# BFQ, Multiqueue-Deadline, or Kyber? Performance Characterization of Linux Storage Schedulers in the NVMe Era

**Zebin Ren**[1], Krijn Doekemeijer[1], Nick Tehrany[2], and Animesh Trivedi[1]

[1]VU Amsterdam
[2]BlueOne Business Software LLC

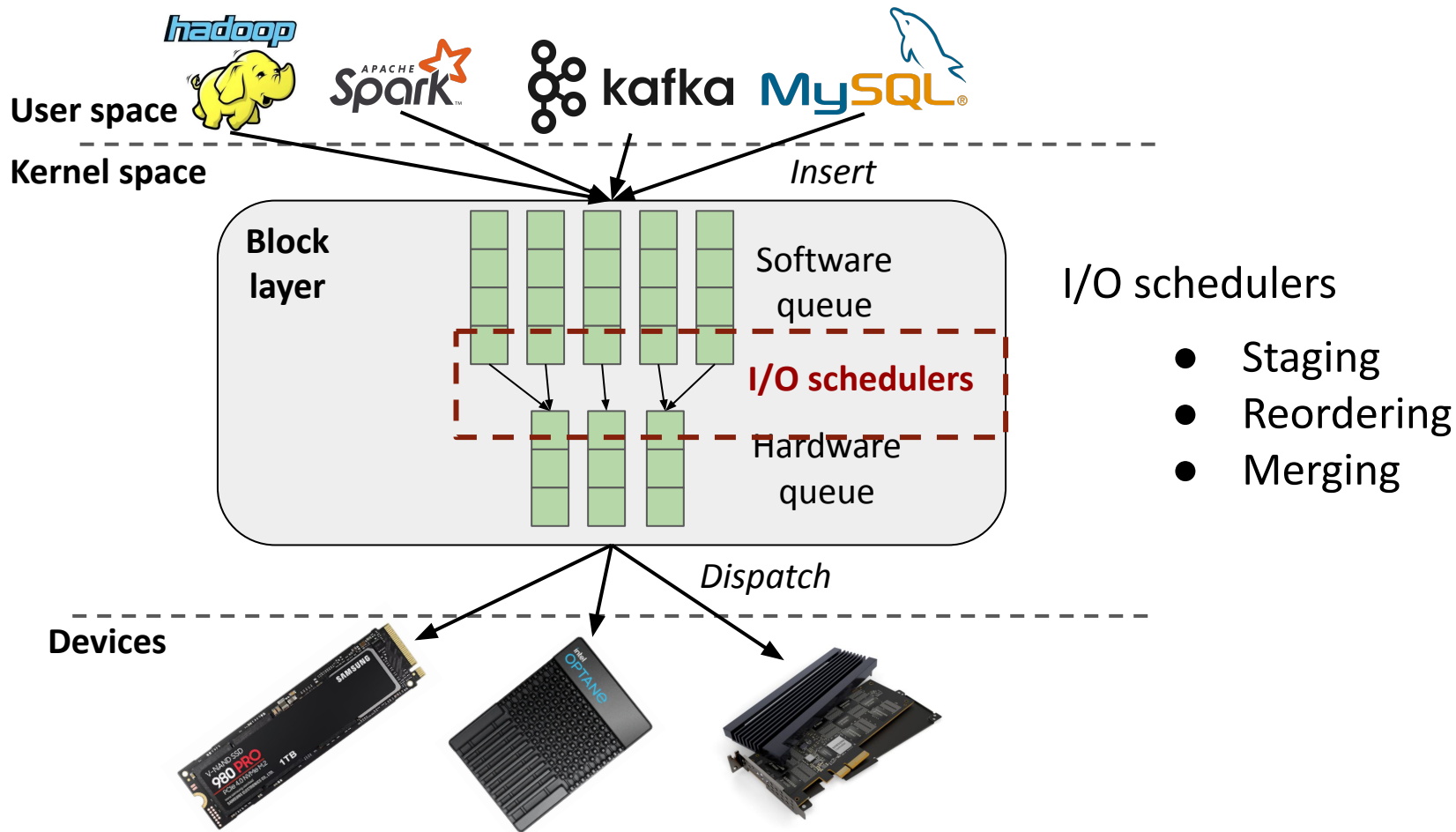This paper won the **best paper award** in ICPE'24.



