



# Solving the Problem of Persistent Data Retention

Overcoming the limitations of both conventional  
RAID arrays and Tape Libraries in assuring affordable  
long-term data integrity

## Abstract

Most people assign a value to digital data based on their perception of its immediate importance rather than when they might need it at some future date. As a result, the vast majority of data is stored in a state of limited usefulness due to how and where it is housed, how it was initially identified or tagged, and how easily it can be retrieved.

For the purposes of this paper, we will distinguish between “transactional” data, which is data that is actively being created, changed, and accessed by an application, and “persistent” data, which has been created and is being retained in an unchanged state for potential future use. When an event happens whereby users have a critical need to locate and retrieve specific portions of this persistent data, they may be severely limited in how quickly and accurately they can access it. Yet, the value of the data may be even more important at that point than when it was originally created.

With the onset of government-mandated regulations surrounding the retention and retrieval of corporate data, it has become critical to know where this persistent data is, to be able to immediately retrieve precisely what is needed, and to be completely sure that the data is still intact even as it is locked away in the archive. The ability to do this can make the difference between winning or losing lawsuits, saving lives, and ensuring that years of irreplaceable work is protected.

Because corporations are increasingly aware of the role and value of persistent data in the enterprise, they are taking a more careful look at storage technologies and methodologies to ensure it is protected. This includes expanding the definitions of data availability and data integrity. These were formerly centered around “mission critical” online transactional data and were generally stored on the most expensive and reliable storage systems.

In this paper, we will discuss the techniques for ensuring that persistent data in the enterprise today will really be available when it is needed in the future.

## The Role of Persistent Data in the Enterprise

- In corporations today, the four basic roles and usage models for persistent data are:
- Backup/Recovery;
- Business Continuity/Disaster Recovery (DR);
- Active Archive; and
- Passive Archive, including Compliance and Deep Archive.

These roles and usage models have both common and conflicting requirements, as detailed in Table I:

| Requirements       | Backup / Recovery | Business Continuity / DR | Active Archive | Passive Archive  |
|--------------------|-------------------|--------------------------|----------------|------------------|
| Application Path   | Backup Software   | Must Be Restored         | Direct         | Must Be Restored |
| Access Frequency   | Seldom            | Rarely                   | Seldom         | Rarely           |
| Retention Period   | Days – Months     | Days – Months            | Years          | Decades          |
| Retrieval Time     | Minutes – Hours   | Hours                    | Seconds        | Days             |
| Retrieval Objects  | Files – Volumes   | Files – Volumes          | Files          | Files            |
| Data Integrity     | Very High         | Very High                | Very High      | High – Absolute  |
| Retrieval Accuracy | Last Good Copy    | Process Recoverable      | Very High      | Exact, Inclusive |
| Other Copies       | Multiple          | Few                      | Few            | Few – None       |
| Storage Cost       | Low               | Very Low                 | Low            | Low – Very Low   |

Table I: Four Basic Usage Models for Persistent Data

In the past, companies used data tape to store almost all of their persistent enterprise data, due to its low capital cost. While tape may serve well for some environments and helps for portability off site, it can cause problems with long retrieval times and difficulty with the data integrity requirements of all four usage models above. As a result, the industry has been moving to low-cost disk solutions. The problem is that conventional disk has limitations for long-term storage, in particular with the requirements for longer retention periods and low costs outlined in Table I.

This is especially true for compliance data, for which guaranteed integrity is not only a business necessity but a legal requirement. In that situation it is critical that companies be able to find the exact copies of specifically-requested data as well as to verify that it has not been changed from its original state.

Tape is still appropriate for storing most types of Passive Archive data, especially when it needs to be stored offsite (Deep Archive). This is because tape is inherently portable, tape is inexpensive, and access time is often not critical. However, enterprises are finding out the hard way that because tape is so portable, it can easily be misplaced, stolen, damaged or lost.

On an industrywide level, what has been desperately needed is a new persistent data storage platform that combines the best attributes of both disk and tape.

