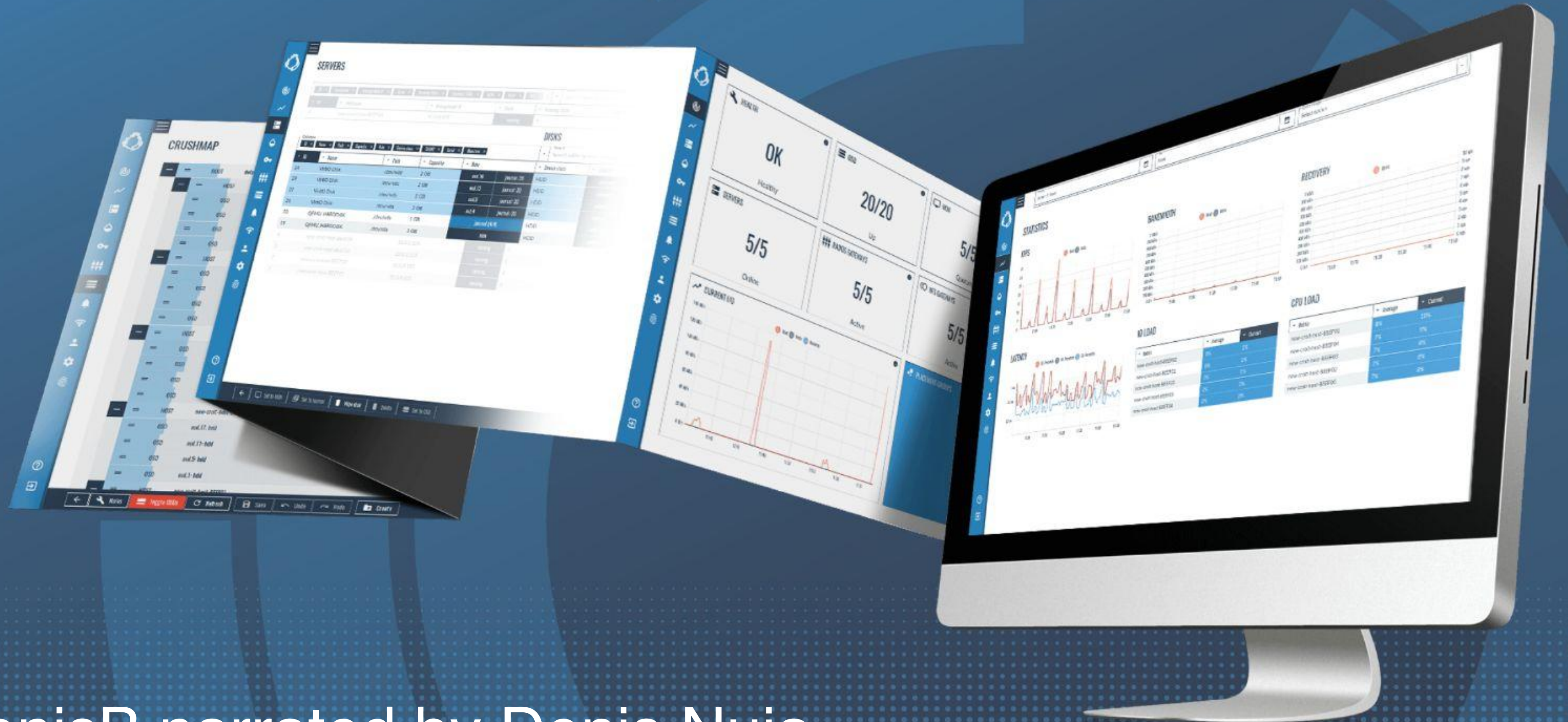


Implementing SPDK DAOS bdev module



A journey by DenisB narrated by Denis Nuja

Why

At the moment DAOS provides just a few connectivity options:

- Direct Key/Value and Array APIs - highly app specific.
- Direct DFS API (partially POSIX complaint) - also app level.
- DFuse (mounting like a volume) - system wide, but with it's own limitations.

What if my application works only with block devices ?

NVMeOF client is already built in into many platforms.

Why

Implementing SPDK bdev module allows to benefit from:

- NVMe-oF target
- iSCSI target
- vhost server
- NVMe-oF multipath
- BDEVs stacking (raid, lvols ...)
- Dev and perf tooling
- Existing unit tests (who likes to write new ones !)
- Probably almost every upcoming features and improvements

How

SPDK has great documentation and a wide range of existing bdev implementations.

The idea was to grab **malloc** or **aio** bdev sources and replace ***open/read/write*** functions with DAOS API and see what's going to happen.

How

SPDK has great documentation and a wide range of existing bdev implementations.

