

The International Journal of Digital Curation

Issue 1, Volume 5 | 2010

Bit Preservation: A Solved Problem?

David S. H. Rosenthal,
Stanford University Libraries, CA

Abstract

For years, discussions of digital preservation have routinely featured comments such as “bit preservation is a solved problem; the real issues are ...”. Indeed, current digital storage technologies are not just astoundingly cheap and capacious, they are astonishingly reliable. Unfortunately, these attributes drive a kind of “Parkinson’s Law” of storage, in which demands continually push beyond the capabilities of systems implementable at an affordable price. This paper is in four parts:

Claims, reviewing a typical claim of storage system reliability, showing that it provides no useful information for bit preservation purposes.

Theory, proposing “bit half-life” as an initial, if inadequate, measure of bit preservation performance, expressing bit preservation requirements in terms of it, and showing that the requirements being placed on bit preservation systems are so onerous that the experiments required to prove that a solution exists are not feasible.

Practice, reviewing recent research into how well actual storage systems preserve bits, showing that they fail to meet the requirements by many orders of magnitude.

Policy, suggesting ways of dealing with this unfortunate situation.¹

¹ This article is based on the paper given by the author at iPRES 2008; received April 2009, published June 2010.

The *International Journal of Digital Curation* is an international journal committed to scholarly excellence and dedicated to the advancement of digital curation across a wide range of sectors. ISSN: 1746-8256 The IJDC is published by UKOLN at the University of Bath and is a publication of the Digital Curation Centre.



Introduction

For years, discussions of digital preservation have routinely featured comments such as “bit preservation is a solved problem; the real issues are ...”.² Indeed, current digital storage technologies are not merely astoundingly cheap and capacious, they are astonishingly reliable. Unfortunately, these attributes drive a kind of “Parkinson’s Law” (Parkinson, [1957](#)) of storage, in which demands continually push beyond the capabilities of systems implementable at an affordable price.

This paper is in four parts. The first part examines a typical claim made by a storage system vendor for the reliability of their product. It concludes that these numbers provide no useful information for bit preservation purposes.

The second, theoretical, part asks what characterizes a solution to the bit preservation problem adequate to the large numbers of bits to be stored and the long durations for which these bits are to be preserved. It proposes “bit half-life” as a metric for bit preservation, discusses the requirements being placed upon preservation systems in terms of this metric, and investigates the feasibility of benchmarking systems to see if they meet these requirements. It concludes firstly that the requirements are so onerous that it is not feasible to measure whether systems meet them, and secondly that bit half-life is an inadequate metric.

The third, practical, part reviews recent investigations into the performance of large-scale storage systems and their components. These studies uniformly report that storage reliability actually delivered to applications such as digital preservation systems is much less than that claimed by the manufacturers of systems and components. Tracking these failures to their root causes shows that every single hardware and software component contributes to some extent to the failures the systems experience. It concludes that current storage technologies fall well short of current requirements for bit preservation.

Given that the actual performance of storage systems is much worse than required, and that even if it improves we still won’t be sure that a system will meet its requirements, the fourth part asks what is to be done. As with paper, content in digital archives will inevitably suffer loss and damage. The question is how to invest the limited funds available for preservation to the best effect in terms of improved data survival. There are many ways in which spending more money can reduce (but never completely eliminate) the probability of loss and damage. What is needed to allow informed investment decisions? How can we encourage the development of cost-effective techniques for long-term bit preservation?

Clarification

It is incumbent on those attacking ideas such as the “solvedness” of bit preservation to focus on the strongest version of the idea³. If proponents really believed that bit preservation was solved, they wouldn’t bother with backups. Of

² The prevalence of this meme is aptly illustrated by the letter from the iPres programme committee accepting this paper. It cites the title as “Bit Preservation - A Problem Solved”.

³ “we should always try to clarify and to strengthen our opponent’s position as much as possible before criticising him” (Popper, [1959](#)).

course, they do. What they really mean by bit preservation being solved is that the set of techniques in common use make it so unlikely that bits will be lost that there is no need for concern at the prospect.

The techniques in which they place such faith are backups and checksums. Their real belief is that if they make a few backup copies of their content, and include in them checksums which they occasionally verify, their content will be safe. The goal of this paper is to show that, while backups and checksums may be adequate for relatively short periods and small amounts of preserved data, the scale and duration of current preservation tasks render them inadequate.

The state of our knowledge about preserving bits can be summarized as:

- *The more copies the safer.* As the size of the data increases, the per-copy cost increases, reducing the number of backup copies that can be afforded.
- *The more independent the copies the safer.* As the size of the data increases, there are fewer storage options available. Thus the number of copies in the same storage technology increases, decreasing the average level of independence.
- *The more frequently the copies are audited the safer.* As the size of the data increases, the time and cost needed for each audit increases, reducing their frequency.

