



# dRAID: Parity Declustered RAID for ZFS

Isaac Huang

# The Problem

RAIDz has many great features, but resilver is slow:

- No write hole, self healing, no special hardware

Pool config	Age	Used	Avail	Resilvered	Time	Speed
3x 12-drive RAIDz2	2 Y	120T	23.6T	4.24T	27h2m	45.68 MB/s
1x 4-drive RAIDz1	4 Y	6.42T	1.23T	2.29T	45h45m	14.58 MB/s
3x 6-drive RAIDz1		8.47T		560G	6h32m	24.00 MB/s
5x 6-drive RAIDz1		109T		3.24T	44h52m	21.03 MB/s

RAIDz Resilver Time

# The Causes

## Why is RAIDz resilver slow?

- Random IO due to block pointer tree traversal
  - Some mitigation with “sequential resilver” work
- Write throughput of a single replacement drive is a bottleneck
- A single VDEV doesn't scale to a large number of child VDEVs
  - Aggregated read throughput for reconstruction is limited

# The “Solutions”

## Best “solutions” available today

- Reduce the number of block pointers by using large (1MB+) blocks
- Increase available total IOPS
  - Use narrow RAIDz (for example: 4-wide RAIDz1 or mirroring)
  - Use lots of small disks or SSDs
- None of these solutions address:
  - Single replacement drive bottleneck
  - Limited aggregate read throughput from child VDEVs

# dRAID: No BP Tree Traversal

## RAIDz1

Disk \ LBA	A	B	C	D	E
0	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	D <sub>3</sub>
1	P <sub>1</sub>	D <sub>4</sub>	D <sub>5</sub>	D <sub>6</sub>	D <sub>7</sub>
2	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	P <sub>0</sub>
3	D <sub>0</sub>	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	X



## dRAID1

Disk \ LBA	A	B	C	D	E
0	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	D <sub>3</sub>
1	P <sub>1</sub>	D <sub>4</sub>	D <sub>5</sub>	D <sub>6</sub>	D <sub>7</sub>
2	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	X
3	P <sub>0</sub>	D <sub>0</sub>	X	X	X
4	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	X	X

# dRAID: Sequential Rebuild

- No block pointer tree traversal
  - Completely sequential like traditional RAID-4/5/6
  - Skips free space by scanning spacemap objects
  - Large IO: not limited by block boundary
- No write hole, like RAIDz-1/2/3
  - RAIDz: variable stripe width
  - dRAID: fixed stripe width, but always allocate full stripes

# dRAID: Parity Declustering

RAIDz1-0					RAIDz1-1					Hot Spare
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>
0	1	2	3	4	5	6	7	8	9	<u>10</u>

Declustering →

dRAID1-0										
Group 0					Group 1					Logical Spare
1	4	5	9	3	2	8	10	7	6	<u>0</u>
2	5	6	10	4	3	9	0	8	7	<u>1</u>
3	6	7	0	5	4	10	1	9	8	<u>2</u>
4	7	8	1	6	5	0	2	10	9	<u>3</u>
5	8	9	2	7	6	1	3	0	10	<u>4</u>
6	9	10	3	8	7	2	4	1	0	<u>5</u>
7	10	0	4	9	8	3	5	2	1	<u>6</u>
8	0	1	5	10	9	4	6	3	2	<u>7</u>
9	1	2	6	0	10	5	7	4	3	<u>8</u>
10	2	3	7	1	0	6	8	5	4	<u>9</u>
0	3	4	8	2	1	7	9	6	5	<u>10</u>

# dRAID: Parity Declustering

Rebuild scales to a large number of drives

- Decouple redundancy group (P+D) from the number of child drives
- Write: spare blocks are rotated among all drives
- Read: shared evenly among all drives

Drive \ IO	1	2	3	4	5	6	7	8	9	10
Read	4	4	4	4	4	4	4	4	4	4
Write	1	1	1	1	1	1	1	1	1	1

Rebuild IO Distribution



# dRAID: Demo

Live demo of dRAID rebuild:

- 43-drive dRAID2: 4 x (8D + 2P) groups, 3 distributed spares
  - 601G used out of 46T
  - Each drive capable of 150 MB/s
- Rebuild 1 failed drive to a distributed spare
  - Read 157.2G, write 17.5G

