



# Fast Reconstruction for Large Disk Enclosures Based on RAID2.0

Qiliang Li, Min Lyu, Liangliang Xu, Yinlong Xu and Wei Wang

University of Science and Technology of China

**ICPP 2021** 

#### Large Disk Enclosures and RAID

 Large disk enclosures are widely used to provide large capacity storage



OceanStor 18000F V5



USTC, CHINA

2

Dell PowerVault MD3060e

- Limitations of conventional RAID
  - Slow reconstruction



#### Large Disk Enclosures and RAID

 Large disk enclosures are widely used to provide large capacity storage



OceanStor 18000F V5



Dell PowerVault MD3060e

- Limitations of conventional RAID
  - Slow reconstruction



USTC, CHINA



- Storage space contains data space and hot spare space
  - Data space: storing data in normal state





- Storage space contains data space and hot spare space
  - Data space: storing data in normal state
  - Hot spare space: storing reconstructed data in case of failure



ADSLAB

- Divide storage space into chunks
  - Each chunk usually not smaller than 64MB





- Divide storage space into chunks
- Construct RAID group by randomly selecting chunks
  - E.g., 2+1 RAID group: group 0





- Divide storage space into chunks
- Construct RAID group by randomly selecting chunks
  - E.g., 2+1 RAID group: group 0, group 1





- Divide storage space into chunks
- Construct RAID group by randomly selecting chunks
  - E.g., 2+1 RAID group: group 0, group 1, group 2





- Divide storage space into chunks
- Construct RAID group by randomly selecting chunks
  - E.g., 2+1 RAID group: group 0, group 1, group 2, etc.



#### Reconstruction in RAID2.0



• When a disk fails, all lost chunks denoted as *recovery tasks* are pending for being reconstructed



#### Reconstruction in RAID2.0



• Perform reconstruction in batches

Disks

- Limited memory capacity and CPU resource
- Each batch sequentially reconstructs a limited number of groups
- Select disks randomly for reconstructed chunks



## Imbalanced Data Layout Within a Batch



#### Random Selection of Surviving Disks





**Random selection of disks** Imbalanced write load

Disks





Hot spare space

7/31/2021

#### Local Load Imbalance

7/31/2021



- Definition of load balance rate  $\lambda$ 
  - $\lambda_{r/w}$ : Ratio of the maximum number of chunks read from/written to a disk to the average in a batch



#### Evaluate Local Load Imbalance

- Simulations of RAID<sub>R</sub>
  - Use the system time as random seed
  - 6+1 RAID group, CDF of 100 batches
  - Single disk failure



USTC, CHINA

### Popular Approach: Design Data Layout

- Guarantee uniform data distribution for balanced failure recovery I/O <sup>[1]</sup>
- Limitations
  - Relocation costs
    - Long used storage status Relocation → Designed dedicated data layout
  - Maintenance costs
    - Data layout after recovery or scaling Relocation Normal data layout
  - Heterogeneous cases
    - Existence of hot spot disks with less rebuilding bandwidth

[1] RAID+: Deterministic and balanced data distribution for large disk enclosures - Zhang et al., at FAST'18





Goal 1: Achieving local read and write load balance in reconstruction

Goal 2: Applicable to heterogeneous rebuilding bandwidth

#### Design of DR-RAID



#### Balance read load

• Step 1 – Batching tasks out of order



#### Balance write load

• Step 2 – Selecting surviving disks based on matching theory



- Initialize the batch with successive serial numbers
- Get read load of each disk and calculate  $\lambda_r$
- If  $\lambda_r \neq 1$ 
  - Select a task  $t_r$  from the current batch to be replaced
  - Select a task  $t_l$  from the remaining pending tasks
  - Replace  $t_r$  with  $t_l$
- Keep replacing until  $\lambda_r = 1$  or there are no satisfying replacements.

Note: The replacement rule is explained with the following example.

Dut of Order ADSLAB

Initialize the batch with successive serial numbers



• Get read load of each disk and calculate  $\lambda_r$ 



• Select a task  $t_r$  from the current batch to be replaced



• Select a task  $t_r$  from the current batch to be replaced



• Select a task  $t_r$  from the current batch to be replaced



• Select a task  $t_l$  from the remaining pending tasks



USTC, CHINA

• Replace  $t_r$  with  $t_l$ 



• Keep replacing until  $\lambda_r = 1$  or there are no satisfying replacements





- Construct the bipartite graph based on the selected batch of tasks
- Find a maximum matching
- Distribute the reconstructed chunks to disks based on the maximum matching



• Construct the bipartite graph based on the selected batch of tasks





• Construct the bipartite graph based on the selected batch of tasks



• Find a maximum matching



USTC, CHINA

#### • Find a maximum matching



USTC, CHINA

Distribute the reconstructed chunks to disks based on the maximum matching



USTC, CHINA

#### Heterogeneous Rebuilding Bandwidth

- Get the available rebuilding bandwidth  $b_i$  of disk  $D_i$
- Add weights based on  $b_i$  into Step I and Step 2
  - Substitute the counter  $c_i$  for  $c_i/b_i$  in Step I
  - Change the capacity of edges connecting to sink T in Step 2

#### Evaluation



#### • Hardware setup

- Local cluster of 11 nodes interconnected through a InfiniBand Switch
- Use iSER to implement the connection of large-scale disk pool
- Hardware details of each node
  - 64GB memory, 48 cores, a 7200-RPM ITB SATA hard disk
- Methodology
  - Default configuration
    - 50 disks, 6+1 RAID group, chunk size of 64MB, batch size of 49
  - Single failure reconstruction

#### **Baselines**



- 2 randomized data placement schemes:  $RAID_R$  and  $RAID_H$ 
  - RAID<sub>R</sub>: Use system time as random seed
  - RAID<sub>H</sub>: Randomize based on Jenkins hash function
- 2 conventional organizations: RAID50 and RAID5C



## **Offline Rebuilding**

- Recovery Throughput
  - Varying disk pool sizes
    - 51% improvement of RAID<sub>R</sub> & 6.6x of RAID5C for 50 disks
  - Varying batch sizes
    - 11% improvement of  $RAID_R$  even when batch size is 294





## **Online Rebuilding**

- Request latency
  - The period of performance degradation is shortened by 28%
- Rebuilding throughput with varied rebuilding bandwidth
  - 59% improvement of RAID<sub>R</sub> & 2.0x of RAID50



USTC, CHINA

#### Conclusion



- Local load imbalance within a recovery batch slows down reconstruction process
- DR-RAID selectively packs recovery tasks into a batch to achieve reconstruction load balance
- DR-RAID increases rebuilding throughput compared with the random data placement scheme in offline rebuilding and varied rebuilding bandwidth

# Thanks for your attention!



# Qiliang Li@USTC leeql@mail.ustc.edu.cn