



Data Protection

Adapting to the Sea Change

2005 Special Edition

Mesabi Group



David G. Hill

The Mesabi Group particularly focuses on two “revolutions” that are dramatically changing the ability of storage and information managers to zero in on key data in near-real-time:

- The separation of storage (e.g., via storage area networks) as an equal partner with processing and networking in delivering business value
- Information Lifecycle Management (ILM), the proactive, global management of data from arrival to archiving

The Mesabi Group focuses on how these revolutions enable the storage infrastructure to serve as a vital component of the IT infrastructure not only to support essential business functions, but also — despite the critics — as an enabler in competitive advantage.

The Mesabi Group translates storage technologies into business value; for example, showing how data protection technologies can help business management with the fundamental business function of risk management.

Mesabi Group is managed by **David G. Hill**.

About This Special Edition

This special edition contains the whole independent original research report with the exception of Chapter 9 Representative Companies for which Chapter 9 Tape: An Ongoing Bulwark for Data Protection has been substituted. The report therefore ends with a good illustration of how many of the key principles can actually be applied and accurately reflects the opinions of the Mesabi Group regarding the role of tape in general and LTO technology in particular.

Table of Contents

Executive Summary	1
Key Report Findings	1
Preface.....	5
Business Continuity and the Backup/Restore Process — Never the Twain Shall Meet?	5
The Sea Change in Data Protection	6
What This Report Is — and Is Not — About	7
<i>Chapter One: The Time Has Come for Change</i>	9
What Data Protection Is.....	9
Data Protection Has To Be Placed in the Right Framework	10
Ride the Sea Change in Data Protection.....	11
<i>Chapter Two: Business Continuity: The Framework for Data Protection</i>	12
Business Continuity and Data Protection	13
Business Continuity Is Not Just Disaster Recovery	13
Disaster Recovery: Let’s Get Physical.....	16
Operational Recovery: Think Logically	17
Disaster Recovery Requires Judgment; Operational Recovery Requires Automation	18
Logical Data Protection Gets Short Shift in Business Continuity	19
Logical Problems Feature Prominently in Data Loss or Downtime	20
Logical Data Protection Problems Manifest in a Number of Ways	21
Do Not Neglect Any Facet of Data Protection.....	22
<i>Chapter Three: Data Protection — Where the Problems Lie</i>	23
Data Protection as It Was in the Beginning	23
Typical Data Protection Technology Today Still Leaves a Lot to Be Desired	25
Operational Continuity/Physical: Generally Strong, But Some Improvement Needed	27
Operational Continuity/Logical: More Attention Needs to Be Paid to Logical Data Protection	27
Disaster Continuity/Physical: Well Done, But Cost and Distance Are Issues	28
Disaster Continuity/Logical: The Danger of Being Under Protected May Be Very Real	29
Summing Up Data Protection Challenges by Category	29

Table of Contents

<i>Chapter Four: Data Protection — Setting the Right Objectives</i>	31
How High Is High Enough for Data Availability?	31
SNIA's Data Value Classification: A Point of Departure	32
Do Not Equate Availability with Value	34
Availability Objectives for Operational Recovery and Disaster Recovery Are Not Necessarily the Same.....	35
Availability Is Not the Only Data Protection Objective	36
All Data Protection Objectives Have to Be Met	38
<i>Chapter Five: Data Protection — Getting the Right Degree</i>	39
General Use Classes of Data	39
Tape Is a Special Case	40
Understanding Degrees of Data Protection	40
The Third Degree — Levels of Exposure	41
Mapping Degrees of Protection.....	41
<i>Chapter Six: Information Lifecycle Management Changes the Data Protection Technology Mix</i>	44
Why Data Lifecycle Management Is Not Enough — The Need for Metadata and Management	45
ILM Is Deep Into Logical Pools of Storage	45
Logical Storage Pools at a High Level	46
Moving Information Across Pools — A Distillation Process	47
Archiving Through a New Lens.....	47
Archiving: The Makeover	48
Protecting Archived Data	49
Data Retention.....	50
Disposition of Data	50
Compliance	51
Creating Data Archive Storage Pools by Data Retention Attributes	52
Security and Privacy	54
Active Archiving and Deep Archiving	55
Active Archiving Requires Active Archive Management	55
Long-term Archiving as Part of an Active Archive	57
ILM Changes the Data Protection Technology Mix.....	57
<i>Chapter Seven: Where Data Protection Technologies Play in the New Model</i> .	59
Back to Basics — Extending the Current Model.....	61

Table of Contents

Current RAID Capabilities Are Not Enough.....	61
Evolving Backup/Restore Software.....	61
Better Data Protection through Better Management Reporting and Automation	63
Moving Data Manually and Electronically — the Place of Vaulting	65
At Your Service — the Role of Service Suppliers.....	65
Disk and Tape — Complementing and Competing with One Another	66
Disk-based Backup	66
Virtual Tape	69
Virtual Tape Library	69
MAID	70
Removable Disk Drives and Disk Media	70
Data Protection Appliances.....	71
Tape Automation	72
Getting to the Point.....	74
Point-in-Time Copy	74
Continuous Data Protection	75
Replication Strategies.....	76
Mirroring.....	77
Dated Replication — Pay Close Attention.....	80
Special Requirements for Compliance	81
The Use of WORM Technology.....	82
WORM Tape.....	83
WORM Disk	83
Electronic Locking.....	84
Guaranteeing the Authenticity of Data	85
Privacy and Confidentiality.....	85
Compliance Appliance	85
Mapping the Base Data Protection Technologies to the ILM-Version of the Data Protection Framework	85
<i>Chapter Eight: Summing Up — Redesigning Data Protection.....</i>	<i>89</i>
Data Protection Is Everyone’s Business.....	89
Synthesizing the Data Protection Frameworks.....	89
Guidelines for Data Protection.....	91
The Challenge Ahead and a Call to Action	92
<i>Chapter Nine: Tape: An Ongoing Bulwark for Data Protection.....</i>	<i>93</i>

Table of Contents

The Challenge to Tape	93
Setting the Case for Tape	93
Why High Availability All the Time Is Unrealistic	94
Reconstruction, Not Just Restoration	94
Process and Technology Diversification	95
I/O Isolation — an Extra Measure of Safety	95
Disk-based Backup Is Not a Panacea	96
Removable Disk Is Emulated Tape.....	96
Doing Tape Well	96
Doing Tape Right.....	96
LTO Ultrium 3 — Popeye after Eating His Spinach.....	97
LTO Ultrium — Working to Make Tape Even More Reliable.....	98
LTO Ultrium Is Up to Speed	98
WORM — Getting an Extra Use from Tape.....	98
LTO Ultrium — For All Generations.....	99
Conclusions	99
Mesabi Group Conclusions	100
Glossary.....	101
Author Profile	105

Figures

Figure 2-1: Overview of Business Continuity	14
Figure 2-2: Business Continuity Is More than Data Protection.....	14
Figure 2-3: Business Continuity Keeps Your Business Running	16
Figure 2-4: Causes of Data Loss or Downtime	20
Table 2-1: Logical Data Protection Problems and Sources	21
Table 2-2: Data Protection Category Matrix	22
Figure 3-1: Data Protection: The Way It Was.....	24
Figure 3-2: Typical Data Protection Today	26
Table 3-1: Data Protection Challenges by Category	30
Figure 4-1: High Availability Depends upon the Entire IT Infrastructure.....	32
Table 4-1: SNIA Data Value Classification.....	33
Table 4-2: Operational Recovery and Disaster Recovery Differences	35
Table 4-3: Consequences of Data Loss	38
Table 4-4: Summing Up Key Data Protection Objectives	38
Table 5-1: Sample Degrees of Data Protection for Application <i>n</i>	42
Figure 6-1: The Storage Pyramid — Tiering and Pooling.....	46
Figure 6-2: ILM Changes the Logical Topology Storage Look	48
Figure 6-3: Data Retention Archive Pools	53
Table 6-1: Adding In Archiving to the Data Protection Category Matrix.....	58
Table 7-1: Where Data Protection Technologies Fit in the Data Protection Framework	60
Table 7-2: Base Data Protection Technologies for Active Changeable Data.....	87
Table 7-3: Base Data Protection Technologies for Archived Data	88
Table 8-1: Data Protection Requirements for Application <i>n</i>	90

Executive Summary

Key Report Findings

1. Enterprises may be *over-investing in some areas* of data protection, *while exposing their IT assets to unacceptable risk by under-investing in other areas* of data protection.

For example, an enterprise may not understand the importance of logical data protection. An Ontrack study showed that nearly 40% of the causes of data loss or downtime are logical, not physical, problems. Yet enterprises may not have in place a high availability (defined as seconds or minutes of annual downtime) logical recovery approach for critical applications. (The tendency is to think in terms of physical solutions, such as mirroring, which are not the answer to logical data protection problems.)

2. The report shows enterprises how to think about where they need data protection as well the degrees of data protection that are required to meet those needs. They can then better determine whether or not they are over-investing or under-investing to be able to meet particular data protection needs. The old bromide “one size does not fit all” applies to how enterprises fulfill their data protection requirements, but not for the basic principles of data protection.

Some large enterprises can afford to have a triad of data centers to ensure a high level of availability in the case of a disaster, whereas other enterprises simply use tape vaulting for disaster recovery, trading lengthy application restores for lower cost. However, all enterprises have to take into account both the need for disaster recovery to a remote site and operational recovery for problems that can be corrected at a local site.

Likewise, all enterprises need to take into account physical problems, such as disk failures, and logical problems, such as database corruption or a computer virus, for both disaster and operational recovery situations. The choice of individual data protection technologies is up to the enterprise — but overall data protection should fit within a common framework that applies to all enterprises.

3. Enterprises want “high availability” as part of data protection, yet virtually all use a “low availability” tape solution as part of their data protection strategy. For true data protection, enterprises should use multiple levels of availability in their overall strategies.

In an effort to ensure high availability for critical applications, many enterprises invest in additional, expensive disk storage arrays for an increased degree of physical availability. At the same time, they invest in tape automation solutions that add more levels of data protec-

tion, but that only deliver low availability (defined as hours or days to restore a particular pool of data).

In fact, enterprises can segment applications into ones that require high availability and ones that can function with low availability. Those applications, with their accompanying data that can get by with a tape automation solution alone, allow users to avoid the extra investment costs for additional disk storage. Note also that the most critical applications are (and should be) protected by both additional disk and tape automation.

Thus, enterprises want and need multiple degrees of data protection. RAID on a production data array provides one degree of physical protection (as the failure of one disk drive can be tolerated without loss of data). A remote mirror can provide a second degree of protection. Where information cannot be lost, a tape solution provides a minimum of one (and generally more, through multiple-generation tape copies) additional degree of protection. Disk does not provide logical data protection; tape does (since the tape is outside the I/O “write” stream that can make logical changes to data). Point-in-time copy capability and its derivatives can provide logical data protection on disk, but require understanding, planning, and investment that many IT organizations have yet to make.

4. Enterprises should implement an overall data protection strategy based on a data protection “framework.” Enterprises are *attacking business continuity, backup and recovery process, and compliance as if they are unrelated problems*, but they really all relate to one another in the context of data protection. The data protection “framework” allows the correct allocation of investment and resources to these three areas, as well as other data protection investments.

Business continuity is a key risk management function of any enterprise. Business continuity is about preventing or ameliorating disruptive impacts on the business that range from threats to survival to productivity drains. One of these disruptive impacts is data loss, and data protection avoids data loss.

Temporary loss of data requires that the data be restored before further disruption can occur. Using *backup and recovery* software is one way of restoring the data.

Compliance data is a special case of data protection to prevent disruptions that would be associated with non-compliance.

Relating all three areas — business continuity, the backup and recovery process, and compliance — through their relationships with data protection is just part of understanding the overall framework of data protection. Understanding that framework is important so that an enterprise can put in place a data protection strategy that takes these three aspects — as well as many others — into the proper con-

text so that the proper degrees of data protection and the proper levels of investment are in place.

5. Fixed content stored in active archives has different data protection and data retention requirements than active, frequently-changed data. By implementing Information Lifecycle Management (ILM) and coordinating it with a data protection strategy, enterprises can improve the cost-effectiveness, availability and performance of their storage.

As information in the form of files or records ages, it tends to become fixed data that is unchanging data. That age varies from the time of creation (e.g., a check entered into the system) to a later time (e.g., closing a transaction in an online transaction processing system). When fixed content data is “distilled” from its active changeable counterparts in an application to an “active archive,” the implications for data protection policies and management are significant.

The traditional backup process is not necessary for fixed content data. A piece of fixed content needs to be replicated after it is captured in an active archive, but no traditional backup process is necessary. Copying the data to a full backup on a regular basis is an unnecessary use of resources since the correct number of data protection copies is already available.

The second major change is the ability to put in place strong data retention policies. Although data retention policies can be applied to a pool of storage where active changeable data is commingled with fixed content data, data retention management is most effective with a fixed content pool of storage. That is because data retention applies only to fixed content data. An open transaction cannot be disposed of and cannot be considered (at that stage of its lifecycle) to be compliant data, since all compliant data has to be unchangeable.

The migration of data to an active archive will eventually have a significant impact on the active changeable side of the house as well. Although an enterprise may find it difficult to identify and separate its fixed content data and move it to an active archive, the active changeable side can also benefit when enterprises can find a way to migrate some fixed content data to an active archive. There will be less data to back up (and restore if necessary), so the burden on the overloaded backup/restore process will be reduced. If critical applications need to be remotely mirrored, the disk space for the remote mirror will be reduced. The upper boundary for fixed content could be as high as 80% or more, but even a movement of 20 to 30% of data could very well have a significant payoff.

6. Enterprises should consider their compliance policies in the context of data protection. Compliance is related to data retention, which is part of data protection.

Compliance data is fixed content information in an active archive. Data retention policies can be applied to this active archive. Com-

pliance is simply a more restrictive set of data retention policies, such as chain-of-custody requirements and privacy constraints.

7. Focusing on high availability and neglecting the other key objectives of data protection is dangerous.

Too often high availability and data protection are considered synonymous. Data availability is only one of four key objectives for data protection — data preservation, data responsiveness, and data confidentiality are the others. An overemphasis on high availability could lead to underweighting the other objectives. If the necessary amount of data preservation is not in place, high availability of an application will not matter. If the correct controls for data confidentiality are not in place, serious consequences could result. If data responsiveness is not in place, data will not be usable. A sense that all the objectives have to be balanced properly is necessary.

8. Information Lifecycle Management (ILM) actually will play an important role in the IT infrastructure — and data protection is a key part of that role.

Active archives of fixed content require different data protection strategies than for active changeable data. For example, the fixed-content data will use replication upon capture of the data for data protection purposes, and the active changeable data will use a backup/restore process. Moreover, implementation of active archiving requires migration of data between an active changeable pool of storage and an active archive pool of storage. ILM supports both of these tasks, among others.

Preface

It's well-known that data protection is a business necessity — yet few agree on exactly what data protection is. And failure to appreciate the full dimensions of the data protection challenge can lead to poor data protection management and costly resource allocation issues. The following example shows some of the difficulties that can arise when enterprises do not have a clear data protection strategy.

Business Continuity and the Backup/Restore Process — Never the Twain Shall Meet?

When asked what words most readily come to mind for “data protection,” the terms “backup/restore” and “business continuity” are likely to top the list. Enterprises clearly understand that all three relate to risk management and that risk management is an essential business task. Very few enterprises, however, understand that improving backup/restore may not improve business continuity. In fact, failure to understand the relationship between the ongoing down-in-the-details task of backup/restore and the global strategy of business continuity may result in unnecessary exposure to risk, under- or over-spending on data protection funding, and wasting of scarce IT administrator resources.

Let's start with **business continuity**. Business continuity attempts to prevent any major disruptions to business processes. Thus, business recovery is clearly different from disaster recovery — a concept with which it is often confused. *Disaster recovery* focuses on minimizing the effects of disaster, while *business continuity* focuses both on avoiding unplanned outages (due to either a disaster or an operational problem) in the first place and on minimizing the effects of unplanned outages. Specifically, business continuity emphasizes high availability — defined as restoration of access to applications within seconds or minutes — and resiliency — the ability of applications to continue running despite outages in systems, storage, or underlying software.

Now let's consider **backup/restore**. While backup is performed routinely, restore is only performed when systems are down as a result of an unplanned outage. Inevitably, the focus of backup/restore, like disaster recovery but unlike business continuity, is to minimize the effects of unplanned outages.

Now consider the practical effects of an over-focus on backup/restore rather than business continuity. Any low-to-high availability continuum clearly shows that the backup/restore process with tape is low availability (where low-availability is defined as restoration of data access to applications within hours or days), while technologies such as remote mirroring are high availability.

The continued high investment in low availability backup/restore process in conjunction with tape automation solutions is clearly inconsistent with the desire to move to the higher availability side of the continuum. Moreover, hours of downtime while a restore is taking place can cost customers and threaten the existence of a company. Take, for example, a recent outage of a European discount retailer: had it run longer than two hours, it could have resulted in the loss of millions of euros—a “business-critical situation” (*Progress Fathom: Business Continuity Down to the Details*, June 2005, www.valleyviewventures.com). Yet IT organizations are not likely to replace their current backup/restore processes anytime soon.

The way to avoid the costs and risks of an over-focus on backup/restore is to better understand an enterprise’s overall requirements for data protection. Even though a rip and replace strategy is typically unthinkable, enterprises need to be aware how much and where to place their bets on the data protection roulette wheel today — and those bets will definitely change tomorrow.

The Sea Change in Data Protection

In the last three years, the technology landscape of data protection has fundamentally changed — a true “sea change.” Disk-based backup, compliance technologies, and information lifecycle management (ILM) are examples of the technologies that are affecting how data protection bets should be placed and by how much. The net result is a sea change, a marked transformation.

These new technologies typically reflect new business processes as well. *Disk-based backup* reflects an increasing appreciation of the importance of a business continuity process. *Compliance technologies* reflect the increased importance of meeting regulatory requirements such as Sarbanes-Oxley. *ILM technologies* reflect a new process that is enabling finer-grained, more cost-efficient control over an enterprise’s data.

The sea change results in a number of questions for which IT organizations must have answers — about their current data protection infrastructure, and about the direction in which that infrastructure needs to evolve. Among these questions are:

1. What is the right target and what are the right objectives for a comprehensive data protection strategy?
2. How are data protection infrastructure holes identified and —if any exist — how are they filled?
3. How are low availability and high availability data protection technologies layered within an overall data protection framework to give sufficient degrees of data protection?
4. How will ILM lead to changes in data protection technologies and strategies?

5. How do all the existing and emerging technologies of the data protection puzzle fit together to help build a roadmap for evolving the data protection infrastructure?

What This Report Is — and Is Not — About

The purpose of this report is to serve as a guide for IT organizations so they are able to more clearly answer these and related questions. This report re-examines the basic principles of data protection in light of all the new demands that are being placed upon the IT infrastructure, and it also looks at how both maturing and emerging data protection hardware and software technologies affect those changes. The framework that arises from these basic principles helps put data protection in context to the overall IT infrastructure and helps IT organizations clarify the choices and options that are available to them for data protection.

However, this report is not a buyer's guide — that would require a never-ending encyclopedia! Although representative companies that offer data protection technologies are listed, the suitability — i.e., applying the criteria of scalability, interoperability, resource use, cost, maturity, vendor acceptability, etc. — of each of their products separately and in concert is dependant on the situation and therefore is unique to each reader. What the report does do, is to identify and examine the key decisions that should be made and strategies that should be implemented before evaluation of products and services can begin.

Moreover, the report is not a deep dive into the various data protection technologies. It examines current and emerging technologies in relation to an overall framework or approach to data protection. Readers can then better understand how to fit technology options into their overall data protection schema.

Where possible, conformity to the terminology used and directions charted by the Storage Networking Industry Association (SNIA) have been used so that users do not have to learn new concepts. However, since SNIA's work and perspectives on data protection are still evolving, there are situations in which this report diverges from the current direction that SNIA is taking.

This report is not the final word about data protection, but rather is intended to arm readers with information that enables them to act more effectively to achieve data protection.

Here is a brief exercise for the reader: before reading further, prepare a short list of questions:

- What is your view of data protection?
- What are you doing now for data protection?
- What are the issues you currently face regarding data protection?

- What actions, if any, are you planning to improve your data protection processes and infrastructure?

After reading the report, compare your answers to these questions before and after. This report will make what is potentially unclear about data protection now, as you start to read the report, obvious after you have finished it.

Chapter One:

The Time Has Come for Change

Studies reveal that data protection — in one form or another — is at the top or near the top of any list of issues facing the management of storage. In the short term, this importance is due to immediate concerns such as “how do I meet regulatory requirements right now.” In the long term, data protection aims to protect the information without which the business cannot function, and which is now a primary source of many enterprises’ competitive advantage. Data protection is therefore a cornerstone of any organization’s management of risk, and risk management is now recognized as one of the fundamental tasks of any enterprise.

Today, data protection is associated primarily with a wide spectrum of IT and business issues:

- Backup and restoration
- Disaster recovery
- Business continuity
- High availability
- Data asset preservation
- Compliance
- Data privacy
- Data security

Yet today’s IT organizations still tend to focus simply on improving backup/restore processes.

What Data Protection Is

Data protection is the mitigation of the risk of loss of or damage to an enterprise’s data on either a temporary or permanent basis.

Data protection is insurance. Therefore, the aim of data protection is not to maximize profits or revenues, or minimize costs, but to minimize worst-case losses. Like regular insurance, data protection insurance is a necessary cost of the prudent business, and balances the costs of unplanned outages against the costs of the insurance policy. A side-effect of data protection may be more cost-effective use of information assets; but users should not require profits from their data protection solutions, any more than from their life insurance policies on key executives.

Unlike the traditional insurance markets, the data protection market offers no “third-party” insurers (with the possible exception of Lloyd’s of London). Enterprises are “self-insured” today, and should

expect to be self-insured tomorrow. Insurance “premiums” are paid internally, in the form of additional hardware, software, and people. One principle remains the same, however — when you pay for data protection insurance, you want to minimize its cost and maximize its value.

“One principle remains the same, however —when you pay for data protection insurance, you want to minimize its cost and maximize its value.”

As we have noted above, data protection seeks not only to ensure the availability of data, but also its confidentiality, privacy, and availability to regulators. This is still insurance — the legal costs of failure to protect confidentiality and privacy, or to fail to supply appropriate information to regulators are high, as are the competitive disadvantages of leaking proprietary information. For example, the attention that is now being turned to “business compliance” has at its heart appropriate protection of data such as e-mails.

Data Protection Has To Be Placed in the Right Framework

IT organizations are actively examining how to improve the data protection function, as shown by an increased interest in disk-based data protection strategies and a number of new replication technologies. Trying to sort through the myriad of choices can be difficult.

The key to choosing any of these strategies and technologies is understanding the overall context, the overall “data protection infrastructure portfolio,” into which individual data protection technologies should fit. Otherwise, what appear to be individually sound decisions may not lead to offering the necessary levels of data protection. Among the problems that can occur are:

- Failure to protect data adequately
- Making the wrong allocation decision (spending too much on areas that do not really require a level of protection and too little on areas that require greater protection)
- Straining the IT administrative resources assigned to data protection even further and with less results than necessary

Without the right framework, enterprises cannot know where to place their longer-term data protection technology investment bets or how much they should place on each bet. And that means that any framework has to take into account the changing world of data protection technology.

Ride the Sea Change in Data Protection

Change that affects the requirements for data protection is coming from several directions. One of the directions is extending and improving what is already being done. An example of this is disk-to-disk backup.

