

# Using a Shared Storage Class Memory Device to Improve the Reliability of RAID Arrays

Sara Chaarawi<sup>1</sup>, Jehan-François Pâris<sup>1</sup>, Ahmed Amer<sup>2</sup>, T. J. Schwarz, S. J.<sup>3</sup>, Darrell D. E. Long<sup>4</sup>

<sup>1</sup>Department of Computer Science, University of Houston, Houston, TX 77204-3010, USA

<sup>2</sup>Department of Computer Engineering, Santa Clara University, Santa Clara, CA 9505, USA

<sup>3</sup>Depto. de Informática y Ciencias de la Computación, Universidad Católica del Uruguay, 11600 Montevideo, Uruguay

<sup>4</sup>Department of Computer Science, University of California, Santa Cruz, CA 95064, USA

schaarawi@gmail.com, paris@cs.uh.edu, a.amer@acm.org, tschwarz@ucu.edu.uy, darrell@cs.uh.edu

## I. INTRODUCTION

While magnetic disks have dramatically increased over the last twenty to thirty years, their access times are still measured in milliseconds and their failure rates still make them one of the least reliable parts of computer systems. In addition, their power consumption limits their usage in portable devices. Storage class memories (SCMs) constitute an emerging class of non-volatile storage systems that address these three issues. First, they promise much lower access times than magnetic disks. Second, they are expected to be much more reliable as they have no moving parts. Third, they will have much lower power requirements. Their sole major drawback is a higher price per Gigabit.

Given these characteristics, one of the first expected applications for SCMs will be intermediary caches for conventional disks. Currently accessed data would be stored in SCMs while dormant data would remain stored on magnetic disks. SCMs are also likely to displace flash memory in portable applications thanks to their higher write endurance.

We propose here another application for SCMs, namely enhancing the reliability of conventional disk arrays. The idea is not new: some of us recently proposed to increase the reliability of two-dimensional RAID arrays by replacing some of their parity disks by SCM devices [PA09]. The main limitation of this approach was its cost as the required number of SCM devices grew as the square root of the number of disks. The solution we propose here is quite different: we propose to group conventional RAID level 5 arrays into small groups of two to three RAID level 5 arrays a single “capstone” SCM device that will allow the set of disks to tolerate all double failures and most triple failures. As a result, the additional SCM device increases the mean time to data loss of the arrays in the group it protects by at least 20,000 percent.

The remainder of the paper is organized as follows: Section 2 discusses the performances of storage class memories and introduces shared parity devices and Section 3 introduces our organization. Section 4 discusses its performance and Section 5 presents our conclusions.

TABLE I. EXPECTED SPECIFICATIONS OF PCM DEVICES. [PA09]

<i>Parameter</i>	<i>Expected Value (2012)</i>
Access time	100 ns
Data Rate	200–1000 MB/s
Write Endurance	10 <sup>9</sup> write cycles
Read Endurance	no upper limit
Capacity	16 GB
Capacity growth	> 40% per year
Mean Time to Failure	10–50 million hours
Ratio of random to sequential access times	1
Active Power	100 mW
Standby Power	1 mW
Shock and Vibration resistance	> 15 g
Cost	< \$2/GB
Cost reduction rate	40 percent/year

## II. PREVIOUS WORK

In this section we briefly discuss the performances of storage class memories. And review previous work on shared parity devices and discuss the performances of storage class memories.

### A. Storage class memories

Storage class memories (SCMs) [N07] constitute a new class of non-volatile storage systems that are both cheaper than volatile main memory, and much faster than conventional disks. Unlike magnetic disk and MEMS [CG+00] technologies, SCMs have no moving parts. In addition, they do not suffer from the potential write-speed limitations of flash memory.

We will focus here on phase-change memories (PCMs) as an exemplar of this new class of storage devices. While it is not yet clear which type of SCMs will eventually succeed on the marketplace, most of our conclusions are likely to hold for any type of SCMs. The most promising PCM technology

relies on the physical properties of chalcogenide materials. At room temperature, these materials can exist in two stable states, namely an amorphous state exhibiting a high resistivity and a crystalline state characterized by a much lower resistivity. Quickly heating the material above its melting temperature and then letting it quickly cool will leave the material in an amorphous state, characterized by a high resistivity. Similarly, heating the material above its crystallization temperature and then letting it cool at a relatively slower rate will leave it in a crystalline, more conductive state.

