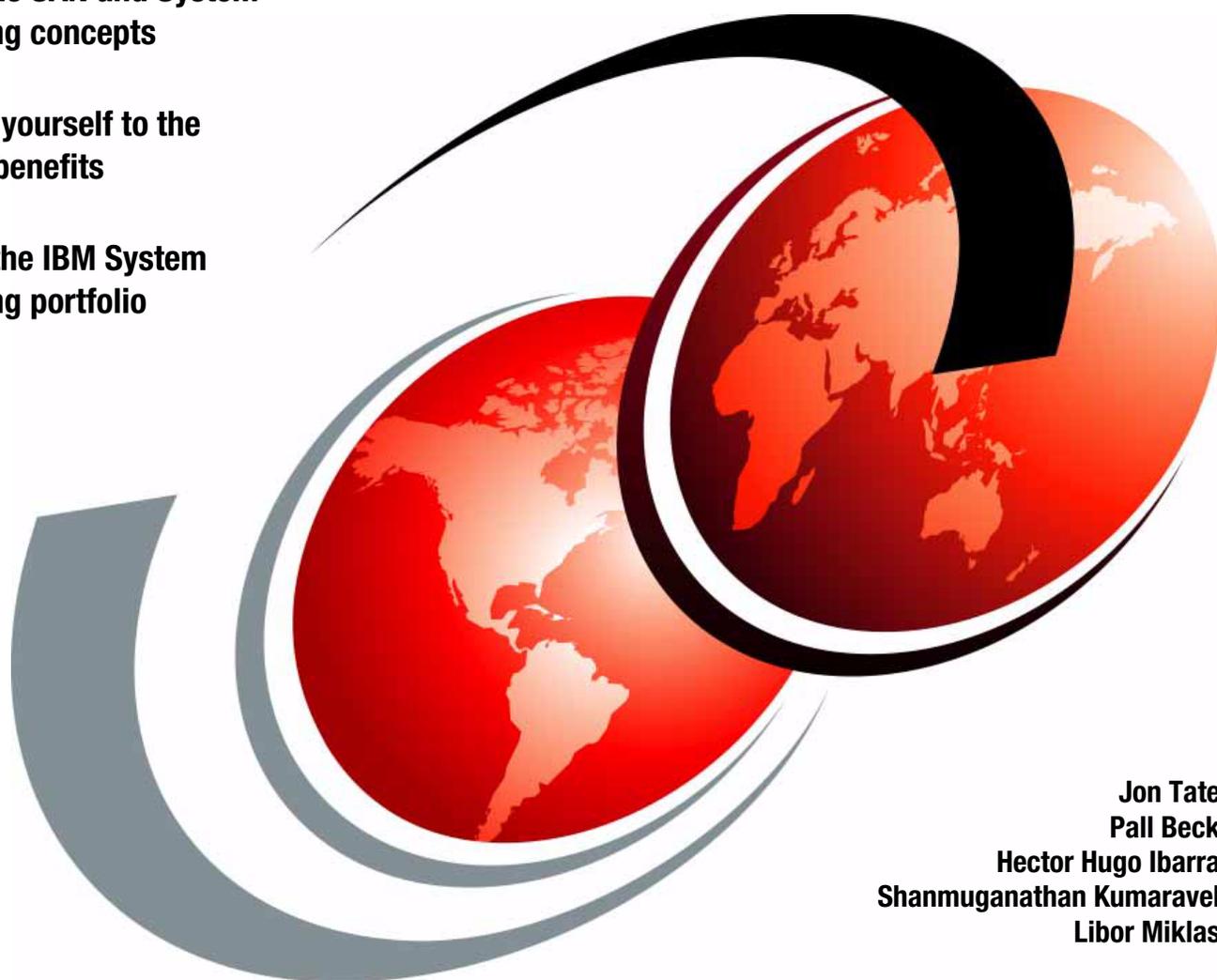


Introduction to Storage Area Networks and System Networking

Learn basic SAN and System
Networking concepts

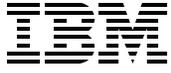
Introduce yourself to the
business benefits

Discover the IBM System
Networking portfolio



Jon Tate
Pall Beck
Hector Hugo Ibarra
Shanmuganathan Kumaravel
Libor Miklas

Redbooks



International Technical Support Organization

**Introduction to Storage Area Networks and System
Networking**

November 2012

Note: Before using this information and the product it supports, read the information in “Notices” on page xi.