A second direction is change in the basic way that the movement and storage of information is carried out in an organization. For example, ILM is not only about moving information from one tier of storage to another, but also about managing stored information differently — and a major effect of the difference in information management is in better data protection. Moreover, ILM leads to an overall change in the mix of data protection technologies (e.g., data replication vs. data backup) that are used within an enterprise.

A third direction of change comes about from changing business requirements. A key illustration is a new business emphasis on compliance. IT business-compliance policies require understanding and implementation of new policies, practices, and procedures as well as possible new hardware and software data protection technologies.

The rest of this report examines the basic principles of data protection in light of these changing business requirements and in light of existing and emerging data protection technologies. The key take-aways that should be kept in mind when reading the rest of this report are:

1. Determine where over-investment and under-investment in data protection technology is taking place, so that your IT organization can direct future investments to shore up the weak spots.
2. Determine what the effects of changing business requirements and technology advances on your enterprise's data protection investment are.
3. Gain a sense of how the major categories of data protection technologies interact, so that you can determine the proper mix and deliver the proper level of service.

Chapter Two:

Business Continuity: The Framework for Data Protection

Business continuity is the mitigation of risk caused by interruption to normal enterprise activities and processes. More specifically, business continuity imposes a software and hardware superstructure on key IT systems and networks that aims to ensure that (a) business-critical applications are available to all end users all of the time despite failures of individual components (resiliency) and (b) when these applications are not available, the outage time is as short as possible (high availability).

Effective business continuity protects key stakeholders' interests, brand reputation, the good-will of customers, and the value-creating activities of the enterprise. If a business-continuity strategy fails, the consequences range from undesirable or unacceptable (customer dissatisfaction or loss of productivity), to severe (economic loss of market valuation/revenue or loss of public or customer confidence), to the most severe (business failure).

For an application to work, all of its components — hardware, software, information storage, and networks — need to work. Of these, information storage is typically the most important, both because its loss can be irreparable and because, over time, it becomes the key bottleneck to recovery. In a recent disaster, a law firm lost all data on its customers, as well as its legal and office applications, network, and PCs. Replacing the software, hardware, and network was the matter of a business day. Replacing the data was effectively impossible, and the lack of data on the law firm's clients forced the firm to fold. Likewise, reloading an application, hot swapping a server, or rerouting messages along a network, today, is practically speaking a matter of minutes; reloading data into a terabyte-sized database for an operational system is a matter of hours or days.

Therefore, the key task of any business continuity strategy is data protection; and, by the same token, a key aim of data protection is business continuity. Furthermore, a business continuity strategy and architecture can serve as a good framework into which to fit data protection technologies and strategies. It is comprehensive; it ensures that the needs of other parts of the architecture aside from storage and the business as well as IT are taken into account; and it fully recognizes the crucial role of information storage. The rest of this Chapter considers how a business-continuity framework can enlighten and improve a data protection strategy.

Business Continuity and Data Protection

To understand why enterprises may not be receiving the level of data protection that they think they are requires an understanding that business continuity is not only about *disaster continuity* (more familiarly, thought of as disaster recovery), but also *operational continuity* — the ability to deal with day-to-day operational problems. The right amount of attention for data protection has to be given to each — and that may not always be the case.

If IT organizations do not understand how day-to-day operations and disaster recovery planning have different requirements for both physical and logical data protection, they may not have the right technology mix — and may make wrong investments — for data protection.

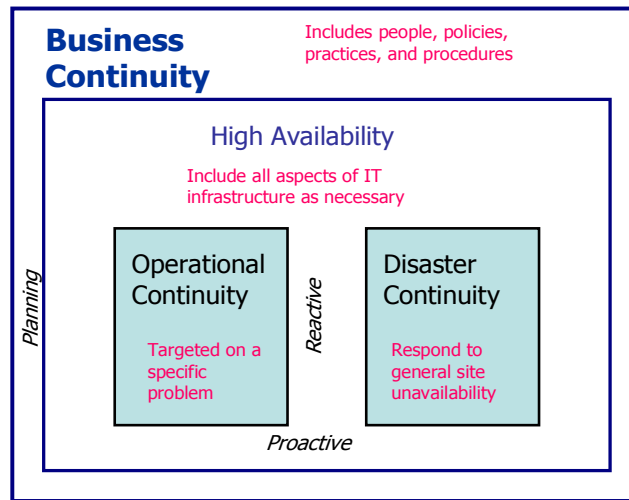
Both operational and disaster continuity require the proper level of both *physical* (storage device level) and *logical* (the data itself) data protection. A data item may be flawed although the disk is functioning perfectly; a disk may crash but the same data on a different disk be preserved.

Although there are ways to ensure logical data protection, the primary emphasis tends, traditionally, to be on the physical side. This can be a problem — database corruption that occurs in the middle of a vital business-intelligence query or customer order is not in the best interests of the business. Operational continuity requires an emphasis on logical data protection after the basic physical data protection requirements have been met.

Business Continuity Is Not Just Disaster Recovery

Business continuity tends to have an information technology flavor, but it is (or should be) an enterprise-wide activity that includes manual practices, processes, and procedures as well. Likewise, business continuity spans both disaster recovery and operational continuity (see Figure 2-1).

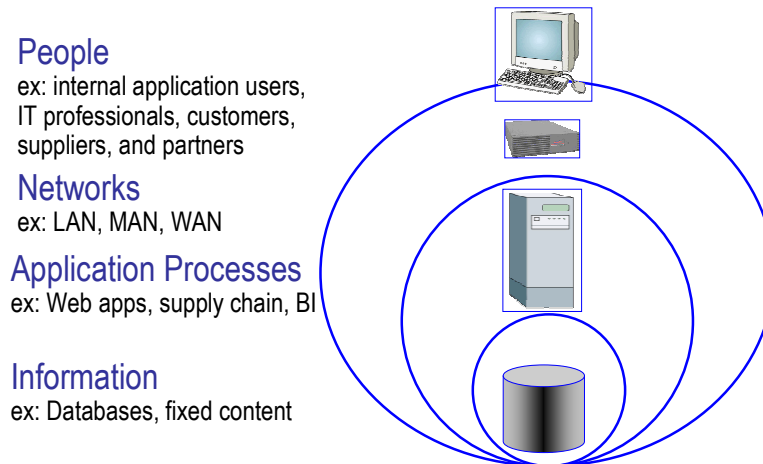
Figure 2-1: Overview of Business Continuity



Source: Mesabi Group, January 2005

On the IT side, data protection is a necessary, but not a sufficient condition for business continuity (Figure 2-2).

Figure 2-2: Business Continuity Is More than Data Protection



Source: Mesabi Group, January 2005

As noted above, information (i.e., useful and usable data) is at the heart of business continuity. Depending on the level of severity, without data protection the ability of processes and people to work successfully is jeopardized, if not impossible.

Yet data protection alone is not enough. Without the applications, processes, and networks working properly—and without people with the right skills in the right places who are able and willing to use them properly—the data cannot be accessed and used.

Likewise, IT may have responsibility for achieving high availability and resiliency of the computer systems — but business continuity involves more than this. Business continuity requires planning, so that the right people with the right skills will have the right tools and right knowledge at the right time in order to respond to a threatened or actual negative service-level-impacting event. Thus, IT typically cannot have full responsibility for overall business continuity.

Business continuity requires planning, so that the *right* people with the *right* skills will have the *right* tools and *right* knowledge at the *right* time in order to respond to a threatened or actual negative service-level-impacting event.

IT's responsibility for the availability of the IT infrastructure builds upon two pillars — operational continuity and disaster continuity. The word “continuity” indicates that proactive actions are the responsibilities of these two functions, for example, provisioning a disaster recovery facility to enhance the prospects of minimizing the impacts of a disaster should one occur.

Note that operational continuity focuses on targeting individual problems that hopefully have a limited scope, whereas disaster continuity has to focus on what would need to be done in the event that the entire IT infrastructure, including all applications and their supporting server, storage, and network services, has to be replicated at a site other than the original home of the applications.

Proactive activities of planning, provisioning, monitoring, and preventive maintenance prepare an enterprise as best as possible against the storms of service-level-threatening or devastating events. If and when such an event occurs, the two pillars turn into operational *recovery* and disaster *recovery* (Figure 2-3).

Both operation recovery and disaster recovery face danger from natural causes (including the inherent limitations in technology) and human-related causes, whether that be inadvertent human error or intent-based malign attacks. Understanding the difference between the two types of recovery is essential for understanding the type of data protection that is best for each.

Figure 2-3: Business Continuity Keeps Your Business Running

	 Operational Recovery	 Disaster Recovery
<u>Time Down</u>	• Minutes/year	• Hours to day per event
<u>Infrastructure Threat</u>	<ul style="list-style-type: none"> • Disk failure • Network congestion • Application performance degradation 	<ul style="list-style-type: none"> • Earthquake • Blizzard • Fire • Flood
<u>Willful Threat</u>	<ul style="list-style-type: none"> • Computer viruses • Hacking 	<ul style="list-style-type: none"> • Disgruntled employee • Terrorist attack

Source: Mesabi Group, January 2005

Disaster Recovery: Let's Get Physical

SNIA's ILM definition of Disaster Recovery is:

Disaster Recovery (DR): The recovery of data access to data and associated processing through a comprehensive process of setting up a redundant site (equipment and work space) with recovery of operational data to continue business operations after a loss of use of all or part of a data center. This involves not only an essential set of data but also an essential set of all the hardware and software to continue processing of that data and business. This may involve some amount of down time to perform the recovery.

Note that an event is considered a "disaster" only when data processing has to be moved from a primary to a secondary site and when that processing is carried out using a different set of computer hardware (including both servers and storage).

Physical data protection is properly the first focus of disaster continuity, but logical data protection needs to be taken into account (as well, of course, as reconstitution of applications, networks, and people resources). If there is a failover to a second site, that site now assumes an operational role, and so the data must not only physically exist at the second site but also be usable immediately.

Problems that can be fixed at the primary site without requiring the movement of data to a second site are not a disaster in a business continuity sense. However, operational problems, such as the long-term loss of a critical database, could still be catastrophic to a business. (While IT would understand the difference from a planning perspective, business management may not understand, and therefore care about, the distinction!)

The recovery from a disaster may take hours to days. That may seem contrary to intention when remote synchronous mirroring can lead to a nearly instantaneous restart for storage. However, even if all the storage for an enterprise is mirrored remotely (which may very well not be the case), storage is only one aspect. The rest of the hardware infrastructure — servers and networks — and software infrastructure — applications, databases, and operating systems — also have to be in place. Additionally, people need to be in place. An assessment process determines when to declare an emergency that results in a total transfer to a disaster recovery site, and that may take time.

Interestingly, operational recovery can often take more time than disaster recovery. While the secondary site may benefit from synchronous mirroring, recovering a primary site either from the secondary site or from a tape backup can take hours — or days.

Most IT organizations have never experienced a true disaster and hopefully never will. That does not mean that the second site is not necessary. In order to get the most out of the investment that needs to be made anyway, one objective might be to make the second site as useful as possible under normal business conditions. For example, failing over to a second site temporarily for planned maintenance or equipment upgrades at the first site and workload balancing are a couple of reasons that the second site might be useful outside of disaster recovery. But, since all that could probably be done at a single site with better economies of scale, the real reason for the second site is “distance separation,” to minimize the risk of both sites being impacted at once. And even though disasters are relatively rare, that “insurance premium” is probably well worth the incremental cost.

Operational Recovery: Think Logically

SNIA’s ILM definition of Operational Recovery is:

Operational Recovery (OR): Recovery of one or more applications and associated data to correct operational problems such as a corrupt database, user error or hardware failure. May use point in time copy or other techniques to create a consistent set of recoverable data.

Note that an operational recovery can be due to either a logical problem (say virus or accidental file deletion) or physical problem (say two drives failing in the same disk array where there was protection against only one failure before the drive that initially failed was rebuilt). That said, operational recovery has a strong emphasis on logical data protection once the basic physical data protection technologies are put in place.

Operational recovery is within the control of the IT organization and is the responsibility of the IT organization. Operational recovery assumes that all (or nearly all) of the users are able and willing to use the applications; i.e., normal working conditions prevail.

Disaster Recovery Requires Judgment; Operational Recovery Requires Automation

Disaster recovery responds to a systemic event that affects all applications either simultaneously or in a rolling manner. Disaster recovery requires a triage approach to recovery, where the most time-sensitive and business-critical applications are restored first. Operational recovery is a response to a non-systemic single event. A logical event typically affects a single application (although dependent applications may also be affected). In the case of a hardware failure, such as a disk array, all applications that use the affected disks would be involved, but other applications (unless there are dependencies) would not be. If disaster recovery is like responding to an epidemic, then operational recovery is like responding to events in an emergency room.

If disaster recovery is like responding to an epidemic, then operational recovery is like responding to events in an emergency room.

An operational recovery may involve other people in the IT infrastructure other than storage personnel (such as a database administrator), but should be able to be handled within the IT organization itself. The same is not true with disaster recovery.

Disaster recovery is really the responsibility of the entire enterprise, although the leadership role for the process may be assigned to the IT organization. (In many cases, IT may be in the unpalatable position of being charged with the total responsibility, but not given sufficient authority or resources to ensure the disaster recovery plan).

Disaster recovery can be a people-intensive process. Although key aspects of the process can be automated (such as failover to a remote site), disaster recovery requires professional management. The range of possible disasters is so wide that human beings have to be able to adjust and compensate for the nature of the particular type of disaster. Although some degree of automation may be possible, adjustment to unplanned situations requires human judgment.

Disaster recovery is always about a service-level-impacting event. Operational recovery may be able to respond to a service-level-threatening event before it becomes a service-level-impacting event.

Whenever feasible, IT professionals want to delegate continuity and recovery responsibility to processes with the highest degree of automation available. The reason is simple: Automation can sense potential SLA-impacting events before they happen (such as an out-of-space condition on a disk array) and may be able to take corrective action based upon policies (hence why there is so much discussion about policy-driven management).

Automation that can take corrective action to SLA-impacting or-threatening events (i.e., self-healing) will be welcomed by IT administrators when it arrives. Software that helps by monitoring, alerting, and advising is available and welcome now.

That does not mean that all problems can be solved through policy-based automation, or that IT administrators will not have a concern about possible loss of control. IT operations always run into surprises that require judgment to resolve. We anticipate, however, that concern over loss of control will probably give way to recognition of the benefits of automation in minimizing SLA-commitment-breaking unavailability.

Logical Data Protection Gets Short Shift in Business Continuity

Operational continuity and disaster continuity each need a different mix of data protection technologies to achieve the planned levels of data protection that an enterprise requires. Yet enterprises may not have a clear understanding of the differences between physical and logical data protection — and that may result in a dangerous lack of attention to logical data protection.

Let us start by examining both physical and logical data protection more closely. As already pointed out, data protection divides into two classes: *physical* data protection and *logical* data protection. To provide full data protection, both are mandatory. Physical data protection focuses on *storage* devices themselves to recover from dysfunction, failure, or destruction of one or more physical components of a storage system. Logical data protection focuses on the protection of the *data* itself; the bit patterns must retain their designated order and completeness. (In other words, a user must get back exactly the data that was put in — reordered bits or missing bits would destroy the integrity of the data and could lead to potentially serious consequences.)

Physical data protection is built upon redundancy (through expansion of storage requirements beyond actual usage, using parity protection schemes and full copies) and locality (separating an “original” copy from a second copy through geographical separation). Logical data protection may use redundancy, but it is built primarily upon isolation (taking a set of data out of the I/O path so nothing can be changed), locking (using software to prevent I/Os from changing a particular piece of data), and fixation (using hardware or software for write once, read many capability to a piece of storage media).

The first focus of *physical* data protection is on data availability, while the first focus of *logical* data protection is on data preservation. Physical data protection can do error correcting and consistency checks to determine that the data is the same as was written. However, that does not mean that the data is correct. (The infamous GIGO — garbage in, garbage out — applies.)

Physical data protection cannot prevent I/O processes from changing data because that is what those I/O processes are put in place to do. Yet those processes, say a virus or data corruption, may change bit patterns so that the data is unusable. Therefore making an exact copy of the data (say local or synchronous remote mirroring) provides physical data protection. The exact copy does not provide any logical data protection because any change in the original — whether right or wrong — is reflected in the copy.

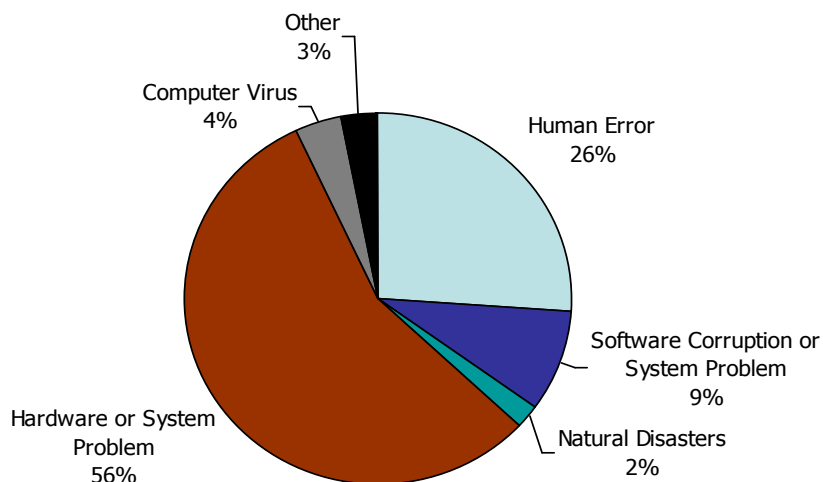
Logical Problems Feature Prominently in Data Loss or Downtime

Ontrack, a data availability service provider, has created a compelling chart detailing the causes of data loss or downtime (Figure 2-4).

Note that site disasters account for only 2% of data loss or downtime events. This low percentage would not excuse any underinvestment in disaster recovery implementations as part of the overall business continuity process, as the expected value of loss is high if the probability of a disaster is multiplied by the large magnitude of loss that could result from a disaster.

Note also that 56% of the problems are caused by hardware or system problems. Although that seems to be a high percentage, since a number of reliable physical data protection technologies exist, that result could be due to an underinvestment in the necessary technologies or some other reason.

Figure 2-4: Causes of Data Loss or Downtime



Source: Ontrack, 2004

Logical problems account for 39% of the causes of data loss or downtime, but the biggest chunk is human error (note that automa-

tion may pay large dividends by avoiding human error). Software program malfunctions can cover a wide range of problems, including database corruption. And even though computer viruses were only 4% of the problem, online attacks on data integrity are only likely to increase, and with them the possibility of a devastating impact.

Why is this important? Because much planning focuses on delivering disaster continuity, but the most likely threats to continuity are likely to come from the operational side. And on the operational continuity side, logical data protection problems shout for attention. But can logical data protection problems really have a significant impact upon an enterprise?

Logical Data Protection Problems Manifest in a Number of Ways

Hardware, software, and people are all extremely complex, so together they represent a combustible mixture in a data center environment. In complex IT environments, what is perhaps surprising is not that problems arise, but that there are not more of them. Among the myriad of potential logical data protection problems, the following represent just some of the possible sources that can result in data loss or downtime.

Table 2-1: Logical Data Protection Problems and Sources

Problem	Source
Data corruption through loss or alteration of data without the application's knowledge and consent	<ul style="list-style-type: none"> • Faulty hardware (bit loss or incorrect ordering) • Software bugs (unexpected conditions reached and responded to incorrectly) • User or IT administrator error (accidental file deletion)
Downtime and/or data corruption through application errors	<ul style="list-style-type: none"> • Faulty application version (introducing new software without sufficient testing) • New system interfaces (semantic interpretation errors) • Database errors (out of order transaction commits, accidental deletion of rows, dropped tables, etc.)
Data corruption through willful action	<ul style="list-style-type: none"> • Externally from viruses and worms • Internally from deliberate tampering
Downtime through unintended human error	<ul style="list-style-type: none"> • Storage system configuration error • Allowing out of space conditions to occur

Source: Mesabi Group, January 2005

Problem diagnosis and assessment is not always easy, especially if the problem is intermittent, manifests itself in only small ways from which it has to be deduced that a larger problem is likely to develop, or is buried as a time bomb set to explode under certain conditions. Early detection to prevent a service-level-impact-threat from becoming a service-level-impacting reality requires eternal vigilance on the part of IT administrators.

A logical data protection problem can affect a key application, whether the application crashes or not. The inability to dispense cash from an automated teller machine or the inability to correctly deliver the right goods to a customer in a timely fashion can affect an enterprise's credibility (and market valuation).

Whether the logical data protection problem is pernicious and persistent or simply quickly diagnosed and corrected, logical data protection must get its full measure of attention.

Do Not Neglect Any Facet of Data Protection

No aspect of data protection can afford **not** to be protected. The target for the IT organization starts with four simple boxes (Table 2-2). Both operational continuity and disaster continuity have a physical *and* a logical component to them. Each box has to be considered individually and all four boxes together have to be considered collectively to devise a data protection solution that meets an enterprise's needs.

Table 2-2: Data Protection Category Matrix

	Operational Continuity	Disaster Continuity
Physical		
Logical		

Source: Mesabi Group, January 2005

Although it seems simple, filling in the matrix is not that easy. The first challenge is in knowing when the levels of data protection are enough. The second challenge is in understanding that the target is moving and knowing how that will affect what needs to go into the matrix to get the right levels of data protection.

Chapter Three:

Data Protection — Where the Problems Lie

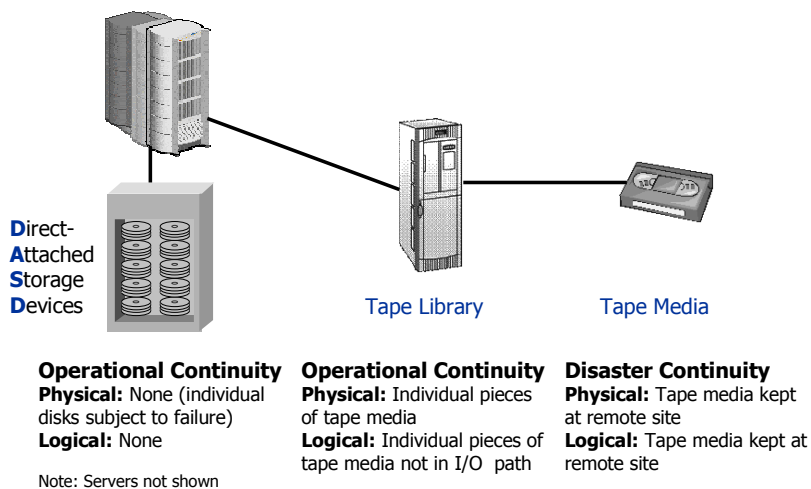
Understanding IT's heritage of data protection technologies is essential to understanding the thinking that still permeates the IT community regarding the nature of data protection. The genesis of that "thinking" is from the limited choice for data protection technologies at the time, as well as cost considerations. Not that many years ago, data protection and the backup/restore process using tape were synonymous. All the other numerous data protection technologies in use today were not only unavailable; they were unthinkable.

Data Protection as It Was in the Beginning

Just prior to the introduction of RAID technology, the primary storage media in use was, as today, Winchester disk and tape (Figure 3-1). (For simplicity's sake, because they were a minor part of the overall market, optical storage technologies and microfiche will be ignored.)