Thus techniques that might be adequate at a small scale will break down as the scale increases.

Claims

How would we know if bit preservation were a solved problem? I suggest that proponents of this claim must feel confident that they could, at a minimum, preserve a petabyte of data undamaged for a century. Petabyte-scale data collections with long-term value, such as the Sloan Digital Sky Survey (SDSS)⁴ and the Protein Data Bank⁵ already exist, so this is asking them to surmount a rather low bar. How confident should proponents feel in their ability to keep a petabyte for a century? I suggest that they should have at least a 50% chance of success. Again, this is a rather low bar.

Proponents might bolster their case that these bars can easily be surmounted by pointing to claims such as: “ST5800 has a MTDDL (Mean Time To Data Loss) of 2.4×10^6 years.”⁶ (Sun Microsystems, 2008), or: “a Pergamum system capable of storing 10^{16} bytes of user data [will have] an MTDDL of 1.25×10^7 hours, or about 1,400 years.” (Storer, Greenan, Miller & Voruganti, 2008). These, and similar claims by other vendors, at first glance make it appear that bit preservation is indeed solved. Off-the-shelf solutions are ready to hand with performance so good that backups and checksums are quite superfluous. But do these claims stand up to examination?

⁴ The Sloan Digital Sky Survey: <http://www.sdss.org/>

⁵ Worldwide Protein Data Bank (WWPDB): <http://www.wwpdb.org/>

⁶ Numbers are expressed in powers-of-ten notation to help readers focus on the scale of the problems and the extraordinary level of reliability required.

Before using Sun's claim for its ST5800 as an example, I should stipulate that the ST5800 is an excellent product. It represents the state of the art in storage technology, and Sun's marketing claims represent the state of the art in storage marketing. Nevertheless, Sun does not guarantee that data in the ST5800 will last 2.4×10^6 years. Sun's terms and conditions explicitly disclaim any liability whatsoever for loss of, or damage to, the data the ST5800 stores (Sun Microsystems, [2006](#)) whenever it occurs.

All that the claim says is that if you watched a large number of ST5800 systems for a long time, recorded the time at which each of them first suffered a data loss, and then averaged these times, the result would be 2.4×10^6 years. Suppose Sun watched 10 ST5800s and noticed that three of them lost data during the first year, four of them lost data after 2.4×10^6 years, and the remaining three lost data after 4.8×10^6 years, they would be correct that the MTDL was 2.4×10^6 years. But we would not consider that a system with a 30% chance of data loss in the first year had solved the bit preservation problem. A single MTDL number isn't a useful characterization of a solution.

Consider the slightly more scientific claim made at the recent launch of the SC5800 by the marketing department of Sirius Cybernetics⁷: "SC5800 has a MTDL of $(2.4 \pm 0.4) \times 10^6$ years". Sirius thus claims that about 2/3 of the failures occurred between 2.0×10^6 and 2.8×10^6 years after the start of the experiment. They didn't start watching 10 SC5800s 2.8 million years ago. So how would they know?

Perhaps, instead of watching, say ten systems for 2.4×10^6 years they watched more systems for a shorter time. Sirius says they will sell 2×10^4 SC5800s per year at $\$5 \times 10^4$ each (a billion-a-year business), and they expect the product to be in the market for 10 years. The SC5800 has a service life of 10 years. So if Sirius watched their entire production of SC5800s ($\$10^{10}$ worth of storage systems) over their entire service life the experiment would end 20 years from now after accumulating about 2×10^6 system-years of data. If their claim is correct they would have about a 17% chance of seeing a single data loss event.

In other words, Sirius Cybernetics claims that the probability that no SC5800 will ever lose any data is over 80%. Or, since each SC5800 stores 5×10^{13} bytes, that there is an 80% probability that 10^{19} bytes of data will survive 10 years undamaged.

If one could believe the Sirius Cybernetics claim, the petabyte would look pretty safe for a century. But the claim clearly isn't based on an experiment that won't provide results until 2028 and even when it does will not validate the number in question. In fact, numbers like these are not the result of experiment at all. No feasible experiment could validate them. They are *projections*, based on models of how components of the system such as disks and software behave.