Fifth Edition (November 2012)

This edition applies to the products in the IBM System Networking portfolio.

© Copyright International Business Machines Corporation 2012. All rights reserved.

Note to U.S. Government Users Restricted Rights -- Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM Corp.

Contents

Notices	xi
Trademarks	xii
Summary of changes	xiii
November 2012, Fifth Edition	xiii
Preface	xv
The team who wrote this book	xvi
Now you can become a published author, too!	xviii
Comments welcome	xviii
Stay connected to IBM Redbooks	xviii
Chapter 1. Introduction	1
1.1 What is a network	2
1.1.1 The importance of communication	2
1.2 Interconnection models	2
1.2.1 The open systems interconnection model	2
1.2.2 Translating the OSI model to the physical world	4
1.3 What do we mean by storage?	5
1.3.1 Storing data	5
1.3.2 Redundant Array of Independent Disks	5
1.4 What is a storage area network?	11
1.5 Storage area network components	12
1.5.1 Storage area network connectivity	13
1.5.2 Storage area network storage	13
1.5.3 Storage area network servers	14
1.6 The importance of standards or models	14
Chapter 2. Why, and how, can we use a storage area network?	15
2.1 Why use a storage area network?	16
2.1.1 The problem	16
2.1.2 The requirements	17
2.2 How can we use a storage area network?	17
2.2.1 Infrastructure simplification	18
2.2.2 Information lifecycle management	19
2.2.3 Business continuity	19
2.3 Using the storage area network components	20
2.3.1 Storage	20
2.3.2 Storage area network connectivity	21
2.3.3 Servers	26
2.3.4 Putting the components together	28
Chapter 3. Fibre Channel internals	31
3.1 First, why the Fibre Channel architecture?	32
3.1.1 The Small Computer Systems Interface legacy	32
3.1.2 Limitations of the Small Computer System Interface	32
3.1.3 Why Fibre Channel?	35
3.2 Layers	37
3.3 Optical cables	40

3.3.1	Attenuation	41
3.3.2	Maximum power	41
3.3.3	Fiber in the storage area network	42
3.3.4	Dark fiber	46
3.4	Classes of service	46
3.4.1	Class 1	47
3.4.2	Class 2	47
3.4.3	Class 3	47
3.4.4	Class 4	47
3.4.5	Class 5	48
3.4.6	Class 6	48
3.4.7	Class F	48
3.5	Fibre Channel data movement	48
3.5.1	Byte encoding schemes	49
3.6	Data transport	51
3.6.1	Ordered set	52
3.6.2	Frames	53
3.6.3	Sequences	55
3.6.4	Exchanges	55
3.6.5	In order and out of order	56
3.6.6	Latency	56
3.6.7	Open fiber control	56
3.7	Flow control	57
3.7.1	Buffer to buffer	57
3.7.2	End to end	57
3.7.3	Controlling the flow	57
3.7.4	Performance	58
	Chapter 4. Ethernet and system networking concepts	59
4.1	Ethernet	60
4.1.1	Shared media	60
4.1.2	Ethernet frame	61
4.1.3	How Ethernet works	62
4.1.4	Speed and bandwidth	63
4.1.5	10 GbE	64
4.1.6	10 GbE copper versus fiber	65
4.1.7	Virtual local area network	67
4.1.8	Interface virtual local area network operation modes	69
4.1.9	Link aggregation	71
4.1.10	Spanning Tree Protocol	71
4.1.11	Link Layer Discovery Protocol	74
4.1.12	Link Layer Discovery Protocol Type Length Values (LLDP TLVs)	74
4.2	Storage area network IP networking	76
4.2.1	The multiprotocol environment	76
4.2.2	Fibre Channel switching	76
4.2.3	Fibre Channel routing	76
4.2.4	Tunneling	76
4.2.5	Routers and gateways	77
4.2.6	Internet Storage Name Service	77
4.3	Delving deeper into the protocols	77
4.3.1	Fibre Channel over Internet Protocol (FCIP)	77
4.3.2	Internet Fibre Channel Protocol (iFCP)	78
4.3.3	Internet Small Computer System Interface (iSCSI)	79

