

Dynamic Redundancy based on Media Reliability

Shehbaz Jaffer, Ashvin Goel, Angela Demke Brown and Bianca Schroeder
University of Toronto

Abstract

Existing drive reliability techniques such as RAID [5] have a fixed redundancy scheme. These techniques do not take the health of the underlying media into consideration while creating redundant data. Also, different drives like SSDs and conventional disks exhibit varying failure patterns [1]. Hence a generic scheme may not work for all media types. Additionally, such schemes assume disks follow a fail-stop model [7]. Moreover, the rebuild times for complete drive failure schemes are high and drastically affect the active I/O workload. With increased drive capacities and a constant throughput, rebuild times continue increasing exponentially [3]. Furthermore, the existing reliability techniques do not account for correlated failures. [6]

We propose DyRe, a Dynamic Redundancy technique that varies redundancy associated with data based on the *health* of the storage media below the data. More redundant information is stored for data on unreliable media. For data on reliable media, less redundancy is stored to optimize storage capacity. DyRe improves failure resiliency of drives within a node or a rack. We assume an overall global resiliency scheme for failure resiliency across racks [2]. With a more failure resilient rack, we can decrease the frequency of expensive cross rack rebuilds. Also, decreasing intra-rack failures will bring down operational costs for drive configuration and replacement. We provide a certain reserved storage for keeping storage health aware redundancy data. A health monitor gives us information about an impending disk sector or block failure. We incrementally create redundant information for data that lies on media impending a failure.

We propose using a number of techniques to detect impending storage failures. For example, we use latent sector errors reported by S.M.A.R.T. counters to detect impending failures [4]. We also use dtrace to monitor increased drive I/O latencies for certain blocks. If the I/O

latency associated with a block is high, a possible reason is the retries made by drive to perform internal ECC and retrieve inaccessible data [8]. We replicate such blocks at file system layer using ZFS. Such feedback based redundancy approach should give stronger resiliency against partial or complete drive failures within a rack.

References

- [1] BALAKRISHNAN, M., KADAV, A., PRABHAKARAN, V., AND MALKHI, D. Differential raid: Rethinking raid for ssd reliability. *Trans. Storage* 6, 2 (July 2010), 4:1–4:22.
- [2] CIDON, A., ESCRIVA, R., KATTI, S., ROSENBLUM, M., AND SIRER, E. G. Tiered replication: A cost-effective alternative to full cluster geo-replication. In *2015 USENIX Annual Technical Conference (USENIX ATC 15)* (Santa Clara, CA, July 2015), USENIX Association, pp. 31–43.
- [3] LEVENTHAL, A. Triple-parity raid and beyond. *ACMQueue* 7, 11 (December 2009).
- [4] MA, A., DOUGLIS, F., LU, G., SAWYER, D., CHANDRA, S., AND HSU, W. Raidshield: Characterizing, monitoring, and proactively protecting against disk failures. In *13th USENIX Conference on File and Storage Technologies (FAST 15)* (Santa Clara, CA, Feb. 2015), USENIX Association, pp. 241–256.
- [5] PATTERSON, D. A., GIBSON, G., AND KATZ, R. H. A case for redundant arrays of inexpensive disks (raid). In *Proceedings of the 1988 ACM SIGMOD International Conference on Management of Data* (New York, NY, USA, 1988), SIGMOD '88, ACM, pp. 109–116.
- [6] PRABHAKARAN, V., BAIRAVASUNDARAM, L. N., AGRAWAL, N., GUNAWI, H. S., ARPACI-DUSSEAU, A. C., AND ARPACI-DUSSEAU, R. H. Iron file systems. In *Proceedings of the Twentieth ACM Symposium on Operating Systems Principles* (New York, NY, USA, 2005), SOSP '05, ACM, pp. 206–220.
- [7] SCHROEDER, B., AND GIBSON, G. A. Disk failures in the real world: What does an mttf of 1,000,000 hours mean to you? In *Proceedings of the 5th USENIX Conference on File and Storage Technologies* (Berkeley, CA, USA, 2007), FAST '07, USENIX Association.
- [8] TSAI, T., THEERA-AMPORN PUNT, N., AND BAGCHI, S. A study of soft error consequences in hard disk drives. In *Proceedings of the 2012 42nd Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)* (Washington, DC, USA, 2012), DSN '12, IEEE Computer Society, pp. 1–8.