The state of the art in this kind of modeling is exemplified by the Pergamum Project at UC Santa Cruz (Storer et al., [2008](#)). Their model includes disk failures at rates derived from (Pinheiro, Weber & Barroso, [2007](#); Schroeder & Gibson, [2007](#)) and sector failures at rates derived from disk vendor specifications. Their system attempts

⁷ Purveyors of chatty doors, existential elevators and paranoid androids to the nobility and gentry of this galaxy (Adams, [1978](#)).

to conserve power by spinning the disks down whenever possible; they make an allowance for the effect of doing so on disk lifetime but it isn't clear upon what they base this. They report that the simulations were difficult:

“This lack of data is due to the extremely high reliability of these configurations - the simulator modeled many failures, but so few caused data loss that the simulation ran very slowly. This behavior is precisely what we want from an archival storage system: It can gracefully handle many failure events without losing data. Even though we captured fewer data points for the triple inter-parity configuration, we believe the reported MTTDL is a reasonable approximation.”

Although the Pergamum team's effort to obtain “a reasonable approximation” to the MTTDL of their system is praiseworthy, there are a number of reasons to believe that it overestimates the reliability of the system in practice:

- The model draws its failures from exponential distributions. They thus assume that both disk and sector failures are uncorrelated, although all measurements of actual failures (Bairavasundaram, Goodson, Schroeder, Arpaci-Dusseau & Arpaci-Dusseau, [2008](#); Talagala, [1999](#)) report significant correlations. Correlated failures greatly increase the probability of data loss (Baker et al., [2006](#); Elerath & Pecht, [2007](#)); this is the reason RAID systems fail to achieve the reliability predicted by the published models, which are based on exponential distributions (Patterson et al., [1988](#)).
- Other than a small reduction in disk lifetime from each power-on event, they assume that failure rates observed in always-on disk usage translate to their mostly-off environment. A study (Williams, Rosenthal, Roussopoulos & Georgis, [2008](#)) published after their paper reports a quantitative accelerated life test of data retention in almost-always-off disks. It shows that the 3.5” disks anticipated by the Pergamum team have data life dramatically worse in this usage mode than 2.5” disks using the same underlying technology.
- They assume that disk and sector failures are the only failures contributing to the system's failures, although a study (Krioukov et al., [2008](#)) shows that other hardware components contribute significantly.
- They assume that their software is bug-free, despite several studies of file and storage implementations (Engler, [2007](#); Jiang, Hu, Zhou & Kanevsky, [2008](#); Prabhakaran et al., [2005](#)) that uniformly report finding bugs capable of causing data loss in all systems studied.
- They also ignore all other threats to stored data (Rosenthal, Robertson, Lipkis, Reich & Morabito, [2005](#)) as possible causes of data loss. Among these threats are operator error, insider abuse and external attack. Each of these has been the subject of anecdotal reports of actual loss of preserved data. In particular, over 20 years at the San Diego Supercomputer Center, operator error is said to have been the cause of three-quarters of all data loss incidents (Moore, [2008](#)).

What can models like this tell us? Their results depend on both:

- the details of the simulation of the system being studied which, one hopes, accurately reflect its behavior, and
- the data used to drive the simulation which, one again hopes, accurately reflect the behavior of the system's components.

Under certain conditions, it is reasonable to use these models to compare different storage system technologies. The most important condition is that the models of the two systems use the same data. A claim that modeling showed system *A* to be more reliable than system *B* when the data used to model system *A* had much lower failure rates for components such as disk drives would not be credible.

These models may well be the best tools available to evaluate different techniques for preventing data loss, but they aren't adequate to determine whether bit preservation is a solved problem. We need to know the *maximum* rate at which data will be lost. The models assume things, such as uncorrelated errors and bug-free software, that all experimental studies show are false. The models exclude most of the threats to which stored data is subject. And in those cases where similar claims, such as those for disk reliability (Schroeder & Gibson, 2007; Pinheiro et al., 2007), have been tested they have been shown to be optimistic. It is not reasonable to assume that these factors are negligible, nor that they affect all systems equally; the models thus provide an estimate of the *minimum* data loss rate to be expected.

Even if we believed the models, the MTDDL number doesn't tell us how much data was lost in the average data loss event. Is petabyte system *A* with a MTDDL of 10^6 years better than a similar size system *B* with a MTDDL of 10^3 years? If the average data loss event in system *A* loses the entire petabyte, where the average data loss event in system *B* loses a kilobyte, it would be easy to argue that system *B* was 10^9 times better.

It is clear that we need a better way to define and measure bit preservation performance. Mean time to data loss is not a useful characterization of how well a system stores bits through time.

Theory

In order to claim that “bit preservation is a solved problem” we would need three things we currently don't have:

- A specific requirement as to how well bits need to be preserved.
- A technique for measuring whether actual systems achieve the required level of bit preservation.
- Measurements of an actual system using the technique that confirm it meets or exceeds the requirement.

