

# Reliability and Security of RAID Storage Systems and D2D Archives Using SATA Disk Drives

GORDON F. HUGHES and JOSEPH F. MURRAY  
University of California San Diego

---

Information storage reliability and security is addressed by using personal computer disk drives in enterprise-class nearline and archival storage systems. The low cost of these serial ATA (SATA) PC drives is a tradeoff against drive reliability design and demonstration test levels, which are higher in the more expensive SCSI and Fibre Channel drives. This article discusses the tradeoff between SATA which has the advantage that fewer higher capacity drives are needed for a given system storage capacity, which further reduces cost and allows higher drive failure rates, and the use of additional storage system redundancy and drive failure prediction to maintain system data integrity using less reliable drives. RAID stripe failure probability is calculated using typical ATA and SCSI drive failure rates, for single and double parity data reconstruction failure, and failure due to drive unrecoverable block errors. Reliability improvement from drive failure prediction is also calculated, and can be significant. Today's SATA drive specifications for unrecoverable block errors appear to allow stripe reconstruction failure, and additional in-drive parity blocks are suggested as a solution. The possibility of using low cost disks data for backup and archiving is discussed, replacing higher cost magnetic tape. This requires significantly better RAID stripe failure probability, and suitable drive technology alternatives are discussed. The failure rate of nonoperating drives is estimated using failure analysis results from  $\approx 4000$  drives. Nonoperating RAID stripe failure rates are thereby estimated. User data security needs to be assured in addition to reliability, and to extend past the point where physical control of drives is lost, such as when drives are removed from systems for data vaulting, repair, sale, or discard. Today, over a third of resold drives contain unerased user data. Security is proposed via the existing SATA drive secure-erase command, or via the existing SATA drive password commands, or by data encryption. Finally, backup and archival disc storage is compared to magnetic tape, a technology with a proven reliability record over the full half-century of digital data storage. In contrast, tape archives are not vulnerable to tape transport failure modes. Only failure modes in the archived tapes and reels will make data unrecoverable.

Categories and Subject Descriptors: B.8.1 [**Performance and Reliability**]: Reliability, Testing, and Fault-Tolerance; C.4 [**Computer Systems Organization**]: Performance of Systems—*Fault tolerance*

---

This research was supported by the Information Storage Industry Center at the University of California at San Diego, sponsored by the Alfred P. Sloan Foundation.

Authors' addresses: G. F. Hughes, Center for Magnetic Recording Research, University of California, San Diego, CA 92093-0401; email: gfhughes@ucsd.edu; J. F. Murray, ECE Dept. University of California, San Diego, CA 92093-0409.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or direct commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 1515 Broadway, New York, NY 10036 USA, fax: +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2004 ACM 1553-3077/04/1200-0095 \$5.00

General Terms: Reliability

Additional Key Words and Phrases: Disk drive, failure prediction, SMART, secure erase, storage systems architecture, storage resource management, SATA, archival storage

---

## 1. INTRODUCTION

Personal computer disk drives (serial ATA computer interface, SATA) are now being used in enterprise-class storage systems, primarily in disk-to-disk backup (D2D) of RAID arrays to nearline SATA arrays. (See Table I glossary). This allows faster RAID backup than tape, simultaneously with RAID user data access, and can be comparable in cost to tape backup. Disk-to-disk-to-tape (D2D2T), which offloads the tape backup operation overhead from the primary RAID user array, is most common. The PC drive marketplace demands that drive costs be balanced against drive reliability design and reliability demonstration tests. SATA drives offer the low costs of ATA drives along with the highest areal bit densities to provide maximum capacity per-drive (up to 400 GB today). A SATA drive can be an ATA drive with a modified circuit board (PCB) to provide the SATA interface changes, or it can be a newly-designed low-cost drive. This is similar to Fibre Channel drives, which can be modified SCSI drives (FC uses the SCSI command set, with a serial physical transport). In this article, a SATA drive will be assumed to be a high-volume ATA drive with a SATA PCB.

This article discusses disk drive reliability of ATA, SATA, and SCSI drives, and ways to obtain high storage system reliability with less reliable individual drives. The authors *do not* advocate one drive interface over others. Drives with *any* interface can be designed and validated at similar reliability levels. The article *does* propose that RAID parity architecture with drive monitoring can provide adequate RAID system reliability using drives designed for lower cost and lower stress PC use. The increasing number of available drive interfaces (ATA, SATA, SCSI, FC, iFC, iSCSI) are simply opportunities for more flexible storage system design choices. The marketplace will sort out their relative merits.

PC ATA drives are designed for daytime office use, not  $24 \times 7$ . The market lifetime of PC drive products is closer to one year than five years, limiting the availability of long-term field reliability data. Few PC drives are returned to their manufacturer for failure analysis, and their service environment is unknown (temperature, vibration, and electrical power). Some ATA drives have been only partially flaw-marked during manufacturing final test, because the necessary test hours are not available in the available manufacturing test floor space and test time. Reliable flaw-marking requires many writes and reads of the entire drive, but each full write/read cycle of a 250GB, 7200rpm drive will take about 80 minutes. Internal self-test commands in these drives can be used for subsequent flaw-marking in storage systems.