Winchester disk technology itself offered no extra measure of data protection. (Winchester disks are disks where the disk media itself and the disk drive are sealed into a single unit. Prior to the introduction of Winchester technology, disk pack media could be removed. The removability of disks is being reintroduced with the removability of RAID groups — which include both drives and media). Each Winchester disk (hereafter referred to as simply "disk") had to stand on its own, so that the mean time between failures (MTBF) for multiple disks was far less than for one disk. Although the technology for disk copies existed, the cost — except for extremely critical online transaction processing systems — was prohibitive. Practically, neither physical nor logical data protection existed.

Magnetic tape solutions provided not only the first line of defense against data problems, but also the last (and any intermediate) line of defense as well. A tape solution consists of tape media, tape drives, and tape automation. Tape media has evolved from reels to more-easily-manipulatable cartridges, and tape drives have shrunk dramatically as a consequence. Furthermore, tape automation (such as a tape library) has improved the flexibility of management of a large number of pieces of tape media, but tape continues to be slower than disk, and the process of transfer a lengthy one.

Figure 3-1: Data Protection: The Way It Was

Source: Mesabi Group, January 2005

Unlike disk, individual pieces of tape media are easily removed from a tape drive and can run in any compatible tape drive. This capability is important because tape media can be transported to a remote site independent of a tape drive, and then put to use in a suitable tape drive at the remote site. Thus, movement of data is dependent upon the availability of transportation, but not upon the availability of a network. Tape drives can operate independently, but are often embedded in tape automation solutions (such as an autoloader or tape library).

The copying of data from disks to tape media is done through the use of backup/restore software. This is the traditional backup/restore process, and it was essentially the only software for data protection.

This process actually provides a great deal of both physical and logical data protection for both operational and disaster continuity. Since each tape copy is on a copy of physical media other than the primary disk copy, tape delivers physical data protection. Since any tape cartridge that is not in a tape drive is not in the I/O path, tape media also deliver logical data protection. Since a tape copy can be physically transported to a disaster recovery site, tape provides both physical and logical data protection for disaster continuity.

Since several generations of tape copy are available, no tape copy represents the only copy, and no single point of failure exists—thus, tape is both the “front line” and “last line of defense” for data protection.

Apart from scalability, reliability, and manageability issues, the key concern with a sole tape-based solution is lack of “high” availability, where high availability might be defined as minutes per year (and surely no more than hours per year).

A major recovery using tape may require hours at best, a day or more as a likely occurrence, and a week or more in extreme circumstances. The first reason is that the recovery process, which is called a restore process, is actually a *rebuild* process. Tape, as a sequential medium, is too slow to handle random online processing. Therefore, before the data can be used for online processing, the contents of tape have to be copied to disk — and that can take a long time. If a mirrored disk copy were already available, the process would be a *restart* using the mirrored disk, and could take seconds to minutes.

The second reason that tapes are not always an optimal solution for recovery is that a number of tapes are likely to be needed to restore a disk system, and that can lead to problems if a part of the physical tape system (such as a piece of tape media) fails or requires significant resuscitation work (such as in the case of an intermittent error condition or constant read retries), or if there is a mistake in the sequencing of the tapes.

Considering the above examples, tape by itself is not sufficient to meet the demands of modern IT organizations.

Typical Data Protection Technology Today Still Leaves a Lot to Be Desired

A lot of change has taken place in data protection since the early days (Figure 3-2). One of the primary improvements in data protection technology was the introduction of redundant array of independent disks (RAID), which provided physical data protection for the price of one or more additional disk drives. (RAID originally stood for redundant array of *inexpensive* disks, but the word *independent* was substituted for inexpensive as the price of disks fell dramatically and the relationship of disks to one another became more important than their cost.)

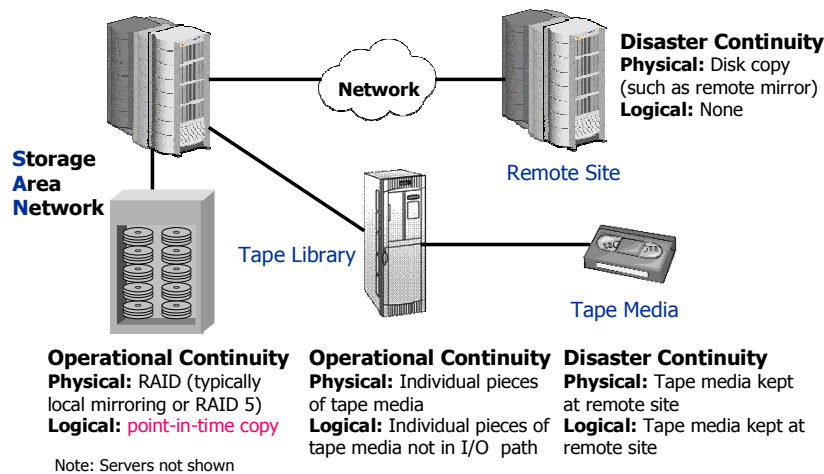
A number of RAID levels exist, but only a few are in common use. RAID 1 is a synonym for mirroring — for every disk that contains “original” data, a corresponding disk exists that has a “copy” of that data. This means that the usable disk space for RAID 1 is only 50% of the total available disk space. Parity RAID levels enable the recalculation of “lost” data in the event of a single failed disk through the use of parity check data. A RAID 5 group requires only one more disk than the number of disks required to hold the working data, although the parity check data and the working data are actually spread across all the disks.

RAID technology delivers dramatically improved availability over an array of unprotected individual disks, and RAID technology typically now forms the first line of physical defense.

Remote mirroring, a variant of RAID 1, delivers fast restart physical disk protection at remote sites to aid in disaster continuity.

A second major advance is point-in-time copy capability, which is a fixed view of the data and therefore is not subject to change from I/O processes. This independence enables a point-in-time copy to deliver logical data protection as of that instant. Since that instant is unlikely to be the instant when a logical failure occurs, the aim of point-in-time copying is to be able to recover with a minimal loss of data. And even though point-in-time capability is available, some organizations do not necessarily use it for logical data protection on the disk array itself, but rather as a vehicle to provide a consistent point for invoking the standard backup/restore process.

Figure 3-2: Typical Data Protection Today



Source: Mesabi Group, January 2005

Although data protection has significantly improved from the pre-RAID days, improvements over the typical data protection configuration are still necessary for a large percentage of enterprises in all four boxes in the data protection category matrix. This necessary improvement may not be a technology issue, but rather an education, adoption, and cost issue.

The myriad of new data protection products that have become available over the past few years, as well as the continued evolution of data protection products and services expected in the near future, indicates that the availability of the necessary technology is probably not the primary inhibitor to implementing effective data protection for business continuity. Understanding, finding, affording, and implementing the appropriate *mix* of data protection technologies is more likely to be the key issue.

The following sections address the state of the art in each of the four boxes in the data protection category matrix.

Operational Continuity/Physical: Generally Strong, But Some Improvement Needed

The addition of RAID technology has made this a strong area and, with concepts like triple mirroring, an enterprise can buy its way to a desired level of physical availability. The operative word is “buy,” as incremental changes in availability can become very expensive. The Achilles heel in RAID technology is that a typical RAID array can only allow one disk failure and still protect the data. During the period in which the RAID group is being rebuilt, the data in the array is exposed to the risk of data loss if a second failure should occur. And the chance for a bit error in rebuilding a large disk drive as compared to a smaller drive is by no means insignificant.

With that said, would not an advance in RAID technology to allow more than one failure in a RAID group be useful? The answer is yes; and one small company has had such a technology for several years, and at least two large storage suppliers are now considering offering the same type of technology. The cost for doing so could be low, as the hot spare that is typically found in RAID arrays could be put to active use. Although there would be a slight performance penalty, such multiple-failure-tolerating RAID technology would be the closest thing to a higher availability “free lunch” that is likely to come along soon.

Operational Continuity/Logical: More Attention Needs to Be Paid to Logical Data Protection

Point-in-time copy capabilities, including snapshot copy capability, have proven to be helpful for logical data protection. A powerful use of point-in-time copy derivative capability, called continuous data protection, is now becoming generally available. A number of other technologies, including replication technologies, virtual tape libraries, and write-once-read-many (WORM) technologies, are also available to aid with logical data protection. In short, tape now has a number of allies — in addition to basic point-in-time copy capability — to help with logical data protection.

Many of the technologies, such as continuous data protection and virtual tape libraries, are still relatively new, so IT organizations may be either unfamiliar with the technologies or still in some stage of evaluating the technologies. However, point-in-time technology has been around for quite awhile and has been used successfully by a large number of organizations. Nevertheless, point-in-time functionality not yet been adopted to the extent that it needs to be in order to provide the right level of logical operational continuity. The lack of adoption may be due to an IT tendency to focus on disaster recovery in general and the physical side of recovery for both operational continuity and disaster continuity, rather than the logical side; but logical operational continuity needs its fair share of attention as part of a comprehensive data protection strategy.

Disaster Continuity/Physical: Well Done, But Cost and Distance Are Issues

Remote mirroring has proven its worth, and has been justifiably successful in the data protection marketplace as a result. However, unless an enterprise already has a data center that can serve as a secondary disaster recovery site for the enterprise's primary site, the cost for establishing a disaster-specific site can be quite expensive.

Cost — network equipment, software, and remote disk array cost — is a barrier to synchronous remote mirroring implementation for many organizations. One reason is that many of the original synchronous remote mirroring products required that the disk array at the second site be the same model as the disk array at the first site. However, more cost-effective remote replication technologies are now available for organizations that are willing to make some concessions in exchange for cost savings. For example, if an organization can tolerate the performance loss penalty in case of a disaster, the ability to use less expensive disk arrays as the target is an option they might find attractive.

The second issue is that the distance between a primary site and a secondary site should be targeted at 300 miles or more. Although 300 miles is an arbitrary figure, it is a distance being mandated for certain compliance activities. Even if an organization is not subject to compliance restrictions, common sense says that if you are planning a long-distance data center, then there is no sense in choosing one 250 miles away if there is any possibility of not falling into compliance at a later date.

However, synchronous remote mirroring is typically used between two data centers that are no more than 100 miles apart. Once again, this is an arbitrary limit, but one that is based upon experience with acceptable response time latency for the valuable online transaction processing (OLTP) applications that can justify the expense of synchronous remote mirroring. The good news is that OLTP data probably is not compliance data (and even if it were, a copy could be sent to a third data center). The bad news is that there is often a large volume of file data — notably e-mails — that might need to be sent to a distant data center.

In summary, asynchronous remote mirroring and other remote replication technologies are available to accommodate the needs of physical (and in some cases logical) data protection at a distant site. The challenge to IT organizations is how to meet the necessary data protection requirements while not having to use more remote sites than is absolutely necessary.

Disaster Continuity/Logical: The Danger of Being Under Protected May Be Very Real

If primary site processing has to move to a secondary site because of a disaster, the former secondary site has to assume the mantle of being the primary site. One of the first questions that needs to be asked is about the length of time that the original primary site will be out of service. If the answer is a week or more, the enterprise may want to shoulder the burden of implementing a complete logical data protection solution if one is not already built in — and it may not be.

Even if an outage may last a week or more, additional logical data protection may not be needed if the disaster recovery site replicates not only disk storage, but tape storage as well. (A subset of the original tape solution may be enough in a pinch if the data center environment's configuration, e.g., its space and power, can accommodate expansion.) If the disaster site does not replicate tape storage, a third-party disaster accommodation arrangement could suffice.

A point-in-time copy (or equivalent) capability might serve as a stopgap measure, but tape would provide secondary physical data protection as well as logical data protection.

In any event, a strategy for logical data protection needs to be put in place now, if the organization has not already done so. There is no sense in making a large investment in a disaster recovery site if you are not protected from permanent loss of data due to database corruption, accidental file deletion, virus, or other logical data protection problems.

Summing Up Data Protection Challenges by Category

IT organizations need to examine the data protection technology challenges to determine how they affect the data protection planning process within their enterprise (Table 3-1). These should be kept in mind when setting the objectives for data protection for their enterprise.

Table 3-1: Data Protection Challenges by Category

	Operational Continuity	Disaster Continuity
Physical	<p>Key Available Technology: RAID 1, RAID 5, and variants</p> <p>Key Challenge: Relatively inexpensive and low performance impact multiple parity RAID</p>	<p>Key Available Technology: Synchronous and asynchronous remote mirroring</p> <p>Key Challenge: Getting RPO as close to zero as possible over long distances</p>
Logical	<p>Key Available Technologies: Point-in-time copy capability and tape</p> <p>Key Challenge: Acceptance of continuous data protection technology</p>	<p>Key Available Technologies: Vaulting and electronic vaulting</p> <p>Key Challenge: Acceptance of dated replication technology as a complement to existing technologies</p>

Source: Mesabi Group, January 2005

Chapter Four: Data Protection — Setting the Right Objectives

Recognizing *where* there may be problems in your data protection strategy is not enough. As an IT strategist, you should also understand *what* the right objectives for a data protection strategy should be. Setting those objectives is critical, but not necessarily easy.

The objectives that the IT strategist should consider are high availability, data preservation, data responsiveness (i.e., getting the data to the user in a reasonable time), and confidentiality. These must be both maximized and balanced, because changes in one area affect others, and therefore too much focus on one objective can lead to problems meeting another objective.

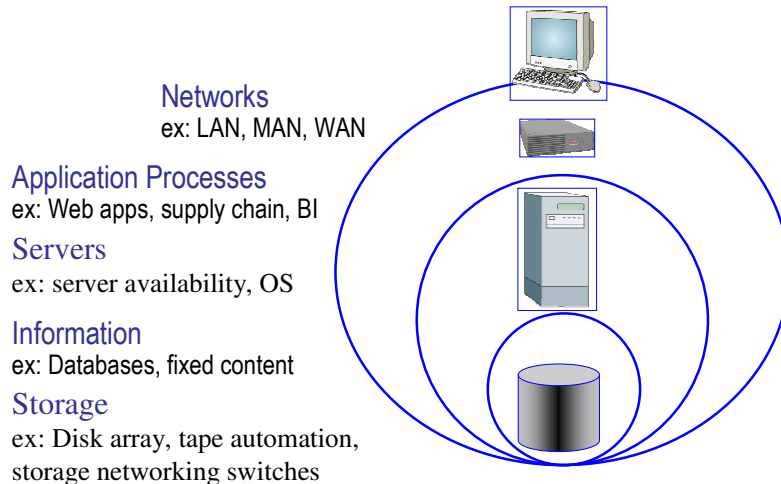
How High Is High Enough for Data Availability?

High availability is a key objective in a data protection strategy, and one of the keystones of business continuity. However, an overemphasis on high availability can lead to problems with data preservation (all the money goes into keeping the systems up, and very little goes into preventing data loss when they do go down), data responsiveness (fault-resilient storage often does not restore as quickly), and confidentiality (all the money goes to keeping the systems up, and very little to protecting the data from hackers). As a result, an organization may not meet its real data protection goals and probably will spend more than necessary for data protection.

As business continuity is more than data protection, so does high availability depend upon the entire IT infrastructure (Figure 4-1). That means that the entire IT infrastructure has to be tuned to the same relative level of protection. However, if any part of the infrastructure can be significantly improved for a relatively small increase in cost, the investment is probably worthwhile. Unfortunately, as availability increases incrementally, costs rise quickly.

However, all else being equal, incremental investment to increase the availability of storage is preferable to incremental investment to increase the availability of other parts of the IT infrastructure, such as servers or databases. If carefully done, investment in storage availability typically can improve data preservation and data responsiveness as well.

Figure 4-1: High Availability Depends upon the Entire IT Infrastructure



Source: Mesabi Group, January 2005

SNIA's Data Value Classification: A Point of Departure

To determine objectively what the level of data protection for each application should be, the Data Protection Initiative (DPI) within the Data Management Forum (DMF) of the Storage Networking Industry Association (SNIA) defines three key terms: (1) Recovery Point Objective, (2) Recovery Time Objective, and (3) Data Protection Window.

- 1) The DMF definition of Recovery Point Objective (RPO) is:

The maximum desired time period prior to a failure or disaster during which changes to data may be lost as a consequence of recovery. Data changes preceding the failure or disaster by at least this time period are preserved by recovery. Zero is a valid value and is equivalent to a "zero data loss" requirement.

This is a definition of the amount of permanent data loss. Permanent data loss means that the data cannot be restored through use of IT. In some cases, manual reentry of data may be possible, but that ability may be infrequent.

- 2) The DMF definition of Recovery Time Objective (RTO) is:

The maximum desired time period required to bring one or more applications and associated data back to a correct operational state.

This is the maximum downtime that an application should suffer for a single failure event.

- 3) The DMF defines Data Protection Window (DPW) as:

Like the backup window, this is the available time during which a system can be quiesced and the data can be copied to a redundant repository without impacting business operations.

This is critical because data has to be protected. The question is, how can that protection be provided without unnecessary “unavailability” of systems?

Building upon these definitions, the DMF goes on to define five classes for data value classification and the resulting RPO, RTO, and DPW for each data value class (Table 4-1).

Table 4-1: SNIA Data Value Classification

Data Value Class	Data Availability	RPO (Data Loss Risk)	RTO (Max. Recovery Time)	DPW (Copy Data Time)
1 – Not Important to operations	90%	1 week	7 days	Days
2 – Important for productivity	99%	1 day	1 day	12 hours
3 – Business important information	99.9%	2 hours	2 hours	10 minutes
4 – Business vital information	99.99%	10 minutes	15 minutes	None
5 – Mission critical information	99.999%	1 minute	1.5 minutes	None

Source: Derived from SNIA “Implementation Guide for Data Protection,” March 2004

The DMF then defines a five step implementation guide for data protection: identify data value class, define best solution, select specific components, check system cost, and confirm decision or change. The five step process seems reasonable, but then the implementation guide adds that: “if cost is too high, change data value class or specific components.” In other words, if you can’t afford it, change your mind about how important it is!

Although the DMF should be given a lot of credit for the work that it is doing — and this work is only a draft in an attempt to tackle a very difficult issue — accepting the data value classification and implementation strategy at face value might not be the best strategy for IT management. An enterprise attempting to use the DMF’s approach should think clearly about its applicability to the enterprise’s circumstances. The reasons for looking closely at the pertinence of the table are as follows:

- *Value is not the same as availability* — making the assumption that they are equivalent can lead to a misallocation of data protection investment dollars.
- *The RPO and RTO for operational recovery and disaster recovery are not necessarily the same* — making the assumption that they are the same can lead to a misapplication of resources when making a recovery.
- *Availability is only one data protection objective* — giving availability excessive weight versus the other objectives can lead to mistakes in protecting data.

The following sections elaborate on these findings.

Do Not Equate Availability with Value

No statistical correlation has been proven between the value of data and the need for the availability of that data, yet that is the fundamental assumption of the data classification scheme.

It is likely that if an application requires high availability, it is mission-critical. Everyone can point out applications that need (i.e., require) 99.999% availability at least because of the severe financial penalties (such as lost revenues) for minimal downtime and therefore can be considered mission-critical applications. A telephone billing system comes to mind. With the increase of importance of Web applications, the number of applications requiring extreme availability can only go up. Although there are counter examples (applications that require very high availability, but do not have high business value to the enterprise itself), the exceptions themselves might not be enough to override the assertion.

However, the converse — if an application is mission-critical, it requires high availability — is not true. First, value may come from several sources. The intellectual property value of digital assets, such as CAD/CAM designs or a film library may actually have a quantifiable value. Accounts receivable is the source of customer revenue. A data warehouse may serve as the basis for a business intelligence analysis that changes pricing strategy or decides where to best invest funds in a marketing campaign. An e-mail system may be used to facilitate customer service. Yet none of these applications are likely to require that downtime is limited only to seconds or minutes per year.

Equating mission-criticality and high availability can lead to “straightjacketed” IT. The revenue generation, operational processes, and decision-making processes may run only in normal or extended business hours or at specified times known in advance. Even if the applications are scheduled to run 24 hours a day by 7 days a week, some reasonable amount of downtime, say even hours in a year if absolutely necessary, is desirable.

Why? Unavailability as described makes no distinction between planned and unplanned downtime. Yet planned downtime can be very beneficial to the enterprise and to its customers — an application upgrade that provides additional functionality, an operating system upgrade that enhances security and reliability, and a server upgrade that improves performance are all examples of when downtime becomes beneficial. IT needs flexibility to improve services — not a straitjacket that requires a Houdini-like performance to make the simplest improvement.

IT needs flexibility to improve services — not a straitjacket that requires a Houdini-like performance to make the simplest improvement.

IT applications that have critical business applications that require extreme availability typically have to live with them. However, IT should not go about creating the need for extreme availability unless absolutely necessary. An application that has significantly serious consequences for minor downtime is something to be avoided if at all possible. Unless the reward is high, IT organizations should not take on both increased risk and expense.

Availability Objectives for Operational Recovery and Disaster Recovery Are Not Necessarily the Same

The response to an operational problem or a disaster situation is different (Table 4-2) because each has different characteristics. RPO and RTO may be conditional and contingent depending on circumstances.

Table 4-2:
Operational Recovery and Disaster Recovery Differences

	Operational Recovery	Disaster Recovery
Problem Locus	Concentrated (typically one application)	General (typically all applications)
Mindset	Fire drill	Disaster relief
Management control	“Emergency room”	“Control center”
Response Team	Ad hoc as particular problem requires	Designated DR team
Resolution Strategy	Intensive care	Triage

Source: Mesabi Group, January 2005

An operational problem is typically isolated to an individual IT infrastructure component, such as an application, database, disk array, or tape library. There may be dependencies (say, among applications) that can cause the impact to magnify over time. Generally, however, there is a single triggering event and a single solution (however complex).

By contrast, a disaster affects all applications. Applications have to be recovered and restarted based upon a set of priorities, i.e., a triage. In an actual disaster recovery situation, the people who are responsible for restarting the applications may not be the ones who had operational responsibility for them.

A mission-critical application that automatically fails over to a remote site may have the same RPO and RTO for both operational and disaster recovery. But an enterprise typically has many important applications in its portfolio and many of these do not require that same level of protection.

An accounts payable (A/P) application can serve as an illustration. If an operational problem occurs with the A/P system during normal business hours and there are no other major problems, the application recovery focus will be on the A/P problem, with an intention to restore as quickly as possible. An RPO and RTO serve as the upper bounds of reasonableness for the recovery process. Since the problem is an operational problem, zero data loss is a feasible goal, but a recovery from tape, if necessary, could take hours.

However, a major disaster changes the rules. Enterprises have to pay their bills, but the A/P application is likely to be pushed to *the* bottom of the priority stack and could take from a few days to a week or more to restore. The chief financial officer (CFO) has a reasonable defense for paying creditors late. (However, the accounts receivable (A/R) application is likely to have a higher priority for restoration in case of a disaster, as enterprises want to maximize cash inflow and minimize cash outflow!)

Availability Is Not the Only Data Protection Objective

The driving goal is to have *data always available securely, with optimal performance, to authorized users anywhere via any connection on any device*. Availability is certainly critical to obtaining that goal, but from a data protection perspective there are really four objectives that are part of that goal and have to be met in working toward it:

- *Data preservation* — data must be consistent and accurate all the time, and also must be complete within acceptable limits.
- *Data availability* — the ability of I/O requests to reach a storage device and take the appropriate action.

- *Data responsiveness* — the ability of I/Os to deliver data to an authorized user according to measures of timeliness that are deemed appropriate for an application.
- *Data confidentiality* — data is available only to those authorized.

Note that data availability is not the same as data preservation. Not all preserved data needs to be immediately accessible. It may take a month to get some historical records back from the tape warehouse for discovery during a legal proceeding, but a month is adequate time. Not all data that needs to be accessed quickly for business intelligence needs to be preserved — in some cases, financials can be quickly reconstituted from sales and other data if the financial spreadsheet is lost.