## Why MAID Is the Right Platform for Storing Persistent Data

While this strategy has limitations for online transactional data where data access requirements are measured in milliseconds, as part of a tiered architecture or an Active Archive of persistent data, properly implemented, MAID technology overcomes the biggest liabilities of both disk and tape-based systems.

A key aim with any long-term storage strategy, however, is to ensure that the data remains intact on whatever storage medium it is placed. When a book or file is put on the shelf for long-term archive, the environment is carefully maintained to ensure it is protected. Too much humidity and the pages can deteriorate and even rot. Magnetic disk or tape media has a similar problem and is subject to “bit rot” which is when the magnetic substrate deteriorates over time. This can lead to the gradual degradation of data integrity in any storage environment, even when it is just sitting there unused.

This is one of the great failings of tape as a long-term storage medium for the four usage models in Table I. That is, users don’t really know if the data is intact until they try to access it. For most companies, it is too labor intensive and expensive to constantly keep mounting the tapes and checking on them. Doing so would cause additional wear on the magnetic tape itself. For all companies, discovering that data is gone at the very moment when they are trying to access it is not an acceptable result.

Thus, when MAID technology was first commercially introduced in the COPAN Systems platform in 2004, engineers recognized that they had multiple challenges to overcome to make long-term, cost-effective digital data preservation absolutely secure for all of the major usage cases. Specifically, the system had to:

1. Quickly power on and off a selected set of disk drives in a manner that assures only the right set of drives are on when needed, and off when not needed – all within a power budget that satisfies multiple simultaneous requests to a large array of drives;
2. Make sure that the data is still intact when it is needed, even though that data had not been accessed in a long time;
3. Ensure that the drives are healthy when they are actually called on to perform;
4. Have enough intelligence and bandwidth within a large array of mostly idle disks to deliver the needed performance from the particular drives that are actually being accessed;
5. Pull all these requirements together into a robust system architecture that also provides high-level

---

functionality to meet the rest of the primary data preservation requirements in Table I, such as intelligent identification and search capability for data stored on disks that are usually off; and

6. Deliver all of this at the lowest possible cost in an extremely dense, enterprise-class package to be competitive with existing tape solutions that is easy to service and maintain, while capable of providing decades of operational life.

To meet these criteria, COPAN Systems built its MAID platform using three key technologies:

1. Three-tiered architecture, with connectivity and computing intelligence at each tier to dynamically distribute the bandwidth and processing capabilities where and when needed;
2. Power-Managed RAID® software for full RAID 5 data integrity protection and fault tolerance; and
3. Disk Aerobics® software which ensures healthy drives and solid data protection while proactively moving data off faulty or dying drives.

The three-tiered architecture allows Power-Managed RAID, I/O operations, and application processing to all be done in parallel using separate processing resources. As a result, the new SGI® COPAN™ 400 MAID system can deliver wire-speed read/write performance across eight 8 Gbps Fibre Channel connections (6,400 MB/second sustained) with any combination of up to 25% of 896 SATA disk drives operational at one time. With the current 2 TB drives, this provides 1,796 TB of persistent data storage in a single COPAN 400 cabinet. In addition, the COPAN 400 features a second power option, enabling up to 50% of the disk to be spun up at a time for a given shelf. This gives IT managers the option of tuning the balance of performance and power that best suits the requirements of their users.

In day-to-day operations, idle drives are powered on in four-drive RAID 5 sets within seconds of receiving an access request. They are physically arranged so that there is no vibrational interaction or local heating issues. Once on, they respond to all other I/O requests at normal millisecond disk speeds just like conventional RAID arrays. All of the drives are hard-mounted with vibrational damping in easily accessible pull-out canisters for scheduled replacement. They are also actively monitored by individual thermal sensors and kept below 45°C, which assures that they always operate well within the drive's maximum reliability range. SGI purposely kept the number of drives in the RAID sets small in order to reduce the probability of a drive member failing. This also minimizes the time needed to rebuild the RAID set if one drive should suddenly fail.