Enterprise class SCSI/FC drives are designed and tested for high reliability under heavy service duty, and are fully tested and flaw-marked in their manufacturing process, making their cost higher [Anderson et al. 2003]. New models are tested for design reliability evaluation, latent design flaws are found and designed out, and thousands of drives are tested to demonstrate reliability. Such

Table I. Definitions and Acronyms

Acronym	Definition
AFR	Annual drive failure rate per year = 8760/MTBF
ATA	AT Attachment—standard “desktop” PC drive interface
D2D	Disk to Disk backup storage architecture
D2D2T	Disk to Disk to Tape backup storage architecture
ECC	Bit Error Correction Code (by drive PCB electronics and stored with user data)
FC	Fiber Channel drive interface, enterprise storage
iFP	FC storage via TCP/IP network
iSCSI	SCSI storage via TCP/IP network
MAID	Massive Array of Idle Disks (powered down when unused)
Nearline Disk	Disk Array with high serial transfer rate; Allows fast RAID array backup
MTBF	Mean Time to Failure (hours to 50% failing)
PCB	Printed Circuit Board—drive electronics
RAID	Redundant Array of Independent Disks (RAID-5 assumed, with parity ECC)
RPM	Disk revolutions per minute, 5400 to 15000
SATA	Serial evolution of ATA, with some SCSI features such as access re-queuing
SCSI	Small Computer System Interface, standard for enterprise systems
SMART	Drive Self Monitoring and Reporting Technology (drive failure prediction)

tests of perhaps 5000 drives over 30 days, produce millions of hours of total drive operating time, which can demonstrate million-hour MTBF *shortterm* reliability, namely the test time term. Reliability maturity testing of many drives over the manufacturing life of the product will assess longterm reliability. ATA drives undergo a low cost version of this testing. Also, SCSI product life is longer, allowing longterm field reliability failures to be captured and analyzed for design correction and improvement. The argument has been made that the higher manufacturing volume of ATA drives makes them highly reliable at low cost, but this is an argument for *shortterm* reliability, which is not the same as longterm.

SCSI drive capacity today is 18-146GB, set by storage system access speed and redundancy requirements, while SATA drives have up to 400GB. This article will use 36GB SCSI drives and 250GB SATA drives for comparisons. Lower SCSI drive capacity means lower bit areal density, which allows higher drive design operating margins than SATA. SCSI drives use higher cost spin bearings and thermal designs. High rpm SCSI/FC 3.5-inch drives use the 2.5-inch drive disk diameter (65mm) rather than the larger 3.5-inch (95mm) disk diameter of ATA and SATA. This reduces heat from disk spin windage, reduces spin bearing loads, and makes space available in the 3.5-inch form factor for cooling and for a more powerful track servo actuator. The lower bit transfer rate improves read channel signal-to-noise for lower error rate and higher margins. The penalty is half the usable disk area of a 3.5-inch disk and, therefore, half the capacity of a similar areal density SATA drive. One compensation is lower probability of an unrecoverable error in a full SCSI drive read (analyzed in Section 7).

It is reasonable to contend that drives more reliable than PC drives will likely have smaller maximum capacity. As just demonstrated, SCSI drive designers take the reliability advantages from the smaller drive capacities that normal RAID systems want and adjust for high performance by the simultaneous access to many drives.

While SATA drives can certainly be designed and tested to be as rugged and reliable as SCSI/FC drives with similar capacities, their costs would then be comparable, because the same technology basis is used in all disk drives. Their differences lie in engineering and reliability design, test, and performance margins. The design specifics are discussed in [Anderson et al. 2003]. Fibre Channel interface PCBs can equally well be put on low-cost PC drives, and some have already been announced. The important point is that enterprise-class reliability comes from drive design and testing, not from the interface.

This cost vs. reliability dichotomy is seen in many other electronic products. For example, electronics for space and undersea applications have significantly higher investment in reliability design and testing than consumer electronics.

SATA systems do have a reliability advantage, in that enterprise storage systems use modest capacity SCSI/FC drives (36GB is common today), and achieve a desired total system capacity by using many drives. A comparable capacity SATA system might use drives seven times larger, such as 250GB capacity, thus needing seven times fewer drives than SCSI. This would allow SATA drives to have seven times the SCSI failure rate, for a similar drive replacement rate burden. However, unrecoverable error specifications in today's SATA drives will be shown to still limit maximum drive capacity, but the problem can be avoided by double array and in-drive parity. There is also an unquantifiable reliability stress with higher capacity per-disk, but uncertainty regarding even the *mean* reliability of SCSI and SATA drives obscures generalizations about such finer distinctions.