Job one in data protection is the preservation of digital assets. RPO states what the acceptable level of data loss is. RPO should be negotiated between the user and the IT group. Quite frankly, RPO will generally be zero for most applications irrespective of RTO requirements. For examples, to return to the A/P example, most CFOs would not rate the RTO of accounts payable as being very high at all, but, while they might like to, would probably not agree to any permanent loss of data. (In a litigious society, creditors might object to accounts receivable having an RPO of zero [highly likely] and accounts payable having a non-zero RPO, which means that they might not get paid unless they complain.)

But RPO is only one measure of data preservation. RPO is what is planned, but intended and real are two different things. If a RAID 5 or mirrored array has two drive failures before the drives are rebuilt from the first failure, *all* the data is at least temporarily lost. The fallback is to the next layer of data protection, which hopefully can deliver the planned-for RTO.

What keeps IT management up at night is the worry about how much failure can be accepted and still be able to recover all the data (Table 4-3). However painful an hour-long loss of availability might be, even more painful would be an extended (weeks to forever) loss of data, which is why tape still plays a major role in the data protection investments of organizations. Tape is IT's safety net. IT organizations cannot abandon their traditional backup/restore processes until they are sure that the alternatives are equally safe. Once again, low availability is better than no availability.

Data responsiveness is also a key objective. For example, data may be available, but access might be slow — too slow and an application becomes virtually unusable. However, some degradation in performance for some period of time might be acceptable. For example, a remote mirrored array at a disaster recovery site may use slower performance disks in order to be affordable. Degraded performance might be acceptable for the time required to swap in another high performance disk array at the primary site.

Table 4-3: Consequences of Data Loss

	RPO (acceptable data loss)	All (unacceptable data loss)
Temporary (within RTO)	Acceptable	Acceptable
Long-term (over RTO, but not permanent)	Acceptable	Pain level depends upon many factors
Permanent	Acceptable	Devastation

Source: Mesabi Group, January 2005

Confidentiality is also a key objective of data protection. In fact, the Data Protection Act in the United Kingdom equates data protection with confidentiality. The Health Insurance Portability and Accountability Act (HIPAA) in the United States mandates confidentiality for health records. Authorized access must be maintained under all circumstances.

All Data Protection Objectives Have to Be Met

Although availability is critical, the other objectives require the appropriate level of attention as well (Table 4-4). No objective is an absolute — even missing or inaccurate data may be tolerable in some situations, or if some information got out to an unauthorized user, the consequences might not be severe. Still, exceptions have to be thought through very carefully.

Table 4-4: Summing Up Key Data Protection Objectives

Data Protection Objective	Observation
Confidentiality	Confidentiality has to be applied at all times.
Data Preservation	The consequences of not having all data complete and accurate have to be thought through very carefully.
Data Availability	No utility in having data that cannot be accessed.
Data Responsiveness	Slowness kills — if too slow, usefulness of information can go to zero.

Source: Mesabi Group, January 2005

Chapter Five:

Data Protection — Getting the Right Degree

A nightmare that an IT manager does not want to live through in reality is a serious data protection problem that causes major disruption to the business. That means that no manager should be willing to rely upon only one line of defense in data protection. A fallback strategy is necessary in case the first line of defense fails for whatever reason. Lines of defense can be seen as providing degrees or layers of protection. Understanding the process of setting up degrees of protection starts with an examination of the general classes of data in use.

General Use Classes of Data

There are three general classes of data from a use perspective:

1. Production data
2. Data protection data
3. Test data

Test data is actually a special case for application development activities. Although test data is useful, it really does not require data protection, as it can be regenerated from scratch if necessary and does not affect business-critical processing. This report focuses only on production data and data protection data.

The same copy of data can be both production and data protection data. For example, the same physical storage system may contain both production data and data protection data at the same time. With RAID 5 for example, a pool of production data would also have physical data protection. RAID 5 calls for an extra disk, which through extra parity checks information striped across all disks delivers the redundancy so that all the data is usable even though any one of the individual drives in the array may fail. A snapshot point-in-time copy could provide a measure of logical data protection because I/Os cannot write to the snapshot copy to alter or destroy the data. In this case (and even in the case of local mirroring [RAID 1] where there are two separate and distinct physical copies of the data), the production and data protection data are commingled. The data is still production data, since that is the purpose of the copy. Data protection is *added in* (i.e., internal or built-in) rather than *added on* (i.e., external or built-on).

In contrast, a disk array that is a remote mirror serves to provide physical data protection at a distance. In this case, data protection is

added on (in the form of a distinct and separate copy of the data) so its purpose is as data protection data.

Now if the primary array fails, the remote target disk array assumes the mantle of serving as the source of production data. However, the remote array reverts to a data protection role if another primary can be brought up or continues to serve a production role if another array is designated as a target for data protection purposes. Still the reason that the remote array was acquired and implemented was for data protection purposes.

Tape Is a Special Case

If a piece of tape media is a copy of random access disk data, then that copy is pure data protection (both physical because of the physical media and logical if no I/Os can write to it). That tape is always data protection data, unlike a remotely-mirrored disk array, the sequential nature of the tape does not permit it to serve a random access data purpose.

That is not to say that tape cannot serve production purposes. In fact, tape has a long track record of serving for production use. Tape can serve a production role where sequential processing of a whole database is necessary (such as a data mining analysis or batch update of a data warehouse when it is offline). Tape can also be used to retrieve selected files that, in effect, use disk as a cache. The original Hierarchical Storage Management (HSM) system used tape. Very large files and large numbers of infrequently-accessed smaller files may very well find a production home on tape.

Understanding Degrees of Data Protection

Data protection comes in degrees (which also can be thought of as layers). The first degree where data protection can be provided is for the primary copy. The primary copy may or may not have data protection (Figure 3-1). If it does, that is the first line of defense for operational continuity. Built-in data protection to the primary data copy can help prevent service-level threatening events (such as a single disk failure) from becoming service-level-negative-impacting events.

However, add-in data protection cannot provide disaster continuity protection and the risk-protection diversification that is necessary for operational continuity protection. At least one add-on copy — a full copy of the data that is physically separate and distinct from the original — is necessary.

Some Crossover Protection from Disaster Continuity to Operational Continuity

Each of the four boxes in the data protection category matrix (Table 2-2) is separate and distinct, but physical disaster continuity protec-

tion may also serve to benefit operational continuity. For example, say that two disks in a local RAID array fail and a failover to a remote site happens. Failure of individual disk drives is a normal outcome at times when using disk drives, not a situation requiring disaster recovery. For all practical purposes in this case, the array could have been located in the primary data center. The distance separation did not matter, but the presence of the other array matters. Thus, counting the remote array layer in both physical disaster continuity and operational disaster continuity is legitimate as long as it is recognized as only one layer.

The Limits of Data Protection Continuum Charts

Layering is important to understand when looking at a data protection continuum chart. The different versions of data protection continuity spectrum charts that are available tend to show a range of data protection solutions on the basis of availability. On one end of the spectrum, tape solutions are often shown in hours/days to recover and triple mirroring at the other end in seconds. The charts are very useful for viewing a broad spectrum of options at a glance, but they do not typically distinguish physical from logical data protection and imply that availability is the only objective.

Not too many — if any — companies are giving up their tape automation systems despite the fact that they are low availability. Low availability is better than no availability, seems to be the message.

The Third Degree — Levels of Exposure

Call one layer of added-in or added on data protection one degree of protection. One degree of data protection means that one failure is tolerable; data is recoverable. If a failure should occur, data protection is at zero degrees. Zero degrees of data protection means no more failures can be accommodated without total and permanent data loss. This is a level of exposure that IT organizations find unacceptable.

That is why additional degrees of data protection are necessary. The question is how many. The minimum number of layers is two. If one failure occurs, the degrees of protection are down to one. Given that technology is not perfect; having only one extra degree of freedom to fall back upon is not advised. So three degrees of data protection is probably a minimum. Each additional layer adds expense, but one or more additional layers may still justify the expense.

Mapping Degrees of Protection

IT administrators should map out the degrees of data protection for each application (Table 5-1). The degrees of protection have to be split between higher availability and lower availability degrees. Once the higher availability degrees are exhausted, availability depends upon the lower degree availability options. Note that the term

“lower availability” should not be considered a pejorative term, but rather reflect the relative difference between the time-based ability of different technologies to restore information.

Table 5-1, below, does not exhaust all the possible combinations and choices of technologies, but should rather be considered as simply an example.

IT administration has to determine if the levels of protection are adequate. If an application requires extreme availability (99.999% uptime or better), the higher-availability requirements would probably not be met in this example.

On the operational side, physical data protection may be adequately covered for many applications. Note that remote replication is double counted — once on the operational side and once on the disaster side. Remote replication may be synchronous or asynchronous mirroring, which can protect against physical, but not logical, failures. However, remote replication may also be remote dated replication where the replica was done as of a certain time. For example, a snapshot point-in-time copy could be made locally and then replicated (also called copying or duplication) to the remote location.

Table 5-1: Sample Degrees of Data Protection for Application *n*

	Operational Continuity	Disaster Continuity
Physical	<p>Higher Availability Degree 1: Local replication Degree 2: Remote replication</p> <p>Lower Availability Degrees 3-5: Tape/ disk-based backup</p>	<p>Higher Availability Degree 1: Remote replication</p> <p>Lower Availability Degree 2-4: Vaulted tapes</p>
Logical	<p>Higher Availability Degree 1: Point-in-time copy Degree 2: Continuous data protection</p> <p>Lower Availability Degree 3-5: Tape/disk-based backup</p>	<p>Higher Availability: Degree 1: Remote dated replication</p> <p>Lower Availability Degree 2-4: Vaulted tapes</p>

Source: Mesabi Group, January 2005

Note that once again tape plays a big role. Tape typically offers at least three degrees of failure (assuming that at least three generations of tape are produced.). The problem is that it is all lower availability. Moreover, having to use older generations of tape is likely to result in even greater time loss.

The advantage of tape, however, is that the addition of each degree of protection is mainly the cost of a piece of media for each additional generation or copy. While there may be a requirement to add more tape drives or to expand the overall tape library infrastructure, this is often administratively easier and more cost effective than adding disk arrays.

The use of a disk-based backup solution (such as a virtual tape library) may improve availability somewhat, but there is still likely to be a gap between such a solution and higher availability solutions.

On the operational side, logical data protection — if the newer continuous data protection approach is not used — does not have the same level of high availability that is available on the physical side. That could create exposure unless point-in-time copy capability is used and is managed very well for logical operational data protection.

As mentioned, Table 5-1, above, should not be used as a guideline, but rather as an example. The disaster continuity side of the house illustrates why this is the case. For example, a large enterprise may have a triad of fully-stocked data centers for disaster recovery (DR). The first DR site serves as the first line of defense, and the other DR site is in place should the first data center become the production site.

Other companies cannot economically justify the risk/reward for three sites and find that two sites — one production and one DR is fine. (Incidentally, all sites in a multiple data center environment may play both a production and DR role.) Still others cannot justify the cost of a second data center or prefer to outsource at least some of the DR requirements.

But number of sites is not the only issue. What is in a DR site is another issue. IT may have a full replica of its primary site at the remote site, may have a partial replica with the ability to bring in additional equipment when necessary (such as tape automation), or may turn to a third party service provider for recovery. All involve time and cost tradeoffs.

There are also logistical issues. Even if the remote site has a recent tape copy available, it may have only one.

Filling in the degrees of data protection for each application helps identify where the levels of data protection are satisfactory and where improvement might be necessary. But IT organizations first need to understand the impact of information lifecycle management on the requirements mix before they start filling in the boxes.

Chapter Six:
**Information Lifecycle Management
Changes the Data Protection
Technology Mix**

Understanding that information lifecycle management (ILM) is much more than the latest fashion in IT terminology is important to IT management. Along with being a process and a strategy, ILM is a new way of looking at the storage infrastructure and managing data throughout its lifecycle. That has profound implications for data protection. For example, ILM focuses attention on the management of fixed content information as being different from that of dynamic, changing information, which was the traditional view of the world. Fixed content information must be replicated as required for data protection. That replication is a one time process, so the traditional day-in-day-out backup process can be eliminated!

Although ILM has profound implications, it can start with an apparently-simple definition. The information lifecycle is the *policy-driven management* of information as it *changes value* throughout the *full range of its lifecycle* from conception to disposition.

The information lifecycle is the *policy-driven management* of information as it *changes value* throughout the *full range of its lifecycle* from conception to disposition.

What is not included in the definition of ILM, but is critical to understanding ILM, is that every piece of data becomes fixed (i.e., read-only) at some time during its lifecycle — and that time is typically short as compared to the full length of its lifecycle. Active changeable data reflects a creation and change process where viewing the data at different times would reveal that the data had not stayed the same. At some point in time, however, this change ends. Even online transaction processing systems updating customer records create data that must be “frozen” after a certain period of time, say at the end of the month or year. An e-mail is information that is fixed upon capture (as replies do not change the e-mail itself). If an IT organization looks carefully at the data managed under its custodial care, a large percentage of the data will probably be fixed.

Why Data Lifecycle Management Is Not Enough — The Need for Metadata and Management

Storage at the device level is about the management of blocks of data, so the migration of data at the block level from one tier of storage to another is often referred to as data lifecycle management. Claims have been made that data lifecycle management is enough; information lifecycle management is pretentious. That is not true.

Examine what data lifecycle management does. Action about data has to be taken on the basis of metadata (i.e., data about data). Block level metadata has to be simplistic. If actions are taken simply on the basis of age (when a block was created or captured) or if the date of last access exceeds a threshold, that migration is data lifecycle management, since there is no knowledge of the underlying content of the data (which is where the data becomes useful as information).

Migrating data in this manner might result in the use of more cost effective storage for the migrated data and ease the burden of managing the non-migrated data. However, data migration simply on the basis of last access or age does not mean that the migrated data is fixed content data. The metadata does not tell us whether a long period without change is an aberration or an indication that the data has become fixed.

The knowledge that a pool of data belongs to a certain class of data and is fixed is critical because policies for managing and allocating resources — data retention and data protection — are different from those for non-fixed content data. And to support these policies, metadata has to be at a record or file level, not a data item level, and contain information about the data's use, not merely its physical storage characteristics (i.e., data + metadata = information). Therefore, the management process is information lifecycle management; data lifecycle management simply does not cut it. And that is why — paraphrasing Voltaire — if information lifecycle management had not been invented, there would have been a need to invent it.

ILM Is Deep Into Logical Pools of Storage

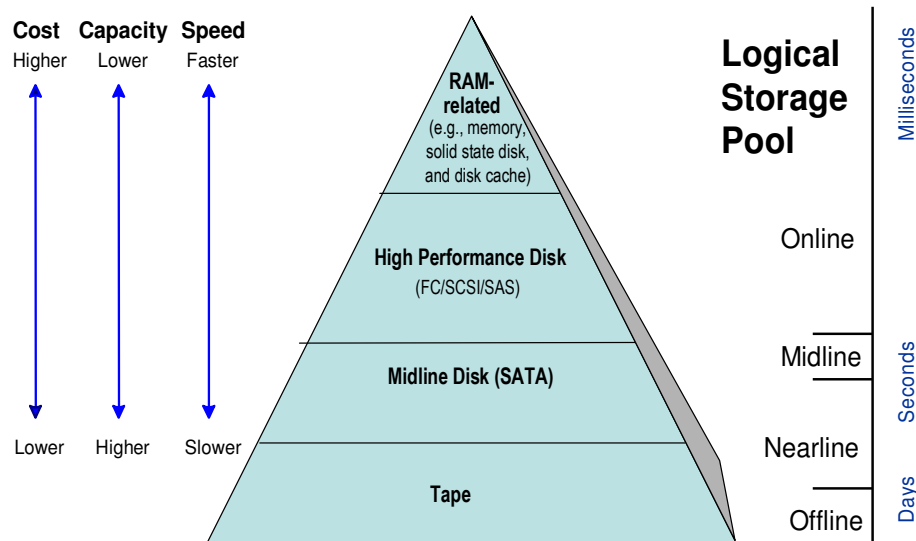
Tiering and pooling are two of the key ideas within information lifecycle management. Figure 6-1 is a simplified storage pyramid that shows the interrelationship between tiering and pooling.

Tiering is the separation of storage into classes by the characteristics of the storage itself: performance (speed and availability), functional capabilities, and cost. As such, tiering is a storage-device-related concept.

Pooling refers to a collection of information that is managed as a homogenous whole for quality of service (QoS) purposes, such as response time and availability. As such, pooling is an information-related process. The objective is to map a pool of information to a choice of storage tier, and the net result is a storage pool.

By knowing the QoS that the information pool requires, a storage administrator can map the pool to the tier that can deliver the appropriate quality of service. The mapping has to take into consideration not only cost, capacity, and speed, but also data protection requirements, such as availability.

Figure 6-1: The Storage Pyramid — Tiering and Pooling



Source: Mesabi Group, January 2005

Logical Storage Pools at a High Level

In the past, whether data was fixed or not, there were primarily only two choices for persistent storage: high-performance disk or tape. The distinction between active changeable and fixed data did not matter. Despite lower performance requirements, fixed data that still had to be available for reading on occasion might not be able to afford the loss of ready access from tape, even though tape was more cost-effective. The introduction of midline disk has changed that. Midline disks are primarily Serial Advanced Technology Attachment (SATA) disks (although parallel ATA [PATA] disk are still sometimes used). Midline disk is suitable for data that does not require the performance (both speed and reliability) of Fibre Channel (FC), SCSI, or Serial SCSI (SAS) disks. That typically includes fixed content and very slowly changing data.

The Nearline and Offline pools focus on data protection (where the sole purpose of the storage pool is on protecting data) whereas the Online and Midline pools focus on using data for production purposes (that is, any business use except data protection). The distinction between production data and data protection data — how they

can be separate and how they can be joined — is important for understanding what mix of data protection technologies can give the intended level of protection.

Moving Information Across Pools — A Distillation Process

The mapping of a pool of information to a tier of storage is an assignment process, so the process could be static and manual. The tiering and pooling assignment process has nothing to do with the movement and migration of data per se. Yet pools of information are not static; inflows to the pool of information and outflows to another pool of storage have to be taken into account. Inflows and outflows are at the “information–object” level. An information-object is the smallest information unit — perhaps a file or a record — that can be differentiated by access only to authorized users, by ownership, by compliance requirements, by identification, and/or by process control.

Migration can be viewed as a distillation process. One key distillation process is to separate information-object “molecules” that are currently active changeable (i.e., information that is likely to change in a foreseeable future) from those that are fixed content (i.e., information that is unlikely to change within the near term). The fixed content information-object “molecules” are distilled from their original storage pool into another storage pool “flask.” And that new “flask” can be called “archiving.”

Archiving Through a New Lens

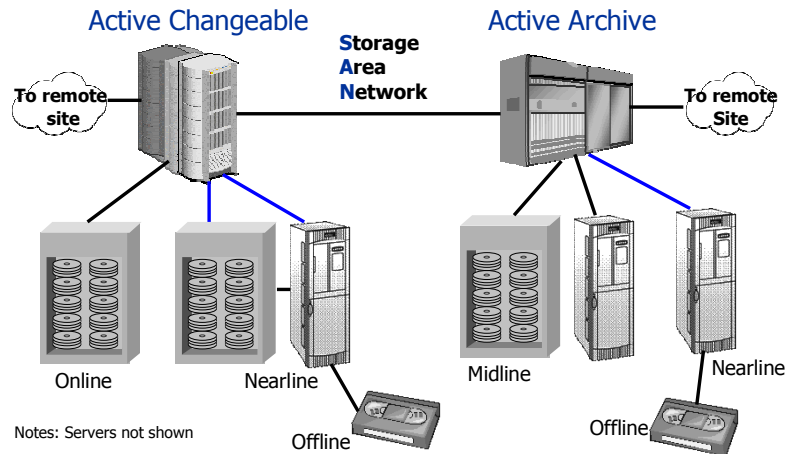
Understanding the concept of the bifurcation of production data into two separate and distinct classes — active changeable data and fixed content data — has achieved some measure of mindshare in IT organizations. A lot of words are being thrown about, such as content addressable storage (CAS), to refer to data with fixed content. These terms all have their place, but the most general term, and the one that is on its way to achieving the greatest popularity and acceptance is “archiving.” In dealing with archiving, ILM fundamentally divides the storage infrastructure into two separate halves: an active changeable side and an active archive side (Figure 6-2).

Note that this is a *logical* arrangement, even though the icons can be misinterpreted as indicating a *physical* arrangement. Mentally separating storage into logical pools is useful whether or not the different pools reside physically on one disk frame and one tape frame or whether or not they are on multiple general purpose and single-purpose appliance oriented frames.

All the logical storage pools are contained in the topology shown in Figure 6-2. Online and Midline pools of storage are for production data (with data protection added-in). Nearline retains its data protection origins from tape automation solutions, but also adds in

Nearline disk (which, while gaining acceptance, is still not typical today). Offline remains data protection data.

Figure 6-2: ILM Changes the Logical Topology Storage Look



Source: Mesabi Group, January 2005

The fundamental changes from the typical data protection configuration today (Figure 3-2) are really on the archive side. The importance of midline disk is not only lower cost and a higher capacity with a performance tradeoff, but also that the management of fixed data through an archive is fundamentally different from the management of active changeable data, even though the two have to work together and exchange information with each other.

For example, the focus on data retention in archiving is much greater than on the active changeable side. On the active changeable side, intermediate changes may (but not always) be discarded because work-in-process does not represent the committed transaction or a usable version of a document. On the archive side, data, whether database or file data, may reach end-of-life and go through a destruction process, but that process should be managed and not ad hoc.

Archiving: The Makeover

Much confusion exists about what archiving really is and what it means to IT organizations. One reason is that the definition of archiving is undergoing a transformation (and really was not well understood anyway!). A starting point for examining archiving is the original SNIA definition of archive:

Archive: A consistent copy of a collection of data, usually taken for the purpose of a business or application state. Archives are normally used for auditing or analysis rather than for application recovery. After files are archived, online

copies of them are typically deleted and must be restored by explicit action.

This is the “classic” definition of archive. A subset of production data has reached the end of its “production” life, but the enterprise has a need to preserve for the data for an additional period of time (up to posterity). One or more tape copies are made and the production copy deleted. The tape(s) are taken to an archive in an underground salt mine (or equivalent) somewhere with the fervent hope that the data never needs to see the light of day again. This is equivalent to data death through suspended animation. The data is not really dead, but it is not really alive either.

That original definition corresponds to the dictionary definition of archive as a place where documents of public or historical interest are preserved.

The new SNIA definition of archive is:

Archive (noun): 1. (noun) a collection of data that is maintained as a long-term record of a business, application, or information state. Archives are typically kept for auditing, regulatory, analysis or reference purposes rather than for application or data recovery. 2. (verb) To copy or move data for purposes of retention; to create an archive.

That new definition reflects the changing times. The reason for the new definition is that archives may have some real “production” purpose instead of just acting as a data mausoleum. There is a lot of read-only data that still can serve real business purposes, such as analysis, reference, and governance.

The makeover has a significant impact. No longer is data, for all realistic purposes, “dead.” Instead, most archived data serves a real business purpose. If archived, “not dead yet” data were small in size, there might not be an impact on data protection, but, these days, fixed content data may be the majority of the data in most enterprises, and the impact on data-protection architecture design is substantial.

Protecting Archived Data

Archived data is fixed data (i.e., the data does not change). Archived data must be protected data; the proper numbers of copies in the proper number of locations have to be prepared. However, archived data does not have to undergo a regular process of backup. Backup is a recurrent cyclical process; if full backups are run weekly, a file that had not changed in a year would be backed up 52 times even although only a limited number of copies of the file are available at any one time.

