit of Order-Preserving Dipspatch

## Eliminating Flush
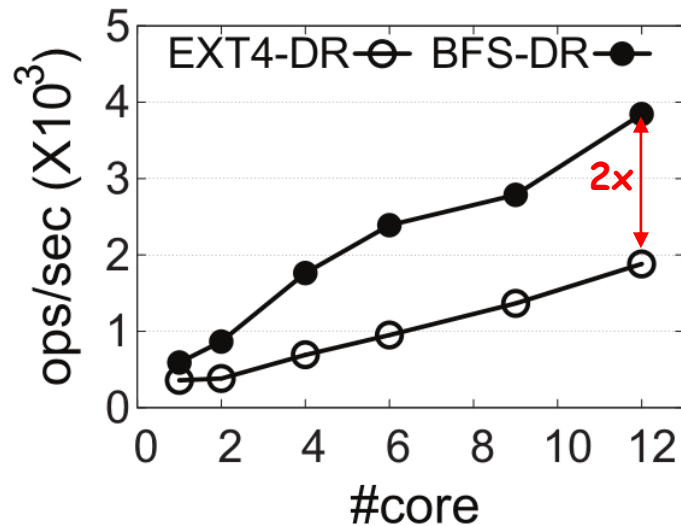
# Eliminating Transfer-and-Flush
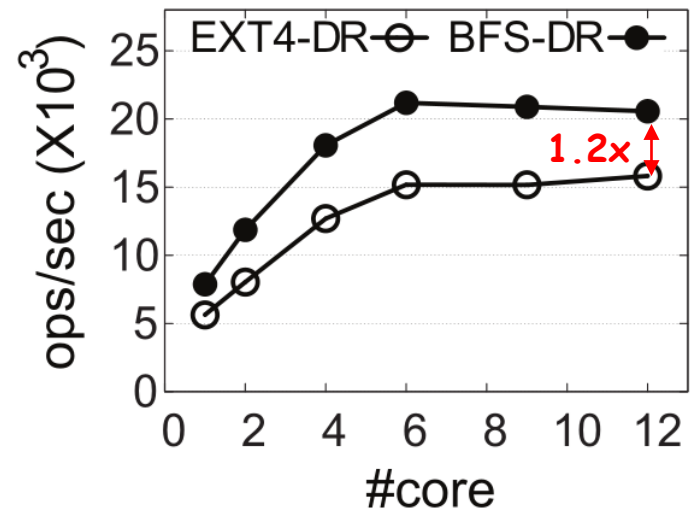


**Eliminating the transfer overhead is critical.**

# Journaling Scalability

- 4 KB Allocating write followed by fsync() [DWSL workload in FxMark]