4.3.4	Routing considerations	81
4.3.5	Packet size	81
4.3.6	TCP congestion control.	81
4.3.7	Round-trip delay	82
4.4	Multiprotocol solution briefs.	83
4.4.1	Dividing a fabric into subfabrics	83
4.4.2	Connecting a remote site over IP	83
4.4.3	Connecting hosts using Internet Small Computer System Interface.	84
Chapter 5.	Topologies and other fabric services	85
5.1	Fibre Channel topologies	86
5.1.1	Point-to-point topology	86
5.1.2	Arbitrated loop topology	87
5.1.3	Switched fabric topology.	88
5.1.4	Single switch topology	89
5.1.5	Cascaded and ring topology	90
5.1.6	Mesh topology.	91
5.1.7	Core-edge topology	92
5.1.8	Edge-core-edge topology	92
5.2	Port types	93
5.2.1	Common port types.	93
5.2.2	Expansion port types	94
5.2.3	Diagnostic port types	95
5.3	Addressing	97
5.3.1	Worldwide name	97
5.3.2	Tape Device WWNN and WWPN	101
5.3.3	Port address	101
5.3.4	The 24-bit port address.	101
5.3.5	Loop address	103
5.3.6	The b-type addressing modes	103
5.3.7	FICON address	104
5.4	Fibre Channel Arbitrated Loop protocols	108
5.4.1	Fairness algorithm	109
5.4.2	Loop addressing	109
5.5	Fibre Channel port initialization and fabric services	110
5.5.1	Fabric login (FLOGI)	110
5.5.2	Port login (PLOGI)	111
5.5.3	Process login (PRLI)	112
5.6	Fabric services	113
5.6.1	Management server	113
5.6.2	Time server.	114
5.6.3	Simple name server	114
5.6.4	Fabric login server	114
5.6.5	Registered state change notification service.	114
5.7	Routing mechanisms.	115
5.7.1	Spanning tree	115
5.7.2	Fabric shortest path first	115
5.8	Zoning	116
5.8.1	Hardware zoning.	117
5.8.2	Software zoning	119
5.8.3	Logical unit number masking	122
Chapter 6.	Storage area network as a service for cloud computing	123

6.1	What is a cloud?	124
6.1.1	Private and public cloud	125
6.1.2	Cloud computing components	125
6.1.3	Cloud computing models	126
6.2	Virtualization and the cloud	129
6.2.1	Cloud infrastructure virtualization	129
6.2.2	Cloud platforms	130
6.2.3	Storage virtualization	132
6.3	SAN virtualization	133
6.3.1	IBM b-type Virtual Fabrics	133
6.3.2	Cisco virtual storage area network	135
6.3.3	N-Port ID Virtualization	137
6.4	Building a smarter cloud	139
6.4.1	Automated tiering	139
6.4.2	Thin provisioning	140
6.4.3	Deduplication	141
6.4.4	New generation management tools	144
6.4.5	Business continuity and disaster recovery	144
6.4.6	Storage on demand	144
	Chapter 7. Fibre Channel products and technology	145
7.1	The environment	146
7.2	Storage area network (SAN) devices	147
7.2.1	Fibre Channel bridges	147
7.2.2	Arbitrated loop hubs and switched hubs	148
7.2.3	Switches and directors	149
7.2.4	Multiprotocol routing	150
7.2.5	Service modules	150
7.2.6	Multiplexers	150
7.3	Componentry	150
7.3.1	Application-specific integrated circuit	151
7.3.2	Fibre Channel transmission rates	151
7.3.3	SerDes	151
7.3.4	Backplane and blades	152
7.4	Gigabit transport technology	152
7.4.1	Fibre Channel cabling	152
7.4.2	Transceivers	157
7.4.3	Host bus adapters	159
7.5	Inter-switch links	161
7.5.1	Cascading	161
7.5.2	Hops	162
7.5.3	Fabric shortest path first	162
7.5.4	Non-blocking architecture	163
7.5.5	Latency	164
7.5.6	Oversubscription	164
7.5.7	Congestion	165
7.5.8	Trunking or port-channeling	165
	Chapter 8. Management	167
8.1	Management principles	168
8.1.1	Management types	168
8.1.2	Connecting to storage area network management tools	170
8.1.3	Storage area network fault isolation and troubleshooting	170