The idea was to grab **malloc** or **aio** bdev sources and replace ***open/read/write*** functions with DAOS API and see what's going to happen.

... it's actually worked !

How

SPDK bdev's bare minimum operation list:

```
switch (bdev_io->type) {  
    // Block IO  
case SPDK_BDEV_IO_TYPE_READ:  
case SPDK_BDEV_IO_TYPE_WRITE:  
    // ...  
case SPDK_BDEV_IO_TYPE_RESET:  
    // Wait until in-flight requests are done  
case SPDK_BDEV_IO_TYPE_FLUSH:  
    // NOOP (DAOS persistent on write)  
case SPDK_BDEV_IO_TYPE_UNMAP:  
    // DAOS supports punching holes  
}
```

How

SPDK_BDEV_IO_TYPE_RESET - there is no reset functionality in DAOS client (and probably could not be).

Current implementation just waits until all in-flight IO requests are completed.

Rest of the command handlers fit naturally in SPDK bdev model.

How

DAOS Array API vs DFS API, which one to use ?:

```
int daos_array_open(...) vs int dfs_open(...);  
int daos_array_read(...) vs int dfs_read(...);  
int daos_array_write(...) vs int dfs_write(...);  
int daos_array_punch(...) vs int dfs_punch(...);
```


How

DAOS Array API vs DFS API, which one to use ?

- The APIs are pretty similar
- Both support batch read and write
- Both provides async mode with DAOS event queue
- DFS API uses Array API underneath, should not be any performance penalties
- DFS API has a slight operational advantage - one could mount DAOS container with **dfuse** and copy the backing file to some other storage.

How

DAOS Array API vs DFS API, which one to use ?

- The APIs are pretty similar
- Both support batch read and write
- Both provides async mode with DAOS event queue
- DFS API uses Array API underneath, should not be any performance penalties
- DFS API has a slight operational advantage - one could mount DAOS container with **dfuse** and copy the backing file to some other storage.
- **The winner is DFS API as it looks familiar to POSIX API.**