@Large Research
Massivizing Computer Systems
https://atlarge-research.com/

GRAPH MASSIVIZER

FUTURE NETWORK SERVICES 6G

CLOUD STARS

NWO

VU VRIJE UNIVERSITEIT AMSTERDAM

# Background



QoS guarantees

- Latency
- Throughput

Up to millions of IOPS
< 10 μs latency

# Background: I/O Schedulers

# Background: What has Changed?

1. Huge improvement of storage performance.



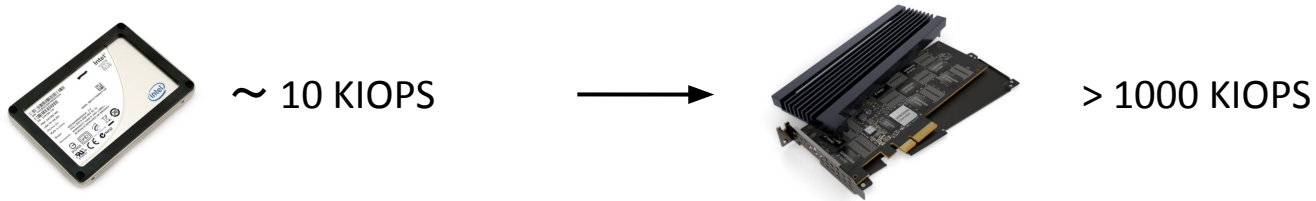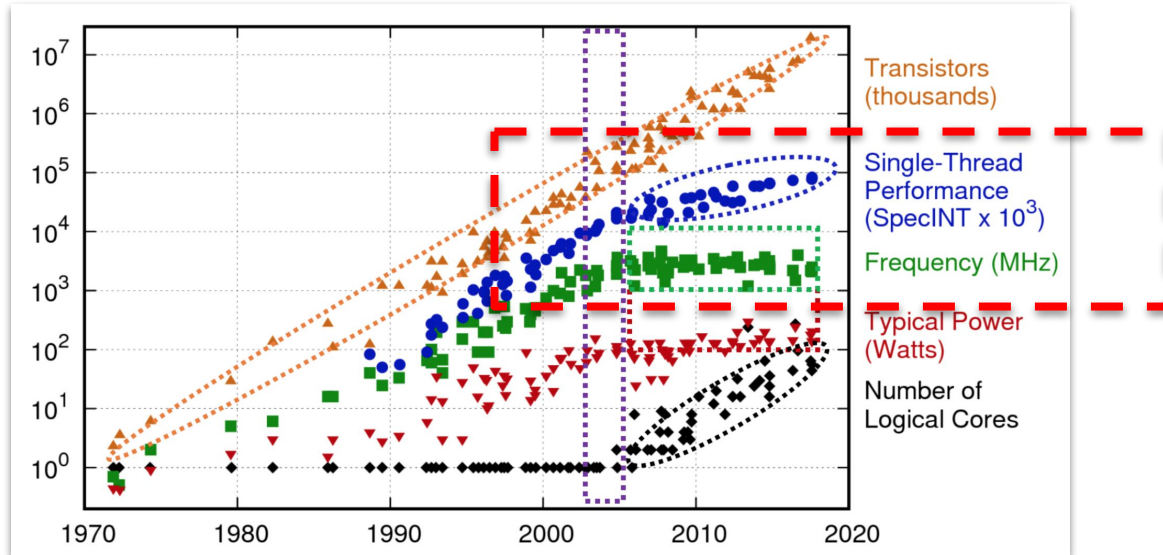~ 10 KIOPS  →  > 1000 KIOPS

2. Improvement of CPU performance stalls.
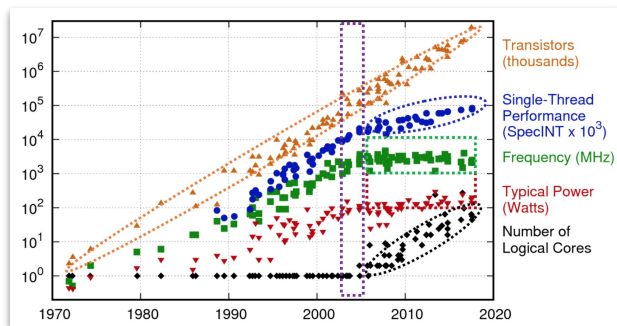
# Background: What has Changed?

1. Huge improvement of storage performance.



~ 10 KIOPS

2. Improvement of CPU performance stalls.



3. Research on I/O schedulers for these SSDs[1-4]

3

# The Linux I/O Schedulers

No plug-and-play implementations.

