

Ensuring Data Survival on Solid-State Storage Devices

Jehan-François Pâris
Dept. of Computer Science
University of Houston
Houston, TX 77204-3010
paris@cs.uh.edu

Darrell D. E. Long
Dept. of Computer Science
University of California
Santa Cruz, CA 95064
darrell@cs.ucsc.edu

Thomas J. E. Schwarz
Dept. of Computer Engineering
Santa Clara University
Santa Clara, CA 95053
tjschwarz@scu.edu

Abstract

The emergence of more reliable, often much faster solid-state storage devices is revolutionizing most aspects of data storage technology. Here we address the impact of these new storage devices on the techniques used to ensure data survival. In particular, we show that the higher bandwidth-to-capacity ratios of the new Storage Class Memory devices will make self-adaptive storage solutions much more attractive.

1 Introduction

Many organizations now maintain very large amounts of data online. This trend is due to many factors; among these are the lower costs of online storage, regulatory requirements (such as the Sarbanes-Oxley Act) and the increasing rate at which digital data are produced.

A critical issue for any large data storage system is how to ensure the survival of the data in the presence of equipment failures. Given the limitations of backup solutions, the best way to achieve this goal is to use redundant storage systems. Two techniques that can be used are mirroring and erasure codes. Mirroring maintains two or more exact copies of the data on distinct disks. Erasure codes, also known as m -out-of- n codes, group disks into sets of n disks that contain enough redundant data to tolerate the loss of $n - m$ disks.

Both techniques have been developed for storage arrays consisting of magnetic disks. While magnetic disks have huge capacities and low storage costs, they also have much higher latencies and are much less reliable than solid state devices. For instance Schroeder and Gibson [SB07] found that “in the field, annual disk replacement rates typically exceed 1%, with 2–4% common and up to 13% observed on some systems.” Pinheiro *et al.* [PW+07] similarly observed annual failure rates varying “from 1.7%, for drives that were in their first year of operation, to over 8.6%, observed in the 3-year old population.”

The high latency of magnetic disks has led to many proposals for faster storage devices. These proposals have in turn resulted in various studies analyzing the performance of file and storage systems incorporating these new technologies [UM03, BB+07 and many

others]. Conversely, the impact of these new technologies on the reliability of storage systems has received insufficient attention.

We propose to fill this gap by investigating the reliability characteristics of the newest and the most promising of these new storage technologies, the so-called *storage class memories*.

The remainder of this paper is organized as follows. Section 2 reviews past and present proposals for alternative storage technologies and introduces storage class memories. Section 3 discusses the reliability characteristics of storage systems using these storage class memories. Finally Section 4 has our conclusions.

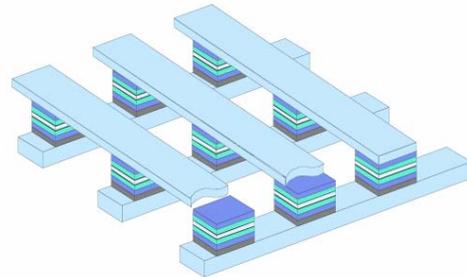


Fig. 1. General organization of a phase change memory [N07].

2 Alternative Storage Technologies

Alternative storage technologies that compete or have competed with magnetic disks include bubble memories [BS71], MEMS-based storage systems [CG+00] and flash memories [PB+97]. Among them, flash memories certainly are the most successful. They are widely used today in digital cameras, cellular phones, digital audio players and various formats of memory cards. Despite their higher cost per byte, they have begun to replace magnetic disks in some portable computers.

There are, however, several factors that limit the applicability of flash drives in general-purpose storage systems. First, they have low write endurance, typically of the order of one million erase/write cycles. Second, they suffer from write-speed limitations. Finally, they do not scale well below 45nm. In the absence of any technological breakthrough, they are not likely to completely replace magnetic disks.

Storage class memories (SCMs) constitute a new class of non-volatile storage systems that are at the same time much cheaper than main memory and much faster than conventional disks. 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.

PCMs contain no moving parts and use cross-bar-type chip structures to access data. As seen on Fig. 1, bits are stored at the intersection of each row and each column of the cross bar structure. Various techniques can be used to encode these data. The most promising approach 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 state.