8.2 Management interfaces and protocols	171
8.2.1 Storage Networking Industry Association initiative	171
8.2.2 Simple Network Management Protocol	173
8.2.3 Service Location Protocol	174
8.2.4 Vendor-specific mechanisms	174
8.3 Management features	177
8.3.1 Operations	177
8.4 IBM Tivoli Storage Productivity Center	178
8.4.1 IBM Tivoli Storage Productivity Center for Data	179
8.4.2 IBM Tivoli Storage Productivity Center for Disk	179
8.4.3 IBM Tivoli Storage Productivity Center for Disk Select	179
8.4.4 IBM Tivoli Storage Productivity Center Basic Edition	180
8.4.5 IBM Tivoli Storage Productivity Center Standard Edition	180
8.4.6 IBM Tivoli Storage Productivity Center for Replication	181
8.4.7 What is IBM System Storage Productivity Center?	181
8.4.8 What can be done from the System Storage Productivity Center?	182
8.5 Vendor management applications	182
8.5.1 b-type	182
8.5.2 Cisco	183
8.6 SAN multipathing software	184
Chapter 9. Security	191
9.1 Security in the storage area network (SAN)	192
9.2 Security principles	193
9.2.1 Access control	193
9.2.2 Auditing and accounting	193
9.2.3 Data security	193
9.2.4 Securing a fabric	194
9.2.5 Zoning, masking, and binding	195
9.3 Data security	195
9.4 Storage area network encryption	196
9.4.1 Basic encryption definition	196
9.4.2 Data-in-flight	198
9.4.3 Data-at-rest	199
9.4.4 Digital certificates	199
9.4.5 Encryption algorithm	199
9.4.6 Key management considerations and security standards	200
9.4.7 b-type encryption methods	201
9.4.8 Cisco encryption methods	203
9.5 Encryption standards and algorithms	205
9.6 Security common practices	206
Chapter 10. Solutions	207
10.1 Introduction	208
10.2 Basic solution principles	208
10.2.1 Connectivity	208
10.2.2 Adding capacity	209
10.2.3 Data movement and copy	209
10.2.4 Upgrading to faster speeds	212
10.3 Infrastructure simplification	213
10.3.1 Where does the complexity come from?	213
10.3.2 Storage pooling	214
10.3.3 Consolidation	217