A key high-availability benefit of the COPAN MAID three-tier architecture is that it offers multiple replication options. Not only is export to tape available for creating a physical tape cartridge from a virtual tape cartridge, but also replication between COPAN 400 VTLs can be done over a WAN to avoid the cost and security risks inherent in the physical transport of tapes. These options offer very fast disaster recovery capability, as well as the ability to consolidate backups from geographically-dispersed sites into a centralized backup facility.

The third key technical innovation, the patent-pending Disk Aerobics software, continuously monitors all operating drives within the system and periodically turns on and monitors each idle drive to ensure that the entire system is healthy and operating within design limits. If Disk Aerobics detects that a particular drive is showing signs of possible failure, it proactively replaces that drive with a spare and retires it from the system. As a result, Disk Aerobics substantially reduces the probability of a RAID set rebuild due to a failed drive, and virtually eliminates the likelihood that a second drive in the RAID set would fail before the RAID set could be restored.

Thus, the combination of these architectural principles in the COPAN 400 overcomes the key issues of using low-cost disk drives to store persistent data for the four key categories of persistent data preservation.

---

## Implementing Disk Aerobics Software

Disk Aerobics is a patent-pending technology which enables the system to:

1. Actively monitor disk health and environmental data;
2. Periodically exercise idle drives to ensure good health;
3. Proactively replace degrading or end-of-life drives; and
4. Retire replaced drives and maintain spares.

The active monitoring of disk health and environmental data involves analyzing a selected subset of the internally-generated SMART (Self-Monitoring, Analysis, and Reporting Technology) data from the drives, as well as checking temperature readings for each drive, canister temperatures, fan temperatures, and voltages. This is all kept in a database, and the system employs heuristic techniques to detect if a drive is becoming marginal or if environmental conditions are moving out of limits.

If any of these parameters falls outside specified limits or if a drive fails to return a SMART status, the data on that drive is immediately backed up to a spare and the drive is proactively replaced in its RAID set and then retired. Because the backup is done directly between the suspect drive and the spare drive without affecting the other members of the RAID set, this is accomplished far more quickly than the time it takes to rebuild an entire RAID 5 set following a failure.

In addition, the SGI COPAN system uses a number of proprietary techniques that anticipate potential data path problems or drive failures before they happen by monitoring I/O timeouts, analyzing CRC errors, and reconstructing bad blocks. Along with SMART data analysis and proactive replacement, the net result is that the vulnerability window for a second drive failure is reduced to milliseconds from what would normally be many hours as is typically needed to rebuild a RAID set if it were allowed to simply fail. Thus, MAID technology, together with Disk Aerobics, anticipates potential failures instead of merely reacting to failures as always-on RAID systems must do.

Disk Aerobics also periodically exercises the idle drives to ensure that they are healthy. To do this, the system sets aside a fixed amount of the power management budget for idle drives so they can each be turned on at least once a month and run through the same set of diagnostics and proactive replacement techniques as the active drives. All of this is accomplished as background tasks with no interruption of normal operational performance.

So how long can disk drives be left powered off before head-disk stiction, corrosion, or some other non-operating problem begins to cause an additional reliability risk? Although industry answers to this question vary, the emerging consensus is that spinning drives up monthly is adequate to keep them healthy. Even with this assurance, the COPAN 400 monitors and analyzes each drive's start up time, the amount of electricity its motor draws and other key indicators as another proactive measure to forestall failure. If any of these parameters begins to increase beyond acceptable limits, the drive is immediately tagged as a candidate for proactive replacement and then retired by the normal procedure. The COPAN 400 also proactively replaces and retires any drive that has exceeded a fixed percentage of the specified power-on hours or start / stops for that drive model, even if all the other operational parameters still look good. When it comes to data integrity, there is zero margin for error.