SCSI failure rates are specified and tested for the worst allowable conditions, namely maximum temperature and worst power quality, and the maximum number of disks in the product model family, typically two to four. So in a two-disk 73GB drive, its one-disk 36GB family member has higher reliability than specified for the overall model family. ATA failure rates are usually nominal PC drive operating conditions and an  $8 \times 5$  duty cycle.

This article proposes ways for SATA storage system designers to maintain high system level reliability using moderate reliability SATA drives, by requiring appropriate SATA drive reliability testing, by increased RAID redundancy, and by storage system management of drive failure warnings.

Suitable SATA hardware architecture is first discussed in Section 2, followed by SATA drive reliability factors that bear on system performance in Section 3. Section 4 presents information on drive operating failure rates. In Section 5, failure mode frequency data is presented from about 4000 drives of a single model and used to estimate nonoperating drive failure rates. Nearline SATA systems, with single and double parity, are analyzed for drive reconstruction failure probability using these failure rates in Section 6. Significant reliability improvement is possible by modestly successful drive failure prediction (SMART), with error counts which can be read from drives by the system software using standard SATA drive commands.

Section 7 calculates the probability of RAID drive stripe failure, if an unrecoverable block error occurs in reading one of the nonfailed drives. This probability is uncomfortably high for large SATA drives, and double stripe parity is

suggested as a solution. Section 8 applies these failure probabilities to nearline D2D2T systems.

D2D archiving is discussed in Section 9, which imposes much stricter disk system reliability requirements. Archiving is assumed to require more than five years of data storage in unpowered drives. This section of the article addresses whether disk drives are suitable to replace tape backup, by comparing operating vs. nonoperating failure modes and failure rates. Reliability is analyzed using the estimated failure rates of nonoperating drives from Section 5. Suitable drive internal technology choices are suggested.

User data security for drives removed from storage systems is addressed in Section 10. This is a serious but largely ignored issue because fully a third of resold disk drives contain unerased user data [Garfinkel and Shelat 2003]. RAID drives are not invulnerable to data security risk.

Finally, Section 11 suggests industry standardization of SATA storage system reliability specifications and reliability test levels to help this industry grow.

## 2. SATA DRIVE STORAGE SYSTEMS

Nearline or offline storage systems should actively manage power similar to tape backup systems, which reduce power and heat load by not powering tape reels in storage. The ATA and SATA specifications already contain appropriate power management commands. MAID system architecture (Massive Array of Idle Disks) appears suitable for nearline D2D systems [Colarelli 2002]. The Storage Networking Industry Association defines MAID as “a storage system comprising an array of disk drives that are powered down individually or in groups when not required. MAID storage systems reduce the power consumed by a storage array” (see [www.snia.org](http://www.snia.org)). Idle drives can be periodically powered up and SATA internal self-test drive diagnostics run. This spends very few system performance resources on drive testing and can initiate copying a failing drives or reconstruction of a defective drive.

An alternative to MAID would be a system similar to a tape library, where drives are removed and remotely stored. User data security is an issue in this case (Section 10).

## 3. SATA RAID DRIVE RELIABILITY REQUIREMENTS

SATA drive customers should specify that drive models being qualified undergo standard SCSI/FC reliability tests, including design verification testing, reliability demonstration testing, and design maturity testing. The specific reliability test requirements should be balanced against the drive price, but reliability MTBF/AFR must either be verified, or be considered as a goal rather than a guarantee. Current SATA and SCSI/FC manufacturers specify similar MTBF and annual failure rates, but do not specify the reliability test levels. SATA unrecoverable error rates are an order of magnitude higher than SCSI (Section 7).

Secondly, SATA drive customers should specify the reliability testing required for each drive in the manufacturing final test. The final test should require flaw-marking of all drive data record locations, with a specified number

of data writes and reads. This is normal for SCSI/FC drives, but PC drives have been sold with only the first 20% of the user data records flaw-tested, and those perhaps only written and read once. Drive customers could safely choose to do part of the flaw-marking in storage systems via existing drive internal test commands to lower the cost of mature SATA drives.

#### 4. DRIVE OPERATING FAILURE RATES

Disk drive longterm annual failure rates data are historically known from field experience in enterprise storage systems. Drive manufacturers analyze returned drives for failure causes and design improvements. Many “failed” SCSI drives are seen, but only a few percent of the ATA field “failures” (the quote refers to the common fact that many “failed” ATA or SCSI drives operate correctly when tested by their manufacturer).