10.3.4	Migration to a converged network	218
10.4	Business continuity and disaster recovery	221
10.4.1	Clustering and high availability	222
10.4.2	LAN-free data movement	224
10.4.3	Disaster backup and recovery	225
10.5	Information lifecycle management	226
10.5.1	Information lifecycle management	227
10.5.2	Tiered storage management	227
10.5.3	Long-term data retention	229
10.5.4	Data lifecycle management	229
10.5.5	Policy-based archive management	230
Chapter 11.	Storage area networks and green data centers	233
11.1	Data center constraints	234
11.1.1	Energy flow in data center	235
11.2	Data center optimization	236
11.2.1	Strategic considerations	236
11.3	Green storage	237
11.3.1	Information lifecycle management	237
11.3.2	Storage consolidation and virtualization	239
11.3.3	On-demand storage provisioning	240
11.3.4	Hierarchical storage and tiering	241
11.3.5	Data compression and deduplication	242
Chapter 12.	The IBM product portfolio	245
12.1	Classification of IBM storage area network products	246
12.2	SAN Fibre Channel networking	248
12.2.1	Entry SAN switches	248
12.2.2	Midrange SAN switches	250
12.2.3	Enterprise SAN directors	255
12.2.4	Multiprotocol routers	260
12.3	IBM System Storage Disk Systems	262
12.3.1	Entry level disk systems	262
12.3.2	Midrange disk systems	264
12.3.3	Enterprise disk systems	267
12.4	IBM Tape Storage Systems	270
12.4.1	Fibre Channel tape drives	271
12.4.2	Autoloaders and entry tape libraries	273
12.4.3	Midrange tape libraries	275
12.4.4	Enterprise tape libraries	276
12.5	Storage virtualization and cloud computing	278
12.5.1	Disk storage virtualization	278
12.5.2	Tape storage virtualization	282
12.5.3	Storage systems for cloud computing	285
12.6	IP-based networking for SAN environments	286
12.7	Hardware solutions for network convergence	287
12.7.1	IBM Virtual Fabric solution	288
12.8	IBM Flex System networking	290
12.8.1	IBM Flex System Fabric EN4093 10Gb Scalable Switch	291
12.8.2	IBM Flex System EN4091 10Gb Ethernet Pass-thru	296
12.8.3	IBM Flex System EN2092 1Gb Ethernet Scalable Switch	298
12.8.4	IBM Flex System FC5022 16Gb SAN Scalable Switch	303
12.8.5	IBM Flex System FC3171 8Gb SAN Switch	309

12.8.6 IBM Flex System FC3171 8Gb SAN Pass-thru.	311
Chapter 13. Certification.	315
13.1 Why certification?	316
13.2 IBM Professional Certification Program	317
13.2.1 About the program	317
13.2.2 Certifications by product	317
13.2.3 Mastery tests.	317
13.3 Storage Networking Industry Association certifications.	318
13.3.1 SNIA Certified Storage Professional (SCSP)	318
13.3.2 SNIA Certified Storage Engineer (SCSE).	318
13.3.3 SNIA Certified Storage Architect (SCSA)	318
13.3.4 SNIA Certified Storage Networking Expert (SCSN-E)	319
13.3.5 SNIA Qualified Data Protection Associate	319
13.3.6 SNIA Qualified Storage Virtualization Associate.	319
13.3.7 SNIA Qualified Storage Sales Professional	319
13.3.8 CompTIA Storage+ Powered by SNIA	319
13.4 Brocade certifications	320
13.4.1 Brocade Accredited Server Connectivity Specialist	320
13.4.2 Brocade Accredited Data Center Specialist	321
13.4.3 Brocade Accredited Fibre Channel connection (FICON) Specialist	321
13.4.4 Brocade Accredited FCoE Specialist	321
13.4.5 Brocade Accredited Internetworking Specialist.	321
13.4.6 Brocade Accredited WLAN Specialist.	321
13.4.7 Brocade Certified Fabric Administrator (BCFA)	321
13.4.8 Brocade Certified Fabric Professional (BCFP)	321
13.4.9 Brocade Certified SAN Manager (BCSM).	322
13.4.10 Brocade Certified Fabric Designer (BCFD).	322
13.4.11 Brocade Certified Architect For FICON (BCAF)	322
13.4.12 Brocade Certified FCoE Professional (BCFCoEP)	322
13.4.13 Brocade Certified Ethernet Fabric Engineer	322
13.4.14 Brocade Certified Network Engineer.	323
13.4.15 Brocade Certified Layer 4-7 Engineer.	323
13.4.16 Brocade Certified Network Professional	323
13.4.17 Brocade Certified Layer 4-7 Professional	323
13.4.18 Brocade Certified Network Designer.	323
13.5 Cisco certification	323
13.5.1 Cisco Certified Entry Networking Technician (CCENT)	324
13.5.2 Cisco Certified Network Associate (CCNA)	324
13.5.3 Cisco Certified Network Associate Security (CCNA Security)	324
13.5.4 Cisco Certified Network Associate Wireless (CCNA Wireless).	324
13.5.5 Cisco Certified Design Associate (CCDA)	324
13.5.6 Cisco Certified Network Professional (CCNP)	325
13.5.7 CCNP Security certification.	325
13.5.8 CCNP Wireless certification	325
13.5.9 Cisco Certified Design Professional (CCDP)	325
13.5.10 Cisco Certified Internetwork Expert (CCIE) - Routing and Switching	326
13.5.11 Cisco Certified Internetwork Expert (CCIE) - Security	326
13.5.12 Cisco Certified Internetwork Expert (CCIE) - Wireless	326
13.5.13 Cisco Certified Design Expert (CCDE)	326
13.5.14 Cisco CCIE Storage Networking.	326
13.5.15 Cisco Certified Architect	326
13.5.16 Cisco specialization tracks	327