The new process involved in archiving fixed data is replication. All the necessary copies can be produced through a replication process at the time of capture in the archive.

Although an individual file or record does not change, a fixed-data archive storage pool is not static in its contents. An archive has both inflows and outflows. Inflows are simply additive — a new piece of fixed data has been added to the archive. Outflows are decremental; data is removed from the archive. If the data is migrated to another piece of storage media, the data is preserved. If the data is not migrated to another piece of media, the archive copy of the data is destroyed. (To be truly destroyed, all other copies would also have to be destroyed.)

The actual archive storage pool is a production pool of storage, and authorized users can access information in the pool for production purposes. The production pool can use data protection technologies for both physical and logical protection, but one or more data protection copies need to be made as well.

Data Retention

When production data was held in only one pool, data retention was the responsibility of storage administrators and storage administration tools for that pool. With the separation into two production pools — an active changeable pool and an archive pool — the responsibility shifts to the archive pool. The active changeable pool does not have to worry about data retention. The active changeable pool saves data, but it does not destroy data. When data moves into the next stage of the lifecycle as fixed data and migrates to an archive pool, data retention becomes an issue.

Data retention now attracts more attention because compliance is a subset of data retention and building and managing an effective compliance process is an important issue for many enterprises. But data retention is also about data destruction. An archive is a place of long-term storage, but long-term does not mean that data lives forever.

Disposition of Data

One of the basic rules of thumb of ILM is the principle of accumulation — enterprises tend to add data to their storage pools faster than they dispose of data. That can lead to problems from increasing the burden upon the data protection processes to overspending on physical storage to overworking of storage administrators. A key focus of the data retention process therefore should be data disposal.

Data disposal may be complex in two ways — the first involves defining policies for data disposal and the second involves actually making sure that all copies of data that is to be disposed is actually destroyed.

Disposal policy can no longer be “laissez faire,” where the original creator (or recipient) of the data may have been entitled to destroy or alter the data (such as certain classes of e-mails) without interference

by administrators. Approval of a disposal policy must take into consideration the inputs from the chief legal counsel, senior management, and business unit management. IT organizations are custodians of data, but IT must act upon direction from the owners of the data — and for this purpose the owners are those who speak officially for the enterprise.

However, IT management is not just an order taker, but a partner in this effort. IT management has the knowledge and experience to put together recommendations that form the context of the discussions. For example, the disposal policy probably should not be benign neglect, where nothing is destroyed — ever. This raises manageability, cost, and, perhaps service-related issues. The recommendations of IT management should contain an impact analysis that discusses the advantages and disadvantages both quantitatively and qualitatively of different choices.

Policy-setting is at the front end of the disposal process; the actual data destruction is at the back end of the disposal process. This presents a dilemma because if all copies of the data are to be destroyed, including data-protection-only copies, the IT organization may lack knowledge or ability to find and destroy copies outside of its control. IT administration can set policies for deletion of the production copy as well as the data protection copies that are within the four-IT “walls” of the enterprise, but cannot control those copies that go to either internal non-IT users or external-to-the-enterprise users. This dilemma is especially acute when discussing the management of compliance data.

Compliance

Compliance data is a subset of data retention management that must adhere to special conditions. In effect, compliance is data retention on steroids. For example, a compliant set of data has at least logical data protection built-in, but the purpose is still to serve a business use. (Data protection in this case is not optional; the data protection has several special conditions built-in to it, including a chain-of-custody process, serious controls on authorization of access, and inalterability of the data.)

Compliance, from a data protection perspective, is putting into place the policies, practices, and procedures to preserve in an unalterable state and to safeguard the confidentiality of data that fall into designated classes for a prescribed (sometimes indefinite) period of time.

Compliance is typically seen as the necessary mandatory response to an authorized third party, such as a government regulator. But it can also serve a number of other purposes as well:

- A voluntary response to a trade association or industry body, in order to increase market size through the use of common practices that make it easier for customers to work with the industry as opposed to a substitute industry.

- A intentional response to protect an enterprise against lawsuits.
- A voluntary good practices response to protect intellectual assets (e.g., patents or trade secrets).
- An effective way to leverage the enterprise's compliance policies, practices, and procedures to provide logical data protection to fixed data (for example, if e-mail has to be saved for two years anyway, why not fossilize it right now?).

Compliance should not be a reactive strategy, only used by enterprises that are subject to mandatory regulation, but rather a proactive strategy that all organizations develop as part of an overall data retention strategy.

Of course, enterprises that are struggling to figure out how to comply with multiple evolving, complex, and apparently conflicting compliance requirements face a challenge to which there are no easy answers, but certain basics apply even to them.

The first issue — chain-of-custody — is both a process and a technology issue. Chain-of-custody is a legal concept that relates to the handling of evidence — in this case, data. Every transaction between the time of creation or capture of data must be completely documented so as to avoid later allegations of tampering or misconduct.

In so far as possible, an IT organization needs to put in place an automated process for controlling all transactions (in the case of data IO requests) to ensure compliance with the chain-of-custody requirement.

Creating Data Archive Storage Pools by Data Retention Attributes

Logically, an archive can be broken into different storage pools by many characteristics — one of the most notable is by data retention characteristics (Figure 6-3).

At first glance, the chart seems shocking; fixed content by definition should mean that the information should be unchangeable, i.e., immutable (and many discussions of archiving affirm the immutability requirement). Unfortunately, that is not quite the case. By and in itself, data is not immutable as long as there is an application change process that is allowed to update the data. An archive management process may close down the ability of an application to change an information object once it is in an archive, but that constraint might not need to be a universal criterion for all information objects.

Figure 6-3: Data Retention Archive Pools

Compliance Not Necessary		Compliance Necessary	
<p>Slowly Changing</p> <p>Any change may result in the replacement of (or appending to) the original information object</p>	<p>Practically Immutable</p> <p>Any change results in an additional and new unique information object.</p>	<p>Guaranteed Immutable</p> <p>Any change results in risk contamination with a new unique information object.</p>	<p>Guaranteed Immutable, but can be Appended</p> <p>Need to be able to handle new information.</p>

Source: Mesabi Group, January 2005

The four basic I/O processes are create, read, update, and delete (a very old acronym for this was CRUD). The create I/O process in an archive that becomes a capture process where information is brought into the archive through a migration process. Read I/Os are allowed to authorized users and often what can be read can be modified even if no update processes may be allowed on information objects already in an archive. Delete I/Os, i.e., destruction of an information object, should be permitted to an authorized user only under the control of an archive management process.

Data in the form of information objects are placed in an archive because there are voluntary reasons for doing so (and therefore compliance is not necessary) or there are mandates for doing so (and therefore compliance is necessary).

Voluntary migration of data to an archive may be useful for data retention, data protection, and storage asset utilization reasons. But some voluntarily-migrated data might be slowly changing. For example, a thirty-year mortgage may be paid off early, or a life insurance policy paid off. These changes are quite infrequent and straightforward, but need to be accounted for. An archiving process that can handle this may be important. The creation of a new, unique information object that reflects changes while keeping the original information object may be a satisfactory solution if additional storage requirements are minimal and if the user application points to the correct version of the information object. Alternatively, the original object could be deleted, since the information object is not subject to compliance requirements.

In the “practically immutable” case, any unanticipated revisions to the original information object would result in both the original and revised versions being stored in the archive.

Information objects that are subject to compliance requirements have greater constraints placed upon them. In the “guaranteed immutable”

pool, information objects cannot be changed at all, but they also cannot be deleted without permission of the archive manager. A potential problem does exist however. A read copy of the information object might be changed without the consent of the archive manager and put in a separate archive as if it were a data protection copy. Then, if the original copy was physically destroyed, the compromised copy could be misconstrued for the original. Although the risk of such contamination is small, the archive manager must have policies in place to deal with this potential type of problem.

“Guaranteed immutable, but can be appended” data is an important data class. Consider medical records, for example. They are subject to privacy and confidentiality requirements, but also to inalterability requirements. As such, they have to be considered fixed content, but, by their nature, they may not only be slowly changing, but in some cases rapidly changing. What has happened is fixed, but there is a need to accommodate change. That can be done through appending new information to the old information. However, appending to existing information and retaining the original may lead to a lot of wasted disk space. Managing this process well will be a requirement of a solution where adding new information over time is a necessity.

Security and Privacy

Data protection and data security are not synonymous terms yet they are not independent terms either. They intersect at both the data preservation and data confidentiality objectives for data protection. For example, data preservation requires prevention of data corruption. A virus is an example of data corruption and an anti-virus program — a security measure — is used to prevent a virus from acting.

Security can take many forms, such as restricting physical access to a data center, but the principal focus of data security is network security. The ability to work online now extends from internal LANs to the ability to access information via the Internet anywhere in the world. The price for the network giving online omniscience is the increased threat of data corruption, information theft, and invasion of privacy. Ensuring data integrity, access only to authenticated users, and adherence to confidentiality and privacy are challenges for data security in support of data protection.

The Trusted Computing Group (TCG) is a vendor-driven standards body that is acting on this challenge. The TCG has a vendor-heavyweight board whose members are IBM, HP, Microsoft, Sony, Sun, Intel, and AMD. Although the challenge is a difficult one, the ability of these companies to “de facto” adopt a standard bodes well for the group. One of their key strategies is to ensure hardware that can be trusted, by inserting a Trusted Platform Module (TPM) chip in the hardware. TCG builds around the use of encryption, which (although it introduces all the problems associated with managing the encryption keys) will be a necessity.

Storage vendors, such as Seagate, are among the companies that are working with TCG. One task is to create a more standardized storage architecture for storage resources in order to put security policies in place that can map down to storage devices. The ability to move data (such as backup data) across a network more securely than today is another key goal.

Providing more solid authentication and stronger data integrity are essential, but do not solve the privacy issue. Data can be easily copied and redistributed improperly. Although no technology is foolproof, work is also being done to improve confidentiality. For example, a user who is authorized to view data only on a monitor will not be allowed to print or retransmit the information.

Enterprises should recognize that, although technology helps to reach data security goals, they still have to put in place the necessary policies, practices, and procedures that limit and restrict access to the information.

Active Archiving and Deep Archiving

The consensus term for the new type of archive is *active archive*. The term for the original archive is now *deep archive*. The distinction is very clear. Dead information that is deep archived is stored on removable media, removed from online files, and transported to an offsite spot for safekeeping. An active archive contains production data, no matter how old or infrequently accessed, that can still be retrieved online.

A deep archive is an electronic landfill, in that the data is likely to never be retrieved. That does not relieve an enterprise from doing the proper media management to ensure that it knows where the data is and how to get it back. However, the first question is why any enterprise would want to deep archive data. If the data has reached the end of its useful life, destruction of the data — rather than deep archiving — should be the choice. If the data still retains some value — such as for regulatory compliance — keeping it in an active archive and sending a replica offline might be a better option. Only if the data must be preserved against significant future risks (e.g., changes in regulations), the volume of data is large, and a long restoration time acceptable, should information be deep archived.

Active Archiving Requires Active Archive Management

Active archiving requires archive management — the umbrella term for managing data retention and data management policies. Active changeable data does not require a manager. Active changeable data is application-controlled; the creation, updating, reading, and deletion functions are all under the aegis of the controlling application, such as an e-mail application. An archive is controlled by a special archive manager application. That application may grant privileges to the originating application, but it does not have to.

Therefore, the archive manager is a key to the best management of an active archive. The production application is no longer always king of the I/O; a storage management application has to kowtow to Microsoft Exchange or Oracle, but an archive manager does not have to. The production application could incorporate the functionality of the archive manager, but the new functionality from the archive manager would still have the overall control.

The reason is simple. Data protection (which includes data retention) is paramount in an archive; the I/O functions are tightly controlled, which was not the case in the production application with active changeable data.

Metadata management is critical for the archive manager. The metadata (other than the basic information of name, size, and date) can be embraced in the following categories:

- *Ownership* — creator, owner, last update, organization, application
- *Access Control* — security clearance; access control list (ACL); browse, read, and write privileges
- *Compliance* — retention policy, earliest deletion date, who has authorization to delete, etc.
- *Identification* — version, identification codes, relationship to other objects
- *Process Control* — workflow information, including approval process

Many IT managers might be concerned about taking on the additional management responsibilities that an active archive requires, so reviewing the benefits would be useful.

Moving data from active changeable disks to the archive side of the IT infrastructure house obviously reduces the amount of data that has to be managed on the active changeable side. This has several benefits for managing the storage on the active changeable side:

- Less data to be backed up on a regular basis minimizes the time needed for backup (and thereby improves either data availability or data-access performance), reduces the need for storage assets to back up the data, and reduces the chances of problems (such as the failure of a backup job to complete or tape media failures).
- Reduces the time required to restore the data from tape should that become necessary.
- Shortens the time that it takes to run a query that spans the entire database.
- Cuts the investment required for high performance disk.

The tradeoff is between these benefits and the additional management burden. If managing fixed content today as inseparable from active changeable data is not a major burden, IT managers may want to delay moving to an active archive and therefore delay taking on the extra burden. However, an honest examination of the whole information infrastructure may very well reveal that the benefits for moving to an active archive in an evolutionary manner (say, by starting with a selected application) will far outweigh the additional responsibilities.

Long-term Archiving as Part of an Active Archive

The adjective “active” in active archive implies some reasonable frequency of access; however much of the data in a large active archive may never (or very infrequently) be accessed. Active also means that the data is online, which implies a “reasonable” response time; however, a reasonable response time to a compliance request may be a day (or more) and not the sub-second or a-few-seconds response time that might be more typical when working “online.”

Those two parameters — expected frequency of access and acceptable response time — can lead to the placement of data in different storage pools, e.g., information pools on different storage tiers.

The age of the data and how long it is expected to be kept is not an issue for that placement. For example, one county keeps well over one hundred years of property deeds online with response time in seconds. That information is likely to be kept as long as the county continues to function.

Long-term archiving raises the specter of technological obsolescence. The active archives themselves are not subject to that problem for the simple reason that they have to be online and therefore the data has to be migrated to storage devices that can be accessed online. The problem could apply to the data protection replicas of the data, however, or to deep archives. The long-term obsolescence problem is being thought about, but is not a major issue of concern today when setting up an active archive.

ILM Changes the Data Protection Technology Mix

The addition of an active archive for fixed content information changes the data protection category matrix (Table 6-1). The reason is that some of the data protection strategies for active archiving are different for active archived information than for active changeable information, such as not necessarily requiring the use of backup/restore software, but rather making dated replicas of the data.

As another example, an application may have data on both the active changeable and the active archive side of the house. That might mean that the RPO and RTO for each side would be different. For example, active transactions in OLTP may require a different (and

probably more stringent) RPO and RTO than closed transactions that are retained for business intelligence purposes.

Table 6-1: Adding In Archiving to the Data Protection Category Matrix

	Operational Continuity		Disaster Continuity	
	Active Changeable	Active Archive	Active Changeable	Active Archive
Physical				
Logical				

Source: Mesabi Group, January 2005

The finer granularity that is expressed in the doubling of the cells in the matrix requires more work on the part of an IT administrator to fill out, but also permits the design of more effective data protection strategies.

Chapter Seven:

Where Data Protection Technologies Play in the New Model

In the preceding chapters, we identified eight data protection categories and developed a framework for characterizing data protection technologies according to the category under which they fall. In order to start deriving practical value from the framework, we will now fill in the eight boxes in the framework with data protection technologies (Table 7-1). Many of the technologies (such as RAID and tape automation) will be familiar, but others (such as continuous data protection) may not be.

Data protection technologies are not always purely for one task or function; there may be a lot of blending, morphing, blurring, and variations in the data protection functionality that any individual product may contain. The focus is therefore on overall technologies and not specific products. IT buyers should “extract the essence” in terms of what function is being performed, where the technology will fit in the framework, and what it adds to the overall degree of protection. Individual products can then be evaluated offline in terms of how they fit one or more of the organization’s data protection needs.

The overview sections that follow for the data protection technologies should serve as a logical sequence for understanding each technology—both standalone and in context with other data protection technologies. Data protection technologies are divided into five large categories for this purpose, as follows:

- *Back to basics* — backup/restore software and RAID are well-known, but what is the impact as today’s backup/restore and RAID technologies morph and change?
- *Disk and tape* as complements to and competitors with one another — the role of disk-based data protection is a current hot topic in data protection. Understanding the interrelationships between disk and tape data protection solutions can yield a better understanding of where each best plays a role.
- *“Getting to the point”* — point-in-time copy capability and its derivatives will play an increasingly important role in logical data protection.
- *Replication strategies* — understanding the difference between replication for physical data protection (e.g., remote mirroring) and replication for logical data protection (e.g., dated-replication that creates time stamped copies) is essen-

tial for matching up the right data protection technology with the right need.

- Special requirements for *compliance* — the additional requirements that compliance puts on data protection technology demand special technologies that organizations have to put in place.

The data protection technologies can then be integrated into a checklist for the production copy and the data protection copies of the data.

Table 7-1: Where Data Protection Technologies Fit in the Data Protection Framework

	Operational Continuity		Disaster Continuity	
	Active Changeable	Active Archive	Active Changeable	Active Archive
Physical	RAID Cloned Point-in-Time Copy Tape Automation* Virtual Tape Library* Continuous Data Protection Data Protection Appliance*	RAID Dated Replication WORM Tape	Synchronous Remote Mirroring Asynchronous Remote Mirroring Dated Replication Vaulting* Electronic Vaulting*	Dated replication Vaulting
Logical	Point-in-Time copy Tape Automation* Virtual Tape library* Continuous Data Protection Data Protection Appliance*	WORM Disk Guaranteed Uniqueness Electronic Locking Dated Replication WORM Tape Compliance Appliance	Dated Replication Vaulting* Electronic Vaulting*	Dated Replication Vaulting

*Backup/restore software is or might be used in conjunction with this technology.

Source: Mesabi Group, January 2005

Back to Basics — Extending the Current Model

Technological change comes about not only by introducing new technologies and services, but also by modifying and morphing current technologies. Moreover, some things may very well stay the same while change comes about in other areas. A good place to start is with current RAID capabilities.

Current RAID Capabilities Are Not Enough

As we noted earlier, current RAID technologies, while quite good, offer unpalatable choices: either hoping that a rebuild of a failed drive will complete before a second disk drive in an array fails or investing in a costly extra mirrored array. Given that the wrong disk drive can be pulled from an array and cause an unexpected second failure, and that the disks in an array may be from the same batch of disks (and is thus more likely to suffer from the same problem that led to the first failure), the organization's comfort level after a single drive failure should not be too high.

A concept called RAID 6 has been discussed in technical circles. The two approaches to RAID 6 — one based upon what is called a Reed-Solomon algorithm and another based upon dual parity — have reportedly not been introduced generally because of efficiency issues.

Typical RAID (excluding mirroring) is based upon a single-parity protection scheme. Multiple-parity RAID would require the equivalent of multiple drives assigned for the parity function. With one company already offering multiple-parity RAID and with at least two major disk array vendors working on the technology, multiple-parity RAID is not an impossibility — and it would be less costly than another mirrored array. Making multiple parity RAID happen faster may result from IT organizations making their desires in favor of multiple-parity RAID known.

Evolving Backup/Restore Software

The one data protection technology common to almost all organizations is the use of backup/restore software to back up and restore sets of data. A backup is what SNIA calls a “dated duplication” of a set of data. Backup/restore software copies the designated set of data from the data source system to a backup target — magnetic tape, magnetic disk, or optical media. Backup is run at discrete intervals. Full backups are a complete copy of a source data set. Incremental backups are backups of only the changes that have been made to the data set since the last backup operation (either full or incremental) was run.

Data on a backup copy is typically not “naturally usable,” which means that the creating application cannot use the data (even on disk) until it has been restored (i.e., put on storage media where the creating application can use it). The restoration requires using a spe-

cial process and software tool (which is typically the backup/restore software that created the backup in the first place). That means that the backup copy is not a replica — the data may be identical, but the format is not. Thus, the copy can only be restored, not restarted as a full replica could be. Failover to a full replica for restarting an application is fast; restoring data to storage media first before an application can be restarted takes more time (and sometimes much more time). Data protection solutions that use backup/restore software therefore inherently have “lower” availability.

Since a backup is typically done once a day, the inherent RPO of a backup copy is one day of potential data loss. RPO therefore is a parameter of the backup scheduling process and not a parameter that an IT administrator sets. Since an RPO of that magnitude may well be unacceptable, IT administrators have to apply other techniques, such as journaling, or employ other data protection technologies, such as remote mirroring, to bring RPO within acceptable levels. Putting multiple degrees of data protection in place, and understanding the implications of what a fallback to each lower degree entails, is essential in the data protection planning process.

The backup/restore process has an implicit assumption that all data changes at some point. The standard backup/restore process is therefore not “fixed-content friendly.” Even though daily incremental backups would not back up unchanged data since the last full backup, each full backup writes out the fixed content information even though it is not going to change.

To get around this problem, a concept called *synthetic full backups* is gaining traction. Changes are only applied (on a “change forever” basis) to the original copy of the data set. If a full copy of the data is required, then the synthetic full backup fulfills the role of a standard full backup. A note of caution is that you have to be careful not to wind up with just one backup instead of a series (i.e., one instead of multiple layers of data protection). The ability to dive deeper into the past is important in case the original backup fails for any reason.

Despite its limitations on availability, RPO, and supporting fixed content data, backup/restore software is still the backbone of most data protection strategies and is likely to remain that way for the following reasons:

- No viable alternative (except for a few emerging solutions) exists that can justify completely doing away with backup/restore processing and still maintain sufficient degrees of data protection.
- Putting in place new policies, procedures, practices, and products is typically very difficult, no matter what the benefits.

- Data protection is so critical that IT organizations must think carefully through the risks of making any significant changes.

A few backup/restore vendors are therefore starting to recast their software in a broader context, acting as the *control center* for the entire data protection process. This strategy will bring these vendors into conflict with other data protection software suppliers because they will be intruding upon the snapshot, archiving, and storage resource management offerings of these other companies.

However, offering a single management interface for an umbrella architectural platform that offers a broad sweep of functionality may be easier for IT organizations to swallow. That way, IT organizations can accept selected functionality as they become comfortable with the functionality and add more functionality as time goes on without having to learn a new interface. The approach also offers the possibility of things such as active archive management, which enables effective “replication management” of fixed content information.

Backup/restore software vendors typically have focused on the scheduling aspects of the backup/restore process, not on management of the overall backup/restore environment. An IT administrator sets in motion one or more backup jobs that may or may not run to completion. Broader management reports for the backup/restore environment, to display information that can lead to a root cause analysis for chronic problems, may not be in the purview of the backup/restore software. For example, backup/restore software may not report information that is important for managing the overall backup/restore environment, such as chronic network congestion and if all critical information on backup clients has been successfully completed on a regular basis. And that leads to the need for better management reporting and automation.

Better Data Protection through Better Management Reporting and Automation

The data protection technologies that have been discussed (RAID and backup/restore software), as well as most of the ones in succeeding sections, are ones that act in a direct fashion; that is, they either protect the production copy better (e.g., RAID) or create a data protection copy (e.g., the backup process). Yet there are management reporting and automation software tools that can play an indirect — but vital — role in data protection. They add value by reducing the chance for human error in very complex environments (such as backup/restore or a storage area network [SAN]) where error can lead to negative SLA-impacting events. For example, it’s better to not have to restore at all than have to worry about how fast you can restore. Preventing a disease is better than having to cure it.