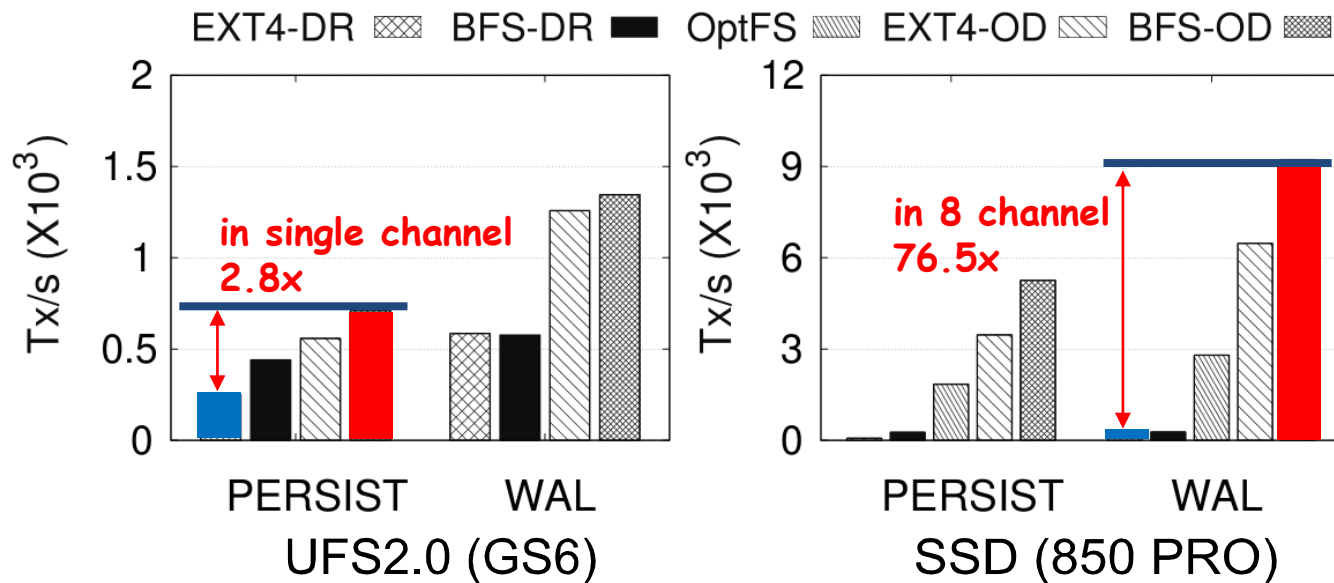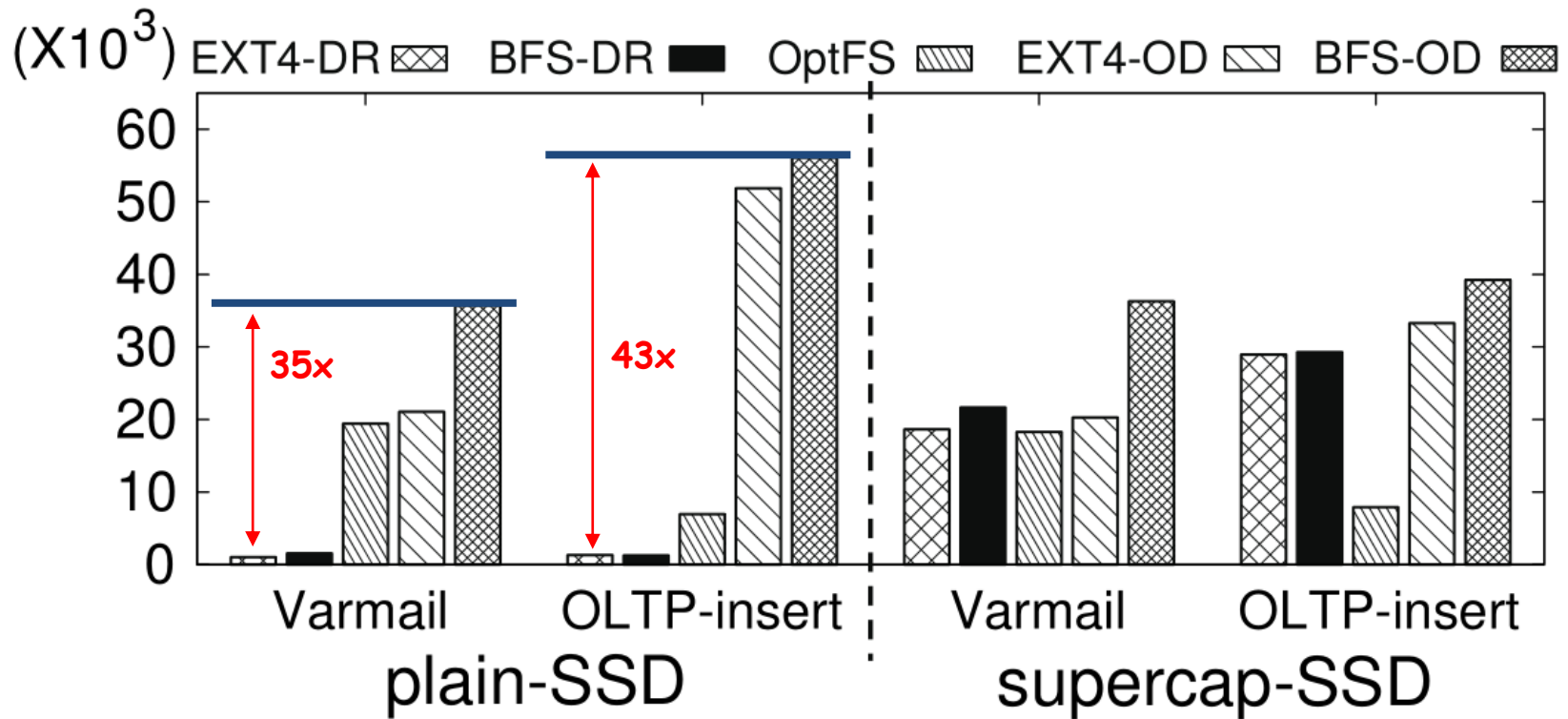**Concurrent Jounrnaling makes Journaling more scalable.**