13.6 The Open Group certifications	327
13.6.1 The Open Group Certified IT Specialists (Open CITS)	327
13.6.2 The Open Group Certified Architect (Open CA)	327
13.6.3 The Open Group certification	328
13.7 Juniper Networks Certification Program (JNCP)	328
13.7.1 JNCP Junos-based certification tracks	328
13.7.2 Service Provider Routing and Switching track	328
13.7.3 Enterprise Routing and Switching track	329
13.7.4 Junos security track	330
13.8 Non-Junos certification tracks	331
13.8.1 E-Series certification track	331
13.8.2 Firewall/VPN certification track	332
13.8.3 SSL certification track	333
13.8.4 Intrusion Detection and Prevention (IDP) Track	333
13.8.5 Unified Access Control (UAC) Track	333
13.8.6 WX certification track	334
Related publications	335
IBM Redbooks	335
IBM Flex System education	336
Referenced websites	337
Help from IBM	338

Notices

This information was developed for products and services offered in the U.S.A.

IBM may not offer the products, services, or features discussed in this document in other countries. Consult your local IBM representative for information on the products and services currently available in your area. Any reference to an IBM product, program, or service is not intended to state or imply that only that IBM product, program, or service may be used. Any functionally equivalent product, program, or service that does not infringe any IBM intellectual property right may be used instead. However, it is the user's responsibility to evaluate and verify the operation of any non-IBM product, program, or service.

IBM may have patents or pending patent applications covering subject matter described in this document. The furnishing of this document does not grant you any license to these patents. You can send license inquiries, in writing, to:

IBM Director of Licensing, IBM Corporation, North Castle Drive, Armonk, NY 10504-1785 U.S.A.

The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law: INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some states do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. IBM may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-IBM websites are provided for convenience only and do not in any manner serve as an endorsement of those websites. The materials at those websites are not part of the materials for this IBM product and use of those websites is at your own risk.

IBM may use or distribute any of the information you supply in any way it believes appropriate without incurring any obligation to you.

Any performance data contained herein was determined in a controlled environment. Therefore, the results obtained in other operating environments may vary significantly. Some measurements may have been made on development-level systems and there is no guarantee that these measurements will be the same on generally available systems. Furthermore, some measurements may have been estimated through extrapolation. Actual results may vary. Users of this document should verify the applicable data for their specific environment.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental.

COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to IBM, for the purposes of developing, using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. IBM, therefore, cannot guarantee or imply reliability, serviceability, or function of these programs.

