



Information Storage and Management

Storing, Managing, and Protecting Digital Information

Managing and securing information is critical to business success. While information storage and management used to be a relatively straightforward and routine operation in the past, today it has developed into a highly mature and sophisticated pillar of information technology. Information storage and management technologies provide a variety of solutions for storing, managing, networking, accessing, protecting, securing, sharing, and optimizing information.

To keep pace with the exponential growth of information and the associated increase in sophistication and complexity of information management technology, there is a growing need for skilled information management professionals. More than ever, IT managers are challenged with employing and developing highly skilled information storage professionals.

This book covers concepts, principles, and deployment considerations across all technologies that are used for storing and managing information. It gives insight into:

- Challenges and solutions for data storage and data management
- Intelligent storage systems
- Storage networking (FC-SAN, IP-SAN, NAS)
- Backup, recovery, and archive (including CAS)
- Business continuity and disaster recovery
- Storage security and virtualization
- Managing and monitoring the storage infrastructure



EMC Proven Professional is the premier certification program in the information storage and management industry. Being proven means investing in yourself and formally validating your knowledge, skills, and expertise by the industry's most comprehensive learning and certification program.

This book helps you prepare for Information Storage and Management exam E20-001 leading to EMC Proven Professional Associate certification. Please visit <http://education.emc.com> for details.

EMC Corporation (NYSE: EMC) is the world's leading developer and provider of information infrastructure technology and solutions that enable organizations of all sizes to transform the way they compete and create value from their information. Information about EMC's products and services can be found at www.EMC.com.

Visit our website at www.wiley.com/compbooks.



COMPUTERS/
Networking/General



Information Storage and Management
Storing, Managing, and Protecting Digital Information

Information Storage and Management

Storing, Managing, and Protecting Digital Information

IS&M

EMC Education Services

Ready to add the ISM book to your reference library?
To order your copy, visit
<http://education.EMC.com/ismbook>



Information Storage and Management

**Storing, Managing, and Protecting
Digital Information**

Edited by
G. Somasundaram
Alok Shrivastava
EMC Education Services



WILEY

Wiley Publishing, Inc.



Foreword

Ralph Waldo Emerson, the great American essayist, philosopher, and poet, once said that the invariable mark of wisdom is seeing the miraculous in the common. Today, common miracles surround us, and it is virtually impossible *not* to see them. Most of us have modern gadgetry such as digital cameras, video camcorders, cell phones, fast computers that can access millions of websites, instant messaging, social networking sites, search engines, music downloads ... the list goes on. All of these examples have one thing in common: they generate huge volumes of data. Not only are we in an information age, we're in an age where information is exploding into a digital universe that requires enhanced technology and a new generation of professionals who are able to manage, leverage, and optimize *storage and information management* solutions.

Just to give you an idea of the challenges we face today, in one year the amount of digital information created, captured, and replicated is millions of times the amount of information in all the books ever written. Information is the most important asset of a business. To realize the inherent power of information, it must be intelligently and efficiently stored, protected, and managed—so that it can be made accessible, searchable, shareable, and, ultimately, actionable.

We are currently in the perfect storm. Everything is increasing: the information, the costs, and the skilled professionals needed to store and manage it—professionals who are not available in sufficient numbers to meet the growing need. The IT manager's number one concern is how to manage this storage growth. Enterprises simply cannot purchase bigger and better “boxes” to store their data. IT managers must not only worry about budgets for storage technology, but also be concerned with energy-efficient, footprint-reducing technology that is easy to install, manage, and use. Although many IT managers intend to

hire more trained staff, they are facing a shortage of skilled, storage-educated professionals who can take control of managing and optimizing the data.

I was unable to find a comprehensive book in the marketplace that provided insight into the various technologies deployed to store and manage information. As an industry leader, we have the subject-matter expertise and practical experience to help fill this gap; and now this book can give you a behind-the-scenes view of the technologies used in information storage and management. You will learn where data goes, how it is managed, and how you can contribute to your company's profitability.

If you've chosen storage and information infrastructure management as your career, you are a pioneer in a profession that is undergoing constant change, but one in which the challenges lead to great rewards.

Regardless of your current role in IT, this book should be a key part of your IT library and professional development.

Thomas P. Clancy
Vice President, Education Services, EMC Corporation
March 2009



Introduction

Information storage is a central pillar of information technology. A large quantity of digital information is being created every moment by individual and corporate consumers of IT. This information needs to be stored, protected, optimized, and managed.

Not long ago, information storage was seen as only a bunch of disks or tapes attached to the back of the computer to store data. Even today, only those in the storage industry understand the critical role that information storage technology plays in the availability, performance, integration, and optimization of the entire IT infrastructure. Over the last two decades, information storage has developed into a highly sophisticated technology, providing a variety of solutions for storing, managing, connecting, protecting, securing, sharing, and optimizing digital information.

With the exponential growth of information and the development of sophisticated products and solutions, there is also a growing need for information storage professionals. IT managers are challenged by the ongoing task of employing and developing highly skilled information storage professionals.

Many leading universities and colleges have started to include storage technology courses in their regular computer technology or information technology curriculum, yet many of today's IT professionals, even those with years of experience, may not have benefited from this formal education, meaning many seasoned professionals—including application, systems, database, and network administrators—do not share a common foundation about how storage technology affects their areas of expertise.

This book is designed and developed to enable professionals and students to achieve a comprehensive understanding of all segments of storage technology. While the product examples used in the book are from EMC Corporation, an

understanding of the technology concepts and principles prepare the reader to easily understand products from various technology vendors.

This book has 16 chapters, organized in four sections. Advanced topics build upon the topics learned in previous chapters.

Part 1, “Information Storage and Management for Today’s World”: These four chapters cover information growth and challenges, define a storage system and its environment, review the evolution of storage technology, and introduce intelligent storage systems.

Part 2, “Storage Options and Protocols”: These six chapters cover the SCSI and Fibre channel architecture, direct-attached storage (DAS), storage area networks (SANs), network-attached storage (NAS), Internet Protocol SAN (IP-SAN), content-addressed storage (CAS), and storage virtualization.

Part 3, “Business Continuity and Replication”: These four chapters introduce business continuity, backup and recovery, local data replication, and remote data replication.

Part 4, “Security and Administration”: These two chapters cover storage security and storage infrastructure monitoring and management.

This book has a supplementary website that provides additional up-to-date learning aids and reading material. Visit <http://education.EMC.com/ismbook> for details.

EMC Academic Alliance

Universities and colleges interested in offering an *information storage and management* curriculum are invited to join the Academic Alliance program. This program provides comprehensive support to institutes, including teaching aids, faculty guides, student projects, and more. Please visit <http://education.EMC.com/academicalliance>.

EMC Proven Professional Certification



This book prepares students and professionals to take the EMC Proven Professional Information Storage and Management exam E20-001. EMC Proven Professional is the premier certification program that validates your knowledge and helps establish your credibility in the information technology industry. For more information on certification as well as to access practice exams, visit <http://education.EMC.com>.