Table 1. Expected specifications of the new storage class memory devices.

<i>Parameter</i>	<i>Expected Value (2012)</i>
Access time	100 ns
Data Rate	200–1000 MB/s
Write Endurance	10^9 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	< \$/GB
Cost reduction rate	40 percent/year

Table 1 displays the most important parameters of the first generation of SCMs. As we can see, they are already almost as fast as main memory and nearly as cheap as magnetic disks. In addition, they have a much better write endurance and better mean times to failures than flash memories.

3 Reliability Analysis

Three major parameters affect the reliability of storage arrays that could be built using first-generation SCMs. These are their high mean times to failure, their relatively low capacities and their very high data

transfer rates. We will discuss the impact of each of these parameters in turn.

3.1 Mean times to failure

As Table 1 shows, future SCMs will have much higher mean times to failure (MTTF) than magnetic disks. While magnetic disks typically have data sheet MTTFs between one and two million hours, SCMs are expected to have MTTFs between ten and fifty million hours. These figures translate into failure rates between 0.00018 and 0.00088 failures per year.

Estimating the reliability of a given storage system means estimating the probability $R(t)$ that the system will operate correctly over the time interval $[0, t]$ given that it operated correctly at time $t = 0$. Unfortunately, $R(t)$ is a function and not a single value. As a result, many studies characterize the reliability of storage systems by their Mean Time To Data Loss (MTTDL).

While MTTDLs are relatively easy to compute, they have their own limitations. First, they assume that storage systems are not replaced until they have reached statistical equilibrium, which is not true as their actual lifetimes are much shorter. Second, they do not take into account variations of device failure rates over their lifetimes. As a result, we prefer to characterize the reliability of a storage system by its *economic life span* $L(r)$, that is, the maximum time interval for which data stored on that system will have a probability r to survive intact [PS08]. To make our figure of merit dimensionless, we will express it in multiples of the MTTFs of the individual storage devices.

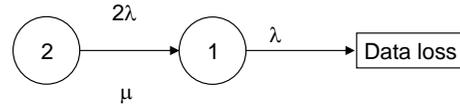


Fig. 2. State probability transition diagram for a mirrored device.

Let us focus for the moment on the performance of mirrored organizations, that is, organizations consisting of two identical devices each holding an identical copy of the data. To further demonstrate the impact of the high reliability of SCMs, let us further assume that the array cannot be repaired during its useful lifetime.

Assuming a constant disk failure rate λ , the survival probability S_1 of a single disk at time t is given by the differential equation

$$S_1' = -\lambda S_1$$

with initial condition

$$S_1(0) = 1,$$

whose solution is the exponential function

$$S_1 = \exp(-\lambda t).$$

Table 2. Economic life span for a single device and a mirrored device.

Number of 9s	Single	Mirrored
1 (90%)	0.10536	0.38013
2 (99%)	0.01005	0.10536
3 (99.9%)	0.00100	0.03213
4 (99.99%)	1.00E-04	0.01005
5 (99.999%)	1.00E-05	0.00317

We model the survival of a mirrored device in the standard Markov model depicted in Fig. 2. We label the non-failure states by the number of existing disks. The initial state is state 2, from which we transition to state 1 at rate 2λ , whenever one of the two disks fails. We can capture the probability p_i of being in state i at time t in a system of ordinary differential equations

$$\begin{aligned} p_2' &= -2\lambda p_2 \\ p_1' &= 2\lambda p_2 - \lambda p_1 \end{aligned}$$

with initial conditions

$$\begin{aligned} p_2(0) &= 1 \\ p_1(0) &= 0 \end{aligned}$$

whose solution is

$$\begin{aligned} p_2(t) &= e^{-2\lambda t}, \\ p_1(t) &= 2e^{-2\lambda t}(e^{\lambda t} - 1). \end{aligned}$$

The data survival probability for our mirrored device is

$$S_2(t) = p_2(t) + p_1(t).$$

To obtain the economic life span $L(r)$ of each organization, we solve observe that $L(r)$ is the solution of the equation $S_n(t) = r$. After setting $\lambda = 1$, the economic life span is expressed in multiples of the device MTTF.