Table 1 displays the most important parameters of the first generation of SCMs. As we can see, they have both much faster and much moiré reliable than magnetic disks

### B. Shared parity devices

There are many applications that require better data survival warranties than those afforded by RAID level 5 arrays. A first option is to switch to RAID level 6 arrays. A cheaper solution is to add to each group of two to three RAID array an extra parity disk [PA09].

Consider, for instance, the disk array displayed in Fig. 1. It consists of two conventional RAID arrays sharing an additional disk  $Q$  containing yet to be specified parity data. For the sake of simplicity, we have represented the two arrays as having separate parity disks while we expect their parity blocks to be distributed among the seven disks forming each RAID array. As Fig. 2 shows, we can define a *virtual parity disk*  $P'$  whose contents are the *exclusive or* (XOR) of the contents of parity disks  $P_0$  and  $P_1$ . (Had the parity blocks been equally distributed among the seven disks of each array, we would have defined a *virtual set of parity blocks*.)

The virtual array consisting of 12 data disks and the virtual parity disk  $P'$  forms a conventional RAID array that protects its contents against any single disk failure. We then define the contents of disk  $Q$  in a way that ensures that the 12 data disks, the virtual parity disk  $P'$  and disk  $Q$  form a RAID level 6 array. This can be done by using an EvenOdd scheme [BB+95], a Row Diagonal Parity (RDP) scheme [CE+04, GL+08] or any other RAID level 6 organization [P08].

We observe that the two parity disks  $P_0$  and  $P_1$  effectively protect all stored data against all single disk failures and all double failures that do not affect two disks in the same array. When combined with the shared parity disk  $Q$ , they protect the same data against any double disk failure and most, but not all, triple disk failures. The triple failures that will result in a data loss are the failures of:

1. three disks in the same RAID array, or
2. two disks in the same RAID array *plus* the shared parity disk  $Q$ .

Since our disk organization comprises 15 disks, it can experience  $\binom{15}{3}$  distinct triple failures. In addition, there are

$\binom{7}{2}$  distinct double and  $\binom{7}{3}$  distinct triple failures for each

of the two RAID arrays. As a result, our disk organization

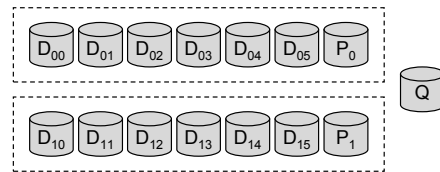


Figure 1. A pair of RAID arrays with an shared parity disk.

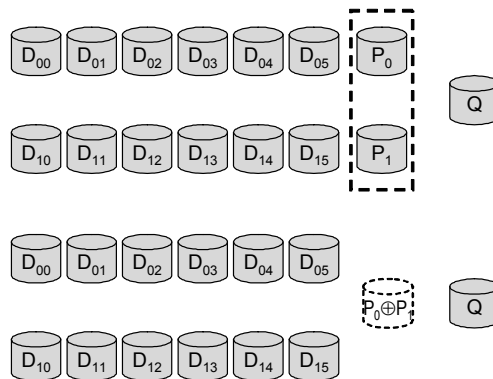


Figure 2. An alternate view of the previous array.

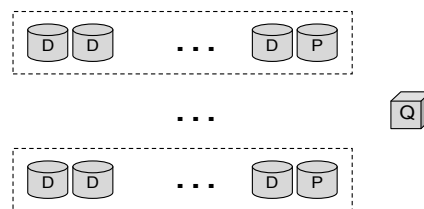


Figure 3. Our model.

will be able to tolerate exactly  $\binom{15}{3} - 2\binom{7}{3} - 2\binom{7}{2} = 343$  of the 455 possible triple disk failures, that is, slightly more than 75% of them.

As one can expect, avoiding any data loss in the presence of all double failures and three quarters of triple failures has a dramatic impact on the mean time to data loss of the two RAID arrays: Pâris and Amer [PA09] found out that adding a shared parity disk to a set of two of three small RAID arrays would increase their MTTDL by at least 14,000 percent.