Contents

Foreword	xix
Introduction	xxi
Section I Storage System	1
Chapter 1 Introduction to Information Storage and Management	3
1.1 Information Storage	5
1.1.1 Data	5
1.1.2 Types of Data	7
1.1.3 Information	7
1.1.4 Storage	8
1.2 Evolution of Storage Technology and Architecture	9
1.3 Data Center Infrastructure	10
1.3.1 Core Elements	10
1.3.2 Key Requirements for Data Center Elements	11
1.3.3 Managing Storage Infrastructure	13
1.4 Key Challenges in Managing Information	14
1.5 Information Lifecycle	14
1.5.1 Information Lifecycle Management	15
1.5.2 ILM Implementation	16
1.5.3 ILM Benefits	17
Summary	18
Chapter 2 Storage System Environment	21
2.1 Components of a Storage System Environment	21
2.1.1 Host	22
2.1.2 Connectivity	24
2.1.3 Storage	26

2.2 Disk Drive Components	27
2.2.1 Platter	28
2.2.2 Spindle	28
2.2.3 Read/Write Head	28
2.2.4 Actuator Arm Assembly	29
2.2.5 Controller	29
2.2.6 Physical Disk Structure	30
2.2.7 Zoned Bit Recording	31
2.2.8 Logical Block Addressing	32
2.3 Disk Drive Performance	33
2.3.1 Disk Service Time	33
2.4 Fundamental Laws Governing Disk Performance	35
2.5 Logical Components of the Host	38
2.5.1 Operating System	39
2.5.2 Device Driver	39
2.5.3 Volume Manager	39
2.5.4 File System	41
2.5.5 Application	44
2.6 Application Requirements and Disk Performance	45
Summary	48
Chapter 3 Data Protection: RAID	51
3.1 Implementation of RAID	52
3.1.1 Software RAID	52
3.1.2 Hardware RAID	52
3.2 RAID Array Components	53
3.3 RAID Levels	54
3.3.1 Striping	54
3.3.2 Mirroring	55
3.3.3 Parity	56
3.3.4 RAID 0	57
3.3.5 RAID 1	57
3.3.6 Nested RAID	59
3.3.7 RAID 3	60
3.3.8 RAID 4	61
3.3.9 RAID 5	61
3.3.10 RAID 6	62
3.4 RAID Comparison	62
3.5 RAID Impact on Disk Performance	65
3.5.1 Application IOPS and RAID Configurations	66
3.6 Hot Spares	67
Summary	67
Chapter 4 Intelligent Storage System	69
4.1 Components of an Intelligent Storage System	70
4.1.1 Front End	70
4.1.2 Cache	72
4.1.3 Back End	77
4.1.4 Physical Disk	77

4.2 Intelligent Storage Array	80
4.2.1 High-End Storage Systems	80
4.2.2 Midrange Storage System	81
4.3 Concepts in Practice: EMC CLARiiON and Symmetrix	82
4.3.1 CLARiiON Storage Array	83
4.3.2 CLARiiON CX4 Architecture	84
4.3.3 Managing the CLARiiON	86
4.3.4 Symmetrix Storage Array	87
4.3.5 Symmetrix Component Overview	89
4.3.6 Direct Matrix Architecture	91
Summary	93
Section II Storage Networking Technologies and Virtualization	95
Chapter 5 Direct-Attached Storage and Introduction to SCSI	97
5.1 Types of DAS	97
5.1.1 Internal DAS	98
5.1.2 External DAS	98
5.2 DAS Benefits and Limitations	99
5.3 Disk Drive Interfaces	99
5.3.1 IDE/ATA	99
5.3.2 SATA	100
5.3.3 Parallel SCSI	101
5.4 Introduction to Parallel SCSI	102
5.4.1 Evolution of SCSI	102
5.4.2 SCSI Interfaces	103
5.4.3 SCSI-3 Architecture	105
5.4.4 Parallel SCSI Addressing	109
5.5 SCSI Command Model	110
5.5.1 CDB Structure	110
5.5.2 Operation Code	110
5.5.3 Control Field	112
5.5.4 Status	112
Summary	113
Chapter 6 Storage Area Networks	115
6.1 Fibre Channel: Overview	116
6.2 The SAN and Its Evolution	117
6.3 Components of SAN	118
6.3.1 Node Ports	118
6.3.2 Cabling	120
6.3.3 Interconnect Devices	121
6.3.4 Storage Arrays	122
6.3.5 SAN Management Software	122
6.4 FC Connectivity	123
6.4.1 Point-to-Point	123
6.4.2 Fibre Channel Arbitrated Loop	124
6.4.3 Fibre Channel Switched Fabric	126

6.5 Fibre Channel Ports	128
6.6 Fibre Channel Architecture	130
6.6.1 Fibre Channel Protocol Stack	131
6.6.2 Fibre Channel Addressing	131
6.6.3 FC Frame	133
6.6.4. Structure and Organization of FC Data	135
6.6.5 Flow Control	135
6.6.6 Classes of Service	136
6.7 Zoning	136
6.8 Fibre Channel Login Types	139
6.9 FC Topologies	139
6.9.1 Core-Edge Fabric	140
6.9.2 Mesh Topology	142
6.10 Concepts in Practice: EMC Connectrix	143
Summary	146
Chapter 7 Network-Attached Storage	147
7.1 General-Purpose Servers vs. NAS Devices	148
7.2 Benefits of NAS	148
7.3 NAS File I/O	149
7.3.1 File Systems and Remote File Sharing	150
7.3.2 Accessing a File System	150
7.3.3 File Sharing	150
7.4 Components of NAS	151
7.5 NAS Implementations	152
7.5.1 Integrated NAS	152
7.5.2 Gateway NAS	153
7.5.3 Integrated NAS Connectivity	153
7.5.4 Gateway NAS Connectivity	154
7.6 NAS File-Sharing Protocols	155
7.6.1 NFS	156
7.6.2 CIFS	156
7.7 NAS I/O Operations	157
7.7.1 Hosting and Accessing Files on NAS	158
7.8 Factors Affecting NAS Performance and Availability	158
7.9 Concepts in Practice: EMC Celerra	162
7.9.1 Architecture	162
7.9.2 Celerra Product Family	165
Summary	166
Chapter 8 IP SAN	169
8.1 iSCSI	171
8.1.1 Components of iSCSI	171
8.1.2 iSCSI Host Connectivity	172
8.1.3 Topologies for iSCSI Connectivity	173
8.1.4 iSCSI Protocol Stack	174
8.1.5 iSCSI Discovery	175

	8.1.6 iSCSI Names	176
	8.1.7 iSCSI Session	177
	8.1.8 iSCSI PDU	178
	8.1.9 Ordering and Numbering	179
	8.1.10 iSCSI Error Handling and Security	180
	8.2 FCIP	181
	8.2.1 FCIP Topology	182
	8.2.2 FCIP Performance and Security	183
	Summary	184
Chapter 9	Content-Addressed Storage	187
	9.1 Fixed Content and Archives	188
	9.2 Types of Archives	189
	9.3 Features and Benefits of CAS	190
	9.4 CAS Architecture	191
	9.5 Object Storage and Retrieval in CAS	194
	9.6 CAS Examples	196
	9.6.1 Health Care Solution: Storing Patient Studies	196
	9.6.2 Finance Solution: Storing Financial Records	197
	9.7 Concepts in Practice: EMC Centera	198
	9.7.1 EMC Centera Models	199
	9.7.2 EMC Centera Architecture	199
	9.7.3 Centera Tools	201
	9.7.4 EMC Centera Universal Access	202
	Summary	203
Chapter 10	Storage Virtualization	205
	10.1 Forms of Virtualization	205
	10.1.1 Memory Virtualization	206
	10.1.2 Network Virtualization	206
	Virtual SAN (VSAN)	207
	10.1.3 Server Virtualization	207
	10.1.4 Storage Virtualization	208
	10.2 SNIA Storage Virtualization Taxonomy	210
	10.3 Storage Virtualization Configurations	211
	10.4 Storage Virtualization Challenges	212
	10.4.1 Scalability	213
	10.4.2 Functionality	213
	10.4.3 Manageability	213
	10.4.4 Support	214
	10.5 Types of Storage Virtualization	214
	10.5.1 Block-Level Storage Virtualization	214
	10.5.2 File-Level Virtualization	215
	10.6 Concepts in Practice	217
	10.6.1 EMC Invista	217
	10.6.2 Rainfinity	220
	Summary	223

