

OpenZFS

dRAID, Finally!



Mark Maybee



Hewlett Packard
Enterprise



- What is dRAID?
- Issues and Answers
 - Permutation arrays
 - Group size
 - Stranded space
 - Space inflation
- Drive Replacement Performance
- When is dRAID going to be available?
- Next steps?

This feature was originally developed at Intel by Isaac Huang (see his OpenZFS Summit talks in 2015, 2016, and 2017). It languished for a time before being rebased by Don Brady. The feature was “picked up” by Cray in 2018 to be used as a key capability in the ZFS version of its distributed storage product. Brian Behlendorf adopted the Cray version of the feature in 2020 and created the current PR. He has worked with Mark Maybee from Cray to keep it rebased and address a number of outstanding issues which needed to be addressed prior to integration.

- **Group Size** - the number of pieces the data is partitioned into plus the amount of parity
 - The amount of parity determines the redundancy
 - The number of data pieces determines the overhead
- **dRAID Size** - the number of drives used for data
 - Does not include spare drives
- **dRaid Row** - a 16M chunk from each drive at the same offset
 - E.g., Row 0 starts at offset 0 of each drive
- **Permutation Slice** - 1 or more rows which are permuted based on the dRAID permutation array
 - Number of rows per slice is derived from LCM of Group Size and dRAID Size

What is dRAID?



OpenZFS

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 vs RAIDz



OpenZFS

- dRAID implementation is layered on RAIDz code
 - Different map_alloc function, same IO pipeline
- RAIDz layout
 - Each group is constructed from a set of physical drives
 - Group columns extend full length of drive
 - Allocations must be multiple of parity + 1
 - Example: 1K data requires 2K space with 1 parity
- dRAID layout
 - Groups divided into “rows” (16MB)
 - Group rows are non-contiguous on physical drives
 - Allocations must be multiple of parity + data
 - Example: 1K data requires 2.5K space in 5 disk group

RAIDz1

Disk LBA	A	B	C	D	E
0	P ₀	D ₀	D ₁	D ₂	D ₃
1	P ₁	D ₄	D ₅	D ₆	D ₇
2	P ₀	D ₀	D ₁	D ₂	P ₀
3	D ₀	P ₀	D ₀	D ₁	X

dRAID1

Disk LBA	A	B	C	D	E
0	P ₀	D ₀	D ₁	D ₂	D ₃
1	P ₁	D ₄	D ₅	D ₆	D ₇
2	P ₀	D ₀	D ₁	D ₂	X
3	P ₀	D ₀	X	X	X
4	P ₀	D ₀	D ₁	X	X

Why Decluster?



- Spare drives are leveraged
 - Spare capacity is distributed across all spindles
- Each group's data is randomly distributed using all drives
 - Decouples redundancy group from number of data drives
- Data reconstruction leverages all drives
 - Reads shared evenly among all drives
 - Writes go to spare blocks distributed across all drives
- Sequential Resilver works!
 - Any row can be reconstructed without block pointer data

[illegible]

Creating a dRAID



OpenZFS

- New top-level vdev type: draid
 - Similar to raidz
 - Draid1, draid2, and draid3 supported
- Specify redundancy group size
 - Ex: draid1:5d -> 5+1 group size
 - Defaults to 8
 - Must be less than number of data drives in draid config
- Specify spare capacity (drive count)
 - Ex: draid1:5d:2s -> 2 drives worth of spare
- Config exposed in zpool status output
 - Spare capacity exposed as pseudo drives

```
# zpool create tank draid2:4d:2s L0 ... L52
```

```
# zpool status
```

```
pool: tank
state: ONLINE
config:
```