Disk storage systems usually operate  $24 \times 7$ , so the historical failure data are operating (power-on-hours) failure rates. These range from 0.3–3% per year [Hughes 2002]. A typical drive specification goal today for both SATA or SCSI/Fibre channel is 0.7% annual failure rate ( $=24 \times 365 / 1,200,000$  hours MTBF). Another manufacturer specified reliability factor is “five-year drive life” This has recently become a drive manufacturer warranty guarantee. There is a contradiction here, since 1,200,000 hour MTBF means that 50% of the drives will survive for 137 years, far more than their design life of 5. The actual reliability goal and its test validation is to meet the 0.7% AFR for 3–5 years. Lower failure rates are expected in enterprise-class drives and systems with professional handling, cooling, vibration, and power. Personal (PC) drives can fail at higher rates, especially if mishandled by users (drive bearings have higher precision and mechanical shock vulnerability than professional film cameras or aircraft gyroscopic instruments).

#### 5. DRIVE NONOPERATING FAILURE RATES

Less is known about nonoperating failure rates, for example, in MAID backup systems where drives spend appreciable time unpowered. Archival MAID system failure rates are especially uncertain. While nonoperating failure rates are expected to be lower than operating rates, they will not be zero. Drives are complex devices with hundreds of possible failure modes, and although fewer modes are relevant in nonoperation, some are more likely—such as head/disk stiction and corrosion. Existing drive design options can ameliorate these nonoperating risks, for example, by ramp loading heads onto spinning disks, instead of having heads sitting on disk surfaces and sliding on the disk surface while the disk starts or stops spinning, risking stiction and wear.

Little is known quantitatively about nonoperating failure rates because drives have not been historically tested and used in backup service. These rates could be  $1/10$  to  $1/2$  of operating failure rates. This estimate comes from Table II, a Pareto list of typical drive failure modes from about 4000 drives returned for test from one drive manufacturer’s failure analysis categorization into 200 failure modes. The failure modes are listed in their frequency of occurrence and are typical of all drives—however, their frequency of occurrence will vary by manufacturer and class of drives. “No problem found” is high on any failure

Table II. Failure Types, Occurrence Frequency, and Primary Stress Condition

Failure mode	Description	Frequency	Stress Condition
Head-Disk interference	Head crash	15.5%	Operating
No problem found	Returned drive tests ok	15.0%	—
Recording heads	Complex nanotech devices	14.5%	Operating
Post manufacture	Drive handling damage	10.1%	Non-Op
Circuit board “PCB”	Many IC components	8.5%	Operating
Head or disk corrosion	Causes disk HDI or defects	7.7%	Non Op
Head assembly “E-Block”	Wires, preamp, coil	6.8%	Operating
Head Disk Assembly	Mechanics, Electrical, Voice coil	3.9%	Operating
Disk defects	Causes HDI or read errors	2.6%	Operating
Drive firmware	Internal operating system	1.9%	Operating
Head-disk stiction	Disk won’t spin up in drive	1.3%	Non-Op
Spindle bearing	Disk spin bearing	1.1%	Operating
Contamination inside drive	Foreign gases or chemicals	0.7%	Op/Non-Op

frequency list. The third column shows the failure frequencies (the least frequent failures have been omitted). The fourth column estimates whether the mode is more likely to occur in an operating or a nonoperating drive.

About 20% of the 4000 failures occurred in nonoperating drive stress modes in Table II, leading to the estimate that the failure rate of a nonoperating drive is 1/10 to 1/2 that of an operating drive. These numbers will be used later for RAID reconstruction failure reliability calculations.

In remote vault backup systems, drive mechanical disconnect and reconnect will induce another failure mode—signal and power connector failure. Unless drive connectors are specifically designed for frequent cycling, this mode can be severe. Manufacturers today typically warn against repeated connector cycling in ATA drives.

## 6. SATA NEARLINE RAID STORAGE SYSTEMS

Nearline systems that offer reliability requirements similar to online RAID systems and drive failure rates of a fractional percent per-year may be bearable if parity reconstruction does not impose significant system performance burdens [Schwarz et al. 1995]. A RAID-5 reconstruction failure (too many drives failing before detection and reconstruction) can be recovered by a tape backup restore, although this is often an emergency situation in enterprise storage.

Lower drive reliability can be countered by RAID designs that insure data reconstruction with two drives failed (double parity, or parity-2). This lowers the reconstruction failure probability since three or more drives must fail before it fails. Double parity is used today in some RAID systems [Lueth 2004].

Table III shows a simple example of RAID 5 double parity for data reconstruction with two disk failures, by adding column parity per-stripe-block to the normal row parity per-RAID-stripe. (This is intended as a simple example, not an efficient design). If drives D1 and D2 fail, reconstruct block S1, B1, using block B1 (column) parity with drives D5, D4, and D3 onto a new hot spare drive D1; then reconstruct block S1, B2 using the reconstructed S1, B1, and stripe S1 row parity onto new D2; then reconstruct S2, B2 onto new D1 using B2 column parity. System performance can be further enhanced by more complex double

Table III. Simple Example of RAID Double Parity

RAID Array Stripes S1, S2, S3	RAID Stripe Data Blocks B1-B4 (3 Stripes + 2 Parity on 5 Drives)				
	S1, B1 (D1)	S1, B2 (D2)	S1, B3 (D3)	S1, B4 (D4)	Parity S1 (D5)
	S2, B1 (D5)	S2, B2 (D1)	S2, B3 (D2)	S2, B4 (D3)	Parity S2 (D4)
	S3, B1 (D4)	S3, B2 (D5)	S3, B3 (D1)	S3, B4 (D2)	Parity S3 (D3)
	Parity B1 (D3)	Parity B2 (D4)	Parity B3 (D5)	Parity B4 (D1)	Parity S4 (D2)

parity schemes [Lee and Park 1996; Lueth 2004]. Unrecoverable block errors must be considered (Section 7).