Section III	Business Continuity	225
Chapter 11	Introduction to Business Continuity	227
	11.1 Information Availability	228
	11.1.1 Causes of Information Unavailability	228
	11.1.2 Measuring Information Availability	229
	11.1.4 Consequences of Downtime	230
	11.2 BC Terminology	231
	11.3 BC Planning Lifecycle	233
	11.4 Failure Analysis	236
	11.4.1 Single Point of Failure	236
	11.4.2 Fault Tolerance	236
	11.4.3 Multipathing Software	238
	11.5 Business Impact Analysis	238
	11.6 BC Technology Solutions	239
	11.7 Concept in Practice: EMC PowerPath	239
	11.7.1 PowerPath Features	240
	11.7.2 Dynamic Load Balancing	240
	11.7.3 Automatic Path Failover	242
	Summary	245
Chapter 12	Backup and Recovery	249
	12.1 Backup Purpose	250
	12.1.1 Disaster Recovery	250
	12.1.2 Operational Backup	250
	12.1.3 Archival	250
	12.2 Backup Considerations	251
	12.3 Backup Granularity	252
	12.4 Recovery Considerations	255
	12.5 Backup Methods	256
	12.6 Backup Process	257
	12.7 Backup and Restore Operations	258
	12.8 Backup Topologies	260
	12.8.1 Serverless Backup	263
	12.9 Backup in NAS Environments	263
	12.10 Backup Technologies	267
	12.10.1 Backup to Tape	267
	12.10.2 Physical Tape Library	268
	12.10.3 Backup to Disk	270
	12.10.4 Virtual Tape Library	271
	12.11 Concepts in Practice: EMC NetWorker	274
	12.11.1 NetWorker Backup Operation	275
	12.11.2 NetWorker Recovery	276
	Summary	278
Chapter 13	Local Replication	281
	13.1 Source and Target	282
	13.2 Uses of Local Replicas	282

13.3 Data Consistency	283
13.3.1 Consistency of a Replicated File System	283
13.3.2 Consistency of a Replicated Database	284
13.4 Local Replication Technologies	286
13.4.1 Host-Based Local Replication	286
13.4.2 Storage Array–Based Replication	288
13.5 Restore and Restart Considerations	295
13.5.1 Tracking Changes to Source and Target	296
13.6 Creating Multiple Replicas	298
13.8 Management Interface	299
13.9 Concepts in Practice: EMC TimeFinder and EMC SnapView	299
13.9.1 TimeFinder/Clone	300
13.9.2 TimeFinder/Mirror	300
13.9.3 EMC SnapView	302
Summary	304
Chapter 14 Remote Replication	307
14.1 Modes of Remote Replication	307
14.2 Remote Replication Technologies	309
14.2.1. Host-Based Remote Replication	309
14.2.2 Storage Array–Based Remote Replication	312
Three-Site Replication	316
14.2.3 SAN-Based Remote Replication	319
14.3 Network Infrastructure	322
14.3.1 DWDM	322
14.3.2 SONET	322
14.4 Concepts in Practice: EMC SRDF, EMC SAN Copy, and EMC MirrorView	323
14.6.1 SRDF Family	323
14.6.2 Disaster Recovery with SRDF	324
14.6.3 SRDF Operations for Concurrent Access	325
14.6.4 EMC SAN Copy	326
14.6.5 EMC MirrorView	327
Summary	328
Section IV Storage Security and Management	331
Chapter 15 Securing the Storage Infrastructure	333
15.1 Storage Security Framework	333
15.2 Risk Triad	334
15.2.1 Assets	335
15.2.2 Threats	336
15.2.3 Vulnerability	337
15.3 Storage Security Domains	338
15.3.1 Securing the Application Access Domain	339
15.3.2 Securing the Management Access Domain	342
15.3.3 Securing Backup, Recovery, and Archive (BURA)	345

15.4 Security Implementations in Storage Networking	346
15.4.1 SAN	346
15.4.2 NAS	351
15.4.3 IP SAN	357
Summary	358
Chapter 16 Managing the Storage Infrastructure	361
16.1 Monitoring the Storage Infrastructure	362
16.1.1 Parameters Monitored	362
16.1.2 Components Monitored	363
16.1.3 Monitoring Examples	366
16.1.4 Alerts	372
16.2 Storage Management Activities	373
16.2.1 Availability management	373
16.2.2 Capacity management	373
16.2.3 Performance management	374
16.2.4 Security Management	374
16.2.5 Reporting	374
16.2.6 Storage Management Examples	375
16.3 Storage Infrastructure Management Challenges	380
16.4 Developing an Ideal Solution	380
16.4.1 Storage Management Initiative	381
16.4.2 Enterprise Management Platforms	383
16.5 Concepts in Practice: EMC ControlCenter Navisphere Manager	384
Summary	390
Index	437

Section

Storage System

In This Section

Chapter 1: Introduction to Information Storage and Management

Chapter 2: Storage System Environment

Chapter 3: Data Protection: RAID

Chapter 4: Intelligent Storage Systems

Chapter 1

Introduction to Information Storage and Management

Information is increasingly important in our daily lives. We have become information dependents of the twenty-first century, living in an on-command, on-demand world that means we need information when and where it is required. We access the Internet every day to perform searches, participate in social networking, send and receive e-mails, take pictures and videos through digital cameras, and satisfy many other personal and professional needs. Equipped with a growing number of content-generating devices, more information is being created by individuals than by businesses. Information created by individuals gains value when shared with others. When created, information resides locally on devices such as cell phones, cameras, and laptops. To share this information, it needs to be uploaded via networks to data centers. It is interesting to note that while the majority of information is created by individuals, it is stored and managed by a relatively small number of organizations. Figure 1-1 depicts this virtuous cycle of information.

The importance, dependency, and volume of information for the business world also continue to grow at astounding rates. Businesses depend on fast and reliable access to information critical to their success. Some of the business applications that process information include airline reservations, telephone billing systems, e-commerce, ATMs, product designs, inventory management, e-mail archives, Web portals, patient records, credit cards, life sciences, and global capital markets.

KEY CONCEPTS

Data and Information

Structured and Unstructured Data

Storage Technology Architectures

Core Elements of a Data Center

Information Management

Information Lifecycle Management

The increasing criticality of information to the businesses has amplified the challenges in protecting and managing the data. The volume of data that business must manage has driven strategies to classify data according to its value and create rules for the treatment of this data over its life cycle. These strategies not only provide financial and regulatory benefits at the business level, but also manageability benefits at operational levels to the organization.

Data centers now view information storage as one of their core elements, along with applications, databases, operating systems, and networks. Storage technology continues to evolve with technical advancements offering increasingly higher levels of availability, security, scalability, performance, integrity, capacity, and manageability.