NAME	STATE	READ	WRITE	CKSUM
tank	ONLINE	0	0	0
draid2:4d:53c:2s-0	ONLINE	0	0	0
L0	ONLINE	0	0	0
L1	ONLINE	0	0	0
L2	ONLINE	0	0	0
L3	ONLINE	0	0	0
...				
L50	ONLINE	0	0	0
L51	ONLINE	0	0	0
L52	ONLINE	0	0	0
spares				
draid2-0-0	AVAIL			
draid2-0-1	AVAIL			

```
errors: No known data errors
```




- Need to provide the drive permutation for each slice
- Generated at pool creation
- If lost, all data access is lost
- Stored in the label
 - Takes up lots of space
- **Answer: Pre-define the permutations**
 - No need to figure it out at pool creation
 - No need to store in label
 - No risk of “losing” this information

OpenZFS

- ...

11	11	11	11	11
12	12	12	12	12
13	13	13	13	13
14	14	14	14	14
15	15	15	15	15
16	16	16	16	16
17	17	17	17	17
18	18	18	18	18
19	19	19	19	19
20	20	20	20	20
21	21	21	21	21
22	22	22	22	22
23	23	23	23	23
24	24	24	24	24
25	25	25	25	25
26	26	26	26	26
27	27	27	27	27
28	28	28	28	28
29	29	29	29	29
30	30	30	30	30
31	31	31	31	31
32	32	32	32	32
33	33	33	33	33
34	34	34	34	34
35	35	35	35	35
36	36	36	36	36

→

0	0	0	0	0	1	1	1	1	1	2	2	2	2	2	3
3	3	3	3	4	4	4	4	4	5	5	5	5	5	6	6
6	6	6	7	7	7	7	7	8	8	8	8	8	9	9	9
9	9	10	10	10	10	10	11	11	11	11	11	12	12	12	12
12	13	13	13	13	13	14	14	14	14	14	15	15	15	15	15
16	16	16	16	16	17	17	17	17	17	18	18	18	18	18	19
19	19	19	19	20	20	20	20	20	21	21	21	21	21	22	22
22	22	22	23	23	23	23	23	24	24	24	24	24	25	25	25
25	25	26	26	26	26	26	27	27	27	27	27	28	28	28	28
28	29	29	29	29	29	30	30	30	30	30	31	31	31	31	31
32	32	32	32	32	33	33	33	33	33	34	34	34	34	34	35
35	35	35	35	36	36	36	36	36	37	37	37	37	37	38	38
38	38	38	39	39	39	39	39	40	40	40	40	40	41	41	41
41	41	42	42	42	42	42	42	43	43	43	43	43	44	44	44
44	45	45	45	45	45	46	46	46	46	46	47	47	47	47	47

...



- Block allocation must fit within a group slice
 - Group slice size = group size * 16M
 - If block does not fit into space available in group slice it is allocated in next available group slice
 - Space can remain unallocated at end of group slice
 - There are many group slices in a dRAID
- Answer: Blocks spanning groups
 - Now can split a block across two groups
 - Leverage multi-row allocation maps developed for RAIDz expansion project

Space Inflation

- Unlike RAIDz, all columns are filled in an allocation
 - Particularly significant with small blocks
 - Any block that is not a multiple of *group size* * *sector size* will end with pad sectors
 - Pad sectors are explicitly zero-filled for sequential resilver to work

		RAIDz1				
		Disk				
LBA		A	B	C	D	E
0		P ₀	D ₀	D ₁	D ₂	D ₃
1		P ₁	D ₄	D ₅	D ₆	D ₇
2		P ₀	D ₀	D ₁	D ₂	P ₀
3		D ₀	P ₀	D ₀	D ₁	X

		dRAID1				
		Disk				
LBA		A	B	C	D	E
0		P ₀	D ₀	D ₁	D ₂	D ₃
1		P ₁	D ₄	D ₅	D ₆	D ₇
2		P ₀	D ₀	D ₁	D ₂	X
3		P ₀	D ₀	X	X	X
4		P ₀	D ₀	D ₁	X	X