The most available I/O schedulers?

Linux I/O schedulers!

4

# The Linux I/O Schedulers

No plug-and-play implementations.

The most available I/O schedulers?

Linux I/O schedulers!

**None**
- Least overhead.
- No performance guarantees.

**BFQ**
- Fair-share between apps.
- Complex, high overhead.

**Kyber**
- Designed for fast SSDs.
- Balancing between read and write.

**MQ-Deadline**
- Issues request with increasing sector order.
- Soft latency deadlines.

4

# Setup



fio workloads

CPU
10 cores

io_uring interface

Linux I/O schedulers

8x

6.2 MIOPS

Latency-sensitive application (L-app)

Application

8x

Throughput-bound application (T-app)

Application

128x

...

...

8x

5

# Research Questions

## RQ1: Overhead



**I/O schedulers**

- Latency?
- Throughput?

# Research Questions

**RQ1: Overhead**

- Latency?
- Throughput?

**RQ2: Scalability**

# Research Questions

**RQ1: Overhead**

- Latency?
- Throughput?

**RQ2: Scalability**

- L-apps.
- T-apps.
- SSDs.

**RQ3: Interference**

- Read L-app + increasing write T-apps.

# RQ1. Overheads

# Overhead

**Latency of L-apps**

2.7% higher



**Throughput of T-apps**

36.7% lower



Slightly higher latency, up to 2.7% higher latency.

Significantly lower throughput, up to 36.7% lower.

I/O schedulers → Significantly higher throughput overhead.

# RQ2. Scalability

# Scalability of L-apps

## QD=16



Cumulative prob vs Latency (μs)

Legend:
- None
- BFQ
- Kyber
- MQ-DL

**Better** ←

## QD=32



Cumulative prob vs Latency (μs)

**Better** ←

Higher workload →higher overhead.     Why?

# Scalability of L-apps: CPU Usage



A single CPU core is 100% utilized.

When CPU bottlenecked → higher latency overheads.

# Scalability of T-apps: 1 SSD vs. 8 SSDs

## 1 SSD



Big gap of scalability on throughput between different I/O schedulers.

# Scalability of T-apps: 1 SSD vs. 8 SSDs



Big gap of scalability on throughput between different I/O schedulers.

More devices → better scalability.

# Scalability of T-apps: 1 SSD vs. 8 SSDs



What cause this scalability issue?

# Scalability of T-apps: CPU Usage



The scalability issues are caused by high CPU contention.

# Lock Overhead of I/O Schedulers



BFQ and MQ-Deadline → high CPU lock overhead.

Adding devices mitigates the lock overhead.

**January, 2024**: Identified by the Linux kernel developers[5][6].

# RQ3. Taming I/O Interference

# 1 L-app (R) + Increasing T-apps (W)



**Foreground**

**Background**

Read L-app — 4 KB

+

Write T-app
...
64 KB — 128x

Increasing

# 1 L-app (R) + Increasing T-apps (W)

**Foreground**

Read
L-app

4
KB

**+**

**Background**

Write T-app
...
64
KB    128x

Increasing



BFQ and Kyber → low latency for the foreground L-app.

# Conclusions

RQ1: What is the overhead of Linux I/O schedulers?

- Minor latency overhead.
- Significantly throughput overhead.

RQ2: What is the scalability of Linux I/O schedulers?

- Latency → depends on CPU.
- Throughput, BFQ and MQ-DL → high lock contention.
- Throughput, Kyber → good, similar to None.

RQ3: Can the Linux I/O schedulers tame I/O inference?

- Only BFQ and Kyber can
  provide bounded performance.

**Throughput of T-apps**

36.7% lower

QD=32

# Take-Home Messages

1. *I/O Schedulers can influence the performance significantly*.
   None has the lowest overhead and highest scalability.
   BFQ has the highest overhead and lowest scalability.

2. *Different schedulers have different locking and scaling overheads.*
   BFQ = MQ-Deadline > Kyber > None.