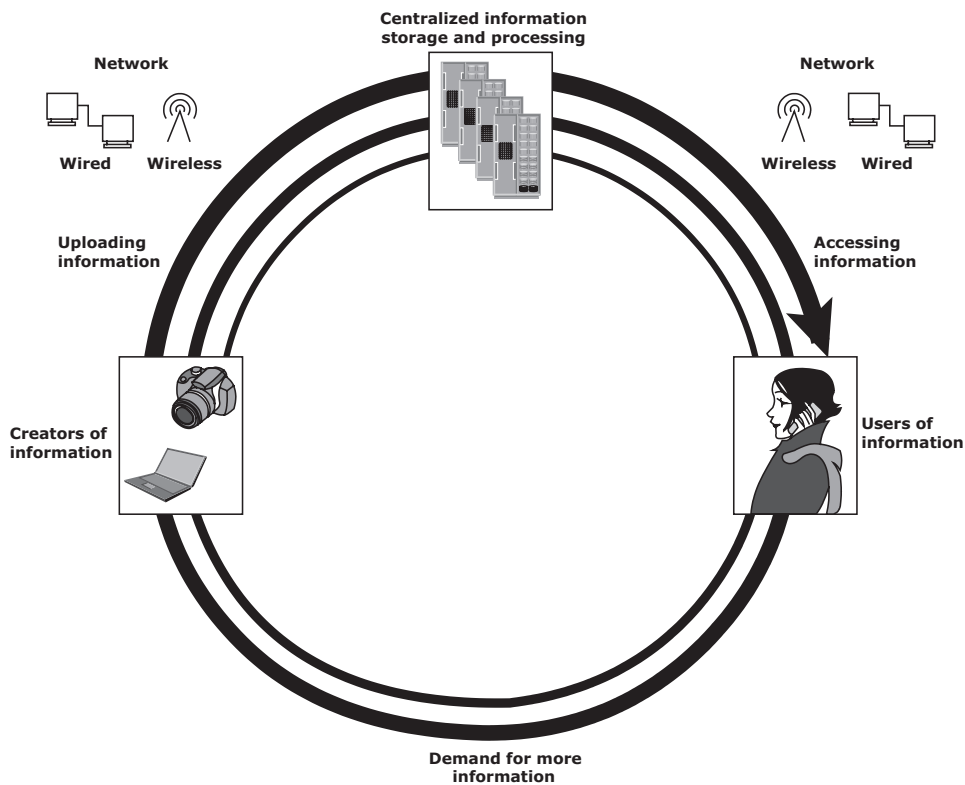


Figure 1-1: Virtuous cycle of information

This chapter describes the evolution of information storage architecture from simple direct-attached models to complex networked topologies. It introduces the information lifecycle management (ILM) strategy, which aligns the information technology (IT) infrastructure with business priorities.

1.1 Information Storage

Businesses use data to derive information that is critical to their day-to-day operations. Storage is a repository that enables users to store and retrieve this digital data.

1.1.1 Data

Data is a collection of raw facts from which conclusions may be drawn. Handwritten letters, a printed book, a family photograph, a movie on video tape, printed and duly signed copies of mortgage papers, a bank's ledgers, and an account holder's passbooks are all examples of data.

Before the advent of computers, the procedures and methods adopted for data creation and sharing were limited to fewer forms, such as paper and film. Today, the same data can be converted into more convenient forms such as an e-mail message, an e-book, a bitmapped image, or a digital movie. This data can be generated using a computer and stored in strings of 0s and 1s, as shown in Figure 1-2. Data in this form is called *digital data* and is accessible by the user only after it is processed by a computer.

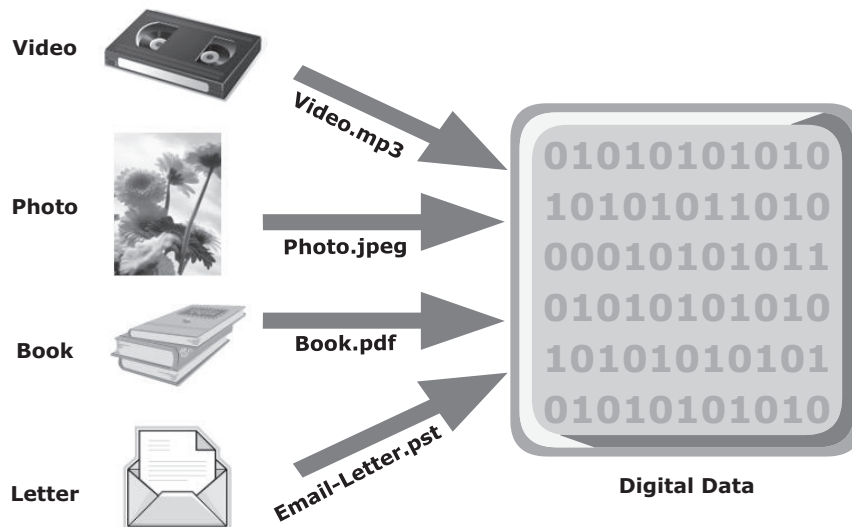


Figure 1-2: Digital data

With the advancement of computer and communication technologies, the rate of data generation and sharing has increased exponentially. The following is a list of some of the factors that have contributed to the growth of digital data:

- **Increase in data processing capabilities:** Modern-day computers provide a significant increase in processing and storage capabilities. This enables the conversion of various types of content and media from conventional forms to digital formats.
- **Lower cost of digital storage:** Technological advances and decrease in the cost of storage devices have provided low-cost solutions and encouraged the development of less expensive data storage devices. This cost benefit has increased the rate at which data is being generated and stored.
- **Affordable and faster communication technology:** The rate of sharing digital data is now much faster than traditional approaches. A handwritten letter may take a week to reach its destination, whereas it only takes a few seconds for an e-mail message to reach its recipient.

Inexpensive and easier ways to create, collect, and store all types of data, coupled with increasing individual and business needs, have led to accelerated data growth, popularly termed the *data explosion*. Data has different purposes and criticality, so both individuals and businesses have contributed in varied proportions to this data explosion.

The importance and the criticality of data vary with time. Most of the data created holds significance in the short-term but becomes less valuable over time. This governs the type of data storage solutions used. Individuals store data on a variety of storage devices, such as hard disks, CDs, DVDs, or Universal Serial Bus (USB) flash drives.

EXAMPLE OF RESEARCH AND BUSINESS DATA



- **Seismology:** Involves collecting data related to various sources and parameters of earthquakes, and other relevant data that needs to be processed to derive meaningful information.
- **Product data:** Includes data related to various aspects of a product, such as inventory, description, pricing, availability, and sales.
- **Customer data:** A combination of data related to a company's customers, such as order details, shipping addresses, and purchase history.
- **Medical data:** Data related to the health care industry, such as patient history, radiological images, details of medication and other treatment, and insurance information.

Businesses generate vast amounts of data and then extract meaningful information from this data to derive economic benefits. Therefore, businesses need to maintain data and ensure its availability over a longer period.

Furthermore, the data can vary in criticality and may require special handling. For example, legal and regulatory requirements mandate that banks maintain account information for their customers accurately and securely. Some businesses handle data for millions of customers, and ensures the security and integrity of data over a long period of time. This requires high-capacity storage devices with enhanced security features that can retain data for a long period.

1.1.2 Types of Data

Data can be classified as structured or unstructured (see Figure 1-3) based on how it is stored and managed. Structured data is organized in rows and columns in a rigidly defined format so that applications can retrieve and process it efficiently. Structured data is typically stored using a database management system (DBMS).

Data is unstructured if its elements cannot be stored in rows and columns, and is therefore difficult to query and retrieve by business applications. For example, customer contacts may be stored in various forms such as sticky notes, e-mail messages, business cards, or even digital format files such as .doc, .txt, and .pdf. Due its unstructured nature, it is difficult to retrieve using a customer relationship management application. Unstructured data may not have the required components to identify itself uniquely for any type of processing or interpretation. Businesses are primarily concerned with managing unstructured data because over 80 percent of enterprise data is unstructured and requires significant storage space and effort to manage.

1.1.3 Information

Data, whether structured or unstructured, does not fulfill any purpose for individuals or businesses unless it is presented in a meaningful form. Businesses need to analyze data for it to be of value. *Information* is the intelligence and knowledge derived from data.

Businesses analyze raw data in order to identify meaningful trends. On the basis of these trends, a company can plan or modify its strategy. For example, a retailer identifies customers' preferred products and brand names by analyzing their purchase patterns and maintaining an inventory of those products.

Effective data analysis not only extends its benefits to existing businesses, but also creates the potential for new business opportunities by using the information in creative ways. Job portal is an example. In order to reach a wider set of prospective employers, job seekers post their résumés on various websites offering job search facilities. These websites collect the resumes and post them on centrally accessible locations for prospective employers. In addition, companies post available positions on job search sites. Job-matching software matches keywords from

résumés to keywords in job postings. In this manner, the job search engine uses data and turns it into information for employers and job seekers.