To handle the proactive replacement needs of any SGI COPAN MAID product, the system uses 40 spare drives (5 per shelf), which should result in scheduled service or maintenance action to replace retired drives with fresh spares no more often than every 12 to 18 months, even with 896 drives in a fully-loaded COPAN 400 MAID system. Allocating spares and keeping track of the number remaining is also a key function of the Disk Aerobics software. If the number of spares has been reduced to less than a pre-determined amount, the COPAN 400 will "phone home" and request an immediate service call to replenish the full inventory of spares.

## Reliability of SATA Drives

A comprehensive study on SATA disk drive reliability was published in the ACM Transactions on Storage. Gordon Hughes and John Murray of the University of California at San Diego looked at the failure modes of 4,000 SATA disk drives, and showed that there was an Annual Failure Rate (AFR) of ~ 2.1% (21 per 1,000 drives). This implies an MTBF (Mean Time Between Failures) of ~400k hours on a 24 x 7 x 365 basis. A typical Fibre Channel or SCSI drive, by comparison, has an AFR of ~0.3%, or an MTBF of ~3M hours.

Table II below shows the breakdown of the failure modes for the SATA drives found in this study (approximately 10% of the failure modes could not be identified).

| Failure Mode                 | Description  | Frequency    | Stress Condition          | AFR per 1000 drives |
|------------------------------|--|--------------|---------------------------|---------------------|
| Head-Disk Interference (HDI) | Head Touch-Down or Crash                           | 15.5%        | Operating                 | 3.3                 |
| No Problem Found             | Drive Returned, but Tests OK                       | 15.0%        | N/A                       | 3.2                 |
| Recording Heads              | Failure of Complex Nano-Technology Devices         | 14.5%        | Operating                 | 3.0                 |
| Post Manufacture             | Drive Handling Damage                              | 10.1%        | N/A                       | 2.1                 |
| Circuit Board (PCB)          | IC Component or Board Failure                      | 8.5%         | Operating                 | 1.8                 |
| Head or Disk Corrosion       | Causes HDI or Disk Defects                         | 7.7%         | Non-Operating             | 1.6                 |
| Head Assembly (E-Block)      | Wires, Preamp, or Coil Fail                        | 6.8%         | Operating                 | 1.4                 |
| Head-Disk Assembly           | Mechanics, Electric, or Voice Coil Fail            | 3.9%         | Operating                 | 0.8                 |
| Disk Defects                 | Causes HDI or Read Errors                          | 2.6%         | Operating                 | 0.5                 |
| Drive Hardware               | Internal Operating System                          | 1.9%         | Operating                 | 0.4                 |
| Head-Disk Stiction           | Disk Won't Spin Up Due to Head Adhesion to Surface | 1.3%         | Non-Operating             | 0.2                 |
| Spindle Bearing              | Disk Spin Bearing Fails                            | 1.1%         | Operating                 | 0.2                 |
| Contamination Inside Drive   | Foreign Materials or Gases Cause Failure           | 0.7%         | Operating / Non-Operating | 0.1                 |
|                              | <b>Total</b>                                       | <b>89.6%</b> |                           | <b>18.8</b>         |

Table II: Failure Analysis of SATA Drives.

Hughes and Murray noted that little is known about non-operating failure rates for idle drives since that has historically not been tested. They expect this number to be between 1/10th and 1/2 of operating failure rates. If we eliminate the “No Trouble Found” and “Post Manufacture (Handling Damage)” categories from Table II, and assume that the AFR for operating failure modes would increase proportionately with power-on time and the AFR for non-operating failure modes would decrease proportionately with power-on time, then the projected AFRs per 1,000 drives, versus power-on duty cycle, would be as shown in Figure 1.