Keeping the Production Site and Disaster Recovery Site Synchronized

From a storage perspective, changes in the storage configuration of a SAN at the primary site must be reflected at the remote disaster site. The user should put a change management process in place that ensures changes at one site are validated and synchronized at the other site. Remember that, if a disaster should occur, the staff that is in place to handle the disaster may be different from the staff that managed the original production site. They will be under a tremendous amount of emotional stress and time pressure. Disaster recovery time is not the right time for them to have to figure out what should have been done to synchronize the storage requirements at both sites. That is the responsibility of the original production site staff.

Production-site staff can do site-storage synchronization manually or with the aid of a software tool, but keeping configurations up-to-date through a change management process is essential. And with the complexity of modern SANs and the criticality of the data, use of a software tool may be essential.

Improving the Backup/Restore Process

The time that IT organizations spend on improving their backup/restore processes probably rivals the time that millions of individuals spend on improving their golf game — and with the same result. Although some improvement may be possible, a limit on improvement is quickly approached, to the point where further investment in improvement is cost-ineffective.

Given the number of applications and a complex IT infrastructure, the problems associated with getting the backup/restore process to work effectively are enormous, including:

- Not having time to run a backup job
- Failure to schedule backup jobs to protect all the data that needs to be protected
- Either not noticing that backup jobs failed or not being able to take corrective actions to rerun the jobs
- Not noticing that the output from the backup job that appears to have succeeded is either too much or too little compared to what it should be
- Not noticing that the file system is going to run out of space shortly if no corrective action is taken
- Not detecting mechanical or other physical errors that will cause a restore to fail, *before* attempting the restore

No technological *deus ex machina* exists to solve all these problems, but non-backup/restore software for monitoring and reporting is available that can help. IT organizations may expect their own

backup/restore software to help with these problems (and in some cases some assistance is available). However, typically, solving these problems is not the focus of the backup/restore software, and so IT organizations may want to examine other software products that can do monitoring and reporting to provide the information that can help address these problems.

Moving Data Manually and Electronically — the Place of Vaulting

The old saw “the more things change, the more they stay the same” applies to vaulting. For both operational and disaster recovery reasons, tape cartridges that are exported (which means that they are physically removed) from tape libraries have to be removed to a remote site. The transportation and storage of these pieces of tape media could be done within the resources of an enterprise, but typically they are outsourced to a third party for which vaulting is a core competency.

In an electronic age, physical transportation of a logical commodity — data — seems somewhat of an anachronism. However, the information packaged in a physically-small (about 4” x 4” x 1”) tape cartridge is already hundreds of gigabytes native (which means uncompressed) and trying to move that amount of data regularly could be a network bandwidth or cost challenge.

Nevertheless, electronic vaulting is making inroads. SMB enterprises and branches of larger companies are among those that are able to take advantage of electronic vaulting, and especially because the incremental amount of their data that needs to be transferred each day is small. And with the increased adoption of synthetic full backups, large IT organizations may be able to take advantage of them for a wider range of purposes. Moreover, electronic vaulting does not just have to be for traditional backing up of data. Instant replication of changes to an active archive could be a potential use as well.

At Your Service — the Role of Service Suppliers

Data protection services have been around for a long time. Tape vaulting and recovery services are two examples. Although still small, the number of professional services organizations that are dedicated to storage in general — and data protection in particular — are starting to grow. Of course, the large professional services organizations have skills in data protection (such as the design of disaster recovery sites) and so do storage vendors (such as planning and implementation of their data-protection-related products).

The number of professional service opportunities for consulting, integration, project management, and knowledge transfer seem endless. Among them are site assessments, architectural planning, product and technology selection, project planning, installation, training, and troubleshooting.

Although the extended data protection category matrix is conceptually simple, applying the principles across a broad application portfolio, complex IT infrastructure, and a large number of existing and emerging data protection technologies can be a daunting task. No surprise, then, that many IT organizations are turning to third party help for expertise that is not within their current skill sets.

Disk and Tape — Complementing and Competing with One Another

The “flashpoint” for how disk and tape complement and compete with each other is in how they interact with traditional backup/restore software. Disk and tape can interact with the traditional backup/restore process in a number of ways.

- *Disk-based backup* substitutes disk for tape as the target for a backup or data restoration job. Although the possible consolidation of tape automation systems as a result means that disk and tape compete with each other, the fact that the backup jobs on disk have to be copied to tape also means that disk and tape complement each other.
- *Virtual tape*, which is often confused with a virtual tape library, is about the more efficient use of tape with the help of a front-end disk cache. Disk is a necessary additive to the solution, so tape and disk are strictly complementary.
- A *virtual tape library* is a particular disk-based backup strategy, so disk and tape have the same relationship as with the overall disk-based backup strategy.
- *Massive Array of Idle Disks (MAID) and removable disk drives and media* are strictly competitive with tape as they represent direct replacements.
- *Data protection appliances* are dedicated hardware/software combinations whose sole purpose is data protection. The appliance may be a disk array (which means possible competition with tape) or an integrated disk array and tape automation solution (which indicates that disk and tape complement each other).
- *Tape automation* is the traditional starting point for discussing backup/restore applications. If there are any substitution effects between disk and tape, disk tends to replace tape, so what the role of tape will be after all the assaults from disk is the key question.

Disk-based Backup

Using disk-based backup in conjunction with traditional backup software, backup jobs copy data to a disk array for data protection rather than to a tape drive. A set of disk drives has to be reserved for this process. The cost of the disk system as well as any software that

is necessary to process the data is an incremental cost to an IT organization since the existing tape automation infrastructure is typically not replaced. Moreover, some change in operational procedures as well as retraining of staff might be necessary. The question then may very well be why so many IT organizations are so interested in disk-based backup.

The answer lies in two words — reliability and speed. Reliability refers to improving the reliability of the data restoration process, and speed refers to shortening the length of time that a backup or data restoration job takes.

Speeding up the Backup/Restore Process — Your Mileage May Vary

A key justification for inserting disk as an additional layer in the backup/restore process is to reduce the time to create a backup copy and the time to restore a given set of data. Although there are ways to do backups at any time (such as from a point-in-time copy of the data), many backup jobs are still run after a production application has been shut down at night. The problem is that the ever-increasing amount of data with which many enterprises have to deal takes longer to back up, but the hours in a night have not changed. This is the “running out of night” (a.k.a. “shrinking backup window”) problem. On the restore side, improving the time to restore data in order to meet quality of service objectives may be equally important.

So how much faster is using disk instead of tape? A publicly-available major-storage-vendor claims to cut backup time by 30 to 60% using a single process, and cut restore time by 90%, using a virtual tape library (VTL), one of the two approaches for disk-based backup. (A VTL is generally regarded as having better performance characteristics than straight disk-based backup.) Another reputable VTL vendor claims a doubling of both backup (which is a write-only process) and data restoration (which is a read-only process). The results held true for both large and small files. (The sources for these claims are available upon request.)

These results are useful for IT organizations trying to relieve the pressure of a shrinking backup window problem for, say, one to three years. On the restoration side, disk can help with partial restores where only selected files (down to individual files) have to be restored. And since partial restores are much more common than full restores, this can be very helpful indeed.

However, while disk-based backup and restore lead to *higher* availability, the result is still low availability, not high availability. Hours may be cut in half, but hours are not minutes or seconds. Keep in mind that tape emulation means that the disk cannot be used natively to run an application. The data has to be copied from one set of disks to another set of disks. That process is a restore, not a

restart. IT organizations have to set their expectations for disk-based backup and restore accordingly.

Improving Restore Reliability

Another justification for disk-based backup that many IT organizations offer is improving the reliability of restore. Many IT organizations have a concern that the potential failure rate on data restorations with their current tape automation infrastructure is higher than they find acceptable. Physically, a RAID-protected disk array can survive the failure of a single disk without loss of data, so the array has a much greater MTBF than a single disk. Generally, tape does not have this advantage (a mirroring technique for tape does not seem to have attracted a great following), as each piece of tape media has to stand on its own MTBF. (Although the aggregated set of tapes through using multiple generations of tape is greater than an individual tape, each deeper dive into the past to restore from tape takes an additional amount of time.)

Keep in mind the caveat that the data must be available on disk (not staged off to tape) for the restoration process. Disk space should be able to accommodate, say, a weekly full backup as well as all the daily incremental backups for a week. This should suffice for most circumstances.

Moreover, recall that if the backup copy was not created in the first place, no data restoration process can take place. Note though that not creating the backup may be a process problem and not a physical problem. This process problem can occur for a number of reasons, including scheduling errors, failure of backup jobs to run to completion (e.g., due to network congestion), and failure to notice that not all the critical data is being backed up. Disk-based restorations today cannot rectify policy and process errors.

Keep in Mind

Note that the disk-based backup and restore process is for operational continuity, not disaster continuity. For operational continuity, the disk-based process adds a layer of physical *and* logical data protection.

Eventually data has to be moved from disk to tape. There are two primary methods for doing this. The first is to send the data from disk through the server that ran the backup job in the first place (the server can be called the media server or the backup server) and then to tape. The reason for doing this is that the media server can keep track of where the data is so that the media server (which is also responsible for the restoration process) can restore data from either disk or tape.

The second approach is to allow the VTL to write directly to tape as a secondary media manager without going through the original me-

dia server. If this happens, a risk exists that the original media server would not be able to restore from tape. (Restoring to the backup disk first and then to the target array for the restoration would add an unnecessary time-consuming step.) Vendors get around this potential problem by using the common technique of writing barcodes on pieces of tape media that can be read by the original media server.

Virtual Tape

The terms “virtual tape” and “virtual tape library” are frequently bandied about as if they were the same, and that can cause confusion, because they are separate and distinct terms. Virtualization makes disk appear as something that it is not naturally. Virtual tape refers to the virtualization of a piece of media rather than the virtualization of the tape drives that go into a tape library.

Virtual tape has a longer history than virtual tape libraries and the use of virtual tape is a commonplace practice on mainframe systems. On the mainframe, the process of writing datasets to tape often left the tapes with a lot of empty space. With virtual tape, multiple datasets are concatenated on disk and then written to tape. Open systems typically have not had the same issue with empty space, but some efficiencies can still be achieved with open systems tape, so virtual tape is now available for open systems as well.

Virtual tape is primarily an asset utilization and ease of management benefit play. Virtual tape achieves indirect benefits for data protection by minimizing the number of tapes that have to be restored, which leads to fewer chances for restoration problems.

Virtual Tape Library

Today’s backup/restore software is designed to minimize the impact upon the existing policies, practices, and procedures of an IT organization. Standard backup/restore software packages can target disk as well as tape. A virtual tape library is software that runs on a disk array to emulate a tape library.

A VTL adds in the cost of the virtual tape library software in addition to the cost of the standard backup/restore software. However, simply retargeting standard backup/restore software from disk to tape requires that each backup job be manually retargeted to disk. That is not true of a virtual tape library. If the number of backup jobs that have to be changed is manageable, straight disk-based backup may be a feasible alternative. A second concern is that there might be a two-terabyte (TB) file system limitation, which would apply to straight disk-based backup, but not to a VTL. The two primary issues are integration and scaling.

A more complex backup environment and/or large capacity backup requirements (say, six TB or greater, as a rough measure) would tend to favor a VTL. Otherwise, straight disk-based backup might be a reasonable choice.

MAID

If nothing else, the acronym MAID catches the eye. MAID stands for Massive Array of Idle Disks. In a MAID, the disk drives are powered down, individually or in groups, when not needed. The premise of MAID is very simple: why spin disks continually when access to the data on those disks is very infrequent? By not spinning disks, savings are accrued on environmental costs (air conditioning and electricity), even though the disk drives may be packed more densely. In addition, reducing the power-on time for the disks means that lower-cost disks targeted for lower duty-cycle applications can be used, such as SATA disks. The net result is that these lower cost disks can have a longer service life than the higher cost disks used in high performance, always-on arrays

MAID is a middle ground between “online” disk and tape. “Online” disks not only have random access, which means that they can access specific information quickly when needed, but are always spinning, whereas tapes are idle until accessed. MAID is a middle ground where the benefits of random access can be used when an I/O request is actually made, but the benefits of idleness are also taken into account when the information is not needed.

MAID occupies the ground in active archives where the archive is still “active” in the sense that the data has to be “online,” but the need for that access is infrequent. MAID therefore occupies the long-term archive ground before deep archiving. That long-term ground could contain pools of infrequently accessed production data, such as old customer histories, old CAD/CAM files, and old camera surveillance data, but could also contain compliance data, such as medical test data. This is the type of data that is written once and only occasionally ever read. When it is needed, however, the expectation is that it will be easy to find and can be accessed relatively quickly, at least when compared to the same data archived on tape.

From a data protection perspective, MAID uses idleness to extend the operational reliability of an array. MAID is therefore suitable for use as a bulk compliance information repository, extending its usefulness to the logical operational side, if the proper software is used.

In the long-term archive space, MAID competes with tape on a bulk cost space basis and with “active” disk on a cost basis. MAID is an emerging technology that so far is supported only by a few smaller companies.

Removable Disk Drives and Disk Media

A few vendors are now reintroducing removable disks. One approach is to bundle a RAID group of Winchester disk drives and associated disk media into a removable magazine that is comparable in form factor to a similar magazine of tape cartridges. The second is to actually decouple the disk media from the disk drive. Both ap-

proaches enable the transportability of disk-stored information for offsite storage. Since the first approach embeds the drive along with the media, the key transportability issue is shock resistance; no one wants to lose data because the magazine was dropped! In the second case, each piece of media has to be put in a protective case to prevent environmental damage.

The direct removable disk media approach may be useful for small businesses (or units of large businesses, such as branch offices) for whom having even a single tape drive introduces a level of expense and complexity that they have a difficult time managing. Moreover, these organizations have to manage their own data protection. The challenge becomes to provide media management without having the formal tools — for example, managing movement of data to and from an offsite location while at the same time ensuring that no disk is lost, misplaced, damaged, or has the wrong version of information.

The bundled group of removable disk drives approach also requires media management, but, since this approach is more likely to be used in larger IT environments, the media management approach that is already used with tape may be employed.

The magazine could represent a backup copy in tape format or a replica in disk format. If the magazine represents a backup, this is an example of a disk-based backup, and using the magazine for restoration would constitute a disk-based restoration. If the data were in disk format, the data would be for a single point in time. Although restoration could be fast (especially important at a disaster recovery site), the problem would be in figuring out how to find and apply all the changes since the disk copy was made.

Removable disk drive and disk drive media approaches are likely to be useful in selected applications, but are unlikely to unseat established mainstream tape solution infrastructures in the near future.

Data Protection Appliances

A data protection appliance is a dedicated, self-contained bundle of software and hardware that serves a specific data protection function, such as acting as a VTL. A standalone VTL appliance would consist of a disk array whose disks serve only to support the VTL and a VTL software package. An integrated VTL appliance would couple a tape library to the disk array either logically through software or physically.

Three interrelated questions arise when discussing data protection appliances:

1. Where should the “intelligence” for a data protection function reside?
2. When should an appliance solution be used in place of a general-purpose solution?

3. If an appliance is used, should it be standalone or integrated?

The first debate is where software intelligence for data protection functions should reside — on an application or database server, in the storage network or on its edge, or on the array — either a general-purpose array or a dedicated appliance array.

The answer should not depend upon philosophical arguments, but rather on business needs. If an enterprise expects to have to scale a data protection function beyond the capabilities of what an individual array might be expected to provide or wants the ability to shift between heterogeneous storage platforms over time, intelligence in the network might make sense when it is generally available. Otherwise, having the data protection capability at the array level can make sense, as that is a familiar level for managing intelligence.

The second debate, between special-purpose appliance servers and general purpose computing servers, goes well beyond their applicability for data protection functions. The answer might well be that both will continue to flourish. An appliance is a black box in which the inputs and outputs are well-defined, but, if there is a problem, only vendor technical experts can solve the problem. Maintenance support is therefore critical. General-purpose solutions can also have problems, and although the problems may not be related to the data protection function, they still may impact it.

The third debate, standalone vs. integrated for a data protection appliance, really comes down to whether the value of improved ease of use through integration outweighs the proprietary lock-in for all parts of an integrated solution.

In summary, IT organizations have to decide how they want their data protection functional intelligence to act, either alone as software that is independent of a particular physical implementation or embedded in a specific physical implementation. Solving an immediate problem with improved manageability over what was done before has to be weighed against whether that solution will scale to meet future demands.

Tape Automation

Frequently, IT organizations are not in love with large, complex electromechanical tape automation systems both for perceived reliability reasons and because of these systems' drain on administrative resources. Moreover, having a very large number of pieces of tape media increases the risk of not being able to easily restore data when most needed. When you add in the inherent lower availability of tape than disk, and top it off with declining disk prices and the rise of a large number of disk-based backup alternatives, it may cause you to leap to the conclusion that tape's days are numbered.

However, that is very unlikely to happen. To paraphrase Mark Twain, reports on tape's death are greatly exaggerated. To see why, look carefully at Table 7-1. Tape is solidly entrenched in each of the eight boxes in the table (as vaulting typically refers to tape). Although disk can have a role in each of the eight boxes, disk is not as established as tape in all eight boxes, as the use of disk depends upon emerging software technologies in many cases. Moreover, a single piece of tape media could fit into each of the eight boxes as necessary. That is not true with non-removable disk. The removability and transportability of tape creates flexibility that might prove vital. For example, tape could be used to recreate data at a third site if the planned disaster recovery site should also fail.

Recall the discussion on degrees of protection. Three degrees of protection are probably the minimum number of layers of protection. For a time-and-revenue-sensitive application, three layers of disk might very well be necessary, but even here tape is probably sensible. The reason is that should all the disks fail there would be a significant revenue loss (and consequent loss of market valuation), but, if the enterprise could never recover its data, it could be out of business — period.

Three sets of disk arrays are extremely unlikely to physically fail. Logically, deliberate external or internal threats or even inadvertent human error (such as pulling out the wrong disk in an array, compounded by cascading procedural errors) might cause unexpected problems. Tape automation delivers “biological diversity” for extra degrees of protection.

Tape should continue to provide a relative cost advantage via tape. This is not true in all situations (as it depends upon how many pieces of tape media are used for each tape drive in a tape automation system). However, the decline in absolute cost of storage (due to continuing price/performance improvement of more than 30% a year for disk drives) means that using disks for disk-based backup, for example, becomes more affordable. That absolute drop in cost also applies to tape media, so strict head-to-head comparisons between tape and disk will tend to favor tape.

A tape automation system consists of pieces of tape media, tape drives, and a robotics automation system. Significant advances have been made in all three areas over the last several years in reliability, manageability, and capacity; and the roadmaps of leading vendors indicate that these trends will continue.

The introduction of new technologies does not always lead to displacement of existing technologies, but rather may lead to a change in the portfolio of functions that they perform. Tape has a long history of adjusting to disk — tape's primary role in batch processing (with extensive sorting of tapes to produce reports) has been replaced by its role in data protection.

Tape will continue to serve as the last line of data protection defense for time-sensitive critical systems and will be employed closer to the front line for not-as-time-sensitive applications. On the active changeable side, these applications will include both the traditional backup/restore processes (although perhaps on the back end of disk-based backup) and continuous data processing applications (where tape will contain copies of the data). On the active archiving side, replicas of ingested fixed content may not be done by the traditional backup/restore process, but copies will need to be made.

Tape may very well play a role in active archives for storing primary copies of very large amounts of compliance data, such as medical records, rights management data, such as videos and music, and bulk data, such as seismic data. The key determinant of disk vs. tape should be the frequency of access and the required response time once the data needs to be retrieved.

We conclude that IT organizations should focus on how disk and tape can best complement each other, as tape is here for the foreseeable future.

Getting to the Point

The ability to create a point-in-time copy of a pool of data is having — and will continue to have, through ever more ingenious uses — a major impact on data protection. A point-in-time (PIT) copy is a “copy” of a pool of data at a chosen instant in time. The advantage is that the copy is frozen in time. Although there is no guarantee that the copy itself does not suffer from data corruption, there is a guarantee that changes after the time of the copy will not change the original pool of data; thus a point-in-time copy provides logical data protection.

Since a PIT copy is considered to be a fully usable collection of data, some vendors offer the capability of writing to a PIT copy. Although a PIT copy can be a starting point for adding in changes, the moment that changes are made it ceases to be a PIT copy. The ability to use the PIT copy as a foundation for change is valuable, but the protection offered for logical data protection requires that a PIT copy be read-only.

Point-in-Time Copy

The two basic “flavors” of PIT copies are PIT clones and snapshots. PIT clones are an exact *physical* copy of a pool of storage. A PIT clone delivers both physical (as of the time of the cloning) and logical (from the time after the cloning) data protection. The price that is paid (other than the cost of the software) is a duplication of the disk storage required by the original pool of storage. The cost (as well as manageability) severely limits the number of clones to only one or a few. The advantage is that the use of the clone for such purposes as

serving as the basis for backup to tape or for production testing does not impact the performance of the production disk array.

A snapshot is a software image of data as of a predefined instant. A snapshot is taken on the original disks where the data is stored and at the time the snapshot is taken the original production data and the snapshot are identical. That means that no additional physical space is required at that instant. The original production data and the snapshot data diverge as writes change the original production data. The approach taken is typically a copy-on-write technique that creates temporary blocks and updated blocks of data. When changes are made, additional physical disk space has to be available and allocated. The process requires the management of index tables.

The key differences between a clone and a snapshot are:

- A clone is an offshoot of mirroring technology, whereas a snapshot uses an indexing strategy
- A clone requires a full physical copy of the original data, whereas a snapshot requires only enough additional space to accommodate all the changes
- A clone is a real *hardware* duplicate or replica; a snapshot is only a *virtual* duplicate or replica

PIT copies can serve many roles, including serving as a starting point for making a backup copy or for application production testing, but a key role is in providing high availability logical data protection.

Continuous Data Protection

With continuous data protection (CDP), an enterprise can create a data protection copy (typically on a disk-array-based data protection appliance) that can recover to *any* point-in-time. Typically, changes are recorded continually by the CDP appliance (using a non-invasive journaling technique that does not require even the momentary halting of an application's I/O processing that has to take place when creating a snapshot point-in-time copy). The journal can be rewound to any point-in-time as the basis for creating an any-point-in-time (APIT) copy of the data without having to know at what point-in-time a copy should have been taken.

When offered on a data protection appliance, CDP offers today's only up-to-the-moment logical data protection and physical data protection with high availability. CDP is not a new backup approach; CDP is an alternative (or, more likely, a complement) to the traditional backup software approach. CDP provides fine granularity over the data restoration process in that logical unit volumes (LUNs) or individual files can be restored.

CDP can be used natively as a temporary alternative to the original array (although with performance degradation, due to using midline disks instead of performance disks). CDP could also serve as the

basis for business intelligence analyses or production application testing without disturbing the performance of the production disk array.

CDP is an operational continuity approach when implemented in a data protection appliance. However, CDP may be made available over a distance as a disaster continuity approach. In this case, it would fall into the dated replication class, where the changes that are continually made all have a time stamp associated with them (i.e., the system knows the time each change was made).

One issue that has been raised is how well a CDP system can handle consistency groups. An application may use data that is spread across multiple physical disks. All that data make up a consistency group; a set of data that has to be synchronized for restoration purposes. This issue may be a non-issue (because frequently data is not spread across multiple disks), but it is a question that IT managers should raise when reviewing CDP products.

The concept of CDP is still sinking into the collective consciousnesses of IT organizations, but it is likely to be one of the technologies that make a major difference in how enterprises will change their data protection infrastructures in the future.