(a) plain-SSD

(b) supercap-SSD

# Mobile DBMS: SQLite

**Barrier enabled IO stack gets more effective as the parallelism of the Flash storage increases.**



EXT4-DR   BFS-DR   OptFS   EXT4-OD   BFS-OD

**in single channel 2.8x**

**in 8 channel 76.5x**
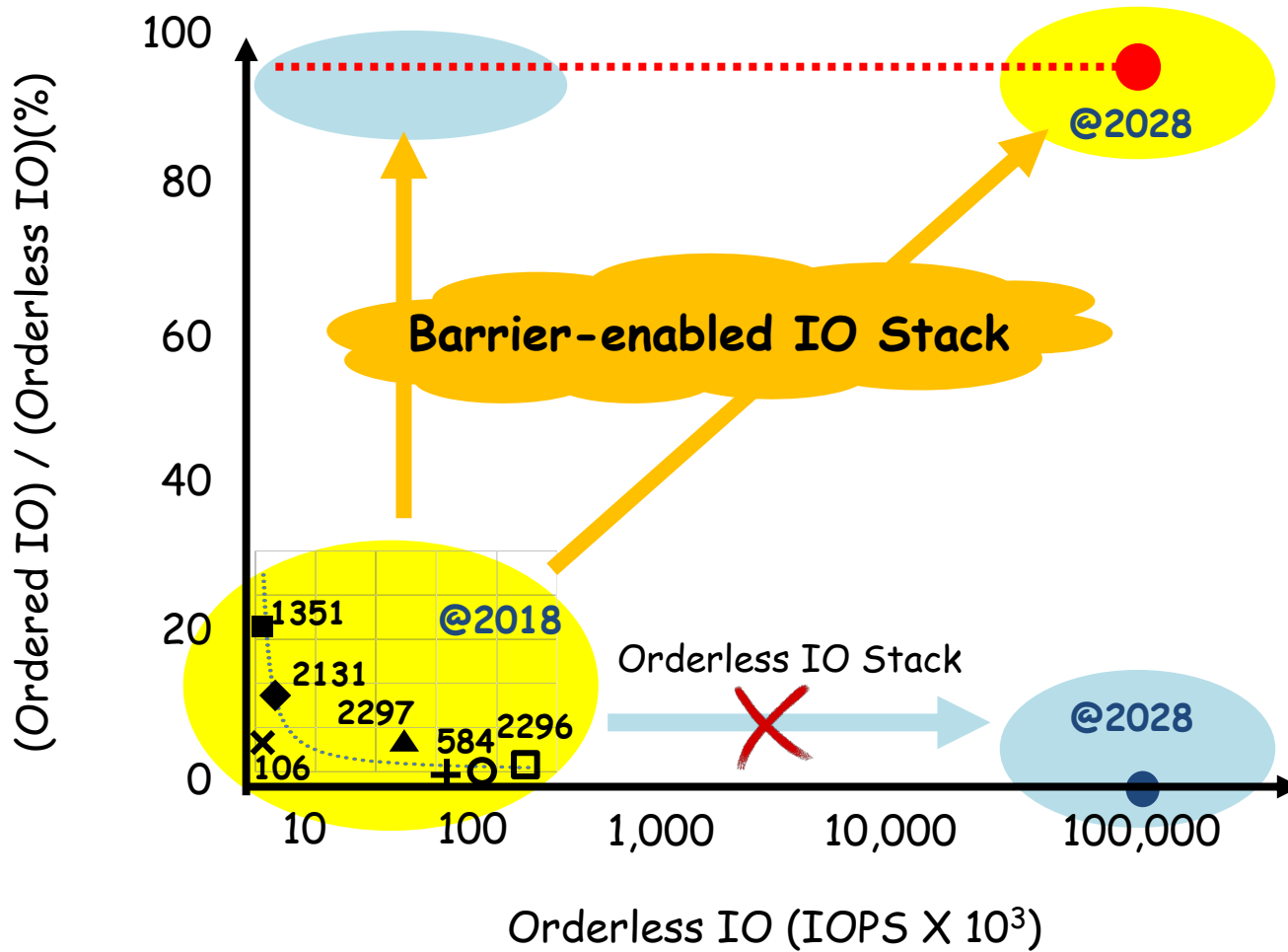
UFS2.0 (GS6)

SSD (850 PRO)

# Server Workload: varmail / Insert(MySQL)

# Conclusion

- Modern IO stack is fundamentally driven by the legacy of rotating media.