Trademarks

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both. These and other IBM trademarked terms are marked on their first occurrence in this information with the appropriate symbol (® or ™), indicating US registered or common law trademarks owned by IBM at the time this information was published. Such trademarks may also be registered or common law trademarks in other countries. A current list of IBM trademarks is available on the Web at <http://www.ibm.com/legal/copytrade.shtml>

The following terms are trademarks of the International Business Machines Corporation in the United States, other countries, or both:

AFS®	i5/OS™	Redbooks (logo)  ®
AIX®	IBM Flex System™	RMF™
AS/400®	IBM SmartCloud™	RS/6000®
BladeCenter®	IBM®	ServerProven®
DB2®	Informix®	Storwize®
Domino®	iSeries®	System i®
DS4000®	Lotus®	System p®
DS6000™	OS/390®	System Storage®
DS8000®	OS/400®	System x®
Easy Tier®	Power Systems™	System z9®
ECKD™	POWER6®	System z®
Enterprise Storage Server®	PowerHA®	Tivoli®
ESCON®	PowerPC®	VMready®
Express Storage™	POWER®	XIV®
FICON®	ProtecTIER®	xSeries®
FlashCopy®	pSeries®	z/OS®
GPFS™	PureFlex™	z/VM®
HACMP™	RackSwitch™	z9®
HyperFactor®	Redbooks®	zSeries®

The following terms are trademarks of other companies:

Intel Xeon, Intel, Intel logo, Intel Inside logo, and Intel Centrino logo are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

Linux is a trademark of Linus Torvalds in the United States, other countries, or both.

LTO, Ultrium, the LTO Logo and the Ultrium logo are trademarks of HP, IBM Corp. and Quantum in the U.S. and other countries.

Microsoft, Windows NT, Windows, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

UNIX is a registered trademark of The Open Group in the United States and other countries.

Other company, product, or service names may be trademarks or service marks of others.

Summary of changes

This section describes the technical changes made in this edition of the book and in previous editions. This edition might also include minor corrections and editorial changes that are not identified.

Summary of Changes
for SG24-5470-04
for Introduction to Storage Area Networks and System Networking
as created or updated on November 17, 2012.

November 2012, Fifth Edition

This revision includes amendments, deletions, and additions to support IBM® strategies and initiatives, and to add an introduction to IBM System Networking, as well as update those areas of Storage Area Networking as appropriate.

For any omissions or inaccuracies, contact Jon Tate (tatej@uk.ibm.com).

Preface

The plethora of data that is created by the businesses of today is making storage a strategic investment priority for companies of all sizes. As storage takes precedence, three major initiatives emerge:

- ▶ Flatten and converge your network

IBM takes an open, standards-based approach to implement the latest advances in the flat, converged data center network designs of today. IBM System Networking solutions enable clients to deploy a high-speed, low-latency Unified Fabric Architecture.

- ▶ Optimize and automate virtualization

Advanced virtualization awareness reduces the cost and complexity of deploying physical and virtual data center infrastructure.

- ▶ Simplify management

IBM data center networks are easy to deploy, maintain, scale, and virtualize, delivering the foundation of consolidated operations for dynamic infrastructure management.

Storage is no longer an afterthought. Too much is at stake. Companies are searching for more ways to efficiently manage expanding volumes of data, and to make that data accessible throughout the enterprise. This demand is propelling the move of storage into the network. Also, the increasing complexity of managing a large numbers of storage devices and vast amounts of data is driving greater business value into software and services.

With current estimates of the amount of data to be managed and made available increasing at 60 percent each year, this outlook is where a storage area network (SAN) enters the arena. SANs are the leading storage infrastructure for the global economy of today. SANs offer simplified storage management, scalability, flexibility, and availability; and improved data access, movement, and backup.

Welcome to the era of *Smarter Networking for Smarter Data Centers*.

The smarter data center with improved economics of IT can be achieved by connecting servers and storage with a high-speed and intelligent network fabric. A smarter data center that hosts IBM System Networking solutions can provide an environment that is smarter, faster, greener, open, and easy to manage.

This IBM Redbooks® publication provides an introduction to the SAN and Ethernet networking, and how these networks help to achieve a smarter data center. This book is intended for people who are not very familiar with IT, or who are just starting out in the IT world.