Replication Strategies

Replication technology is one of the strong suits in an IT data protection strategy, because of its promise of “high” availability for both operational and disaster continuity.

A replica is a copy, a duplicate, or a reproduction. For data protection purposes, a replica has to be a separate physical copy. A replica also has to be naturally usable — which means that the application that uses the information has to be able to use it directly. A replica can either be a dated duplication, which means that it has a stamped time of creation, or it can be undated.

A backup/restore copy is not a replica, since the copy is not naturally usable without undergoing a transformation using the restore functionality of the backup/restore software.

A PIT clone is a replica, since the clone is a separate physical copy of the data. A snapshot on a production copy is not a replica, since there is no additional physical copy. However, snapshots that are copied to a target either individually or in the context of continuous data protection result in a replica.

One of the big discussion topics in replication is where the software intelligence to manage the replication process should be located. The three choices are host-based, storage-network-based, and disk-array-based. The different choices can be examined on the basis of cost, scalability, manageability, performance, and use of IT resources. The decision in favor of a particular product depends upon the applica-

tion requirements and budget of an IT buyer, but the basic principles of replication apply to all three.

One key distinction in replication is between mirroring, which is undated replication, and dated replication, which is a simpler way of saying naturally-usable dated duplication. Mirroring is valuable for physical data protection for disaster recovery. Dated replication is useful for data protection for all aspects of disaster continuity as well as the active archiving side of operational continuity. If CDP is also considered a replication technology (and there is no reason that it shouldn't be), then dated replication is the only type of data protection technology that covers the entire data protection category matrix.

Mirroring

The job of mirroring is to create an exact copy (also replica or duplicate) of data on a source disk to a target disk. Mirroring is a continuous process. That means that the mirroring process does not take time off—whenever the source data is online; the mirrored copy should be as well. This also means that mirroring provides physical data protection, but not logical data protection, since data corruption would be copied to the target disk.

Mirroring is a process for the active changeable data side, and not the active archive data side. The active changeable data side will typically have from a little to a lot of fixed content data in it; data that is treated as if it were active changeable data. From a mirroring perspective, the fixed data would be copied only once (unlike making full backups of data, where fixed data is copied each time a full backup is made), so there is no real overhead. Thus, after the initial copy has been made, there would be no network demands from the fixed data. The only burden of mirroring would be that the remote array would be weighed down with the cost burden of perhaps more expensive disks than would otherwise have to be the case.

However, mirroring would be inappropriate for an active archive. Granted, inflows to the archive are changes, but other replication techniques are sufficient for one-time changes without introducing the costs and management requirements of mirroring.

One exception to these rules is when the mirrored copy is split off, which means that updates from the source no longer take place. At that point, the mirrored copy is now a point-in-time clone. The clone would offer both logical and physical data protection and could be used as when doing a backup. When the clone is put back into service as a mirror, a resynchronization process has to take place.

Although mirroring can be done locally, local mirroring is typically synchronous and goes under the name of RAID 1 (and variants). When mirroring is typically mentioned from a replication perspective, the discussion is really about remote mirroring. The two primary flavors of mirroring are *synchronous remote mirroring* and

asynchronous remote mirroring. A compromise between these two — *semi-synchronous mirroring* — has not received a lot of attention.

Synchronous Remote Mirroring

A *synchronous remote mirror* maintains an exact up-to-date copy of the data located on part or all of a local (also called primary or source) disk array with that of a remote (also called secondary or target) disk array. Every write I/O on the local array is immediately sent to the remote array. No further I/O write actions are performed by an application until the remote array acknowledges that it has also written the I/O to one of its disks. Thus the source and target are always identical, which is why this approach is called synchronous.

The big advantage of synchronous remote mirroring is that its RPO and RTO are (or can be made to be) zero. An IT organization does not have to worry about loss of data, so the data preservation objective is met, and failover can avoid loss of availability, so data availability SLA requirements are met. This seems like the best of all possible worlds, so it is no wonder that synchronous remote mirroring has done so well.

However, synchronous remote mirroring does not have all the answers. It provides physical data protection, but no logical data protection. That is fine for its basic purpose of providing physical data protection in case of a disaster, or for a secondary purpose of helping recover from hardware failures at the primary site. IT simply has to not ask synchronous remote mirroring to do a task (logical data protection) that it was not designed to do.

A second issue has been cost — for software, for storage network hardware to connect both the local and remote arrays to a WAN, for a separate remote array and surrounding IT infrastructure, and for having a private dedicated network line with sufficient bandwidth. Depending upon the nature of sunk fixed costs (such as for storage networking hardware) and variable costs (such as the size of an array), an IT organization may find that it cannot economically justify synchronously mirroring applications for other than time- and revenue-sensitive, mission-critical applications. This means that non-mission-critical applications (which could be the bulk of the application portfolio in terms of storage requirements) are not protected by synchronous mirroring.

A third issue is the latency inherent in remote communications due to the fact that the speed of light is finite. (Other latencies other than the speed of light are involved in the process, but the other latencies are fixed, whereas the latency due to light is variable based upon distance.) Latency in acknowledgement of writes at the target site to the source site introduces delays in the ability of an application to continue to do new transactions. This problem is not noticeable at “short” distances, but becomes a problem at “long” distances. Al-

though the longest distance is entirely arbitrary (as the impact depends not only on the latency, but on whether or not that latency actually noticeably degrades the performance of a particular application), a common rule of thumb is a *maximum* distance of 100 km (about 60 miles).

This distance seems reasonable for disaster recovery purposes until the organization realizes that these distances may still put an enterprise in line for the same disaster event (a hurricane, for example). Although arguably somewhat arbitrary, the *minimum* distance between disaster recovery sites is likely to be 500 km (roughly 300 miles). And that argues for the need to have another replication technology that either complements (for those enterprises that can afford it) or supplants (for those that must have the longer distance and cannot afford both) synchronous remote mirroring.

Asynchronous Remote Mirroring

The purpose of an *asynchronous remote mirror* is to maintain a copy of data on a source data array at a distant target disk array. In an asynchronous remote mirror, an application does not wait for an acknowledgement of an I/O write request from a remote target array before moving on to its next task. That means that the source and target arrays are not necessarily identical, since there is no guarantee that the remote site has actually received and written the I/O request successfully.

Asynchronous remote mirroring can take advantage of lower-speed (and, if necessary, less reliable) networks, in contrast to synchronous remote mirroring, which requires a high-speed, highly reliable network. The reason in the case of synchronous remote mirroring is that an application is tightly coupled with the remote site from a performance perspective, so a high-speed network is needed even if the volume of information sent is relatively low. With asynchronous remote mirroring, an application's processing performance is independent of the target site, so it does not depend as heavily upon the speed of the network.

An asynchronous remote mirror runs the risk of having a non-zero RPO. Depending upon a particular implementation, an RPO may be 15 seconds, 30 seconds, or minutes. An IT organization may face another unpalatable choice — business conditions (such as financial processing) require a zero RPO, but the necessity to have a disaster recovery site at a distance greater than necessary for guaranteed synchronicity physically forces a non-zero RPO.

A possible resolution to the dilemma is to accept an emergency non-zero RPO in case of a disaster. Disasters do not occur very frequently, and for such a situation a non-zero RPO might be tolerable although not desirable.

Another possible resolution is a “work-around” — for example, having three sites where one site is production, one for asynchronous

mirroring, and one for synchronous mirroring. The assumption is that the synchronous-mirrored site would have time to transmit the necessary changes to the asynchronously-mirrored site before it too went down. (In the case of an instantaneous disaster that affected both of the synchronously-linked sites, the loss of RPO data would probably be of small relative concern.) That solution is expensive and will probably be used only for time- and revenue-sensitive, mission-critical systems.

However, less expensive workarounds that enable recovery from the asynchronously-mirrored target disk array without loss of data may be possible. For example, journaling all unconfirmed transactions to a single disk somewhere within synchronous-mirroring range, even if the technique used to synchronize may be different from mirroring. We recommend that buyers pay attention to how particular data protection suppliers deal with this issue.

Dated Replication — Pay Close Attention

Dated replication — a time-stamped copy of data — is a new term, but the term is necessary to separate it from non-dated replication. Mirroring is non-dated replication. Mirroring in general — and synchronous remote mirroring in particular — has been the glamorous replication techniques, but they provide physical data protection only. Dated replication is important in that it provides both physical and logical data protection.

That does not mean that dated replication is better than non-dated replication, because synchronous remote mirroring — a form of non-dated replication — may provide better RPO and RTO for physical disaster recovery. However, dated replication may provide a competitive alternative to asynchronous remote mirroring for disaster recovery of active changeable data, as well as serving as a technology to be used in conjunction with the inflows and outflows that are associated with active archives.

Consequently, understanding more about dated replication is important, and that starts with understanding what a data replica is. A data replica is a natively usable real copy of data, which means that the application that created the data should be able to read the copy without any intervening transformation or movement of data. A backup copy is a dated duplication, but it is not a dated replica, as typically a backup copy is not natively usable.

A snapshot is not a dated replica, as it is not a physically distinct copy of the data separate from the original pool of data. However, a snapshot can be used as the basis for creating a dated replication on another set of storage media.

A PIT clone is a dated replication, but it is also a full copy, and one of the advantages of many dated replication techniques is that the replica is updated with changes over time. The dated replication

would then be dynamic and not static, which is useful in an active archiving world.

Dated replication is quietly infiltrating the replication market and will continue to do so. There are multiple approaches and strategies for dated replication. Apart from snapshots, the following are some techniques that can be used:

- *I/O journaling* — all I/Os (including both volume and file-level I/Os) are copied and time-stamped, so that one replica of the full data pool can be restored to many different times.
- *Periodic replication* — this is a process where the original data pool is synchronized with the replica data pool on a periodic basis.
- *Copy-on-close* — a copy is made of a file (typically not a database) when an application finishes writing the file.
- *Copy-upon-insertion* — the requisite number of copies are made upon ingestion of data into an active archive; this is a “once and done” approach.

Dated replication can be divided into “local” and “remote,” where local solutions would be on a LAN or SAN (for help with operational continuity), whereas remote solutions could be over a WAN (for help with disaster continuity).

IT organizations can use the term “dated replication” in helping them classify and examine a number of data protection techniques that may have other names. IT organizations should pay close attention to dated replication approaches to see where they fit in the data protection framework, and if they match an organization’s particular requirements.

Special Requirements for Compliance

Naturally, each type of compliance request requires extensive individual attention in terms of policies, practices, and procedures. What should not be lost in all this effort is that the basic principles remain constant.

Enterprises should not look at compliance as just a one-off activity, but rather as part of the overall retention management program for its active archives. A digital rights management program for the control of digital assets should be part of this effort, and as many of the same principles as for regulatory compliance apply, from data preservation to controlling who is allowed to have access to the information. Protecting trade secrets and intellectual property offered for sale over the Internet are important, and this information must comply with internal requirements, even though they are not subject to regulatory compliance requirements.

The Use of WORM Technology

Compliance data must be immutable (i.e., guaranteed inalterable at least until the expiration date for retention has been reached). One way to do this is to write the data on a piece of write once read many (WORM) storage media. Some optical media technologies are physically WORM. In contrast, neither magnetic disk nor magnetic tape are inherently physically WORM media, but can be made so logically — that is through the use of software or firmware (which is software frozen in hardware).

Some enterprises are concerned that they need physically WORM media to ensure regulatory compliance. This is generally not the case. Regulators try to specify functionality, e.g., immutability, and rather than a particular technology that delivers that functionality. Requirements remain, but technology evolves.

Issues with Physically Destroying WORM Media

A physical device can be destroyed, either according to policy or not according to policy. If it is according to policy, destroying a piece of physical media means that all of the data is destroyed, so it's critical that all of the data be expired. If all the data has not expired, policy should not authorize the physical destruction of a piece of storage media.

But this presents a problem. When data reaches its expiration date, the change in status does not necessarily mean that the data has to be destroyed, but rather that it is eligible to be destroyed. In many cases, no problem exists if the organization chooses to retain the data beyond the end of its “freshness date.” This may not be true for all data; policy may require the deletion of some data immediately after the data is eligible to be deleted.

WORM disk may have the option to delete the data through a software process, but WORM tape probably does not have that option. The alternative to destroying the data physically is to have the data encrypted. The data can be logically deleted by the simple expedient of disposing of the encryption keys. However, IT organizations then have the burden of putting into place a key management program for managing encryption keys. An IT organization does not want to lose any encryption keys, but may need to have procedures in place for disposing of certain keys as necessary.

Encrypting creates another management burden besides key management. The content of encrypted data cannot be examined without decrypting the data first, which requires time and resources. A comprehensive metadata repository may provide an index that can facilitate the search process without the burden of opening up all documents, but not in all cases. A business intelligence analysis or compliance request may need to search the detail of each document.

Despite the burdens of managing encryption, IT organizations may have no other choice than to put an encryption key strategy in place for the management of compliance data.

Physical destruction can also be outside of policy, whether unintentional, such as a head crash, or willful, malicious, and illegal destruction of a piece of storage media. To counter that, data protection copies have to be made. Those copies need to be direct copies, and be only WORM-enabled. Otherwise, a non-WORM enabled copy could be altered and then rewritten to a WORM copy. The WORM copy would have the appearance of presenting the data correctly, but would not be correct. As with other areas of compliance, while software and hardware can help, the responsibility is for IT management to put in place the proper policies, practices, and procedures to ensure that the duplication process is done correctly.

WORM Tape

Typically, WORM tape refers to tape cartridges; not to the tape drives in which the tapes are read. Electronic keys or inalterable firmware on the tape cartridge itself turn that piece of media into a WORM tape. Data on a WORM tape cannot be rewritten or reformatted, but can be appended (until the tape runs out).

Depending upon the vendor, a WORM tape cartridge may be a purchasable stock keeping unit (SKU) item, which comes only in WORM format, or initialized as a WORM tape cartridge at the time of first use. In either case, visual identification of a tape cartridge (color of cartridge for a permanently-designed WORM tape cartridge or color of label for a tape cartridge that is initialized as a WORM tape cartridge) can help prevent mishandling of tapes when manually handled after removal from a tape library.

Using WORM tape in conjunction with WORM disk is logical from a data protection perspective, as the data protection copy on tape also has to retain the compliance characteristics that were required on disk. However, remember that moving data from WORM disk to WORM tape is a replication process and not a traditional backup process.

The traditional weekly full backup would not apply because a rotational scheme involving multiple generations of tape implies that the tapes are reusable. (Tapes from an expired generation are put in a pool of “scratch” tapes to be reused.) This is clearly not the case with WORM tape. Although a synthetic full backup approach might be used, data probably would be copied from disk to tape by a periodic replication method.

WORM Disk

The non-erasable, non-rewritable functionality that enables WORM disk capability comes from software at the operating system (OS) level of a storage system — at the network-attached-storage (NAS)

“head” or at the storage controller/server level. When WORM capability is invoked, no one — not even a system administrator with superuser privileges — is allowed to rewrite or modify data. Building in the necessary WORM software functionality requires the storage system vendor to have access to change the OS kernel itself, which means that a proprietary or Linux-based OS would most likely be the chosen OS.

WORM disk is actually a misnomer (but the term can stand) because the actual disk drives themselves have nothing to do with the WORM functionality. Theoretically, “WORM-protected” disk drives could be removed from a system and moved to another system which does not offer WORM protection, but that type of event should be detectable and is extremely unlikely.

One advantage that WORM disk may have, and that the current generation of WORM tape does not have, is the ability to actually delete expired data (assuming that each piece of data is managed on a file basis, although blocks of data on a logical volume could also have expiration dates). This functionality is important because organizations need to be able to reuse disk space as well as to allow the use of encryption as a choice if desired.

WORM disk may present a planning problem that does not affect WORM tape. If compliance data grows much more rapidly than anticipated from a WORM tape perspective, the only requirement is additional WORM-enabled tape media. If the compliance requirements outgrow a disk array, that could present a problem from the WORM disk perspective. In a non-WORM environment, an IT organization might migrate data from an older, smaller capacity disk array to a newer, larger capacity disk array. The older array would either be repurposed or sold. In a WORM case, the data might be migrated to a new array (with suitable precautions), but then the older array would be rendered unusable and unmarketable, as data cannot be deleted until at least their expiration dates. One solution is to select an array that can start small, but can expand (if necessary) to meet future requirements. Another is to select WORM disk functionality, where the disks assigned for the storage of compliance data are virtual disks that can be changed (i.e., migrated) to other physical disk drives as the need arises.

Electronic Locking

The ability to put an electronic “lock” on a piece of data for a prescribed period of time is a key piece of functionality that a WORM-disk can provide. A time-based lock might be used on non-compliance-related data in an active archive to deliver an easy method of logical data protection, since no write I/Os can tamper with or destroy the data. The time lock need not be very long and might be automatically renewable if the decision is made to retain the data.

Guaranteeing the Authenticity of Data

With compliance data, the question of authenticity may come up. Two documents may differ only slightly, and yet those small differences may be significant. The uniqueness of a particular document may be guaranteed by a combination of software and hardware called *content addressable storage* (CAS). Uniqueness does not guarantee authenticity, but, if the creation date and the date put in a protected storage state are incorporated as part of the document, then the genuineness of the document is much closer to being verified.

Privacy and Confidentiality

Authentication is not only for determining the genuineness of a document, but also for ensuring proper access to a piece of data. Privacy is the seclusion of the data from unauthorized users' view, and confidentiality protects exposure of information that should remain secret unless authorized. The focus of this report has been on the storage side of data protection, but confidentiality is a key objective of data protection and must not be neglected.

Compliance Appliance

A compliance appliance is an intersection of a production copy of the data and a data protection appliance. Whereas a data protection copy's sole purpose in life is to serve as the basis for restoring or restarting access to data, a compliance appliance is also a production version of the data, in the sense that applications that need access to the data can access it.

The debate of a general purpose storage system that contains compliance data vs. a dedicated appliance storage system for managing only compliance data does not have the same level of intensity that it has in the case of a data protection appliance. A sense of a need to isolate compliance data from the rest of the information infrastructure would tend to favor the use of a compliance appliance. The software that manages compliance can work in conjunction with the appliance's operating system without conflicting with other requirements. That might be important if two opposing demands — the need to keep a revenue-producing application running around the clock and the need to satisfy a regulatory request — come into conflict on a general-purpose storage system that does not have sufficient performance to satisfy both. A general-purpose system could still be used, but the requirements would have to be thought through carefully.







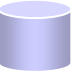
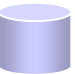


Mapping the Base Data Protection Technologies to the ILM-Version of the Data Protection Framework

The preceding sections described many current and emerging data protection technologies. Table 7-2 (for Active Changeable Data) and

Table 7-3 (for Active Archived Data) on the following pages show where these technologies fit.

Before the boxes in the ILM-version of the data protection framework can be filled in by each application within an application portfolio, these choices have to be examined. Requirements have to jibe with available technologies and available budgets.




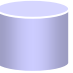




Table 7-2: Base Data Protection Technologies for Active Changeable Data

	Technology	Physical	Logical
Production Copy			
	RAID	X	
	Point-in-Time Copy Snapshot Clone	Secondary	X Primary
Data Protection Copy			
Local (Operational Continuity)			
	Tape Automation*	X	X
	Virtual Tape Library*	X	X
	Continuous Data Protection	Secondary	Primary
	Data Protection Appliance*	X	X
Remote (primarily Disaster Continuity)			
	Mirroring Synchronous Asynchronous	X X	
	Dated Replication	X	X
	Vaulting*	X	X
	Electronic Vaulting*	X	X

* Backup/restore software is or might be used in conjunction with this technology

Source: Mesabi Group, January 2005

Table 7-3: Base Data Protection Technologies for Archived Data

	Technology	Physical	Logical
Production Copy			
	RAID	X	
	WORM Disk		X
	Guaranteed uniqueness		X
	Electronic locking		X
Data Protection Copy			
	Local (Operational Continuity)		
	Dated Replication	X	X
	WORM Tape	X	X
	Remote (primarily Disaster Continuity)		
	Dated Replication	X	X
	Vaulting	X	X

Source: Mesabi Group, January 2005

Replication is repeated at both the local and remote site to indicate that it is the primary technology for both.

Chapter Eight:

Summing Up — Redesigning Data Protection

You have reached this point in one of two ways. One is that you have read the report faithfully from the beginning — a daunting task for which you are to be commended! The second way is that you are looking ahead to see what the answer is without having read the report in detail. If the latter, you are welcome to read ahead, but this chapter is not the “Cliff Notes” for the report. However, it does sum up the report to help you think about redesigning your data protection strategies and infrastructure.

Data Protection Is Everyone’s Business

Data protection is not only the responsibility of the storage administrator and the backup administrator, but also that of the database administrator, the network administrator, and the security officer. The responsibility also extends to IT management and to everyone in the enterprise. And the problem is that too often when a function is the responsibility of everyone, no one is held accountable.

IT management has the lead role, because it has the overall perspective on operations (what is happening today), tactics (what should be done to improve the data protection infrastructure in the upcoming months and year), and strategy (what plans need to be put in place for the long-term betterment of the data protection infrastructure). Note that tactics refer to what an IT organization can realistically do within its current budgeting cycle, and strategy refers to what resources will be necessary to affect change beyond the current budgeting cycle as well..

IT management has to manage the change process, which is beyond the scope of individual contributors. The change is not only in technologies, but also in policies, practices, and procedures. Change affects tasks that individuals perform (so roles and responsibilities change and training may be needed) and may affect the departmental structures (which department is responsible for which functions). Moreover, change requires money. IT management not only has to lead the change within the IT organization, but also work with other organizations within the enterprise, if a data protection change process is to work.

Synthesizing the Data Protection Frameworks

IT organizations are action-oriented. Although the actual work may be time-consuming, the overall actions can be easily specified:

1. *Separate active archive data from active changeable data.* Start with the reordered and revised matrix (Table 8-1). For each application, determine what the mix of active changeable data is versus fixed content data. Is there a “distillation” process that can be put in place to move the fixed content into an active archive? Does it make sense to do so? What does it take to put an active archive in place? Can the active archive serve the needs of multiple applications or do you need to create appliances, such as for a compliance application? (The analysis devolves to the left side of Table 8-1 for those applications for which an active archive is not a consideration.)

Table 8-1: Data Protection Requirements for Application *n*

	Active Changeable		Active Archive	
	Operational Continuity	Disaster Recovery	Operational Continuity	Business Continuity
Physical				
Logical				

Source: Mesabi Group, January 2005

2. *Determine the minimum acceptable requirements for both the left side (Active Changeable) and right side (Active Archive) for all data protection objectives.* No easy way exists to set requirements. What an organization might like (say, 100% availability) may not be affordable, and what an organization can afford from a budget perspective (say 95% availability) may not be acceptable to the business. Start with a realistic number. When the time comes later for analysis, determine the cost. If the cost is too high, opt for a solution that is closest to the desired goal, but still meets budget requirements. Determine the impact that the gap between the two choices has on the business. Perhaps a business case can be put together to justify additional funding. Remember that, in this whole process, all objectives — data preservation, data availability, data responsiveness, and data confidentiality — have to be met. Do not get locked into defining a requirement that is not comfortable to you. For example, as discussed earlier, RPO for operational continuity and disaster continuity may be different.
3. *Determine the degrees of data protection (including both high and low availability alternatives) that can help meet selected objectives.* This report focuses strongly on the data preservation and data availability objectives, which is where the degrees of protection come in. Data responsiveness, and

some aspects of availability, are overall IT infrastructure issues. The confidentiality objective must be worked through in conjunction with the security officer.