In this section we suggest a metric that would be more useful than MTDDL, and ask whether it is possible to characterize actual systems in terms of this metric.

Defining a Solution

The most abstract model of a bit preservation system is as a black box, into which a string of bits $S(0)$ is placed at time $T(0)$ and from which at subsequent times $T(i)$ a string of bits $S(i)$ can be extracted. The system is successful if $S(i) = S(0)$ for all i .

No real-world system can be perfect and eternal, so real systems will fail. The simplest model of these failures is analogous to the decay of radioactive atoms. Each bit in the string independently is subject to a random process that has a constant small probability per unit time of causing its value to flip. The time after which there is a 50% probability that a bit will have flipped is the “bit half-life”.

The requirement of a 50% chance that a petabyte will survive for a century translates into a bit half-life of 8×10^{17} years. The current estimate of the age of the universe U is 1.4×10^{10} years, so this is a bit half-life approximately $6 \times 10^7 U$.

Measuring a Solution

Because current storage systems are extraordinarily reliable, measuring their bit half-life involves observing very large numbers of bits for a very long time. If you wanted to take a year to measure whether a system met the petabyte for-a-century requirement you might watch a thousand such systems, an exabyte of data. If the system were just good enough, you would see a single bit flip in just five of the systems.

Even if one were able to afford this experiment, doing so would be challenging. Data must be read from the system and compared with their expected value. Even if each bit is checked only once at the end of the year, the comparisons have to be performed with less than 1 chance in 10^{19} of any error.

In practice, estimates of bit half-life would have to be based upon the same models as estimates of MTDDL, and would thus share many of the same difficulties.

Of course, no-one expects an individual system of hardware and software to last a century. Preserved data will be migrated regularly to newer hardware and software. It might be thought that a system that migrated data every five years could meet the target with one-twentieth the reliability. There are two fallacies in this:

- The serial migration does not decrease the number of bit-years of exposure to risk.
- The migration process cannot simply be assumed to be error-free.

In the big picture, migration processes are part of the total system whose reliability is to be measured. Errors in these processes will also contribute to reducing the bit half-life.

Assessment

There is no escape from the problem that the size of the data collections to be preserved and the times for which they must be preserved mean that experimental confirmation that the technology chosen is up to the job is not economically feasible. Even if it was, the results would not be available soon enough to be useful. What this argument demonstrates is that, far from bit preservation being a solved problem, it is in a very specific sense an *unsolvable* problem. Even if we believed a system we developed was reliable enough, there are no feasible experiments that could confirm our belief in time to be useful.

Bit half-life is a more informative metric than MTDDL, because it is a measure of the reliability of the data, not a measure of the reliability of the system storing it. The data’s survival is what we care about. It thus captures the fact that the impact of a data

loss event depends not just on when it happens, but also on how much data is lost. It is still far from ideal:

- Bits in real storage systems do not fail independently; they exhibit significant correlations in space and time (Bairavasundaram et al., [2008](#)). These correlations make failure more likely than it otherwise would be. This observation doesn't invalidate the simple "radioactive decay" model; it merely makes adequate bit half-life a necessary but not sufficient condition for a system to meet the requirement.
- Like MTDL, it is a statistical estimate and thus, like MTDL, it is not useful without an uncertainty interval.
- Because storage systems are so reliable, it is just as difficult to measure bit half-life as it is to measure MTDL.

Practice

As enterprises such as Google (Chang et al., [2006](#)) and institutions such as the Sloan Digital Sky Survey and the Large Hadron Collider⁸ collect petabytes of data with long-term value that must remain on-line to be useful, questions of the economics and reliability of storage systems have become the focus of researchers' attention.

Storage Failures

Papers at the 2007 FAST conference used data from NetApp (Schroeder & Gibson, [2007](#)) and Google (Pinheiro et al., [2007](#)) to study disk replacement rates in large storage farms. They showed that the manufacturer's MTTF numbers were optimistic. Subsequent analysis of the NetApp data (Jiang et al., [2008](#)) showed that all other components contributed to the storage system failures, and:

‘Interestingly, [the earlier studies] found disks are replaced much more frequently (2-4 times) than vendor-specified [replacement rates]. But as this study indicates, there are other storage subsystem failures besides disk failures that are treated as disk faults and lead to unnecessary disk replacements.’

Two studies, one at CERN (Kelemen, [2007](#)) and one using data from NetApp (Bairavasundaram et al, [2008](#)), greatly improved on earlier work using data from the Internet Archive (Baker et al., [2006](#); Schwarz et al., [2006](#)). They studied *silent data corruption* in state-of-the-art storage systems; events in which the content of a file in storage changes with no explanation or recorded errors.