### III. OUR PROPOSAL

These excellent results convinced us to consider replacing the shared parity disk by a more reliable device in order to achieve even higher MTTDLs. SCMs constituted a natural choice because they are expected to be much more reliable than magnetic disks and offered fairly high data rates. In addition, they are much less susceptible to be affected by vibrations.

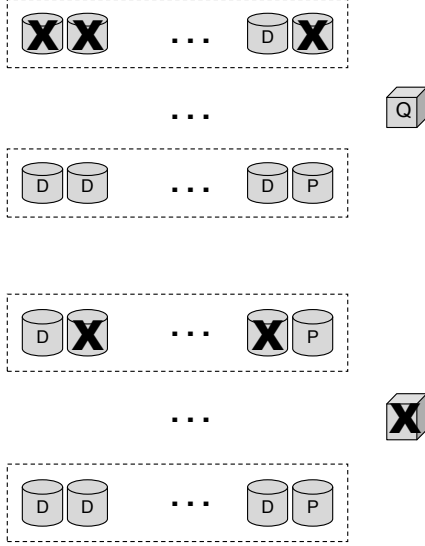


Figure 4. Triple failures resulting in a data loss.

Consider the disk array displayed in Figure 3. It consists of  $m$  RAID level 5 arrays comprising  $n$  disks each plus an additional shared parity SCM device  $Q$ . We characterize the contents of device  $Q$  in the following manner. We first define a virtual parity disk  $P'$  that is formed by XORing the parity blocks of the  $m$  RAID arrays. We then populate device  $Q$  in such a way that it forms a single RAID level 6 array with the data blocks of the original arrays and the virtual parity disk  $P$ .

Since our new organization constitute a RAID level 5 array, it can tolerate all single and all double disk failures. As Figure 4 shows, the triple failures that result in a data loss consist of:

1. A failure of three disks in the same RAID array, and
2. A failure of two disks in the same RAID array combined with a failure of the shared parity device.

As our disk organization comprises  $mn+1$  devices, it is subject to  $\binom{mn+1}{3}$  distinct triple failures. Of these failures,

$\binom{mn}{3}$  are failures of three of the  $mn$  disks and  $\binom{mn}{2}$  are failures of the shared parity device and two of the  $mn$  disks.

Since there are  $\binom{n}{2}$  distinct double and  $\binom{n}{3}$  distinct triple failures for each of the  $m$  RAID arrays, our system will be able to tolerate  $\binom{mn+1}{3} - m\binom{n}{3} - m\binom{n}{2}$  of the  $\binom{2n+1}{3}$  possible triple device failures.

#### IV. PERFORMANCE EVALUATION

Estimating the reliability of a storage system means estimating the probability  $R(t)$  that the system will operate in a correct fashion over the time interval  $[0, t]$  given that it operated correctly at time  $t = 0$ . Computing that function

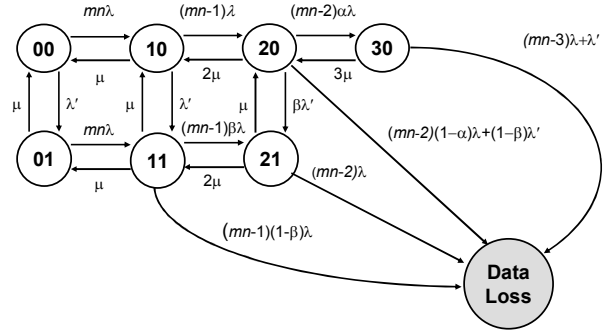


Figure 5. Simplified state transition probability diagram for our system.

requires solving a system of linear differential equations, a task that becomes quickly intractable as the complexity of the system grows. A simpler option is to use instead the mean time to data loss (MTTDL) of the storage system, which is the approach we will take here.

Our system model consists of an array of disks with independent failure modes. When a disk fails, a repair process is immediately initiated for that disk. Should several disks fail, the repair process will be performed in parallel on those disks. We assume that disk failures are independent events and are exponentially distributed with mean  $\lambda$ . In the same way, we assume that SCM device failures are also exponentially distributed, but with mean  $\lambda' < \lambda$ , reflecting the higher reliability of these devices. In addition, we require all repairs to be exponentially distributed with mean  $\mu$ . Both hypotheses are necessary to represent each system by a Markov process with a finite number of states.