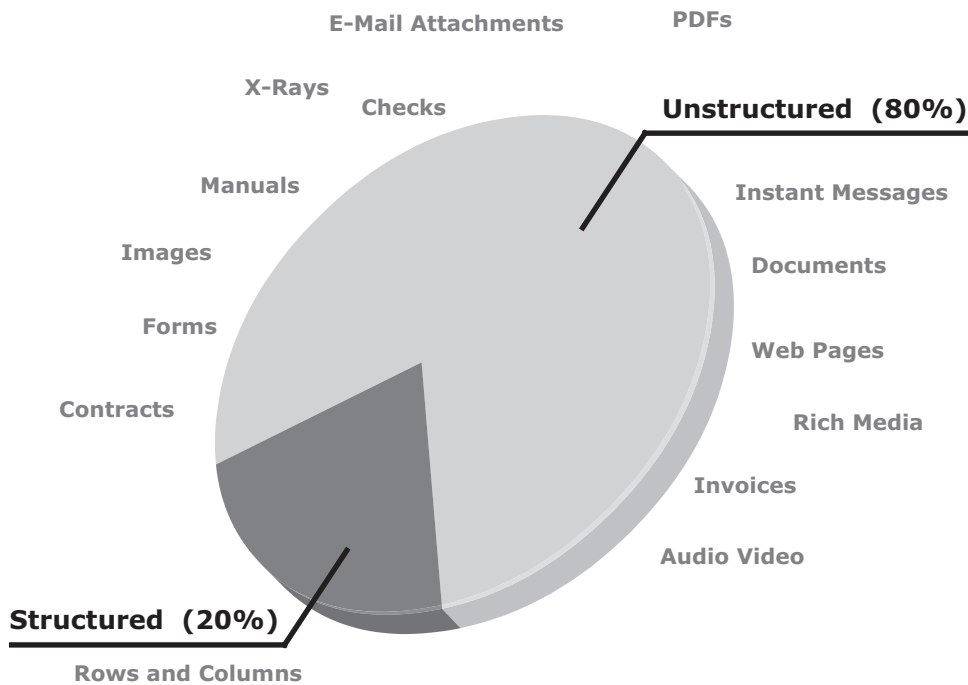


Figure 1-3: Types of data

Because information is critical to the success of a business, there is an ever-present concern about its availability and protection. Legal, regulatory, and contractual obligations regarding the availability and protection of data only add to these concerns. Outages in key industries, such as financial services, telecommunications, manufacturing, retail, and energy cost millions of U.S. dollars per hour.

1.1.4 Storage

Data created by individuals or businesses must be stored so that it is easily accessible for further processing. In a computing environment, devices designed for storing data are termed *storage devices* or simply *storage*. The type of storage used varies based on the type of data and the rate at which it is created and used. Devices such as memory in a cell phone or digital camera, DVDs, CD-ROMs, and hard disks in personal computers are examples of storage devices.

Businesses have several options available for storing data including internal hard disks, external disk arrays and tapes.

1.2 Evolution of Storage Technology and Architecture

Historically, organizations had centralized computers (mainframe) and information storage devices (tape reels and disk packs) in their data center. The evolution of open systems and the affordability and ease of deployment that they offer made it possible for business units/departments to have their own servers and storage. In earlier implementations of open systems, the storage was typically internal to the server.

The proliferation of departmental servers in an enterprise resulted in unprotected, unmanaged, fragmented islands of information and increased operating cost. Originally, there were very limited policies and processes for managing these servers and the data created. To overcome these challenges, storage technology evolved from non-intelligent internal storage to intelligent networked storage (see Figure 1-4). Highlights of this technology evolution include:

- **Redundant Array of Independent Disks (RAID):** This technology was developed to address the performance and availability requirements of data. It continues to evolve today and is used in all storage architectures such as DAS, SAN, and so on.
- **Direct-attached storage (DAS):** This type of storage connects directly to a server (host) or a group of servers in a cluster. Storage can be either internal or external to the server. External DAS alleviated the challenges of limited internal storage capacity.
- **Storage area network (SAN):** This is a dedicated, high-performance *Fibre Channel (FC)* network to facilitate *block-level* communication between servers and storage. Storage is partitioned and assigned to a server for accessing its data. SAN offers scalability, availability, performance, and cost benefits compared to DAS.
- **Network-attached storage (NAS):** This is dedicated storage for *file serving* applications. Unlike a SAN, it connects to an existing communication network (LAN) and provides file access to heterogeneous clients. Because it is purposely built for providing storage to file server applications, it offers higher scalability, availability, performance, and cost benefits compared to general purpose file servers.
- **Internet Protocol SAN (IP-SAN):** One of the latest evolutions in storage architecture, IP-SAN is a convergence of technologies used in SAN and NAS. IP-SAN provides block-level communication across a local or wide area network (LAN or WAN), resulting in greater consolidation and availability of data.

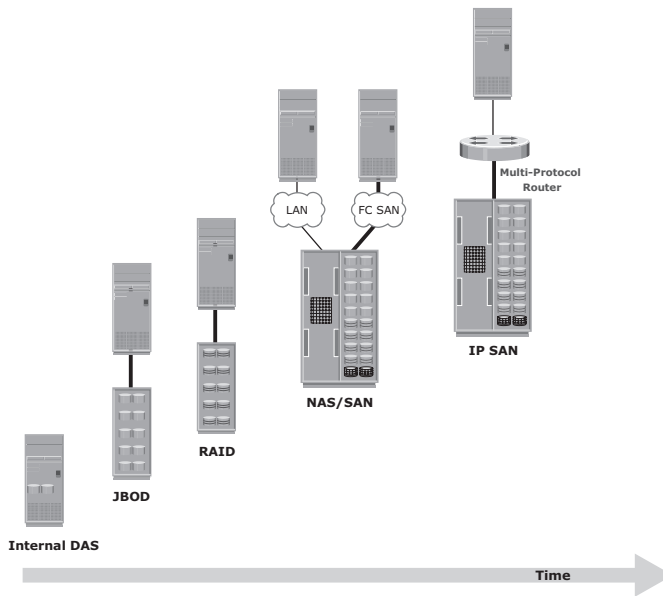


Figure 1-4: Evolution of storage architectures

Storage technology and architecture continues to evolve, which enables organizations to consolidate, protect, optimize, and leverage their data to achieve the highest return on information assets.

1.3 Data Center Infrastructure

Organizations maintain data centers to provide centralized data processing capabilities across the enterprise. Data centers store and manage large amounts of mission-critical data. The data center infrastructure includes computer storage systems, network devices, dedicated power backups, and environmental controls (such as air conditioning and fire suppression).

Large organizations often maintain more than one data center to distribute data processing workloads and provide backups in the event of a disaster. The storage requirements of a data center are met by a combination of various storage architectures.

1.3.1 Core Elements

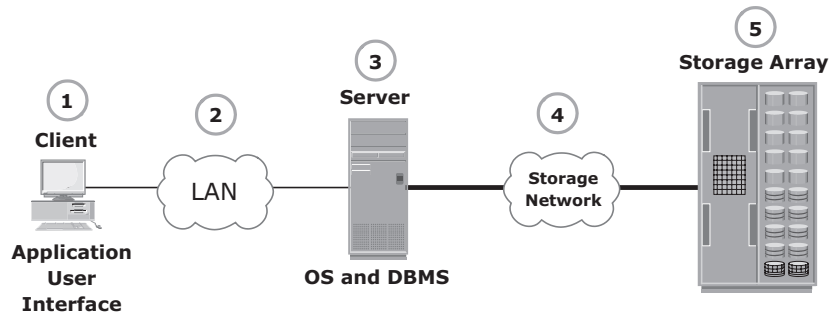
Five core elements are essential for the basic functionality of a data center:

- Application:** An application is a computer program that provides the logic for computing operations. Applications, such as an order processing system, can be layered on a database, which in turn uses operating system services to perform read/write operations to storage devices.