3. Use **Kyber** to prioritize **foreground reads with background writes.**
   HotCloudPerf'24 *A Systematic Configuration Space Exploration of the Linux Kyber I/O Scheduler*

Paper: https://atlarge-research.com/pdfs/2024-io-schedulers.pdf
Source code: https://github.com/ZebinRen/icpe24_io_scheduler_study_artifact

GitHub

PDF

# Thank you!
# Questions?

Paper: https://atlarge-research.com/pdfs/2024-io-schedulers.pdf
Source code: https://github.com/ZebinRen/icpe24_io_scheduler_study_artifact

# Resources

## Images used:

https://www.samsung.com/nl/memory-storage/nvme-ssd/980-pro-pcle-4-0-nvme-m-2-ssd-1tb-mz-v8p1t0bw/
https://www.intel.com/content/www/us/en/products/details/memory-storage/data-center-ssds/optane-dc-ssd-series.html
https://www.anandtech.com/show/12376/samsung-launches-zssd-sz985-up-to-800gb-of-znand
https://www.storagereview.com/review/intel-x25-v-ssd-review-40gb

## References

[1] Till Miemietz, Hannes Weisbach, Michael Roitzsch, Hermann Härtig: K2: Work-Constraining Scheduling of NVMe-Attached Storage. RTSS 2019: 56-68
[2] Mohammad Hedayati, Kai Shen, Michael L. Scott, Mike Marty: Multi-Queue Fair Queuing. USENIX Annual Technical Conference 2019: 301-314 2018
[3] Jaehyun Hwang, Midhul Vuppalapati, Simon Peter, Rachit Agarwal: Rearchitecting Linux Storage Stack for µs Latency and High Throughput. OSDI 2021: 113-128
[4] Jiwon Woo, Minwoo Ahn, Gyusun Lee, Jinkyu Jeong: D2FQ: Device-Direct Fair Queueing for NVMe SSDs. FAST 2021: 403-415
[5] https://www.phoronix.com/news/BFQ-IO-Better-Scalability
[6] https://www.phoronix.com/news/MQ-Deadline-Scalability

# Resources

## Linux I/O schedulers

1. BFQ (Budget Fair Queueing) https://www.kernel.org/doc/html/latest/block/bfq-iosched.html
2. Two new block I/O schedulers for 4.12 https://lwn.net/Articles/720675/
3. Deadline IO scheduler tunables
https://docs.kernel.org/block/deadline-iosched.html#:~:text=The%20goal%20of%20the%20deadline,value%20in%20units%20of%20milliseconds.
4. BFQ I/O Scheduler For Linux Sees Big Scalability Improvement https://www.phoronix.com/news/BFQ-IO-Better-Scalability
5. MQ-Deadline Scheduler Optimized For Much Better Scalability

## New I/O schedulers

1. Myoungsoo Jung, Wonil Choi, Shekhar Srikantaiah, Joonhyuk Yoo, and Mahmut T. Kandemir. HIOS: A Host Interface I/O Scheduler for Solid State Disks. ISCA 2014.
2. Mingyang Wang and Yiming Hu. An I/O Scheduler Based on Fine-Grained Access Patterns to Improve SSD Performance and Lifespan. In Symposium on Applied Computing, SAC 2014.
3. Hui Lu, Brendan Saltaformaggio, Ramana Rao Kompella, and Dongyan Xu. vFair: Latency-Aware Fair Storage Scheduling via per-IO Cost-Based Differentiation. SoCC 2015.
4. Jiayang Guo, Yiming Hu, Bo Mao, and Suzhen Wu. Parallelism and Garbage Collection Aware I/O Scheduler with Improved SSD Performance. IPDPS 2017.
5. Minhoon Yi, Minho Lee, and Young Ik Eom. 2017. CFFQ: I/O Scheduler for Providing Fairness and High Performance in SSD Devices. IMCOM 2017.
6. Mohammad Hedayati, Kai Shen, Michael L. Scott, and Mike Marty. Multi- Queue Fair Queuing. In 2019 USENIX Annual Technical Conference, USENIX ATC 2019.
7. Till Miemietz, Hannes Weisbach, Michael Roitzsch, and Hermann Härtig. K2: Work-Constraining Scheduling of NVMe-Attached Storage. RTSS 2019.
8. Jaehyun Hwang, Midhul Vuppalapati, Simon Peter, and Rachit Agarwal. Rearchitecting Linux Storage Stack for $\mu$s Latency and High Throughput. OSDI 2021.
9. Jiwon Woo, Minwoo Ahn, Gyusun Lee, and Jinkyu Jeong. D2FQ: Device- Direct Fair Queueing for NVMe SSDs. FAST 2021.
10. Jieun Kim, Dohyun Kim, and Youjip Won Fair I/O Scheduler for Alleviating Read/Write Interference by Forced Unit Access in Flash Memory. HotStorage 2022.
11. Caeden Whitaker, Sidharth Sundar, Bryan Harris, and Nihat Altiparmak. Do We Still Need I/O Schedulers for Low-Latency Disks?. HotStorage 2023.