Building an accurate state-transition diagram for our disk organization is a daunting task as we must distinguish between failures of the shared parity device  $Q$  and failures of the other disks as well as between failures of disks belonging to the same disk array and failures of disks belonging to distinct arrays. Instead, we present here a simplified model. Figure 5 displays the simplified state transition probability diagram for a system with  $mn$  disks each and a shared parity device  $Q$ . Each state is identified by a pair  $\langle xy \rangle$  where  $x$  represents the state of the shared SCM device  $Q$  and  $y$  stands for the number of failed disks.

State  $\langle 00 \rangle$  represents the normal state of the system when all its components are all operational. A failure of one of the  $mn$  disks will bring the system to state  $\langle 01 \rangle$ . The transition rate is  $nm\lambda$  reflecting that disk failures are independent processes. A failure of any of the remaining  $(mn-1)$  disks would bring the array into state  $\langle 02 \rangle$ . When the array is in state  $\langle 02 \rangle$ , a failure of any of the remaining  $(mn-2)$  disks will cause a triple disk failure. As we saw before,  $m\binom{n}{3}$  of all  $\binom{mn}{3}$  possible triple disk failures will result in a data loss. Define

$$\alpha = 1 - m \binom{n}{3} \binom{mn}{3}^{-1}$$

as the fraction of triple disk failures that will *not* result in a data loss. The two transitions corresponding to the failure of a third disk can then be expressed as

1. A transition to state <03> with rate  $\alpha(mn - 2)\lambda$
2. A failure transition with rate  $(1 - \alpha)(mn - 2)\lambda$

We assume that a failure of four or more disks will always result in a data loss. This is strictly true when  $m = 2$  and remains a fairly good approximation when  $m > 2$  as long as the device repair rate  $\mu$  remains much higher than the disk failure rate  $\lambda$ .

Let us consider how our model represents failures of the shared parity device. A failure of that device when the array is in state <00> will bring the array in state <10> while the same failure when the system already has a failed disk will bring it in state <11>. State <11> can also be reached from state <10> if one of the  $mn$  operational disks fail. When the array is in state <11>, a failure of any of the remaining  $(mn - 1)$  disks will bring the array in a state where it has two failed disk and a failed shared parity device. As we saw

before,  $m \binom{n}{2}$  of all  $\binom{mn}{2}$  possible double disk failures occurring when the shared parity device is not operational will result in a data loss. Define

$$\beta = 1 - m \binom{n}{2} \binom{mn}{2}^{-1},$$

that is, the fraction of disk failures that will *not* result in a data loss. The two failure transitions leaving state <11> can then be expressed as

1. A transition to state <12> with rate  $\beta(mn - 1)\lambda$
2. A failure transition with rate  $(1 - \beta)(mn - 1)\lambda$

In the same way, the two failure transitions corresponding to a failure of the shared parity device when the array already has two failed disks are

1. A transition to state <12> with rate  $\beta\lambda'$
2. A failure transition with rate  $(1 - \beta)\lambda'$

We assume that a the array is in state <03> will never tolerate the simultaneous failure of three disks and its shared parity device result in a data loss. This is strictly true when  $m = 2$  as a system consisting of two RAID level 5 arrays and a shared parity device cannot tolerate any quadruple disk failure. It remains a fairly good approximation when  $m > 2$  as long as the device repair rate  $\mu$  remains much higher than the disk failure rate  $\lambda$ .

Finally, we assume that the simultaneous failure of three or more data disks and the shared parity device will always result in a data loss. This is strictly true when  $m = 2$  and remains a fairly good approximation when  $m > 2$  as long as the device repair rate  $\mu$  remains much higher than the disk failure rate  $\lambda$ .

Disk repair transitions return the array from state <03> to state <02> then from state <02> to state <01> and, finally, from state <01> to state <00>. Similar transitions return the array from state <12> to state <11> and from state <11> to state <10>. Their rates are equal to the number of failed disks times the disk repair rate  $\mu$ . Repair transitions corresponding to a repair of the shared parity device will

bring the array from state <10> to state <00>, from state <11> to state <01> and from state <12> to state <02>.