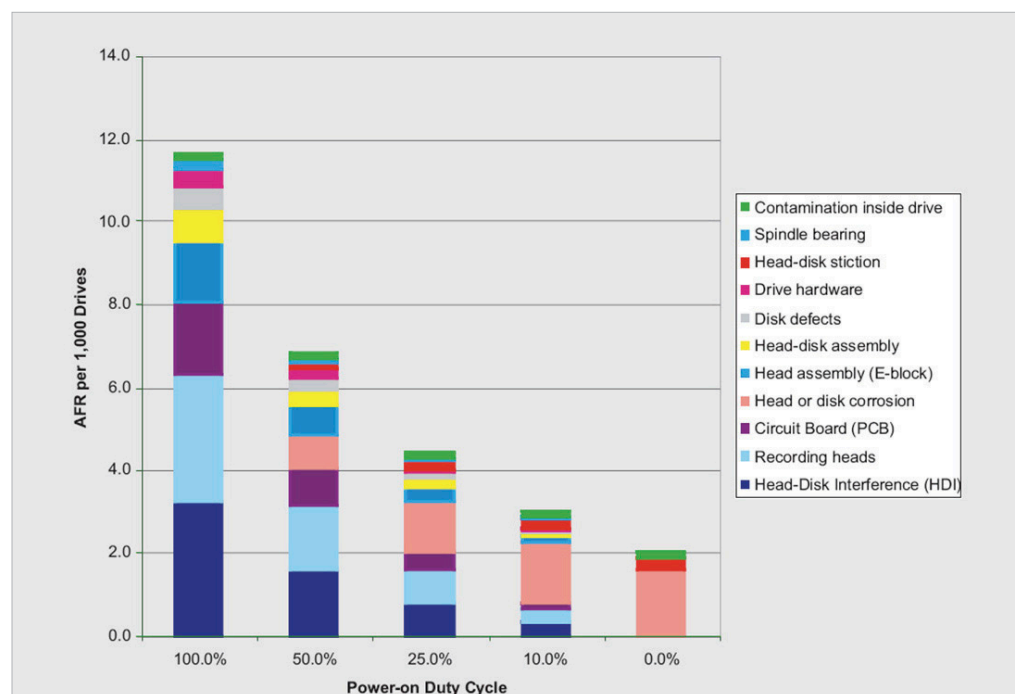


Figure 1: Projected AFR for SATA Drives Versus Power-On Duty Cycle

This would imply that the AFR at a 25% power-on duty cycle should be about 40% of that at 100% power-on duty cycle, or almost equivalent to a Fibre Channel / SCSI drive. Because of Disk Aerobics, the results for the COPAN MAID platform are actually significantly better, as shown below.

## COPAN MAID Dramatically Improves SATA Reliability

We have been tracking the reliability of the MAID products in the field since the first COPAN system was shipped in 2004. The actual AFR of the over 10,000 SATA drives installed in customer environments was at 0.42%. This is equivalent to an MTBF of over 3 million hours. Thus, the measured failures for the drives in the COPAN MAID systems was less than 25% of the failures for conventional SATA drives in the Hughes study; and the COPAN study was tested on a sample size over twice as large. As a result, the projected service life of the drives should improve by at least five times over typical always-spinning drive storage systems. The accumulated field test time for these drives is now over 70 million hours. In addition, there have been no incidents of lost or unavailable data. This is truly enterprise-class reliability and data availability.

The combination of Disk Aerobics and MAID is unquestionably the best solution for storing persistent data. The combination of enterprise-class reliability and data availability, along with superior disk drive performance at a cost comparable to tape, makes the COPAN 400 ideal for all four of the persistent enterprise data usage models.

#### About SGI

SGI is a global leader in large-scale clustered computing, high performance storage, HPC and data center enablement and services. SGI is focused on helping customers solve their most demanding business and technology challenges. Visit [www.sgi.com](http://www.sgi.com) for more information.

**Corporate Office**  
46600 Landing Parkway  
Fremont, CA 94538  
tel 510.933.8300  
fax 408.321.0293  
[www.sgi.com](http://www.sgi.com)

North America +1 800.800.7441  
Latin America +55 11.5185.2860  
Europe +44 118.912.7500  
Asia Pacific +61 2.9448.1463