# Backup Slides

# CPU or NVMe SSD, What is the Bottleneck?



To saturate a SSD:

Multiple I/O requests

Enough CPU resources → 4 processes

4

# SSD Performance



(a) Vary request sizes, 1 SSD.

(b) Vary # processes with 1 SSD.

(c) Vary # processes with 8 SSDs.

# L-app Scalability



**Figure 2: Intra-process scalability latency CDFs with increasing queue depth (QD); Note the different x-axis scale for (e).**

(a) QD = 1 — None: 77.3 μs, BFQ: 79.4 μs, Kyber: 78.3 μs, MQ-DL: 78.3 μs

(b) QD = 16 — None: 136.2 μs, BFQ: 136.2 μs, Kyber: 136.2 μs, MQ-DL: 136.2 μs

(c) QD = 32 — None: 183.3 μs, BFQ: 199.7 μs, Kyber: 189.4 μs, MQ-DL: 185.3 μs

(d) QD = 64 — None: 276.5 μs, BFQ: 329.7 μs, Kyber: 305.2 μs, MQ-DL: 297.0 μs

(e) QD = 256 — None: 806.9 μs, BFQ: 1,155.1 μs, Kyber: 1,003.5 μs, MQ-DL: 946.2 μs

**Figure 3: Inter-process scalability latency CDFs with increasing number of L-apps; Note the different x-axis scale for (e).**

(a) 1 L-app — None: 77.3 μs, BFQ: 84.8 μs, Kyber: 78.3 μs, MQ-DL: 78.3 μs

(b) 16 L-apps — None: 205.8 μs, BFQ: 236.5 μs, Kyber: 226.3 μs, MQ-DL: 218.1 μs

(c) 32 L-apps — None: 419.8 μs, BFQ: 485.4 μs, Kyber: 456.7 μs, MQ-DL: 448.5 μs

(d) 64 L-apps — None: 897.0 μs, BFQ: 1,028.1 μs, Kyber: 962.6 μs, MQ-DL: 938.0 μs

(e) 256 L-apps — None: 3,883.0 μs, BFQ: 4,292.6 μs, Kyber: 4,227.1 μs, MQ-DL: 4,112.4 μs

33

# L-app CPU cost



(a) Intra-process

(b) Inter-process

**Figure 4: CPU usage for intra/inter-process concurrency.**

# L-app Scalability Heatmap



(a) 1 core, 1 process
y-axis: QD

(b) 1 core, QD 1
y-axis: # processes

(c) 10 cores, QD 1
y-axis: # processes

# SSD Scalability



(a) Total throughput    (b) Total CPU usage

**Figure 9: T-app inter-process scalability (10 cores, 10 concurrent T-4KiB-apps) with an increasing number of SSDs.**

# L-app Interference



**(a) T-4KiB-app (reads)**

**(b) T-64KiB-app (reads)**

**(c) T-4KiB-app (writes)**

**(d) T-64KiB-app (writes)**

**Figure 11: L-app tail latency with an increasing number of interfering background applications;** *Note: scales differ on the y-axis and they are in Milliseconds!*

# T-app Interference



(a) T-64KiB-apps (reads)

(b) T-64KiB-apps (writes)

**Figure 12: Read throughput (IOPS) of a T-4KiB-app workload with an increasing number of interfering background T-64KiB-app workload. Note: The y-axis is log-scale.**

# 1 T-app (R) + Increasing T-apps (W)

**Foreground**

Read T-app
...
4 KB    128x

**+**

**Background**

Write T-app
...
64 KB    128x

Increasing



Foreground: Read T-app

Legend:
- None
- Kyber
- BFQ
- MQ-DL

X-axis: Number of T-app (64KiB), values 0, 1, 2, 4, 8, 16, 32
Y-axis: Throughput (KIOPS), $10^{-1}$ to $10^{3}$

BFQ and Kyber → higher bandwidth for the foreground T-app.