- In Flash Storage, the PERSIST order can be controlled while in HDD, it cannot.

- In Barrier-enabled IO stack, we eliminate the Transfer-and-Flush in controlling the storage order.

- To storage vendors,

    "Support for barrier command is a must."

- To service providers,

    "IO stack should eliminate not only the flush overhead

    but also the transfer overhead."

# It is time for a change.

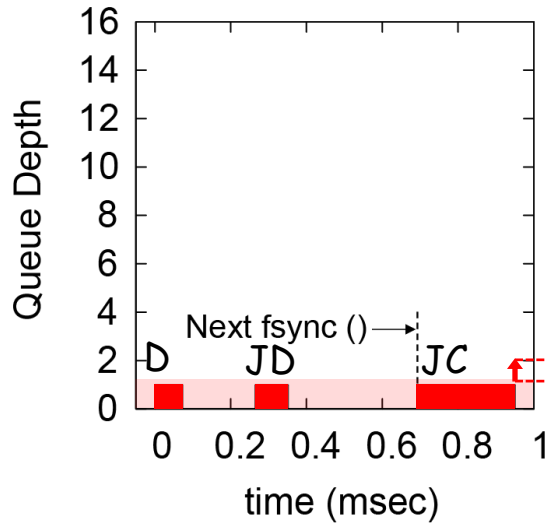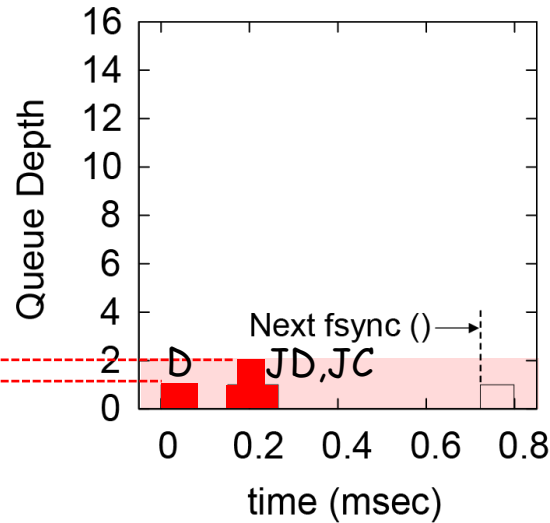**https://github.com/ESOS-Lab/barrieriostack**