SATA storage server software can also periodically read the SMART failure prediction warning flags from SATA drives [Hamerly and Elkan 2001; Hughes 2002; Murray and Hughes 2003; Murray et al. 2004]. A drive which triggers a warning can be copied onto a hot spare, eliminating the RAID drive reconstruction performance penalty and possible unrecoverable errors in reading all other drive data in the failing drive's stripe during a reconstruction (see Section 7). The attainable SMART prediction accuracy is 10–50%, the maximum number reflecting the limited number of failure modes that drives internally monitor out of hundreds of potential modes, and the unpredictability of sudden failures. Although this prediction accuracy appears modest, it can significantly raise SATA system reliability, as will be shown in this section.

SATA system designers should consider reading the drive SMART attribute data directly and implementing the "Smarter SMART" algorithm [Hughes 2002] in storage server software to significantly raise the in-drive SMART warning accuracy to about 50%. This data is a single sector whose format and read command is defined in the ATA specification, and the new algorithm is simple to run. The SMART threshold triggers in today's drives has accuracy as low as 10%, and "Smarter SMART" has 40–60%.

SMART has to have a false alarm rate of around 0.2% per-year to avoid unduly increasing "no problem found" drives, low compared to drive failure rates (Table II). At a 1% drive failure rate, 0.2% is a 20% NPF increase. This small false alarm rate means that the burden of precautionary backup spares in a SATA system is modest.

For analysis of RAID redundancy reliability, let  $p_f$  be the annual probability of a drive failure in a SATA storage system,  $p_{fp}$  the probability that a SMART failure prediction will be made before a drive fails, and  $p_{nfp}$  the probability that a failure prediction is made, but the drive does not fail (within the next few months).  $P_{fp}$  values of 40–60% at 0.2%  $p_{nfp}$  are shown in experimental data from nearly 4000 drives testing the "Smarter SMART" algorithm [Hughes 2002; Murray and Hughes 2003]. Using SATA system SMART, the probability of an unpredicted failure is  $p_f(1-p_{fp})$ , so a 1% failure rate would drop to 0.5% with 50%  $p_{fp}$ .

The probability  $P_{fs2}$  of a parity-2 RAID reconstruction failure of three or more of  $n$  independent drives in a stripe is

$$\begin{aligned}
 P_{fs2} &= 1 - P(0 \text{ fail}) - P(1 \text{ fail}) - P(2 \text{ fail}) \\
 P_{fs2} &= 1 - (1 - p_f)^n - np_f(1 - p_f)^{n-1} - n(n-1)p_f^2(1 - p_f)^{n-2}/2 \\
 &= 1 - (1 - p_f)^{n-2}[1 + (n-2)p_f + (n-1)(n-2)p_f^2/2]. \quad (1)
 \end{aligned}$$



Table IV. RAID-5 Operating Failure Probability, and With 50% SMART

$P_f$ , n Drive Stripe	Parity 1, no Fail Predict	Parity 1, Fail Predict	Parity 2, no Fail Predict	Parity 2, Fail Predict
1%, 4	$5.9 \times 10^{-4}$	$1.5 \times 10^{-4}$	$4.0 \times 10^{-6}$	$5.0 \times 10^{-7}$
3%, 4	$5.2 \times 10^{-3}$	$1.3 \times 10^{-3}$	$1.1 \times 10^{-4}$	$1.3 \times 10^{-5}$
1%, 14	$8.4 \times 10^{-3}$	$2.2 \times 10^{-3}$	$3.4 \times 10^{-4}$	$4.4 \times 10^{-5}$
3%, 14	$6.4 \times 10^{-2}$	$1.8 \times 10^{-2}$	$7.7 \times 10^{-3}$	$1.1 \times 10^{-3}$

The failure probability  $P_{fs1}$  in single parity reconstruction is found by omitting the last term in the middle equation:  $P_{fs1} = 1 - (1 - p_f)^n - n p_f (1 - p_f)^{n-1}$ .

Table IV uses Equation (1) to compare single parity reconstruction fail  $P_{fs1}$ , and double parity fail  $P_{fs2}$ , in 4 and 14 drive stripes, with 1% and 3% drive annual failure probability  $p_f$ , and 50% failure prediction accuracy  $p_{fp}$ .