How

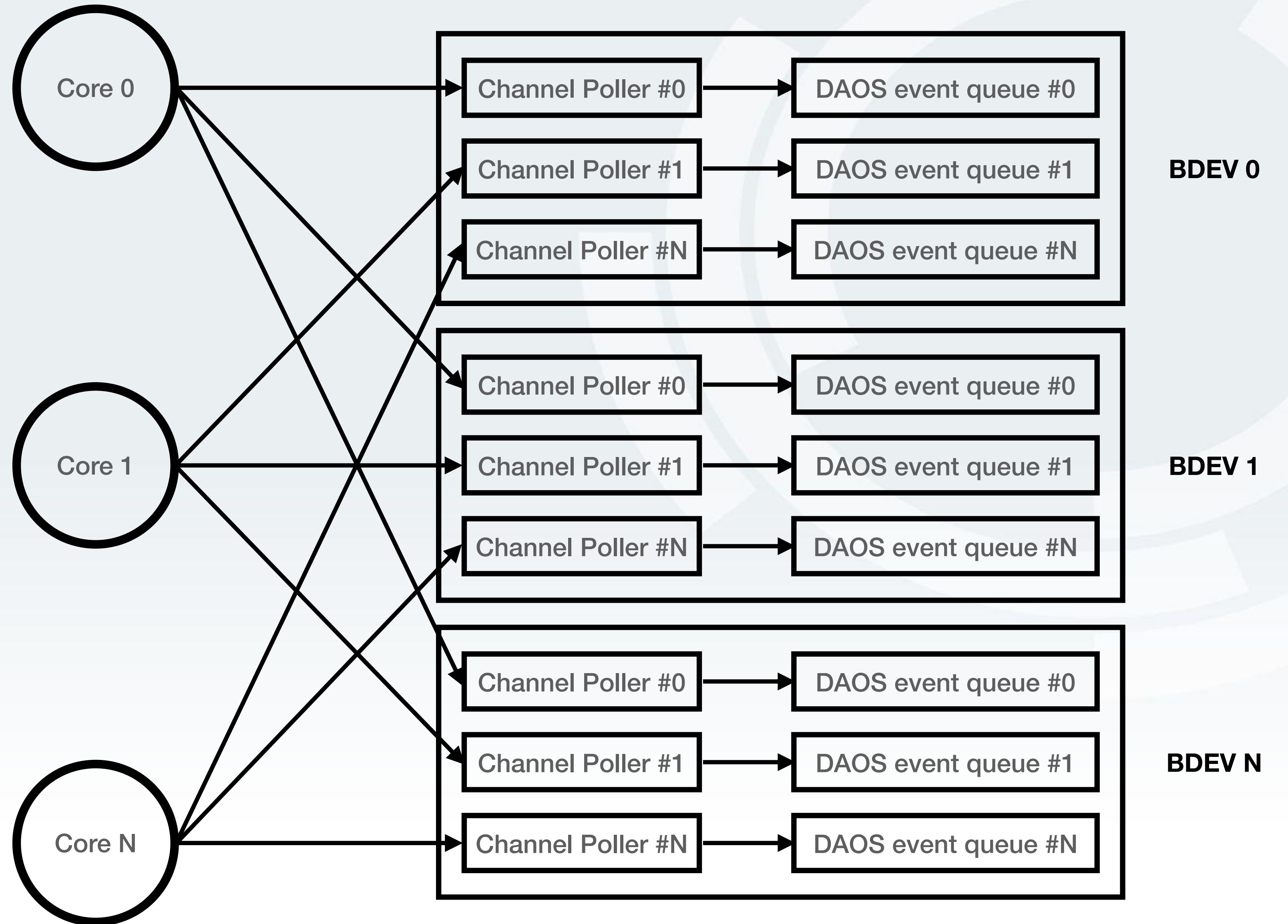
Initial version event handling approach

Every device's channel has its own poller with DAOS event queue attached.

This approach allows almost linear performance scaling (up to certain point) with increasing number of cores allocated to the target.

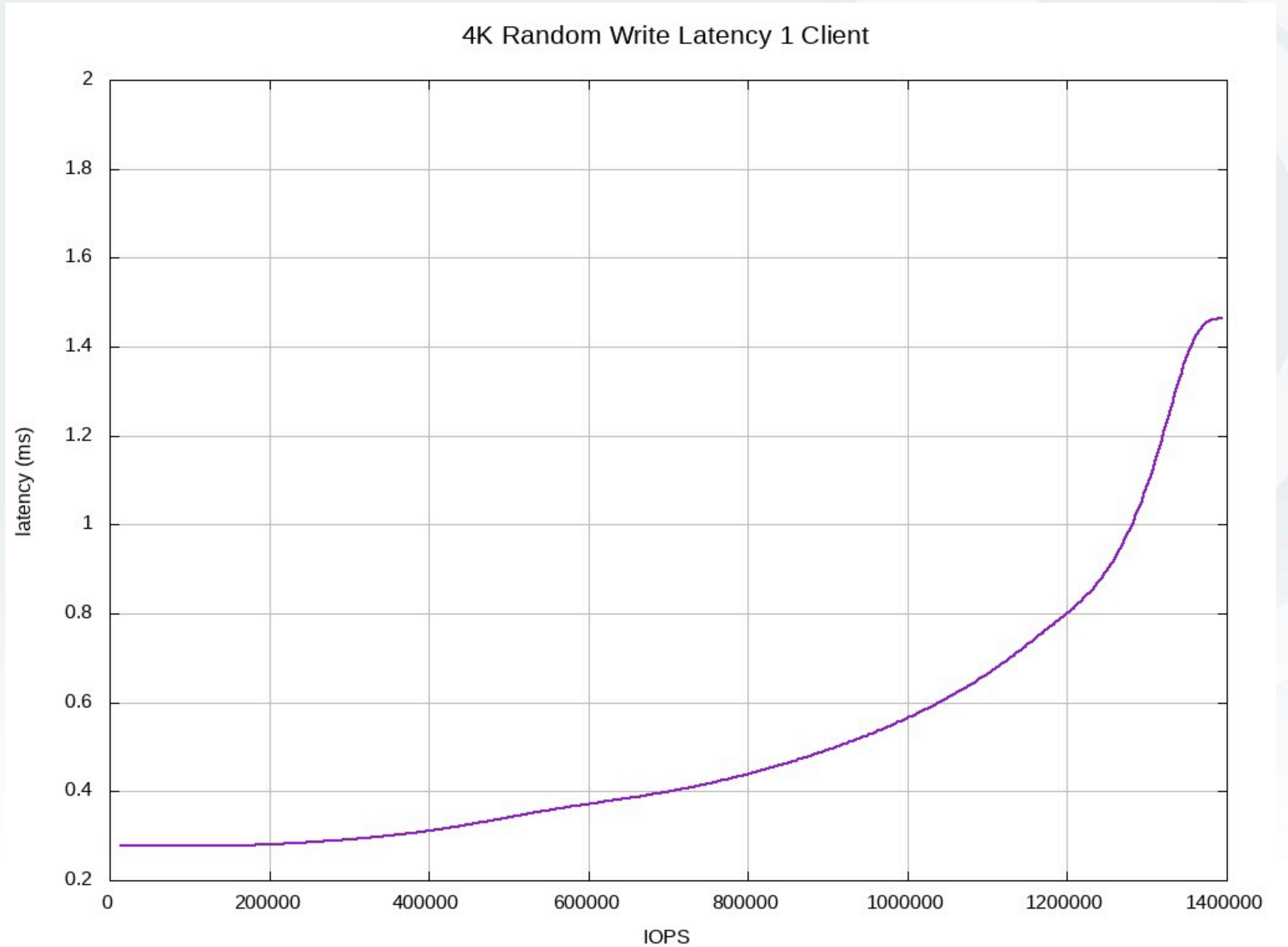
How

It looks like this:



**Latency vs IOPS graph
4k random write
1 client (vdbench)**

**nonRDMA Ethernet TCP
on DAOS cluster and client**



Pitfalls along the way

Prior DAOS v2.2 there wasn't per thread pool/cont cache which significantly affected multithreaded performance (but not multiprocess one).

With v2.2 release cache and more granular locking were introduced resolving the issue and increasing the multithreaded performance.

More details: <https://daosio.atlassian.net/browse/DAOS-11230>

Pitfalls along the way

MxN problem: M devices create N channels/event queue each. As turned out, the number of DAOS event queue:

- fabric provider specific and limited
- does not scale after certain limit

For instance, for TCP provider the limit is 256 queues, which means 8 bdev with 32 cores enabled.

Pitfalls along the way

Over provisioning: there's no way to reserve space inside the DAOS container. This might create an issue when someone creates more bdevs than there are free space and the error would happened on the writes.

Next iteration

Fixed pool of DAOS event queues and possible the fixed number of poller across all DAOS bdevs to distribute load evenly.

Better recovery procedure on errors: in the current version the module treats most of the errors as recoverable, which might not be always the case and proper self shutdown is preferable.

Various QoL improvements.

Acknowledgement

SPDK community in Slack, especially James Harris for patiently answering all silly questions, suggestion and review.

DAOS team, in particularly, Johann Lombardi and Mohamad Chaarawi for helping with never ending stream of questions!

How are actually use it in the Croit Platform

The screenshot shows the 'Create NVMe-oF block device' configuration page in the Croit platform. The page is titled 'Create NVMe-oF block device' and is for a device named 'new-bdev'. The configuration includes:

- Group:** daos
- Type:** DAOS
- Dataset:** Cluster: daos | Pool: benchmark | Dataset: benchmark-replicated
- Size:** 1 TiB
- Blocksize:** 4096
- Limit IOPS:** 10000
- Logical Volume Store:** A toggle switch is currently turned off, with the text 'Whether the device should be used as a logical volume store, or a simple bdev.'

At the bottom of the configuration panel, there are 'SAVE' and 'CLOSE' buttons. The left sidebar shows navigation options: Status, Servers, Config, Gateways, DAOS, Compute, Health, Settings, and Logout.

The screenshot shows the 'SPDK info' page in the Croit platform, displaying the configuration for the 'new-bdev' device as a JSON object. The configuration includes:

```
1 {
2   "name": "613319c4-7994-4de5-b7fe-e5e5fa82b788",
3   "aliases": [
4     "4678d701-ba76-4134-b5ee-6d6a28aef681/416b234a-6f0d-4114-919a-3e906c3c40dd"
5   ],
6   "product_name": "Logical Volume",
7   "block_size": 4096,
8   "num_blocks": 52428800,
9   "uid": "613319c4-7994-4de5-b7fe-e5e5fa82b788",
10  "assigned_rate_limits": {
11    "rw_ios_per_sec": 0,
12    "rw_mbytes_per_sec": 0,
13    "r_mbytes_per_sec": 0,
14    "w_mbytes_per_sec": 0
15  },
16  "claimed": false,
17  "zoned": false,
18  "supported_io_types": {
19    "read": true,
20    "write": true,
21    "unmap": true,
22    "write_zeroes": true,
23    "flush": false,
24    "reset": true,
25    "compare": false,
26    "compare_and_write": false,
27    "abort": false,
28    "nvme_admin": false,
29    "nvme_io": false
30  },
31  "driver_specific": {
32    "lvol": {
33      "lvol_store_uid": "410d97ef-bbfc-4483-b79b-02463cfa6e9f",
34      "base_bdev": "4678d701-ba76-4134-b5ee-6d6a28aef681",
35      "thin_provision": false,
36      "snapshot": false,
37      "clone": false
38    }
39  }
40 }
41
42
43
44
```

The left sidebar is identical to the previous screenshot, showing navigation options: Status, Servers, Config, Gateways, DAOS, Compute, Health, Settings, and Logout.

Thank you for your
attention!

www.croit.io
contact@croit.io