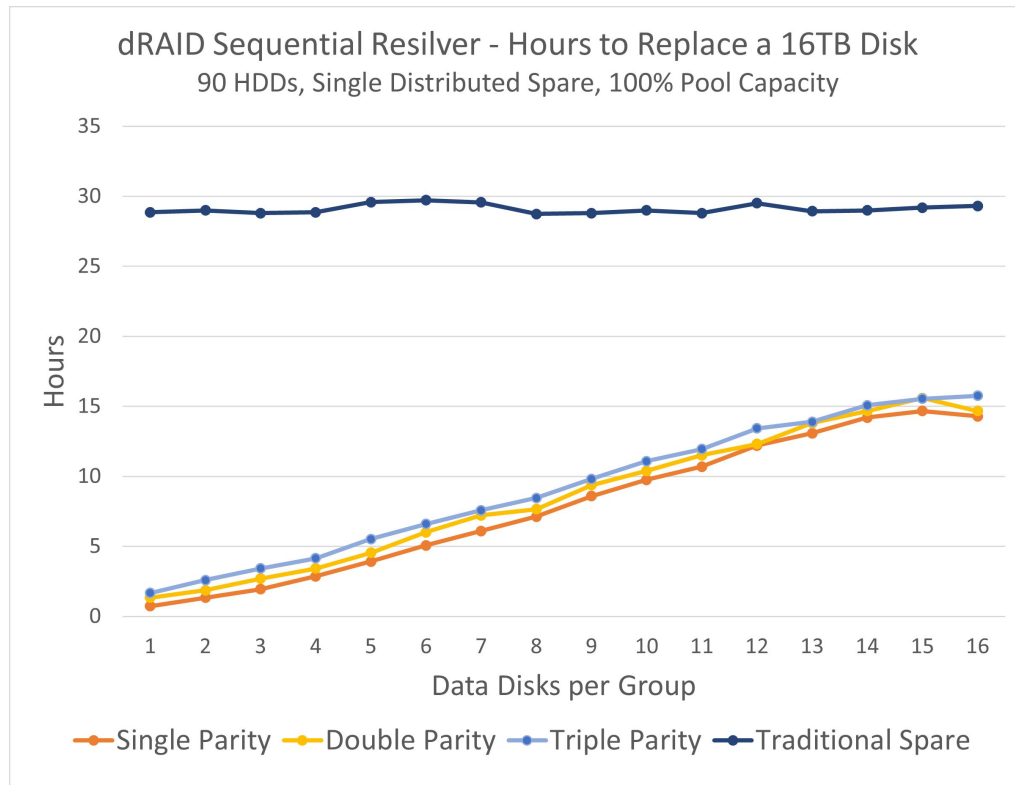
- Answers: Large blocks & Allocation classes

Drive Replacement Performance



OpenZFS

- Uses distributed spare
 - Spare capacity on all drives
- Declustering distributes reads
 - Failed drive impacts all groups
 - Reconstructed using all groups
- Smaller groups == faster rebuild
 - Fewer reads for reconstruction
- Much faster than “hot spare”
 - Small groups 10x faster
 - Large groups 2x faster



When will dRAID be Available?



- Originally targeting OpenZFS 2.0 release
 - Hence the title of this talk
 - Last minute changes/improvements pushed us out
- Now targeting the OpenZFS 2.1 release
 - Late 2020, early 2021 time frame

- Get it integrated!
 - Some small test case issues
 - Final code reviews
- dRAID expansion?
 - See you at the code-a-thon!

Questions?

Performance Tease



OpenZFS

- Preliminary performance results
- *Includes DirectIO support*
- 12 NVMe drives
 - dRAID2:d9:s1 configuration
 - 1 group of 9+2, 1 spare
 - 30G/s hardware write bandwidth
 - 42G/s hardware read bandwidth
- Achieving 80% of max write
 - 24G/s streaming writes
- Achieving 85% of max read
 - 36G/s streaming reads