It can be seen that failure prediction can reduce RAID failure probability an order of magnitude, even with modest 50% prediction accuracy, particularly when low reconstruction failure probability with double parity is required.

The top row parity-1  $5.9 \times 10^{-4} = 0.0006 = 0.06\%$  will be used later as a typical RAID annual failure rate for comparisons.

## 7. RECONSTRUCTION FAILURE PROBABILITY

When a stripe drive fails, reconstruction must successfully read all the data in the remaining stripe drives without an unrecoverable error in any drive block. Large capacity drives bring a significant chance that this reconstruction may fail. A typical SATA drive specification for an unrecoverable read error is  $10^{-14}$  per-bit read. Reconstruction of a failed 250GB SATA drive in a 14 drive stripe can have up to a 26% chance of an unrecoverable read error ( $10^{-14} \cdot 8 \cdot 250 \cdot 10^9 \cdot 13 = 0.26$ ). Consequently, unrecoverable errors may occur in reading the entire contents of a RAID stripe of large capacity SATA drives, all meeting their specifications. A typical enterprise-class SCSI/FC specification is  $10^{-15}$  per-bit read, so reconstruction of a 36GB drive in a 14 drive stripe can have  $10^{-15} \cdot 8 \cdot 36 \cdot 10^9 \cdot 13 = 0.004$  read errors, or a 0.4% chance of reconstruction failure. Improving the SATA unrecoverable read error rate to the SCSI level still leaves a 3% chance of a reconstruction failure. The basic problem is that too big an independent physical unit of storage has failed, requiring full reads of similar size units. One small compensation is that RAID unrecoverable error rate issues can be tested in reasonable time scales, unlike reliability MTBF/AFR.

With parity-2 RAID, unrecoverable block errors can be corrected via the second cross parity stripe (Table III). This will succeed because only the few error blocks have to be read from the second parity stripe. However, the problem returns if two drives fail. A mirrored system could read the few unrecoverable blocks off the mirror array.

To eliminate this risk, in-drive parity stripes can be added—stripes in each drive which protect sequential series of blocks from unrecoverable errors. If a parity block is stored every  $3.6\text{GB}/512 \approx 7,000,000$  blocks, then the unrecoverable error probability per-in-drive stripe would be the same as a SCSI 36GB drive, even with today's SATA unrecoverable error rate ten-times higher than SCSI. The loss of RAID capacity is negligible, and there is negligible

reconstruction overhead since all drive blocks have to be read during a reconstruction. However, this does require the usual RAID parity quadruple access for each write access to an in-drive stripe (read old block contents, read parity block, compute new parity block, write new block and new parity block). This would be a performance limitation for random block I/O, but SATA arrays are best suited for serial access-like backup. This additional in-drive parity would leave Equation (1) valid. The unrecoverable error failure probability is then roughly the SCSI probability above, squared, or  $0.004^2 = 2 \times 10^{-5}$ . That is far smaller than the 0.06% nominal parity-1 RAID in Table IV. Of course in-drive stripes work with any interface drive.

In normal Parity-1 RAID, rereads of an unrecoverable error are unlikely to succeed because the drive will already have unsuccessfully tried many rereads using different drive internal states [Hughes 2002]. Assuming complete drive flaw-marking, unrecoverable errors are due to limitations of the drive error correction code (ECC) in data blocks having minimum readback signal-to-noise ratio. The latter is constrained by physics and drops inexorably as areal density rises. More powerful ECC codes could be designed into SATA drives, with more code interleaves, or a larger symbol bit size used, than the byte-symbols Reed-Solomon ECC now uses. ECC codes also exist which can find the *locations* of bit errors in a data record, even when they are unable to correct them. If this location information could be returned to the user application, it could determine if the desired data is not in the error locations, or if the error locations are in the record slack space, or in disposable data. Drives send no such information today, only the occurrence of an unrecoverable record read error [Hughes 2002]. Object storage drives might allow such features [Storage Networking Industry Association [www.snia.org](http://www.snia.org)].

## 8. D2D2T DRIVE BACKUP SYSTEMS

Magnetic tape has long been used to backup computer information, customarily during dedicated low-use overnight windows when tape backup programs can be run without interfering with user applications. Over time, these windows have shrunk to near nonexistence due to the globalization of business through the Internet and World Wide Web. While tape system speeds have accelerated and tape capacities have increased, they have not keep pace with the demand for shorter backup windows and the rapidly escalating volume of disk drive data being backed up. Nearline disk systems which backup daily to tape have their drives powered for a significant fraction of the time, so Table IV should apply for stripe failure probability. If backup occurs once a day and six-nines reliability per-stripe is required, then any configuration with probability better than  $365 \times 10^{-6} \sim 4 \times 10^{-4}$  is acceptable, and there are several in Table IV.

RAID stripes are vulnerable to drive failure during power-off periods, which can be significant for archival systems. This will be discussed next.