- **Database:** More commonly, a database management system (DBMS) provides a structured way to store data in logically organized tables that are interrelated. A DBMS optimizes the storage and retrieval of data.
- **Server and operating system:** A computing platform that runs applications and databases.
- **Network:** A data path that facilitates communication between clients and servers or between servers and storage.
- **Storage array:** A device that stores data persistently for subsequent use.

These core elements are typically viewed and managed as separate entities, but all the elements must work together to address data processing requirements.

Figure 1-5 shows an example of an order processing system that involves the five core elements of a data center and illustrates their functionality in a business process.



- 1 A customer places an order through the AUI of the order processing application software located on the client computer.
- 2 The client connects to the server over the LAN and accesses the DBMS located on the server to update the relevant information such as the customer name, address, payment method, products ordered, and quantity ordered.
- 3 The DBMS uses the server operating system to read and write this data to the database located on physical disks in the storage array.
- 4 The Storage Network provides the communication link between the server and the storage array and transports the read or write commands between them.
- 5 The storage array, after receiving the read or write commands from the server, performs the necessary operations to store the data on physical disks.

Figure 1-5: Example of an order processing system

1.3.2 Key Requirements for Data Center Elements

Uninterrupted operation of data centers is critical to the survival and success of a business. It is necessary to have a reliable infrastructure that ensures data is accessible at all times. While the requirements, shown in Figure 1-6, are applicable to all elements of the data center infrastructure, our focus here is on storage

systems. The various technologies and solutions to meet these requirements are covered in this book.

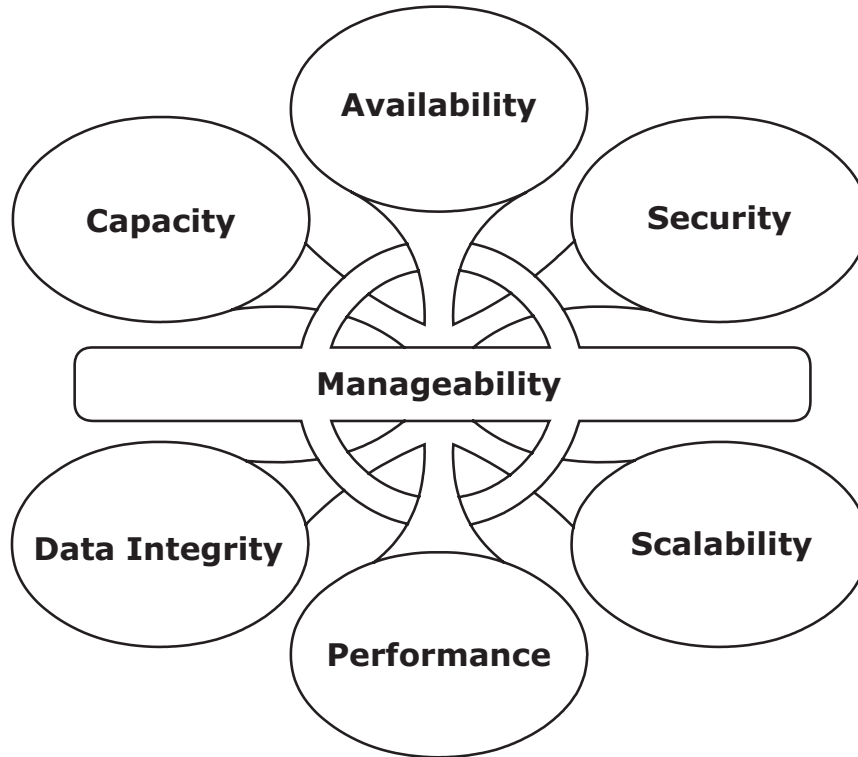


Figure 1-6: Key characteristics of data center elements

- **Availability:** All data center elements should be designed to ensure accessibility. The inability of users to access data can have a significant negative impact on a business.
- **Security:** Policies, procedures, and proper integration of the data center core elements that will prevent unauthorized access to information must be established. In addition to the security measures for client access, specific mechanisms must enable servers to access only their allocated resources on storage arrays.
- **Scalability:** Data center operations should be able to allocate additional processing capabilities or storage on demand, without interrupting business operations. Business growth often requires deploying more servers, new applications, and additional databases. The storage solution should be able to grow with the business.

- **Performance:** All the core elements of the data center should be able to provide optimal performance and service all processing requests at high speed. The infrastructure should be able to support performance requirements.
- **Data integrity:** Data integrity refers to mechanisms such as error correction codes or parity bits which ensure that data is written to disk exactly as it was received. Any variation in data during its retrieval implies corruption, which may affect the operations of the organization.
- **Capacity:** Data center operations require adequate resources to store and process large amounts of data efficiently. When capacity requirements increase, the data center must be able to provide additional capacity without interrupting availability, or, at the very least, with minimal disruption. Capacity may be managed by reallocation of existing resources, rather than by adding new resources.
- **Manageability:** A data center should perform all operations and activities in the most efficient manner. Manageability can be achieved through automation and the reduction of human (manual) intervention in common tasks.

1.3.3 Managing Storage Infrastructure

Managing a modern, complex data center involves many tasks. Key management activities include:

- *Monitoring* is the continuous collection of information and the review of the entire data center infrastructure. The aspects of a data center that are monitored include security, performance, accessibility, and capacity.
- *Reporting* is done periodically on resource performance, capacity, and utilization. Reporting tasks help to establish business justifications and chargeback of costs associated with data center operations.
- *Provisioning* is the process of providing the hardware, software, and other resources needed to run a data center. Provisioning activities include capacity and resource planning. *Capacity planning* ensures that the user's and the application's future needs will be addressed in the most cost-effective and controlled manner. *Resource planning* is the process of evaluating and identifying required resources, such as personnel, the facility (site), and the technology. Resource planning ensures that adequate resources are available to meet user and application requirements.

For example, the utilization of an application's allocated storage capacity may be monitored. As soon as utilization of the storage capacity reaches a critical

value, additional storage capacity may be provisioned to the application. If utilization of the storage capacity is properly monitored and reported, business growth can be understood and future capacity requirements can be anticipated. This helps to frame a proactive data management policy.

1.4 Key Challenges in Managing Information

In order to frame an effective information management policy, businesses need to consider the following key challenges of information management:

- **Exploding digital universe:** The rate of information growth is increasing exponentially. Duplication of data to ensure high availability and repurposing has also contributed to the multifold increase of information growth.
- **Increasing dependency on information:** The strategic use of information plays an important role in determining the success of a business and provides competitive advantages in the marketplace.
- **Changing value of information:** Information that is valuable today may become less important tomorrow. The value of information often changes over time.

Framing a policy to meet these challenges involves understanding the value of information over its lifecycle.

1.5 Information Lifecycle

The *information lifecycle* is the “change in the value of information” over time. When data is first created, it often has the highest value and is used frequently. As data ages, it is accessed less frequently and is of less value to the organization. Understanding the information lifecycle helps to deploy appropriate storage infrastructure, according to the changing value of information.

For example, in a sales order application, the value of the information changes from the time the order is placed until the time that the warranty becomes void (see Figure 1-7). The value of the information is highest when a company receives a new sales order and processes it to deliver the product. After order fulfillment, the customer or order data need not be available for real-time access. The company can transfer this data to less expensive secondary storage with lower accessibility and availability requirements unless or until a warranty claim or another event triggers its need. After the warranty becomes void, the company can archive or dispose of data to create space for other high-value information.