The Kolmogorov system of differential equations describing the behavior of the array is

$$\frac{dp_{00}(t)}{dt} = -(mn\lambda + \lambda')p_{00}(t) + \mu p_{01}(t) + \mu' p_{10}(t)$$

$$\frac{dp_{10}(t)}{dt} = -(mn\lambda + \mu')p_{10}(t) + \mu p_{11}(t) + \lambda' p_{00}(t)$$

$$\frac{dp_{01}(t)}{dt} = -((mn - 1)\lambda + \lambda' + \mu)p_{01}(t) + mn\lambda p_{00}(t) + \mu' p_{11}(t) + 2\mu p_{02}(t)$$

$$\frac{dp_{11}(t)}{dt} = -((mn - 1)\lambda + \mu' + \mu)p_{11}(t) + mn\lambda p_{10}(t) + \lambda' p_{01}(t) + 2\mu p_{12}(t)$$

$$\frac{dp_{02}(t)}{dt} = -((mn - 2)\lambda + \lambda' + 2\mu)p_{02}(t) + (mn - 1)\lambda p_{01}(t) + \mu' p_{12}(t) + 3\mu p_{03}(t)$$

$$\frac{dp_{12}(t)}{dt} = -((mn - 2)\lambda + 2\mu + \mu')p_{12}(t) + \beta(mn - 1)\lambda p_{12}(t) + \beta\lambda' p_{02}(t)$$

$$\frac{dp_{03}(t)}{dt} = -((mn - 3)\lambda + \lambda' + 3\mu)p_{03}(t) + \alpha(mn - 2)\lambda p_{02}(t)$$

where  $p_{ij}(t)$  is the probability that the system is in state <ij> with the initial conditions  $p_{00}(0) = 1$  and 0 otherwise.

Solving the Laplace transforms of these equations we derive from them the mean time to data loss (MTTDL) of the array using the relation

$$MTTDL = \sum_{i,j} p_{ij}^*(0),$$

The result is the quotient of two polynomials too large to be represented here.

Figure 6 displays on a logarithmic scale the MTTDLs achieved by an array organization consisting of two RAID arrays of seven disks each sharing a common parity device. We assumed that the disk failure rate  $\lambda$  was one failure every one hundred thousand hours, that is, slightly less than one failure every eleven years. These values correspond to the high end of the failure rates observed by Pinheiro et al. [PWB07] and Schroeder and Gibson [SG07]. The failure rates for the SCM device were expressed in relation with that of the disk failure rates. Disk repair times are expressed in days and MTTDLs as their  $\log_{10}$  in years.

We can then see that the beneficial effects of replacing the shared parity disk with a more reliable SCM device are already significant when the SCM device is five times more reliable than a regular disk. Conversely, these beneficial effects approach their maximum as soon as the same device becomes ten times more reliable than a regular disk. We also observe that these benefits remain fairly constant over a fairly wide range of repair times: replacing the shared parity disk

by a shared parity device increases by 40 – 60 percent the MTTDL of the array.

While these results are good, we need to keep in mind that SCM devices are likely to remain much more expensive than disks for a long time. As a result, a cheaper option might be to mirror the shared parity disk and/or selecting a more reliable disk make for the parity disks.

**Table 1. Comparing the MTTDLs with our organization consisting of all disks. RAID5 and RAID6 have two stripes of 7 data disks and 1 or 2, respectively, parity disks. The values are for a repair time of 24 hrs.**

Organization	Relative MTTDL
RAID 5	<b>0.00096</b>
All Disks	<b>1</b>
RAID 6	<b>1.0012</b>
SSD 5 × better	<b>1.4274</b>
SSD 10 × better	<b>1.5080</b>
SSD 100 × better	<b>1.5887</b>
SSD never fails	<b>1.5982</b>

### V. CONCLUSION

Storage class memories (SCMs) constitute an emerging class of non-volatile storage devices that promise to be significantly faster and more reliable than magnetic disks. We propose to add one of these devices to each group of two or three RAID level arrays and store on it additional parity data. Our new organization can tolerate all double disk failures, most triple disk failures and most failures involving two disks and the SCM device without incurring any data loss. As a result, the additional parity device increases the mean time to data loss of the arrays in the group it protects by at least 20,000 percent.

More work is still needed to evaluate the impact of irrecoverable read errors and to evaluate cheaper alternatives.