## 9. D2D DRIVE ARCHIVE SYSTEMS

Digital magnetic tape archival reliability is a highly mature technology, proven over a half-century of digital tape archival use. Longterm reliability must be proved, but takes a long time to prove, unlike performance specifications. Fifty-

year-old tape has been successfully read as has hundred-year-old analog voice recordings [Daniel et al. 1999]. No other digital storage technology has existed long enough to prove such archivability. Little justification exists for assuming that disk drive archival reliability is comparable to tape. Tape is designed for archival backup, and has historically validated data preservation standards which test for the failure modes found during tape's century of history, such as dropouts, tape wrap sticking, and data print through. Tape transport devices failures can be repaired and the tape data recovered, unlike disk drive failures.

Today, disk drive costs are dropping far below tape *drive* costs, and near or below tape *cartridge* costs. Consequently, archival backup D2D systems have been proposed, which eliminate conventional magnetic tape backup entirely. A D2D system can run at the full speed of disk, and can efficiently use the higher capacity of SATA disk drives (up to 400GB today), because backup is primarily serial data storage, not needing the high random access speed of the enterprise storage systems being backed up—which get high random access speed by running many smaller capacity drives in parallel.

Disk storage has not historically had such an archival role, and disks are traditionally designed for a five-year service life. Their design considerations include stored data lifetime, as well as electronic and mechanical reliability. Disk magnetic media is designed to retain data against thermal decay for five years (with 100% margin, i.e., ten-year nominal design data life). Thermal decay slowly turns stored magnetic bit states into magnetic noise, and is a serious issue today. The bit size is so small at 60–100GB per-disk that simple Boltzmann  $kT$  thermal energy at room temperature slowly disorders bit “0” vs. “1” magnetization states. It's a hard physics grounded limitation and a subject of major technology conferences [Moser 2002]. The thermal decay problem can be avoided by re-recording the data blocks, or by exchanging data with another drive. This operation can also detect latent block or drive failures. Drive operating firmware is also stored on disk, and archival drives should be designed to store their firmware at lower bit densities that resist thermal decay.

What failure modes can be expected in archival disk drives? Drive failure analysis is a mature engineering discipline with few mysteries, but many failure modes need to be considered (Table II). In contrast, tape archives are not vulnerable to tape transport failure modes. Only failure modes in the archived tapes and reels are fatal. Archival disks might be removed from storage systems and stored unpowered, perhaps in remote sites for disaster protection like tape archive vaults. Meeting drive manufacturer temperature and humidity specifications is important, particularly avoiding low temperature and high humidity, which can cause moisture condensation inside drives.

“Ramp load” drives park their recording heads off the disk surfaces, which may avoid long term head-disk stiction issues that contact stop-start drives can have (whose heads sit on disks when not spinning). Power-off archival storage stresses fewer spin motor bearing failure issues. Disc lube migration and contamination suggest that on-off power cycles not be excessive. Drives ingest air as they cool after a power cycle, along with any contamination or moisture in the air (their air filters remove submicron particles, but water molecules are only an Angstrom in size). Leaving drives spinning for long periods, with the

Table V. Nonoperating RAID Backup Stripe Failure Probabilities

$P_f$ , n Drive Stripe	Parity 1, no Fail Predict	Parity 1, Fail Predict	Parity 2, no Fail Predict	Parity 2, Fail Predict
0.1%, 4 drives	$6.0 \times 10^{-6}$	$1.5 \times 10^{-6}$	$4.0 \times 10^{-9}$	$5.0 \times 10^{-10}$
1.5%, 4 drives	$1.3 \times 10^{-3}$	$3.3 \times 10^{-4}$	$1.3 \times 10^{-5}$	$1.7 \times 10^{-6}$
0.1%, 14 drives	$9.0 \times 10^{-5}$	$2.3 \times 10^{-5}$	$3.6 \times 10^{-7}$	$4.5 \times 10^{-8}$
1.5%, 14 drives	$1.8 \times 10^{-2}$	$4.8 \times 10^{-3}$	$1.1 \times 10^{-3}$	$1.4 \times 10^{-4}$

heads left at one track position, risks failure due to “fly stiction.” Drive power management commands can park heads off the disks in these situations, or power-off drives completely.

Table V shows stripe failure probabilities in nonoperating archive disk systems, calculated using Equation (1) with a best-case nonoperating drive failure probability of 1/10 of 1%, and a worst-case probability of  $\frac{1}{2}$  of 3%, from Section 5. It appears that stripe failure rates approaching  $10^{-8}$  per-year might be attained, and even better if drives are powered and tested more frequently than once a year.

In disk drive archives, each drive should be periodically powered up in a MAID system, or mounted in a drive tester. At least once a year, drive data integrity self tests should be run, including SMART drive internal failure prediction tests, and all drive data exchanged with another drive to avoid the thermal decay problem. The testers should be able to unlock and relock drives, if the SATA password system is used for user security, and is able to mount a full RAID stripe for reconstructing the data from a failed drive. D2D has an advantage over tape here, in that a disk failure can be recovered by reading a specific few drives, while recovering from a tape failure could require reading an entire alternative backup set (assuming it exists). Drives scheduled to be removed from systems, due to apparent failure or upgrade, should have their user data securely erased (Section 10).