4. *Determine if current data protection processes and technologies meet the requirements of data protection to the required degree.* This is a gap analysis to see if there are any differences between what is “required” and what is being delivered today.
5. *If there are any deficiencies (or excesses) in data protection, determine what classes of data protection technology might close the gap and then evaluate those technologies for suitability.* The first goal is to determine the feasibility of the technology from a general sense. The second goal is to examine the particulars and costs of specific vendor implementations of those technologies that can work in conjunction with your IT infrastructure.

Although the steps are discussed in linear fashion, the actual process is likely to be more iterative, may branch in different directions, and may drill down to finer levels of granularity.

Guidelines for Data Protection

The preceding section discussed how to go about using the frameworks. Some general guidelines (or rules of thumb) might be useful from an overall perspective (fully recognizing that prescribing bromides is a lot easier than taking them).

1. Make sure that you have the necessary layers of both physical and logical data protection strategies in place. You do not want to leave gaps in your protection coverage.
2. Take the complexity out of the process wherever you can. One way of doing this is to minimize the number of things that can go wrong (which is why using dated replication instead of backup/restore software for fixed content information is important).
3. An ILM plan that creates a strategy for active archiving requires data classification, a pooling and tiering plan, and data retention policies as part of the process.
4. Have a triage plan in place in case of emergency; but such a plan first requires being able to isolate the interactions among applications (if everything is interrelated, that may be a problem).
5. Remember that physical data protection is often a simple matter of “dialing-up” the level of protection that you need, but that logical data protection has to be thought through very carefully to determine the level of protection that is provided.

6. Know who is going to address a general class of problems, and how, before it occurs. Make sure that you have at least two people (preferably in different sites) who can take action in case one person, for whatever reason, is not able to deal with the problem

The Challenge Ahead and a Call to Action

IT organizations face a dilemma regarding data protection. If they maintain the status quo and try to do only what they are now doing, they face the prospect of increased unmanageability, with the prospect of not protecting all the data all the time as it should be protected. If they make changes in their IT processes and infrastructure, they run the risks inherent in introducing new technologies, including management of the organizational change process

This report will help organizations think about preserving what they can preserve, but also help organizations move forward to adapt to the coming sea change. This means that the status quo is not acceptable. IT organizations have to move forward — and they have to do it now.

Chapter Nine:

Tape: An Ongoing Bulwark for Data Protection

Face up to it and plan accordingly: Tape is here to stay for data protection. Disk will have an expanded role in data protection, but tape will continue to use its strengths to complement disk.

Disk is good for front-line data protection defense for high availability and low short-term data loss risk. Tape supplies long-term data preservation that is a mandatory part of any data protection policy for an enterprise.

To meet those tape needs, the LTO Program offers a well-accepted tape technology with a roadmap for the future. This open tape technology will continue to serve the interests of IT organizations for open systems tape for the foreseeable future.

The Challenge to Tape

Two of the key principles of data protection are data availability and data preservation. Data availability is measured as a recovery time objective (RTO). As part of its ongoing work, the Storage Networking Industry Association (SNIA) recommended that for an RTO to achieve 99.999% annual uptime ('five nines') means a maximum of 1.5 minutes downtime per year and 99.99% uptime means 15 minutes downtime per year ('four nines').

Tape — by its sequential nature and due to the fact that it is not the active production copy — provides lower availability than four or five nines. And many organizations, for good and sufficient reason, have that level of availability burned into their service level agreements (SLAs).

Thus, the requirement of four or five nines of availability (as well as low short-term data loss risk) is leading to the greater use of disk-based data protection both as a first line defense (such as mirroring and continuous data protection) and a second line defense (such as disk-based backup, including virtual tape libraries). This greater attention paid to, and expanded role of, disk has led some pundits (and vendors) to proclaim that disk will replace tape. Nothing could be further from the truth.

Setting the Case for Tape

On the surface, the inevitability of tape is not intuitively obvious. Let's take a zero-based approach to understand why if tape had not been invented, it would have been necessary to invent it.

The discussion starts with data preservation — the second of the key data protection principles. One measure of data preservation is recovery point objective (RPO) — the amount of data that is exposed to permanent loss. SNIA suggests 1 minute in a five nines environment and 10 minutes in a four nines environment.

However, the RPO represents only the first fallback design point to which the data protection infrastructure can respond in the event of a failure. Additional failures may expose all data to the risk of permanent data loss. For example, the RPO for a single failure in a RAID 5 array is zero. However, a second failure in the array before the first failed drive is rebuilt represents the potential loss of all the data (which is why there is backup).

And complete (or nearly so) data preservation, which is tape's strength, is mandatory (i.e., long-term business survival depends upon it), whereas availability is only vital (the risk of revenue and market valuation loss has to be balanced with feasibility and costs to achieve higher availability).

In fact, disk-based data protection cannot achieve all the necessary goals for data preservation technically, economically, and process-wise — and tape can.

Why High Availability All the Time Is Unrealistic

High availability is possible under typical data protection scenarios, such as normal hardware failures. It may not be feasible in 'edge case' scenarios, such as a regional disaster that physically destroys all on-site hardware. Data protection has to shield against uncommon situations where the probability of a potential problem may be low, but the potential loss would be huge.

Take a disaster situation where, in the best case, an enterprise is able to afford a fully-equipped remote spare data center plus hot or warm failover after the predicted lengthy loss of its original production data center. High availability has been preserved, but the enterprise is now at risk during the original data center's shutdown; and although the original data center's applications may continue to run during the shutdown, the remote data center may not preserve all of the key enterprise information that is vital to the business's long-term survival.

Reconstruction, Not Just Restoration

In fact, no organization can afford to keep more than an absolutely necessary number of data centers. Instead, for 'just in case' planning, the organization creates an affordable "virtual" data center — and tape media provides that capability.

When all else fails, tape media provides the data building blocks upon which working applications in an IT infrastructure can be re-

constructed. Thus, tape media allow users to create a “virtual” data center that can be translated into a real data center as needed.

The media properties that are essential for this functionality are removability and portability. That means that the media can be taken from its creation point and moved physically to where it is needed for either a virtual or real data center.

Tape buys both time to reconstruct a real data center and peaceful sleep at night, knowing that tape represents an emergency “spare tire.”

Process and Technology Diversification

No rational investment manager would be caught without a balanced portfolio, as forecasts for any one type of investment vehicle are highly uncertain. By the same logic, even if by chance an organization could afford to have enough degrees of protection through disk alone, it would not be a wise idea.

For example, all the disk mirrors imaginable would not prevent against the propagation of logical corruption swiftly to each and every disk mirror. Continuous data protection systems can defend against logical data protection problems, but are subject to the same physical failure as any other fixed disk system. Tape provides protection against both physical and logical data protection problems in extreme situations where, for whatever reason, disk systems are unable to meet the need.

When one uses a common process and technology, such as disk, a data protection system can be exposed to a systematic or a random problem, such as a human error that is not detected until too late but that still has a far-reaching impact. For example, replication from a full to an empty array can backfire if the replication is accidentally done the wrong way (as has happened). And that error doesn't have to be manual. Automation is the computerization of human processes, and errors can creep into automated processes (e.g., software bugs) as well.

Tape offers a different management process and technology that offers a fallback position in case of emergency. Yes, the backup/restore software management process with tape may not always be a popular one, but it is a safety measure. No, absolute certainty is never possible, but having two separate and distinct processes has been shown in the real world to add a large margin of safety.

I/O Isolation — an Extra Measure of Safety

At least one data protection copy should be safe from inappropriate write I/Os, say from an accidental “delete all” command from a superuser, who thought that the command was being applied locally

and not globally. Tape media outside of tape drives in slots of a tape automation system or stored offsite comply with this dictum.

Disk-based Backup Is Not a Panacea

But what about disk-based backup, such as the use of a VTL? Say that two copies of data are kept on disk, both at the production site, or one at the production site and one at a spare site. Isn't that enough?

No. First, as noted above, backup disk, onsite or offsite, is much more susceptible to the same problems that afflict the primary disk set than is tape — i.e., tape represents “technology diversification.” Second, even backup disk is “online” and in operation, and therefore not isolated from whatever problems afflict the primary hardware; but tape is always clearly “offline” and operating separately — i.e., tape offers “logical isolation.”

Removable Disk Is Emulated Tape

What about removable disk?

Here the questions of scale and costs of change come into play. Tape solutions scale well (and “carbon copies” of a particular backup can be made). Moreover, in order to ease a transition from today's tape technology, removable disk would have to assume the characteristics of tape; including compression and sequential access, and automation using the same management processes (including media management). In that case, enterprises should ask themselves why not use tape, which also provides technology diversification and is already automation-friendly. Some might argue that disk has lower cost; but for scalable solutions tape still has a cost advantage over both removable and fixed disk.

Doing Tape Well

To recap, the purpose of tape is to ensure the long-term data preservation of all of an enterprise's data even in the face of an extreme challenge. For dealing with extremes, fixed disk in fixed locations is more vulnerable and does not provide the requisite variety in both process and technology. Yes, disk works well the vast majority of the time, but tape is there as the last bastion of defense. The basic need — guaranteed long-term full data preservation — has not been met and cannot be met by disk.

Doing Tape Right

But tape cannot be done halfheartedly. Tape management disciplines, such as media management to keep track of what information is available on both onsite and offsite pieces of media, has to be in place. At a minimum two copies have to be kept — one onsite and one offsite. And even starting from a zero-based approach, an enter-

prise may very well come up with the tape rotation strategy that is currently in place (except for perhaps one copy that is held in a disk-based backup).

Enterprises can turn to a number of well-seasoned tape technologies for their tape requirements. The technology towards which many enterprises are turning comes from storage vendors that support the LTO (Linear Tape Open) Program-developed specifications. Those specifications cover both tape drives and the tape media itself for what is called LTO Ultrium technology.

Four companies supply tape drives that meet LTO Ultrium specifications as follows:

- Quantum
- HP
- IBM
- Tandberg

All four companies have a long history of success in the manufacturing of tape drives.

In addition, IT buyers have a choice of buying LTO Ultrium tape media from a number of established suppliers as well including:

- FujiFilm
- Imation
- Maxell
- Sony
- TDK

The objective of any open technology is interchangeability and compatibility of products. With LTO, any certified LTO tape cartridge can work with any certified drive of the same — at least — LTO tape generation.

Although IT buyers typically would buy tape drives from one supplier and tape media from one or two tape cartridge suppliers, those buyers have the peace of mind in knowing that they could switch if they had to for whatever reason. They know that competition among the suppliers of both the drives and the media will continue to drive quality, availability, price, and the development of the next generation of product. LTO technology offers compatibility between generations. It is designed to read back two generations and write back one helping users to ease implementation and protect investments.

LTO Ultrium 3 — Popeye after Eating His Spinach

LTO Ultrium is now on its third generation. LTO Ultrium 3 has a native capacity of 400 GB per tape cartridge and up to 80 GB per second native performance. Compression at an assumed ratio of 2 to 1 is typically used with tape, so IT buyers can plan on actually get-

ting double the native numbers. What that translates into — as an illustration — is that the entire 19 million books in the U.S. Library of Congress could be contained on about 12 Ultrium 3 cartridges and that the movie *Finding Nemo* could be stored in about half a minute.

LTO Ultrium — Working to Make Tape Even More Reliable

Tape drive manufacturers continue to strive to make tape drive technology even more reliable than it has been. One of the key reliability development efforts in LTO tape technology involves servo tracking mechanisms that improve the precision tracking of the head and media in order to help ensure accurate reads and writes of data.

Another is read-after-write verification, where data is immediately read after writing. That not only helps ensure that the data was written correctly in the first place, but also should increase a system administrator's confidence that the data can be restored if and when necessary.

To further increase reliability, LTO technology uses advanced microcode to detect and correct errors. This capability is important for reducing the number of possible failures when trying to restore information from tape.

LTO Ultrium Is Up to Speed

An issue with linear tape technology prior to LTO Ultrium Generation 2 was back hitching. Back hitching, also called “shoe shining,” is the stop and go action of a tape drive when the tape drive cannot match the speed at which it is receiving data to its tape speed. Stopping tape that is moving fast dead in its tracks and then restarting again and again is not good for performance, nor for the reliability of a tape drive and its associated tape cartridges.

Starting with LTO Ultrium 2, all generations of LTO technology have (or will have when available) speed matching, where the drive can sense how fast it is receiving data and take the proper action to avoid back hitching. Speed matching not only improves performance, but also reliability. Another feature added with LTO Gen 3 is the expansion of the buffer from 64MB to 128MB. The larger buffer helps to keep data streaming to the tape.

WORM — Getting an Extra Use from Tape

Using tape technology for storing compliance data is a good idea for a number of reasons including data protection and low archival costs, and especially where tape technology is already in place. LTO technology offers a WORM (write-once, read-many) formatted cartridge. An LTO WORM cartridge helps to address compliance regulations for being non-rewritable. Data can be appended to a partially written WORM tape, but parts of the tape that have been already written to are unalterable. A WORM LTO cartridge is two-tone to

enable easy visual identification to distinguish it from a standard non-WORM LTO cartridge.

LTO Ultrium — For All Generations

The LTO Program has already demonstrated the robustness of its technology, not only in the first and second generations, but now in the third generation of the technology.

The LTO roadmap now extends into the next three generations 4, 5 and 6. Enterprises can be comfortable knowing that LTO Ultrium is a “living, breathing” technology that will continue to evolve as rapid data growth continues.

Conclusions

Tape will continue to be a best practice for enterprises as part of their data protection infrastructure portfolio. Mobility (i.e., the combination of removability from a tape system and portability, which enables physical transportability over long distances) is not an option, but rather a necessity for the virtual data center. In addition, process and technology diversification as well as I/O isolation provide those extra safeguards that may prove all the difference if called upon.

So enterprises should maintain their tape management disciplines and not worry about how to replace tape, but rather how to put together a balanced data protection portfolio where disk and tape complement each other in the appropriate mix. And the LTO Ultrium technology, with shipments of over 1 million tape drives and 30 million tape cartridges, points the way with demonstrated great success in the real world.

The bottom line is that tape is an insurance policy that no other technology can match, covering the enterprise against a risk that it cannot afford to take.

Mesabi Group Conclusions

At the end of the Preface, you were asked to prepare answers to a list of questions: what your view of data protection is, what is being done now on data protection, what the issues are regarding data protection, and what actions, if any, are planned to improve the data protection processes and infrastructure. If you have faithfully read through the report, you should be in a position to write down your answers now and compare them to what you originally wrote. Please do so before continuing. (If you skipped writing down your original comments, you may want to do so now.)

Please compare your before and after list. If they are one and the same, there are at least two possible answers as to why. If you knew the answers from the report upfront and so didn't need to change, congratulations. Hopefully, this report will strengthen your resolve to do what you already knew should be done and you can use the report as ammunition to further your point of view with others.

If your answers before and after are the same, but those answers do not share the directions recommended in this report, there are also at least two possible answers. One is that you read the report and could not accept all the conclusions. A second answer is that you read the report in detail, but perhaps did not understand it. If so, that is the fault of the report, but you may also wish to reread it.

The more important case is where the before and after answers are not the same. The difference in your answers should reflect what you think and not just regurgitate the ideas expounded in the report. The exercise was not to test your knowledge of what the report said, but whether or not the report helped to change your thinking process about data protection and whether or not you are likely to take action in the future on the basis of that knowledge that would be different than what you have done before.

Don't worry if understanding all the needs for change is difficult. One organization is implementing a large fixed content (i.e., active archive) system and has a hard time understanding that when data is ingested into the system for the first time, the process of data protection is replication and not the traditional backup/restore software process that they have been so accustomed to over the years. Take your time in understanding how the changes in data protection might affect you. Reread selected sections, if that would help, in a few weeks or months after you have a chance to think about the impact of change.

The sea change for data protection is here. Here's hoping that you can use this report to ride the waves of the sea change successfully.

Glossary

Active archiving	Data where frequency of access is active rather than inactive, while frequency of updating is non-existent so the data is fixed (i.e., unchanging) and not subject to I/O writes that would change the data
Active changeable	Data where frequency of access is active rather than inactive while frequency of updating leads to changes in the data so that the data is not fixed (i.e., unchanging)
Appliance	A storage system that is dedicated to a specific function, e.g., data protection appliance or compliance appliance
Archive	A long-term collection of data that typically is fixed content data, i.e., no I/O writes are allowed to change the data
Asynchronous remote mirroring	Remote mirroring where the source and the target may not necessarily be identical because of a delay from the target in acknowledging a write
Backup/restore	Backup is a dated (i.e., specified time) duplication of a designated set of data from a data source on one set of media (typically disk) to a backup set of media (either disk or tape) while restore takes the data from a previously created set of data on backup media and copies it to a set of media from which an application that uses the data can access it
Business continuity	A business function that attempts to prevent any major disruptions to business process both through planning, to avoid unplanned outages in the first place, and then on minimizing the effects of unplanned outages if they do occur
Clone copy	A point-in-time copy that also creates an additional physical copy
Compliance	A subset of data retention policies and procedures that must adhere to more rigid and rigorous conditions
Continuous data protection	The ability to create a copy of data that can be restored to any point-in-time
Data availability	The ability of I/O requests to reach a storage device and take the appropriate action
Data confidentiality	Data is available only to those authorized
Data lifecycle management (DLM)	Managing data as blocks without underlying

	knowledge of the content of the blocks based upon limited metadata (i.e., creation date, last accessed)
Data preservation	Data must be consistent and accurate all the time, and also must be complete within acceptable limits
Data protection	The mitigation of the risk or loss of or damage to an enterprise's data on either a temporary or permanent basis
Data responsiveness	The ability of I/Os to deliver data to an authorized user according to measures of timeliness that are deemed appropriate for an application
Data retention	The policies and practices around when specific data should be kept and when it should be disposed of
Data security	Data security shares in common with data protection the requirements for data preservation and data confidentiality, but has a primary focus on the network side more than on the storage side
Dated-replication	A time-stamped new physical copy of data
Deep archiving	The original definition of archiving where production data was written to another set of storage media (typically tape) and moved offsite while the original version was deleted (typically from disk)
Degree of data protection	Each degree of data protection is a layer that can tolerate one point of failure
Disaster continuity	Proactive planning, provisioning, monitoring, and preventive maintenance to minimize the impact of a devastating event upon an enterprise if one should occur
Disaster recovery	The attempt to minimize the impact of a disaster upon business processes if a disaster should occur
Disk-based backup	Using a disk array rather than a tape automation system as the target of a backup process
Electronic locking	Through the use of software, "lock" data from being modified
Electronic vaulting	Moving data for data protection purposes from a source site to a target site over a network
High availability	A relative term to indicate that the unavailability of an IT application to users is measured in terms of seconds or minutes per year
Information lifecycle management	The policy-driven management of information as it changes value throughout the full range of its life-

(ILM)	cycle from conception to disposition
Logical data protection	Protection of the data itself from change through unauthorized or erroneous I/O requests
Long-term archiving	Active archived data where the frequency of access has fallen so low that a tier of more cost-effective storage may be a more appropriate place to house the data
Low availability	A relative term to indicate that the unavailability of an IT application to users is measured in terms of hours or days per year
MAID (Massive Array of Idle Disks)	Spinning down disks when they are not accessed increases lifetime and lowers cooling and electricity costs while at the same time preserves the ability for online access as required
Midline	Online access to production data that is on more cost-effective storage than high-performance disk
Mirroring	Duplicating the data in an array on another array
Nearline	Data protection data that can be accessed “online” typically only by authorized specialists
Offline	Data protection data that requires a manual process to put it back in a network-accessible state (such as inserting into a tape library)
Online	Production data that can be directly accessed by a user over a network; typically considered to have access to high performance disks, but that is not a necessity
Operational continuity	Proactive planning, provisioning, monitoring, and preventive maintenance to prevent a service-level impacting event for specific applications from occurring in the first place or to minimize the impact of such an event if it does occur
Operational recovery	The attempt to minimize the impact of a service-level impacting event for specific applications when one occurs
Physical data protection	Protection of data by preserving the physical integrity and functionality of the physical substrate upon which the data resides, travels, or is processed
Point-in-time copy	A “copy” of a pool of data “frozen” (i.e., made unchangeable) at a chosen instant of time
Pooling	A collection of information that is managed as a homogenous whole for quality of service (QoS)

	purposes
RAID (redundant array of independent disks)	The ability to use one or more disks than is necessary for the actual data itself as a buffer against the failure of one (and possibly more) disks
Removable disk	The ability to remove a RAID group of disk drives (which also contain the disk media) as a whole or the ability to separate the disk platters themselves from the disk drives
Replication	Creating a physical copy of data apart from the original copy
RPO (recovery point objective)	The difference between the time when a failure occurs and the previous time when a set of data was available (such as a tape from the previous day) to which a recovery is made results in a potential permanent loss of all changes to data for the intervening time
RTO (recovery time objective)	The time required to return an application to a working state after a downtime situation occurs
Sea change	A marked transformation over time
Snapshot copy	A point-in-time copy
Storage pool	A mapping of a pool of information to a storage tier
Synchronous remote mirroring	Remote mirroring where the source and target pools of information are identical
Tiering	The separation of storage into classes by the characteristics of the storage itself
Vaulting	Typically the movement of data on tapes from a target site to a protected remote site
Virtual tape	Making disk appear as a piece of virtual tape media (not a tape library) so that data can be more efficiently written to tape media
Virtual tape library	Use of disk as if it were a tape library through a process of creating virtual tape drives on disk
WORM (write once, read many)	The ability to write only once to a piece of media, but the ability to read that media as often as necessary

Author Profile

David Hill, Principal, Mesabi Group

Through his writing, speaking, and research, David Hill has become a recognized thought leader in the field of storage and storage management. He looks at how enterprises can best go about adopting new and improved storage and storage management policies, practices, and technologies that not only meet immediate requirements, but also help position themselves for future growth.

Hill created the Mesabi Group to focus on the underlying dynamics that will drive the evolution of the storage infrastructure to better serve business. Two key areas of focus are storage networking and information lifecycle management. At the core of both is how policy-driven storage management software will evolve the intelligence that will be necessary to manage both data and the storage that data resides upon more efficiently and effectively.

Hill's perspective is how storage intersects with business requirements. For example, the key business function of risk management intersects with information technology through business continuity. Data protection is an essential and key component of business continuity.

As part of his research, Hill has done online presentations for UBS Investment Research on the structural dynamics of the storage market (also available in a Commentary entitled "You Can't Tell the Major Enterprise Storage Players Without a Scorecard"), the role of iSCSI in storage networking, and the direction in which information lifecycle management is heading.

Prior to founding Mesabi Group, Hill was an industry analyst at the Aberdeen Group for a number of years. As the Vice President of Storage Research and founder of the Storage & Storage Management practice, David Hill emphasized how leading enterprises could leverage their enterprise-wide storage investment to derive additional business value that ranges from TCO/ROI advantages to competitive advantage. Hill led both quantitative and qualitative market research studies during his tenure at Aberdeen Group. For example, he authored a report on "Storage Execution in a Time of Scarcity" that was done in conjunction with a survey distributed by InfoStor magazine.

Before Aberdeen, Hill spent many years at Data General where, among other activities, he directed Data General's internal IT data centers as well as managed the introduction of new analytical tools and business systems. While at EMC, he carried out strategic marketing, competitive analysis, sales force planning, and market fore-

casting. He has an advanced degree from the Sloan School at the Massachusetts Institute of Technology.

Hill is frequently quoted in such publications as Business Week, Computerworld, E Commerce Times, InfoStor, Network World, and Storage Networking World Online. In 2005, DStar named him one of the leading storage analysts.

David can be reached at david.hill@valleyviewventures.com.