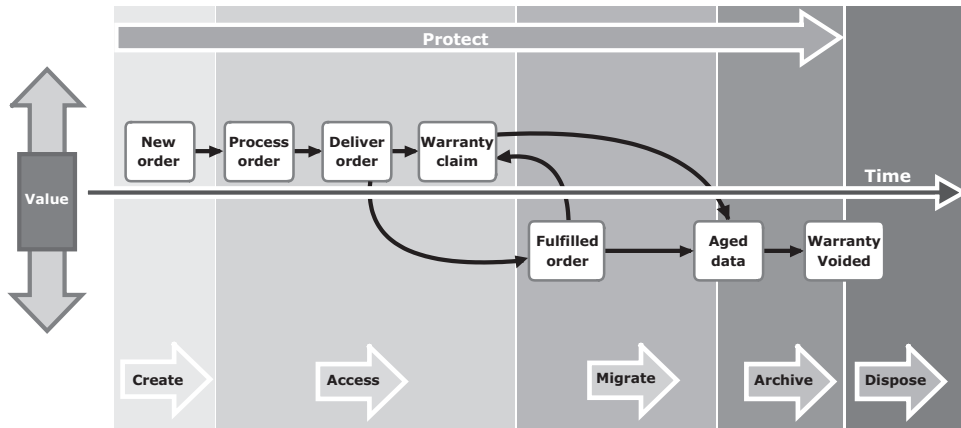


Figure 1-7: Changing value of sales order information

1.5.1 Information Lifecycle Management

Today's business requires data to be protected and available 24×7 . Data centers can accomplish this with the optimal and appropriate use of storage infrastructure. An effective information management policy is required to support this infrastructure and leverage its benefits.

Information lifecycle management (ILM) is a proactive strategy that enables an IT organization to effectively manage the data throughout its lifecycle, based on predefined business policies. This allows an IT organization to optimize the storage infrastructure for maximum return on investment. An ILM strategy should include the following characteristics:

- **Business-centric:** It should be integrated with key processes, applications, and initiatives of the business to meet both current and future growth in information.
- **Centrally managed:** All the information assets of a business should be under the purview of the ILM strategy.
- **Policy-based:** The implementation of ILM should not be restricted to a few departments. ILM should be implemented as a policy and encompass all business applications, processes, and resources.
- **Heterogeneous:** An ILM strategy should take into account all types of storage platforms and operating systems.
- **Optimized:** Because the value of information varies, an ILM strategy should consider the different storage requirements and allocate storage resources based on the information's value to the business.

- **Tiered Storage:** Tiered storage is an approach to define different storage levels in order to reduce total storage cost. Each tier has different levels of protection, performance, data access frequency, and other considerations. Information is stored and moved between different tiers based on its value over time. For example, mission-critical, most accessed information may be stored on Tier 1 storage, which consists of high performance media with a highest level of protection. Medium accessed and other important data is stored on Tier 2 storage, which may be on less expensive media with moderate performance and protection. Rarely accessed or event specific information may be stored on lower tiers of storage.

1.5.2 ILM Implementation

The process of developing an ILM strategy includes four activities—classifying, implementing, managing, and organizing:

- *Classifying* data and applications on the basis of business rules and policies to enable differentiated treatment of information
- *Implementing* policies by using information management tools, starting from the creation of data and ending with its disposal
- *Managing* the environment by using integrated tools to reduce operational complexity
- *Organizing* storage resources in tiers to align the resources with data classes, and storing information in the right type of infrastructure based on the information's current value

Implementing ILM across an enterprise is an ongoing process. Figure 1-8 illustrates a three-step road map to enterprise-wide ILM.

Steps 1 and 2 are aimed at implementing ILM in a limited way across a few enterprise-critical applications. In Step 1, the goal is to implement a storage networking environment. Storage architectures offer varying levels of protection and performance and this acts as a foundation for future policy-based information management in Steps 2 and 3. The value of tiered storage platforms can be exploited by allocating appropriate storage resources to the applications based on the value of the information processed.

Step 2 takes ILM to the next level, with detailed application or data classification and linkage of the storage infrastructure to business policies. These classifications and the resultant policies can be automatically executed using tools for one or more applications, resulting in better management and optimal allocation of storage resources.

Step 3 of the implementation is to automate more of the applications or data classification and policy management activities in order to scale to a wider set of enterprise applications.

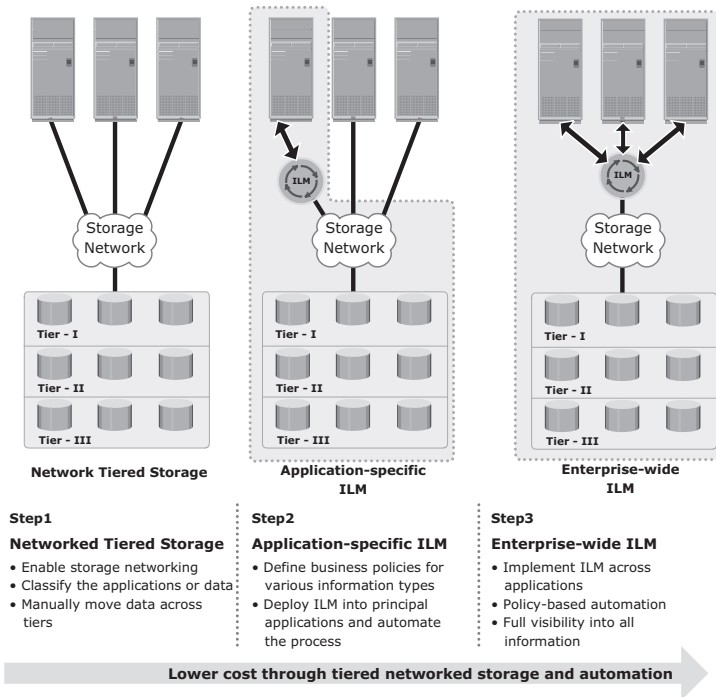


Figure 1-8: Implementation of ILM

1.5.3 ILM Benefits

Implementing an ILM strategy has the following key benefits that directly address the challenges of information management:

- *Improved utilization* by using tiered storage platforms and increased visibility of all enterprise information.
- *Simplified management* by integrating process steps and interfaces with individual tools and by increasing automation.
- *A wider range of options* for backup, and recovery to balance the need for business continuity.
- *Maintaining compliance* by knowing what data needs to be protected for what length of time.
- *Lower Total Cost of Ownership (TCO)* by aligning the infrastructure and management costs with information value. As a result, resources are not wasted, and complexity is not introduced by managing low-value data at the expense of high-value data.

Summary

This chapter described the importance of data, information, and storage infrastructure. Meeting today's storage needs begins with understanding the type of data, its value, and key management requirements of a storage system.