## 10. USER DATA SECURITY

Disk drive data security is controlled when drives are physically inside storage systems. Drives carry well-known file systems making user data potentially less secure than tape backup, where files directories are hidden and complex access security systems are standard. Striped RAID systems also have security vulnerabilities, such as transaction processing systems, where individual user records are often smaller than a single disk block.

Security in tape backup and archiving systems has been historically accomplished by maintaining physical control of the tape media and erasing it (magnetically degaussing) before relinquishing physical control.

Removed backup or archival drives raise the same security issues. For backup, a simple method would be to lock SATA drives using the existing SATA user/master password system, with the password securely controlled inside the storage system. Drives scheduled to be removed from systems due to apparent failure or upgrade should be automatically secure-erased before removal (a standard ATA command, which can take 20–60 minutes), or locked against access with a random 256-bit ATA password (in milliseconds). The ATA and SCSI specifications both contain a secure-erase command. We have tested the

secure-erase command in many drives from 0.2GB to 250GB, and all ATA drives designed since the command was added (10–15GB and greater) have it, while no SCSI drives have implemented this optional command. Detailed test data is available from the authors. The basic rule is not to allow user data to be readable from drives when they leave the physical protection of a storage system. Alternatively, the data can be encrypted before storage.

User data security should not depend on manual erasure policies. Fully a third of after-market used drives contain unerased user data [Garfinkel and Shelat 2003], and RAID striping does not necessarily avoid the security risk.

## 11. A PROPOSAL FOR A SATA STORAGE SYSTEM STUDY

A systems study of design factors for backup and archival disc systems would assist this new storage industry. A class of “Enterprise SATA” drives might emerge, still allowing the high SATA drive capacities at acceptable cost. Standardized SATA storage systems reliability specifications and reliability test levels could significantly help. Such a study could be undertaken by a SATA storage system trade or standards group.

### ACKNOWLEDGMENTS

The authors thank Dave Anderson, the reviewer, Roger Bohn, Walter Burkhard, and John Palmer for valuable discussions and comments.

### REFERENCES

- ANDERSON, D., DYKES, J., AND RIEDEL, E. 2003. More than an interface—SCSI vs. ATA. In *Proceedings of the 2nd Annual Conference on File and Storage Technology (FAST)* (March).
- DANIEL, E., MEE, C. D., AND CLARK, M. C. 1999. Magnetic recording, the first 100 years. IEEE Press, 20.
- COLARELLI, D., GRUNWALD, D., AND NEUFELD, M. 2002. The case for massive arrays of idle disks. *2002 Conference on File and Storage Technologies*.
- GARFINKEL, S. L. AND SHELAT, A. 2003. Remembrance of data past: a study of disk sanitization practices. *IEEE J. Secur. Privacy* (Jan.-Feb.) 17–25.
- HAMELRY, G. AND ELKAN, C. 2001. Bayesian approaches to failure prediction for disk drives. In *The 18th International Conference on Machine Learning*. 1–9.
- HUGHES, G. F. 2002. Improved disk drive failure warnings. *IEEE Trans. Reliab.* 51, (Sept.), 350–357.
- HUGHES, G. F. 2002. Wise Drives, *IEEE Spectrum* (Aug.).
- LEE, M.-Y. AND PARK, M.-S. 1996. Double parity sparing for performance improvement in disk arrays. In *Proceedings of International Conference on Parallel and Distributed Systems—ICPADS 1996*. IEEE, Los Alamitos, CA, 169–174.
- LUETH, C. 2004. NetApp data double parity RAID for enhanced data protection with RAID DP. Network Appliance Report TR3298 (Jan.).
- MOSER, A., TAKANO, K., MARGULIES, D., ALBRECHT, M., SONOBE, Y., IKEDA, Y., SUN, S., AND FULLERTON, E. 2002. Magnetic recording: advancing into the future. *J. Phys. D: Appl. Phys.* 35. R157–67.
- MURRAY, J. F. AND HUGHES, G. F. 2003. Hard drive failure prediction using non-parametric statistical methods. *International Conference on Artificial Neural Networks*. Istanbul.
- MURRAY, J. F., HUGHES, G. F., AND KREUTZ-DELGADO, K. 2004. Comparison of machine learning methods for predicting failures in hard drives. To appear *J. Mach. Learn. Res.*
- SCHWARZ, T. J. E. AND BURKHARD, W. A. 1995. Reliability and performance of RAIDs. *IEEE Trans. Mag.* 31, 2 (March), 1161–1166.
- Storage Networking Industry Association OSD Technical Work Group. [www.snia.org](http://www.snia.org).

Received August 2004; revised September 2004; accepted September 2001