The NetApp study looked at the incidence of silent storage corruption in individual disks in RAID arrays. The data was collected over 41 months from NetApp's filers in the field, covering over 1.5×10^6 drives. They found over 4×10^5 silent corruption incidents. More than 3×10^4 of them were not detected until RAID restoration and could thus have caused data loss despite the replication and auditing provided by NetApp's row-diagonal parity RAID (Corbett et al., [2004](#)).

The CERN study used a program that wrote large files into CERN's various data stores, which represent a broad range of state-of-the-art enterprise storage systems (mostly RAID arrays), and checked them over a period of six months. A total of about 9.7×10^{16} bytes was written and about 1.92×10^8 bytes was found to have suffered

⁸ CERN, Worldwide LHC Computing Grid: <http://lcg.web.cern.ch/LCG/>

silent corruption, of which about 2/3 was persistent; re-reading did not return good data. In other words, about 1.2×10^{-9} of the data written to CERN's storage was permanently corrupted within six months. We can place an upper bound on the bit half-life in this sample of current storage systems by assuming that the data was written instantly at the start of the six months and checked instantly at the end; the result is 2×10^8 or about $10^{-2}U$. Thus to reach the petabyte for a century requirement we would need to improve the performance of current enterprise storage systems by a factor of at least 10^9 . Readers interested in some of the varied and complex causes of silent data corruption are referred to (Elerath, 2009) on disk technology and (Bairavasundaram et al., 2008) on other storage system components.

Surviving Storage Failures

Despite the manufacturer's claims, current research shows that state-of-the-art storage systems fall so many orders of magnitude below our bit preservation requirements that we cannot expect even dramatic improvements in technology to fill the gap. Maintaining a single replica in a single storage system is not an adequate solution to the bit preservation problem.

Practical digital preservation systems must therefore:

- Maintain more than one copy by *replicating* their data on multiple, ideally different, storage systems.
- Audit or (*scrub*) the replicas to detect damage, and repair them by overwriting the known-bad copy with data from another.

The more replicas and the more frequently they are audited and repaired the longer the bit half-life we can expect. This is, after all, the basis for the backups and checksums technique in common use. In fact, current storage systems already use versions of these techniques, for example in the form of RAID (Patterson, Gibson & Katz, 1988). Despite this, the bit half-life they deliver is inadequate. Unfortunately, adding the necessary inter-storage-system replication and scrubbing is expensive.

2007 cost figures from the San Diego Supercomputer Center (Moore et al., 2007) show that maintaining a single on-line copy of a petabyte for a year then cost about $\$1.5 \times 10^6$. A single near-line copy on tape cost about $\$5 \times 10^5$ a year⁹. These costs decrease with time, albeit not as fast as raw disk costs. The British Library estimates a 30% per annum decrease. Assuming that this rate continues for at least a decade, if you can afford about 3.3 times the first year's cost to store an extra replica for a decade, you can afford to store it indefinitely. So, adding a second replica of a petabyte on disk would cost about $\$3.5 \times 10^6$ and on tape would cost about $\$1.4 \times 10^6$. Adding cost to a preservation effort to increase reliability in this way is a two-edged sword; doing so necessarily increases the risk that preservation will fail for economic reasons.

Disk storage has a long history of rapid decline (Christensen, 1997). It is tempting to assume that, as time goes by, it will be economically feasible to add more replicas of given content, and that doing so in effect reduces the problem of a petabyte-for-a-century to a petabyte-for-a-decade. There are three flaws in this argument:

⁹ SDSC reports that the 2008 costs are $\$1.05 \times 10^6$ and $\$4.2 \times 10^5$

- It ignores the fact that raw media cost, which has decreased rapidly, is a minority of the total cost of storage (Moore et al., [2007](#)). The other factors have not decreased exponentially with time.
- It ignores the fact that, just as storage cost per byte has dropped exponentially, the amount of data to be stored has grown exponentially.
- A senior disk drive engineer recently pointed out that although the technology for ever larger consumer 3.5" disk drives is available, the business case for them has collapsed (Anderson, [2009](#)). This makes it likely that in the near future the exponential drop in the cost of large-scale disk storage will cease, at least for a few years.

If the cost-per-byte curve stays flat, even for a few years in the face of exponential growth in storage, demand the result will be a crisis in preservation economics.