### REFERENCES

[BB+95] M. Blaum, J. Brady, J. Bruck, and J. Menon, EvenOdd: An efficient scheme for tolerating double disk failures in RAID architectures, *IEEE Trans. Computers* 44(2):192–202, 1995.

[BM93] W. A. Burkhard and J. Menon. Disk array storage system reliability. In *Proc. 23<sup>rd</sup> International Symposium on Fault-Tolerant Computing*, Toulouse, France, pp. 432–441, June 1993.

[CE+04] P. Corbett, B. English, A. Goel, T. Gracanac, S. Kleiman, J. Leong, and S. Sankar, Row-diagonal parity for double disk failure correction, *Proc. USENIX Conference on File and Storage Technologies (FAST 2004)* San Francisco, CA, pp. 1–14, 2004.

[CG+00] L. R. Carley, G. R. Ganger, and D. F. Nagle, MEMS-based integrated-circuit mass-storage systems. *Communications of the ACM*, 43(11):73–80, Nov. 2000.

[GL+08] W. Gang, L. Xiaoguang, L. Sheng, X. Guangjun, and L. Jing, Generalizing RDP codes using the combinatorial method, *Proc. 7th IEEE International Symposium on Network Computing and Applications (NCA 2008)*, Cambridge, MA, pp. 93–100, July 2008.

[N07] S. Narayan, Storage class memory a disruptive technology, *Presentation at Disruptive Technologies Panel: Memory Systems of SC '07*, Reno, NV, Nov. 2007.

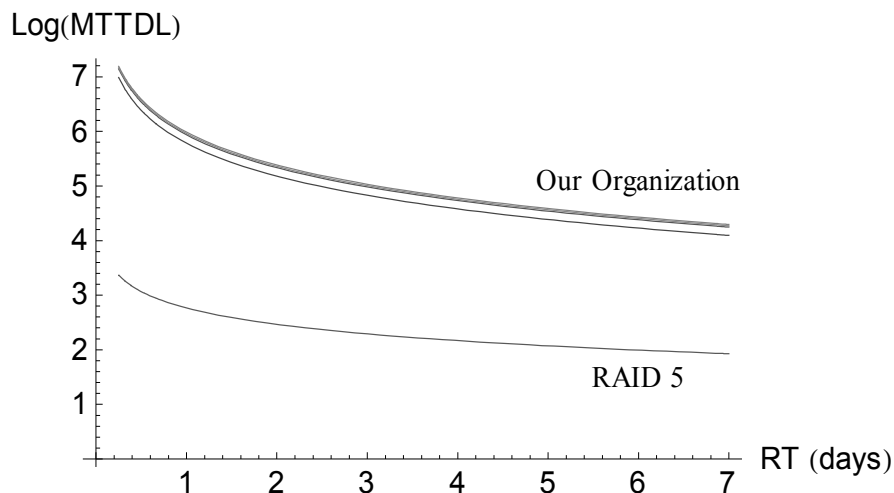
[PA09a] J.-F. Pâris, A. Amer and D. D. E. Long. Using storage class memories to increase the reliability of two-dimensional RAID arrays, *Proc. 17th International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems (MASCOTS 2009)*, pp. 125–132, Sep. 2009.

[PA09b] J.-F. Pâris and A. Amer. Using Shared Parity Disks to Improve the Reliability of RAID Arrays, *Proc. 28<sup>th</sup> International Performance of Computers and Communication Conference*, Phoenix, AZ, pp. 129–136, Dec. 2009.

[P08] J. S. Plank, The RAID-6 liberation codes, *Proc. 6th USENIX Conference on File and Storage Technologies 2* pp. 1–14, Feb. 2008.

[PWB07] E. Pinheiro, W.-D. Weber and L. A. Barroso, Failure trends in a large disk drive population, *Proc. 5<sup>th</sup> USENIX Conference on File and Storage Technologies*, San Jose, CA, pp. 17–28, Feb. 2007.

[SG07] B. Schroeder and G. A. Gibson, Disk failures in the real world: what does an MTTF of 1,000,000 hours mean to you? *Proc. 5<sup>th</sup> USENIX Conference on File and Storage Technologies*, San Jose, CA, pp. 1–16, Feb. 2007.



**Figure 6. Comparing the MTTDLs achieved by our organization with those achieved by set of identical disk arrays lacking a shared parity disk.**