These results are summarized in Table 2 for reliability level r equal to 0.9, 0.99, 0.999, 0.9999, and 0.99999, that is, 1, 2, 3, 4, and 5 nines. As we can see, a single device achieves a life span of 1% of its MTTF with probability 99%. An SCM module with an MTTF of ten million hours would thus have an economic life span of one hundred thousand hours or slightly more than eleven years. More demanding applications could use a mirrored organization to achieve the same economic life span with a reliability level of 99.99%.

This is an excellent result as we considered that the devices could not be repaired during their useful life. SCMs thus appear to become the technology of choice to store a few gigabytes of data at locations where repairs are impossible (satellites) or uneconomical (remote locations).

Let us consider now repairable organizations, starting with repairable mirrored organizations. Figure 3 displays our model. As we are interested in economic life spans, we set $\lambda = 1$ and obtain the following system of differential equations:

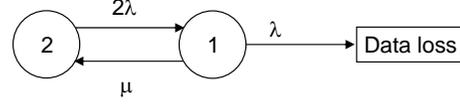


Fig. 3. State probability transition diagram for a mirrored device with repair.

Table 3: Economic life span for repairable mirror devices with variable repair rates μ and failure rate $\lambda = 1$.

Nines	$\mu = 10^3$	$\mu = 10^5$	$\mu = 10^7$
2	5.041230	502.53200	—
3	0.502747	50.0265	—
4	0.051149	5.00041	—
5	0.006010	0.500027	50.1456

Table 4: Economic life span for a repairable RAID Level 5 configuration with variable repair rates μ and failure rate $\lambda = 1$.

Nines	$\mu = 10^3$	$\mu = 10^5$	$\mu = 10^7$
2	0.1148000	11.16920	—
3	0.0123000	1.111890	111.155
4	0.0019840	0.111100	11.1100
5	0.0005124	0.011120	1.11097

$$\begin{aligned} p_2'(t) &= -2p_2(t) + \mu p_1(t) \\ p_1'(t) &= 2p_2(t) - (1 + \mu)p_1(t) \\ p_0'(t) &= p_1(t) \\ p_0(0) &= p_1(0) = 0, p_2(0) = 1 \end{aligned}$$

Its solution is given by

$$\begin{aligned} r(t) &= p_1(t) + p_2(t) \\ &= \frac{\exp(-(1/2)t(3 + \mu + R)) \cdot (-3 - \mu + R + \exp(tR)(3 + \mu + R))}{2 \cdot R} \end{aligned}$$

where R is

$$R = \sqrt{1 + \mu(6 + \mu)}.$$

We then calculated the economic life spans for various levels of nines and selected values of μ/λ between 10^3 and 10^7 . We tabulate the results in Table 3. Unfortunately, the result becomes numerically hard to evaluate when μ is 6 to 7 orders of magnitudes larger than λ . Of course, the economic life span goes towards infinity as μ/λ goes to infinity. A closer look at Table 3 shows that increasing this μ/λ ratio by 100 yields about the same increase in the economic life span.

We can apply the same approach to estimate the reliability of a RAID level 5 consisting of SCM devices. Assuming a 9+1 configuration, our system of ordinary differential equations becomes

$$\begin{aligned} p_2'(t) &= -10p_2(t) + \mu p_1(t) \\ p_1'(t) &= 10p_2(t) - (9 + \mu)p_1(t) \\ p_0'(t) &= p_1(t) \end{aligned}$$

and

$$p_0(0) = p_1(0) = 0, p_2(0) = 1$$

with solution

$$r(t) = \frac{\exp(-\frac{1}{2}t(19 + \mu + R)(-19 - \mu + R + \exp(tR)(19 + \mu + R)))}{2R}$$

with $R = \sqrt{1 + \mu(38 + \mu)}$. These results are summarized in Table 4. As before, we observe that as increasing the μ/λ ratio by 100 yields about the same increase in the economic life span.

These calculations lead us to believe that the potential for fast repair in SCM limits the need for large amounts of parity data. Since repair can only occur after error detection, engineering SCM-based storage systems will have to focus on detecting errors quasi-instantaneously instead of designing more error resilient storage schemes. File systems for SCM need to be able to perform repair efficiently and in particular need to be able to remap failed storage components without incurring access overhead. With the somewhat limited write endurance of SCM, this capability could presumably be part of wear-leveling.