Further, without detailed understanding of the rates at which different mechanisms cause loss and damage, it isn't possible to derive from a desired bit half-life the appropriate number of replicas¹⁰ and thus the cost implication of replication. At small scales the response to this uncertainty is to add more replicas, but as the scale increases this rapidly becomes unaffordable.

Replicating among identical systems is much less effective than replicating among diverse systems. Identical systems are subject to common mode failures, for example caused by a software bug in all the systems damaging the same data in each. On the other hand, purchasing and operating a number of identical systems will be considerably cheaper than operating a set of diverse systems.

Each replica is vulnerable to loss and damage. Unless they are regularly audited they contribute little to increasing bit half-life. The bandwidth and processing capacity needed to scrub the data are both costly; adding these costs increases the risk of failure. Custom hardware (Michail, Kakarountas, Theodoridis & Goutis, [2005](#)) could compute the SHA-1 (National Institute of Standards and Technology [NIST], [1995](#)) checksum of a petabyte of data in a month, but doing so requires impressive bandwidth - the equivalent of three gigabit Ethernet interfaces running at full speed the entire month. User access to data in preservation systems is typically infrequent; they are therefore rarely designed to provide such high-bandwidth read access. System cost increases rapidly with I/O bandwidth, and the additional accesses to the data (whether on disk or on tape) needed for scrubbing themselves potentially increase the risk of failure.

The point of writing software that reads and verifies stored data in this way is to detect damage and exploit replication to repair it, thereby increasing bit half-life. How well can we do this? RAID is an example of a software technique of this type applied to disks. In practice, the CERN study (Kelemen, [2007](#)) looking at real RAID systems from the outside showed a significant rate of silent data corruption, and the NetApp study (Bairavasundaram et al., [2008](#)) looking at them from the inside showed a significant rate of silent disk errors that would lead to silent data corruption. A study (Krioukov et al., [2008](#)) of the full range of current algorithms used to implement RAID

¹⁰ The number can be quite large; a study of paper journals (Yano, [2008](#)) found between 3 and 31 copies were needed to achieve loss probabilities over a century of between 10^{-3} and 10^{-6} given various plausible loss rates of the individual copies. The lower repairability of paper copies inflates these numbers, while their greater durability deflates them, as against digital copies.

found flaws leading to potential data loss in all of them. Both this study, and another from IBM (Hafner, Deenadhayalan, Belluomini & Rao, [2008](#)), propose improvements to these algorithms but neither claim that they can eliminate silent corruption, or even accurately predict its incidence:

“while we attempt to use as realistic probability numbers as possible, the goal is not to provide precise data loss probabilities, but to illustrate the advantage of using a model checker, and discuss potential trade-offs between different protection schemes.” (Krioukov et al., [2008](#))

Thus although replication and scrubbing are capable of decreasing the incidence of data loss in current storage systems, they cannot eliminate it completely. And the replication and scrubbing software itself will contain bugs that can cause data loss. It must be doubtful that we can implement these techniques well enough to increase the bit half-life of systems with an affordable number of replicas by 10^9 .

It takes experiments with petabytes of storage to characterize the performance of current systems accurately. Even if we believed we had implemented replication and audit well enough to improve performance by 10^9 , we could not afford to do the experiments that would be needed to confirm it.

Policy

If bit preservation were a solved problem then it would be reasonable to expect that no bits would be lost. This is not the case; just as in paper archives, preserved content in digital archives will be lost or damaged. Setting unreasonable expectations for the performance of our preservation systems, for example by continually making unsupported claims to have solved the bit preservation problem, is simply setting ourselves up to be perceived as failures.

If preserved bits will be lost, the question becomes how to invest the limited funds available to reduce the rate of loss as much as possible. It is a commonplace that if you can measure something you can improve it. The history of technology markets such as CPUs and graphics chips show that competition between vendors based on widely accepted standard benchmarks can drive rapid improvements in component cost-performance. Alas, although raw storage cost is easily measured and is the subject of effective competition to decrease cost per byte (Christensen, [1997](#)), long-term storage reliability is very hard to measure and the accepted metric for it is not very informative. Competition to reduce the cost of a given level of bit preservation is therefore much less effective.

It is in the interest of the digital preservation community to improve competition in their market. How could this be done?

- Agreement on a metric for bit preservation performance is an essential first step. It would be extremely valuable if it were possible to define one that was easily measurable, but this seems rather unlikely.
- Given this, it seems likely that numbers for bit preservation performance will continue to be generated by models. Achieving consensus on modeling techniques is important, especially as it appears that traditional techniques are running into difficulties (Elerath & Pecht, [2007](#); Storer et al., [2008](#)).