# Queue Depth

Epoch 1: {write (D), write (JD) }
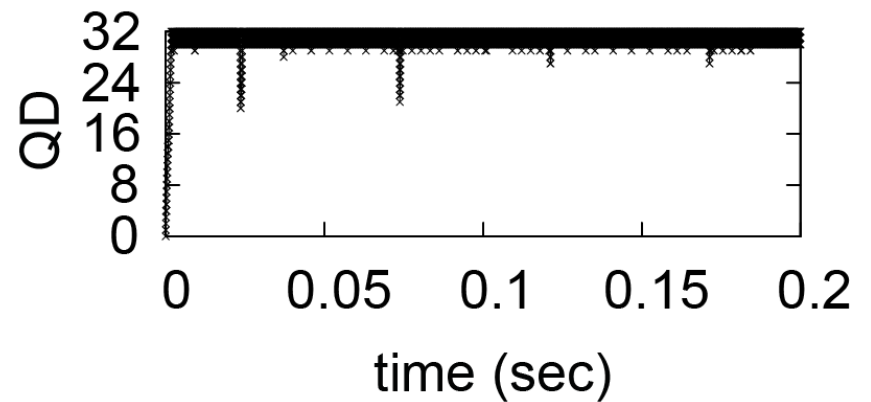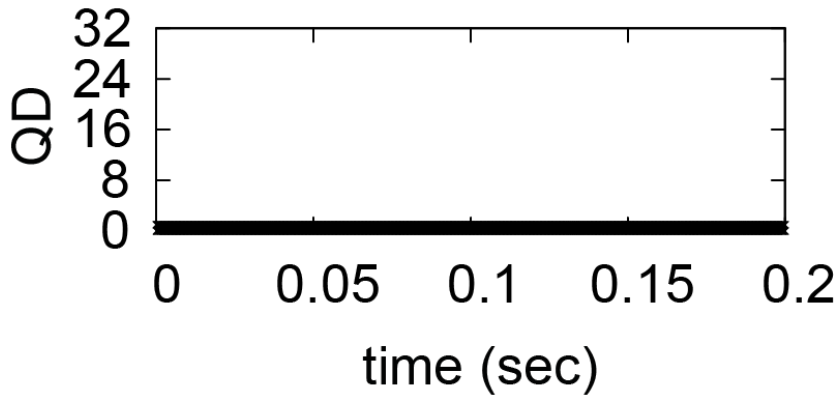
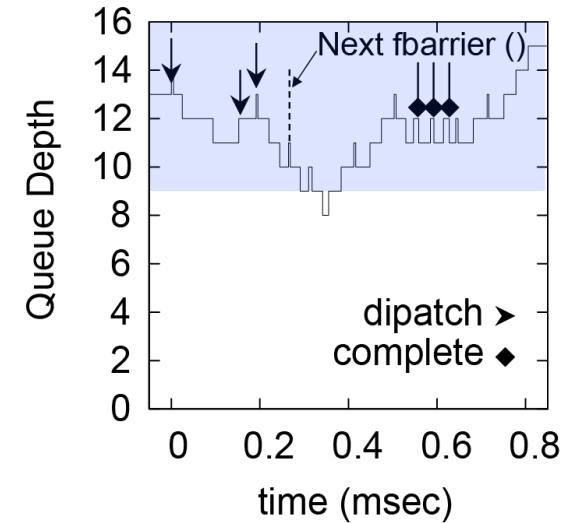Epoch 2: {write (JC)}