# dRAID: Demo Results

Rebuild completed in 37 seconds:

- Aggregate throughput: read 4350.6 MB/s, write 484.3 MB/s
  - 3x faster than RAIDz resilver
  - Scales to more drives
- Combined read/write of a single drive at 115.12 MB/s

# dRAID Downsides: Space Inflation

## RAIDz1

Disk \ LBA	A	B	C	D	E
0	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	D <sub>3</sub>
1	P <sub>1</sub>	D <sub>4</sub>	D <sub>5</sub>	D <sub>6</sub>	D <sub>7</sub>
2	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	P <sub>0</sub>
3	D <sub>0</sub>	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	X

## dRAID1

Disk \ LBA	A	B	C	D	E
0	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	D <sub>3</sub>
1	P <sub>1</sub>	D <sub>4</sub>	D <sub>5</sub>	D <sub>6</sub>	D <sub>7</sub>
2	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	X
3	P <sub>0</sub>	D <sub>0</sub>	X	X	X
4	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	X	X

# dRAID Space Inflation: Mirrored Metaslab

LBA \ Disk	A	B	C	D	E
0	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	D <sub>3</sub>
1	P <sub>1</sub>	D <sub>4</sub>	D <sub>5</sub>	D <sub>6</sub>	D <sub>7</sub>
2	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	X
3	P <sub>0</sub>	D <sub>0</sub>	X	X	X
4	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	X	X

LBA \ Disk	A	B	C	D	E
0	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	D <sub>3</sub>
1	P <sub>1</sub>	D <sub>4</sub>	D <sub>5</sub>	D <sub>6</sub>	D <sub>7</sub>
2	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	D <sub>2</sub>	X
3	P <sub>0</sub>	D <sub>0</sub>	D <sub>1</sub>	X	X
n	.	.	.	.	.
n+1	D <sub>0</sub>	D <sub>0</sub>	.	.	.
n+2	.	.	.	.	.

Regular Metaslab

Mirrored Metaslab

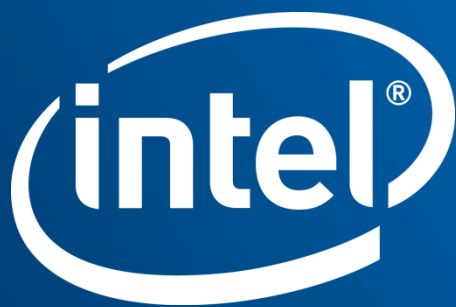
# dRAID: Downsides

- Allocated block size inflation
  - Doesn't work well for small blocks
  - Mitigated by mirrored metaslabs
- Rebuild cannot verify block checksum:
  - Still need to verify checksums by traversing BP tree
- Parity group and spare capacity chosen at VDEV creation

# dRAID: Project Status

Feature complete:

- Code: <https://github.com/zfsonlinux/zfs/pull/5841>
- Document: <https://github.com/zfsonlinux/zfs/wiki/dRAID-HOWTO>
- Need community help: reviewing, testing, patching, porting.



# dRAID: Demo

```
NAME      STATE  READ WRITE CKSUM
tank      ONLINE 0   0   0
  draid2-0 ONLINE 0   0   0
    sdb    ONLINE 0   0   0
    sdd    ONLINE 0   0   0
    sde    ONLINE 0   0   0
    .....
    sdar   ONLINE 0   0   0
spares
$draid2-0-s0 AVAIL
$draid2-0-s1 AVAIL
$draid2-0-s2 AVAIL
```



# dRAID: Demo

```
# zpool offline tank sde
```

```
# zpool replace tank sde '$draid2-0-s1
```

```
# zpool status
```

```
scan: rebuilt 17.4G in 0h0m36s with 0 errors
```

```
tank          DEGRADED  0  0  0
draid2-0      DEGRADED  0  0  0
sdb           ONLINE   0  0  0
sdd           ONLINE   0  0  0
spare-2       DEGRADED  0  0  0
sde           OFFLINE  0  0  0
$draid2-0-s1  ONLINE   0  0  0
```

# dRAID: Demo