- These models will need agreed data. Better and more widely available data about the real world performance storage components are thus important. Realistic studies have only begun to be published, and they aren't yet based on shared metrics. The effort by Usenix and Carnegie-Mellon¹¹ to establish a repository for suitably anonymized data of this kind is to be commended.
- Storage systems are currently designed using completely inadequate models of how components fail. One problem is that these failures are highly correlated, making the models complex and difficult. A shared model of the threats against which bits need to be preserved, models of these threats, and data regarding their incidence is also important.
- Anecdotal evidence suggests that operator error and insider abuse are major causes of data loss in large storage farms; they are difficult to model or characterize. This is in part because sites are very reluctant to admit to data loss incidents. An anonymous incident reporting system modeled on NASA's Aviation Safety Reporting Conclusions System¹² would be very valuable in understanding the mechanisms of, and defending against, these failures.

The fact that it is possible for digital information to be copied perfectly does not mean that it always will be. While perfection is not within the grasp of real-world engineers, improvement is always possible. However, improvement takes money, and without the research outlined above we are unable to make rational trade-offs between the cost of preserving content to a given level of reliability and the cost of the losses implied by the given level.

Conclusions

As we have seen, the case that bit preservation is a solved problem rests on the conviction that the conventional techniques of backups and checksums are more than adequate to the scale of the problem. This conviction is odd. Press accounts (e.g., Brodtkin, 2008) of companies, presumably using the conventional techniques, nevertheless losing essential data are common. Awareness that systems frequently encounter scaling problems is also widespread, as is the expectation that the future demands for preserving digital content will be enormous.

But the case for bit preservation not being solved does not rest on this cognitive dissonance. It rests rather on the many orders of magnitude mismatch between the reliability requirements implied by society's expectations of the amount of data to be preserved and the length of time for which it should be preserved, and the observed performance of current storage hardware and software.

Were every bit to come adequately endowed with capital to provide guaranteed funds through time its preservation would not be a major concern, although it would still not be a solved problem. Like almost all engineering problems, bit preservation is fundamentally a question of budgets. Society's ever-increasing demands for vast amounts of data to be kept for the future are not matched by suitably lavish funds. Thus, absent a technological miracle, bit preservation is a problem with which we are doomed to struggle indefinitely.

¹¹ Usenix, the computer failure data repository (CFDR): <http://cfdrr.usenix.org/>

¹² NASA. Aviation Safety Reporting System: <http://asrs.arc.nasa.gov/>

Acknowledgements

Thanks are due to Michael Bax and the LOCKSS engineering team for critical readings of drafts of this paper, and to the staff of the San Diego Supercomputer Center for the discussions that started me thinking along these lines.

References

- Adams, D. (1978). *The Hitch-Hiker's Guide to the Galaxy*. British Broadcasting Corporation.
- Anderson, D. (2009) Hard drive directions. In *Designing Storage Architectures for Preservation Collections*. Library of Congress, September 22-23, 2009. Retrieved March 30, 2010 from http://digitalpreservation.gov/news/events/other_meetings/storage09/docs/2-4_Anderson-seagate-v3_HDtrends.pdf
- Bairavasundaram, L., Goodson, G., Schroeder, B., Arpaci-Dusseau, A. C., & Arpaci-Dusseau, R. H. (2008). An analysis of data corruption in the storage stack. In *Proceedings of 6th USENIX Conference on File and Storage Technologies*.
- Baker, M., Shah, M., Rosenthal, D. S. H., Roussopoulos, M., Maniatis, P., Giuli, T., et al. (2006). A fresh look at the reliability of long-term digital storage. In *Proceedings of EuroSys2006*.
- Brodkin, J. (2008, August). Loss of customer data spurs closure of online storage service "The Linkup". *Network World*. August 11, 2008.
- Chang, F., Dean, J., Ghemawat, S., Hsieh, W. C., Wallach, D. A., Burrows, M., et al. (2006). Bigtable: A distributed storage system for structured data. In *Proceedings of the 7th Usenix Symposium on Operating System Design and Implementation*, pp. 205–218.
- Christensen, C.M. (1997). *The innovator's dilemma: When new technologies cause great firms to fail*. Harvard Business School Press.
- Corbett, P., English, B., Goel, A., Grcanac, T., Kleiman, S., Leong, J., et al. (2004). Row-diagonal parity for double disk failure correction. In *3rd Usenix Conference on File and Storage Technologies*.
- Elerath, J. G. (2009) Hard-disk drives: The good, the bad, and the ugly. *Comm. ACM* v52 n6. <http://cacm.acm.org/magazines/2009/6/28493-hard-disk-drives-the-good-the-bad-and-the-ugly/fulltext>
- Elerath, J. G., & Pecht, M. (2007). Enhanced reliability modeling of RAID storage systems. In *DSN '07: Proceedings of the 37th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, pp. 175–184. Washington, DC: IEEE Computer Society.