3.2 Storage Capacity Considerations

The relatively small storage capacities of SCMs prevent us from extending these conclusions to larger storage systems. Consider for instance a storage system with a capacity of fifty terabytes. By the expected arrival time of the first generation of SCMs in 2012, we are very likely to have magnetic disks with capacities well exceeding two terabytes. We would use at most fifty of them. To achieve the same total capacity, we would need 6,250 SCM modules. In other words, each magnetic disk would be replaced by 125 SCM modules.

Let λ_m and λ_s respectively denote the failure rates of magnetic disks and SCM. To achieve comparable device failure rates in our hypothetical storage systems, we would need to have $\lambda_s < \lambda_m/125$, which is not likely to be the case.

Our second conclusion is that storage systems built using first-generation SCMs will be inherently less reliable than comparable storage systems using magnetic disks. This situation will remain true as long as the capacity of SCM modules remains less than λ_s/λ_m times that of a magnetic disk.

3.3 Bandwidth Considerations

M -out-of- n codes are widely used to increase the reliability of storage systems because they have a much lower space overhead than mirroring. The most popular of these codes are $n-1$ -out-of- n codes, or RAID5s, which protect a disk array with n disks against any single disk failure [CL+94, PG+88, SG+99]. Despite their higher reliability, $n-2$ -out-of- n codes [BM93, SB92] are much less widely used due to their higher update overhead.

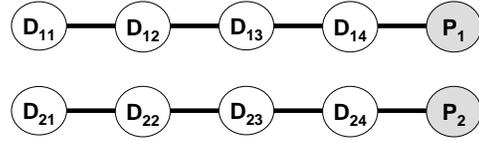


Fig. 4. A storage system consisting of two RAID arrays.

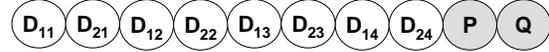


Fig. 5. The same storage system organized as a single 8-out-of-10 array.

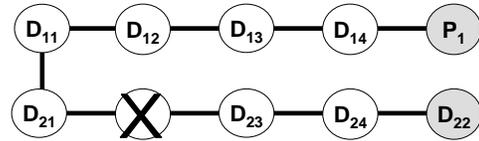


Fig. 6. How the original pair of RAID arrays could be reorganized after the failure of data disk D_{22} .

Consider for instance a small storage system that consists of ten disk drives and assume that we are willing to accept a 20% space overhead. Figure 4 shows the most likely organization for our array: the ten disks are grouped into two RAID level 5 arrays, each having five disks. This organization will protect the data against any single disk failure and the simultaneous failures of one disk in each array.

Grouping the 10 disks into a single 8-out-of-10 array would protect the data against *any* simultaneous failure of two disks without increasing the space overhead. This is the organization that Fig. 5 displays. As we mentioned earlier, this solution is less likely to be adopted due to its high update overhead.

A very attractive option would be to use two RAID level 5 arrays and let the two arrays reorganize themselves in a transparent fashion as soon as one of them has detected the failure of one of its disks. Fig. 6 shows the outcome of such a process: parity disk P_2 now contains the data that were stored on the disk that failed and the other parity disks while parity disk P_1 now contains the parity of the eight other disks.

A key issue for self-reorganizing disk arrays is the duration of the reorganization process [PS+06]. In our example, a quasi-instantaneous reorganization would achieve the same reliability as an 8-out-of-10 code but with a lower update overhead. Two factors dominate the duration of this reorganization process, namely, the time it takes to detect disk failure and the time it takes to perform the array reorganization itself. This latter time is proportional to the time it takes to overwrite a full disk. Assuming a future disk capacity of 2 terabytes and a future disk bandwidth of 100MB/s, this

operation would require slightly more than five hours and a half. The situation is not likely to improve as disk bandwidths have always grown at a lower rate than disk capacities.

With data rates in excess of 200 MB/s, SCMs make this approach much more attractive as the contents of a 16 GB module could be saved in at most 80 seconds. By the time larger SCM modules become available, we can expect their bandwidths to be closer to 1 GB/s than to 10 MB/s. As a result, we can expect self-reorganizing sets of RAID arrays built with SCM modules to be almost as reliable as single m -out-of- n arrays with the same space overhead.