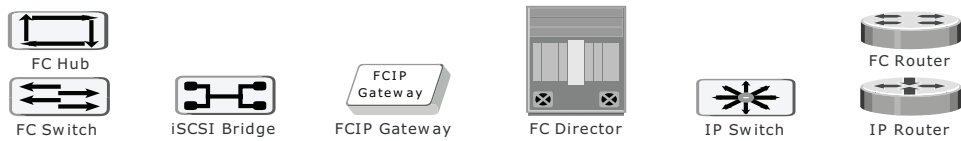
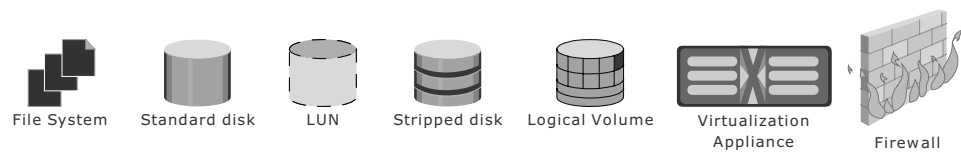
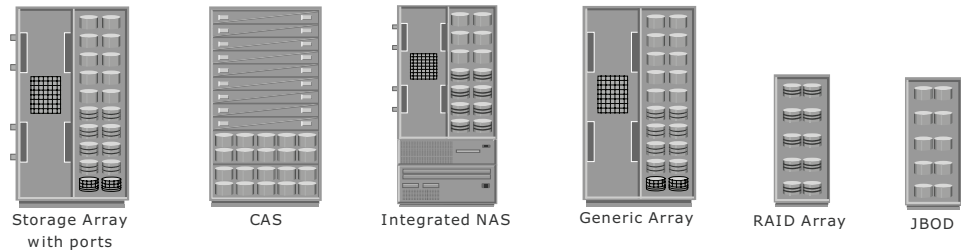
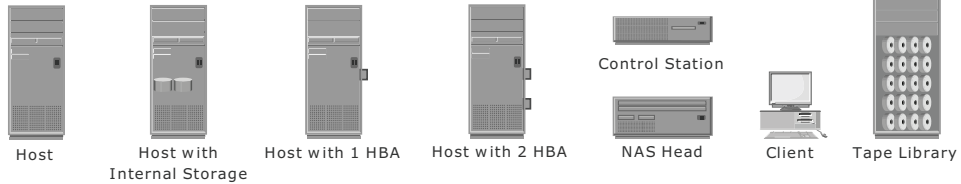
This chapter also emphasized the importance of the ILM strategy, which businesses are adopting to manage information effectively across the enterprise. ILM is enabling businesses to gain competitive advantage by classifying, protecting, and leveraging information.

The evolution of storage architectures and the core elements of a data center covered in this chapter provided the foundation on information storage. The next chapter discusses storage system environment.

EXERCISES

1. A hospital uses an application that stores patient X-ray data in the form of large binary objects in an Oracle database. The application is hosted on a UNIX server, and the hospital staff accesses the X-ray records through a Gigabit Ethernet backbone. Storage array provides storage to the UNIX server, which has 6 terabytes of usable capacity.
 - Explain the core elements of the data center. What are the typical challenges the storage management team may face in meeting the service-level demands of the hospital staff?
 - Describe how the value of this patient data might change over time.
2. An engineering design department of a large company maintains over 600,000 engineering drawings that its designers access and reuse in their current projects, modifying or updating them as required. The design team wants instant access to the drawings for its current projects, but is currently constrained by an infrastructure that is not able to scale to meet the response time requirements. The team has classified the drawings as “most frequently accessed,” “frequently accessed,” “occasionally accessed,” and “archive.”
 - Suggest a strategy for design department that optimizes the storage infrastructure by using ILM.
 - Explain how you will use “tiered storage” based on access frequency.
 - Detail the hardware and software components you will need to implement your strategy.
 - Research products and solutions currently available to meet the solution you are proposing.
3. The marketing department at a mid size firm is expanding. New hires are being added to the department and they are given network access to the department’s files. IT has given marketing a networked drive on the LAN, but it keeps reaching capacity every third week. Current capacity is 500 gigabytes (and growing), with hundreds of files. Users are complaining about LAN response times and capacity. As the IT manager, what could you recommend to improve the situation?
4. A large company is considering a storage infrastructure—one that is scalable and provides high availability. More importantly, the company also needs performance for its mission-critical applications. Which storage topology would you recommend (SAN, NAS, IP SAN) and why?

Icons used in this book



Information Storage and Management

Published by
Wiley Publishing, Inc.
10475 Crosspoint Boulevard
Indianapolis, IN 46256
www.wiley.com

Copyright © 2009 by EMC Corporation

Published by Wiley Publishing, Inc., Indianapolis, Indiana

Published simultaneously in Canada

ISBN: 978-0-470-29421-5

Manufactured in the United States of America

10 9 8 7 6 5 4 3 2 1

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 646-8600. Requests to the Publisher for permission should be addressed to the Legal Department, Wiley Publishing, Inc., 10475 Crosspoint Blvd., Indianapolis, IN 46256, (317) 572-3447, fax (317) 572-4355, or online at www.wiley.com/go/permissions.

Limit of Liability/Disclaimer of Warranty: The publisher and the author make no representations or warranties with respect to the accuracy or completeness of the contents of this work and specifically disclaim all warranties, including without limitation warranties of fitness for a particular purpose. No warranty may be created or extended by sales or promotional materials. The advice and strategies contained herein may not be suitable for every situation. This work is sold with the understanding that the publisher is not engaged in rendering legal, accounting, or other professional services. If professional assistance is required, the services of a competent professional person should be sought. Neither the publisher nor the author shall be liable for damages arising herefrom. The fact that an organization or Web site is referred to in this work as a citation and/or a potential source of further information does not mean that the author or the publisher endorses the information the organization or Web site may provide or recommendations it may make. Further, readers should be aware that Internet Web sites listed in this work may have changed or disappeared between when this work was written and when it is read.

For general information on our other products and services please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Library of Congress Cataloging-in-Publication Data is available from the publisher.

Trademarks: Wiley and the Wiley logo are trademarks or registered trademarks of John Wiley & Sons, Inc. and/or its affiliates, in the United States and other countries, and may not be used without written permission. All other trademarks are the property of their respective owners. Wiley Publishing, Inc. is not associated with any product or vendor mentioned in this book.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

EMC², EMC, EMC Centera, EMC ControlCenter, AdvantEdge, AlphaStor, ApplicationXtender, Avamar, Captiva, Catalog Solution, Celerra, Centera, CentraStar, ClaimPack, ClaimsEditor, ClaimsEditor Professional, CLARAlert, CLARiiON, ClientPak, CodeLink, Connectrix, Co-StandbyServer, Dantz, Direct Matrix Architecture, DiskXtender, DiskXtender 2000, Document Sciences, Documentum, EmailXaminer, EmailXtender, EmailXtract, eRoom, Event Explorer, FLARE, FormWare, HighRoad, InputAccel, Invista, ISIS, Max Retriever, Navisphere, NetWorker, nLayers, OpenScale, PixTools, Powerlink, PowerPath, Rainfinity, RepliStor, ResourcePak, Retrospect, Smarts, SnapshotServer, SnapView/IP, SRDF, Symmetrix, TimeFinder, VisualSAN, Voyence, VSAM-Assist, WebXtender, where information lives, xPression, xPresso, Xtender, and Xtender Solutions are registered trademarks and EMC LifeLine, EMC OnCourse, EMC Proven, EMC Snap, EMC Storage Administrator, Acartus, Access Logix, ArchiveXtender, Atmos, Authentic Problems, Automated Resource Manager, AutoStart, AutoSwap, AVALONidm, C-Clip, Celerra Replicator, CenterStage, CLARevent, Codebook Correlation Technology, Common Information Model, CopyCross, CopyPoint, DatabaseXtender, Digital Mailroom, Direct Matrix, EDM, E-Lab, eInput, Engenuity, FarPoint, FirstPass, Fortress, Global File Virtualization, Graphic Visualization, Infiniflex, InfoMover, Infoscapes, InputAccel Express, MediaStor, MirrorView, Mozy, MozyEnterprise, MozyHome, MozyPro, OnAlert, PowerSnap, QuickScan, RepliCare, SafeLine, SAN Advisor, SAN Copy, SAN Manager, SDMS, SnapImage, SnapSure, SnapView, StorageScope, SupportMate, SymmAPI, SymmEnabler, Symmetrix DMX, UltraFlex, UltraPoint, UltraScale, Viewlets, Virtual Provisioning, and VisualSRM are trademarks of EMC Corporation. All other trademarks used herein are the property of their respective owners. © Copyright 2009 EMC Corporation. All rights reserved. Published in the USA. 01/09