-
- Engler, D. (2007). A system's hackers crash course: Techniques that find lots of bugs in real (storage) system code. In *Proceedings of 5th USENIX Conference on File and Storage Technologies*.
- Hafner, J. L., Deenadhayalan, V., Belluomini, W., & Rao, K. (2008). Undetected disk errors in RAID arrays. *IBM J. Research & Development* 52(4/5).
- Jiang, W., Hu, C., Zhou, Y., & Kanevsky, A. (2008). Are disks the dominant contributor for storage failures? A comprehensive study of storage subsystem failure characteristics. In *Proceedings of 6th USENIX Conference on File and Storage Technologies*.
- Kelemen, P. (2007). Silent corruptions. In *8th Annual Workshop on Linux Clusters for Super Computing*.
- Krioukov, A., Bairavasundaram, L. N., Goodson, G. R., Srinivasan, K., Thelen, R., Arpaci-Dusseau, A. C., et al. (2008). Parity lost and parity regained. In *Proceedings of 6th USENIX Conference on File and Storage Technologies*.
- Michail, H. E., Kakarountas, A. P., Theodoridis, G., & Goutis, C. E. (2005). A low-power and high-throughput implementation of the SHA-1 hash function. In *Proceedings of the 9th WSEAS International Conference on Computers*.
- Moore, R. L., D'Aoust, J., McDonald, R. H., & Minor, D. (2007). *Disk and tape storage cost models*. In *Archiving 2007*.
- Moore, R., (personal communication, 2008)
- National Institute of Standards and Technology. (1995). *Secure hash standard (SHA-1)*. Federal Information Processing Standard Publication 180-1, NIST; Washington, D.C.
- Parkinson, C. N. (1957). *Parkinson's law*. Buccaneer Books.
- Patterson, D. A., Gibson, G., & Katz, R. H. (1988). A case for redundant arrays of inexpensive disks (RAID). In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, pp. 109–116.
- Pinheiro, E., Weber, W.-D., & Barroso, L. A. (2007). Failure trends in a large disk drive population. In *Proceedings of 5th USENIX Conference on File and Storage Technologies*.
- Popper, K. (1959). *Logic of scientific discovery*. (cf. Chapter X, footnote *5.) Hutchinson.

- Prabhakaran, V., Agrawal, N., Bairavasundaram, L., Gunawi, H., Arpaci-Dusseau, A. C., & Arpaci-Dusseau, R. H. (2005). IRON file systems. In *Proceedings of the 20th Symposium on Operating Systems Principles*.
- Rosenthal, D. S. H., Robertson, T. S., Lipkis, T., Reich, V., & Morabito, S. (2005, November) Requirements for digital preservation systems: A bottom-up approach. *D-Lib Magazine 11*(11). Retrieved from <http://www.dlib.org/dlib/november05/rosenthal/11rosenthal.html>
- Schroeder, B., & Gibson, G. (2007). Disk failures in the real world: What does an MTTF of 1,000,000 hours mean to you? In *Proceedings of 5th USENIX Conference on File and Storage Technologies*.
- Schwarz, T., Baker, M., Bassi, S., Baumgart, B., Flagg, W., van Imngen, C., et al. (2006). Disk failure investigations at the Internet Archive. In *Work-in-Progress Session, NASA/IEEE Conference on Mass Storage Systems and Technologies*.
- Storer, M. W., Greenan, K. M., Miller, E. L., & Voruganti, K. (2008). Pergamum: Replacing tape with energy efficient, reliable, disk-based archival storage. In *Proceedings of 6th USENIX Conference on File and Storage Technologies*.
- Sun Microsystems. (2006). *Sales terms and conditions, Section 11.2*. Retrieved from http://store.sun.com/CMTemplate/docs/legal_terms/TnC.jsp#11
- Sun Microsystems. (2008). ST5800 presentation. Sun PASIG Meeting.
- Talagala, N. (1999). *Characterizing large storage systems: Error behavior and performance benchmarks*. (Doctoral Dissertation), Computer Science Division, University of California at Berkeley, Berkeley, CA, USA.
- Williams, P., Rosenthal, D. S. H., Roussopoulos, M., & Georgis, S. (2008). Predicting the archival life of removable hard disk drives. In *Archiving 2008*.
- Yano, C. (2008). *How many journal copies?* A Preliminary Report. Presentation to ALA.