4 Conclusions

SCMs are a new class of more reliable and much faster solid-stage storage devices that are likely to revolutionize most aspects of data storage technology. We have investigated how these new devices will impact the techniques used to ensure data survival and have come with three major conclusions. First, SCMs are likely to become the technology of choice to store a few gigabytes of data at locations where repairs are either impossible (satellites) or uneconomical (remote locations). Second, the low capacities of first-generation SCM modules will negate the benefits of their lower failure rate in most other data storage applications. Finally, the higher bandwidth-to-capacity ratios of SCM devices will make self-adaptive storage solutions much more attractive.

More work is still needed to estimate how the higher shock and vibration resistance of SCMs will affect their field reliability and find the most cost-effective self-reorganizing array solutions for SCMs.

References

- [BB+07] T. Bisson, S. Brandt, and D. D. E. Long, A Hybrid disk-aware spin-down algorithm with I/O subsystem support, In *Proc. International Performance Conference on Computers and Communication (IPCCC '07)*, New Orleans, LA, Apr. 2007
- [BS71] A. H. Bobeck and H. E. D. Scovil, Magnetic bubbles, *Scientific American*, 224(6):78-91, June 1971
- [BM93] W. A. Burkhard and J. Menon. Disk array storage system reliability. In *Proc. 23rd International Symposium on Fault-Tolerant Computing*, Toulouse, France, pp. 432-441, June 1993.
- [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.
- [CL+94] P. M. Chen, E. K. Lee, G. A. Gibson, R. Katz and D. A. Patterson. RAID, High-performance, reliable secondary storage, *ACM Computing Surveys* 26(2):145-185, 1994.
- [E07] E-week, Intel previews potential replacement for flash memory, www.eweek.com/article2/0,1895,2021815,00.asp
- [N07] S. Narayan, Storage class memory a disruptive technology, *Presentation at Disruptive Technologies Panel: Memory Systems of SC '07*, Reno, NV, Nov. 2007
- [PS+06] J.-F. Pâris, T. J. E. Schwarz and D. D. E. Long. Self-adaptive disk arrays. In *Proc. 8th International Symposium on Stabilization, Safety, and Security of Distributed Systems*, Dallas, TX, pp. 469-483, Nov. 2006.
- [PS08] J.-F. Pâris and T. J. Schwarz. On the possibility of small, service-free disk based storage systems, In *Proc. 3rd International Conference on Availability, Reliability and Security (ARES 2008)*, Barcelona, Spain, Mar. 2008, to appear.
- [PG+88] D. A. Patterson, G. A. Gibson and R. Katz. A case for redundant arrays of inexpensive disks (RAID). In *Proc. SIGMOD 1988 International Conference on Data Management*, Chicago, IL, pp. 109-116, June 1988.
- [PB+97] P. Pavan, R. Bez, P. Olivo, E. Zanoni., Flash memory cells-an overview, *Proceedings of the IEEE*, 85(8):1248-1271, Aug 1997.
- [PW07] E. Pinheiro, W.-D. Weber and L. A. Barroso, Failure Trends in a Large Disk Drive Population, In *Proc. 5th 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? In *Proc. 5th USENIX Conference on File and Storage Technologies*, San Jose, CA, pp. 1-16, Feb. 2007.
- [SG99] M. Schulze, G. A. Gibson, R. Katz, R. and D. A. Patterson. How reliable is a RAID? In *Proc. Spring COMPCON 89 Conference*, pp. 118-123, Mar. 1989.
- [SB92] T. J. E. Schwarz and W. A. Burkhard. RAID organization and performance. In *Proc. 12th International Conference on Distributed Computing Systems*, pp. 318-325 June 1992.
- [UM+03] Mustafa Uysal, Arif Merchant, Guillermo A. Alvarez. Using MEMS-based storage in disk arrays, In *Proc. 2nd USENIX Conference on File and Storage Technologies (FAST '03)*, San Francisco, CA, pp. 89-101, Mar.-Apr. 2003.
- [XB99] L. Xu and J. Bruck: X-code: MDS array codes with optimal encoding. *IEEE Trans. on Information Theory*, 45(1):272-276, Jan. 1999.