fsync() in EXT4          fsync() in BarrierFS          fbarrier() in BarrierFS

Intel X25-M
35 K IOPS
2009

830 PRO
80 K IOPS
2012

850 PRO
100 K IOPS
2014
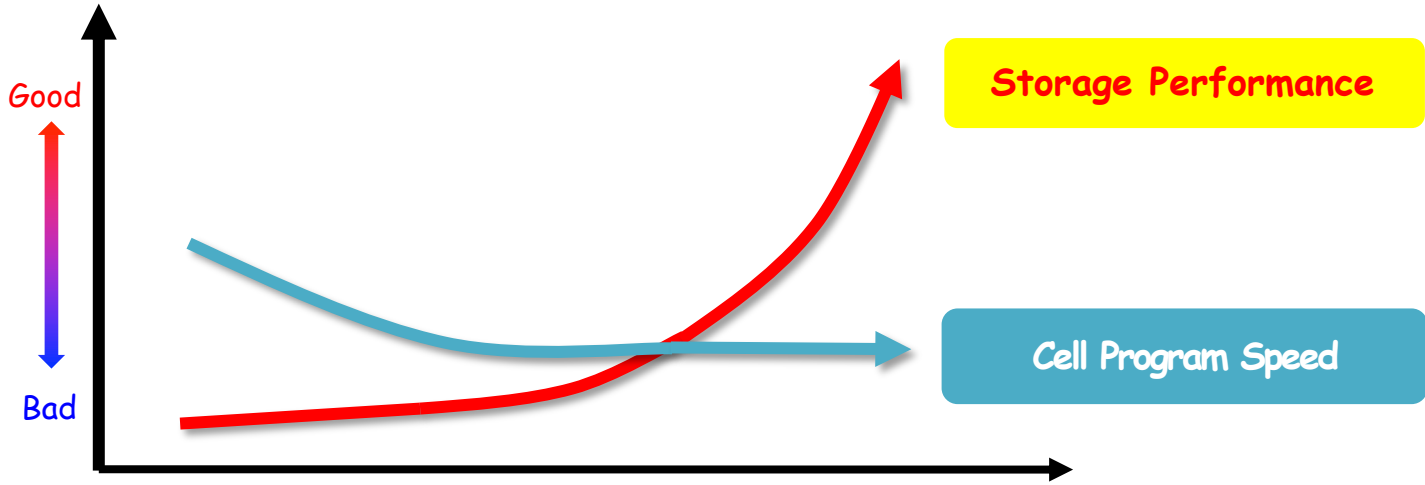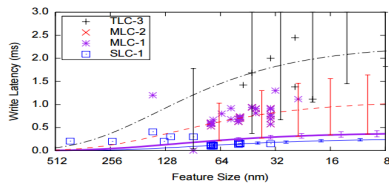
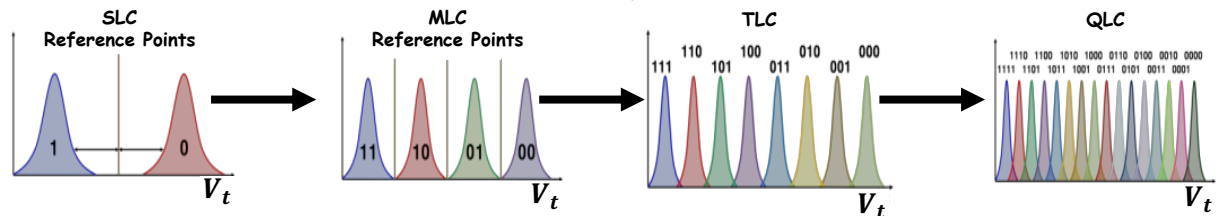Intel 600p
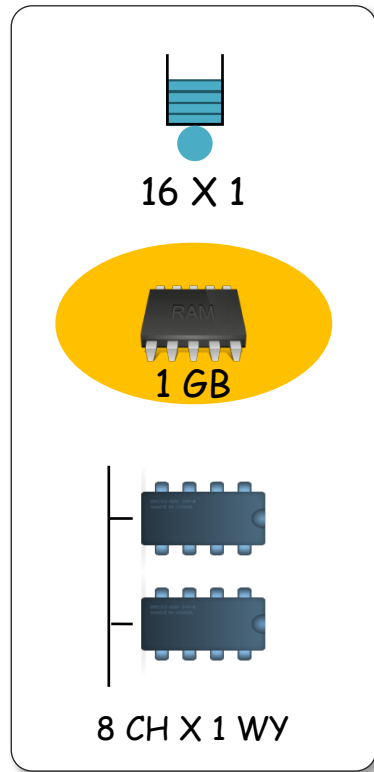155 K IOPS
2016

960 PRO
380 K IOPS
2016

PM1725
1 M IOPS
2015

PCI Gen4
10 M IOPS
2022

Good

Bad

**Storage Performance**

Cell Program Speed

**Finer Process Technology (FAST12)**

**Multi Bits/Cell**

SLC
Reference Points

MLC
Reference Points

TLC

QLC

$V_t$

$V_t$

$V_t$

$V_t$

# Storage Evolution



16 X 1

1 GB

8 CH X 1 WY

**10 K IOPS**

64 K X 64 K

1 TB

64 CH X 2 WY

**1 M IOPS**

# To Mitigate the Transfer-and-Flush overhead

- Eliminate Flush
    - Transactional checksum [IronFS,2005]
    - OptFS [2013], NoFS[2015], FeatherStitch[2007]
    - 'cache barrier'[2005], `nobarrier` option in EXT4[2010]
- Eliminate Transfer

- To reduce frequent fsync() calls
    - Log Structured Merge Tree[1996]
    - Multiple Command Queues [NVMe,2005]

# Dual Mode Journaling: fbarrier()