# Lock in the I/O Schedulers

In block/mq-deadline.c

```c
84   struct deadline_data {
85           /*
86            * run time data
87            */
88
89           struct dd_per_prio per_prio[DD_
90
91           /* Data direction of latest dis
92           enum dd_data_dir last_dir;
93           unsigned int batching;
94           unsigned int starved;
95
96           /*
97            * settings that change how the
98            */
99           int fifo_expire[DD_DIR_COUNT];
100          int fifo_batch;
101          int writes_starved;
102          int front_merges;
103          u32 async_depth;
104          int prio_aging_expire;
105
106          spinlock_t lock;
107          spinlock_t zone_lock;
108  };
109
```

```c
826      */
827  static void dd_insert_requests(struct blk_mq_hw_ctx *hctx,
828                                 struct list_head *list, bool at_head)
829  {
830          struct request_queue *q = hctx->queue;
831          struct deadline_data *dd = q->elevator->elevator_data;
832
833          spin_lock(&dd->lock);
834          while (!list_empty(list)) {
835                  struct request *rq;
836
```

```c
572  static struct request *dd_dispatch_request(struct blk_mq_hw_ctx *hctx)
573  {
574          struct deadline_data *dd = hctx->queue->elevator->elevator_data;
575          const unsigned long now = jiffies;
576          struct request *rq;
577          enum dd_prio prio;
578
579          spin_lock(&dd->lock);
580          rq = dd_dispatch_prio_aged_requests(dd, now);
581          if (rq)
582                  goto unlock;
```

# Lock in the I/O Schedulers   Reduced lock contention

Dispatch

```
      /* Maps an I/O priority class to a deadline scheduler priority. */
@@ -600,6 +607,15 @@ static struct request *dd_dispatch_request(struct blk_mq_hw_ctx *hctx)
        struct request *rq;
        enum dd_prio prio;

+       /*
+        * If someone else is already dispatching, skip this one. This will
+        * defer the next dispatch event to when something completes, and could
+        * potentially lower the queue depth for contended cases.
+        */
+       if (test_bit(DD_DISPATCHING, &dd->run_state) ||
+           test_and_set_bit(DD_DISPATCHING, &dd->run_state))
+               return NULL;
+
        spin_lock(&dd->lock);
        rq = dd_dispatch_prio_aged_requests(dd, now);
        if (rq)
@@ -616,6 +632,7 @@ static struct request *dd_dispatch_request(struct blk_mq_hw_ctx *hctx)
        }
```

# Lock in the I/O Schedulers   Reduced lock contention

Insertion

```
+/*
+ * If we can grab the dd->lock, then just return and do the insertion as per
+ * usual. If not, add to one of our internal buckets, and afterwards recheck
+ * if if we should retry.
+ */
+static bool dd_insert_to_bucket(struct deadline_data *dd,
+                                struct list_head *list, int *seq)
+        __acquires(&dd->lock)
+{
+       struct dd_bucket_list *bucket;
+       int next_seq;
+
+       *seq = atomic_read(&dd->insert_seq);
+
+       if (spin_trylock(&dd->lock))
+               return false;
+       if (!test_bit(DD_INSERTING, &dd->run_state)) {
+               spin_lock(&dd->lock);
+               return false;
+       }
+
+       *seq = atomic_inc_return(&dd->insert_seq);
+
+       bucket = &dd->bucket_lists[get_cpu() & DD_CPU_BUCKETS_MASK];
+       spin_lock(&bucket->lock);
+       list_splice_init(list, &bucket->list);
+       spin_unlock(&bucket->lock);
+       put_cpu();
+
```

https://lore.kernel.org/linux-block/ede4179c-8fa5-4496-ac21-4e3fda41df81@kernel.dk/

# Lock in the I/O Schedulers  Reduced lock contention

Results

```
Device              QD       Jobs      IOPS     Lock contention
========================================================================
null_blk            4        32        1090K    92%
nvme0n1             4        32        1070K    94%
```

```
With that in place, the same test case now does:

Device           QD       Jobs      IOPS     Contention          Diff
========================================================================
null_blk         4        32        2250K    28%                 +106%
nvme0n1          4        32        2560K    23%                 +112%
```

# Who Are We/Am I?