For more information, and a deeper dive into the SAN world, you might find the following Redbooks publications especially useful to expand your SAN knowledge:

Implementing an IBM b-type SAN with 8 Gbps Directors and Switches, SG24-6116

Implementing the IBM System Storage SAN32B-E4 Encryption Switch, SG24-7922

IBM System Storage b-type Multiprotocol Routing: An Introduction and Implementation, SG24-7544

IBM Converged Switch B32, SG24-7935

Also, be sure to see the IBM System Networking Redbooks portal for the latest material from the International Technical Support Organization (ITSO):

<http://www.redbooks.ibm.com/portals/networking>

The team who wrote this book

This book was produced by a team of specialists from around the world working at the IBM International Technical Support Organization (ITSO), Poughkeepsie Center.



Jon Tate is a Project Manager for IBM System Storage® SAN Solutions at the International Technical Support Organization (ITSO), San Jose Center. Before joining the ITSO in 1999, he worked in the IBM Technical Support Center, providing Level 2 support for IBM storage products. Jon has 27 years of experience in storage software and management, services, and support, and is both an IBM Certified IT Specialist and an IBM SAN Certified Specialist. He is also the UK Chairman of the Storage Networking Industry Association.



Pall Beck is a SAN Technical Lead in IBM Nordic. He has 15 years of experience working with storage, both for dedicated clients and for large shared environments. Those environments include clients from the medical and financial sector which includes some of the largest shared SAN environments in Europe. He is a member of the SAN and SAN Volume Controller best practice community in the Nordics and in EMEA. In his current job role, he is a member of the Solutions Consultant Express+ (SCE+) Storage Deployment Specialists, responsible for SCE+ storage deployments around the globe. Pall is also a member of a team that helps in critical situations and does root cause analyzes. He is coauthor of the *Implementing SVC 5.1* and *SVC Advanced Copy Services 4.2* Redbooks publications. Pall has a diploma as an Electronic Technician from Odense Tekniske Skole in Denmark and an IR in Reykjavik, Iceland. He is also an IBM Certified IT Specialist.



Hector Hugo Ibarra is an Infrastructure IT Architect specialized in cloud computing and storage solutions currently working at IBM Argentina. Hector is the ITA Leader for The VMware Center of Competence, designated in 2006. He specializes in virtualization technologies and assisted several global IBM clients in deploying virtualized infrastructures across the world. Since 2009, he has been working as the Leader for the Argentina Delivery Center Strategy and Architecture Services department from where major projects are driven.



Shanmuganathan Kumaravel is an IBM Technical Services Specialist for the ITD-SSO MR Storage team of IBM India. He has supported SAN and disk products from both IBM and Hewlett Packard since August 2008. Before this, he worked for HP product support which provides remote support on HP SAN storage products, servers, and operating systems, including HP UNIX and Linux. Shan is a Brocade Certified SAN Designer (BCSD), Brocade Certified Fabric Administrator (BCFA), and an HP Certified Systems Engineer (CSE).



Libor Miklas is a Team Leader and an experienced IT Specialist working at the IBM Delivery Center Central Europe in Czech Republic. He demonstrates 10 years of practical experience within the IT industry. During the last six years, his main focus continues to be on backup and recovery and on storage management. Libor and his team support midrange and enterprise storage environments for various global and local clients, worldwide. He is an IBM Certified Deployment Professional of the IBM Tivoli® Storage Manager family of products and holds a Masters Degree in Electrical Engineering and Telecommunications.

Thanks to the previous authors of the first, second, third, and fourth editions of this book:

Angelo Bernasconi
Rajani Kanth
Ravi Kumar Khattar
Fabiano Lucchese
Peter Mescher
Richard Moore
Mark S. Murphy
Kjell E. Nyström
Fred Scholten
Giulio John Tarella
Andre Telles

Thanks to the following people for their contributions to this project:

Sangam Racherla
International Technical Support Organization, Poughkeepsie Center

Special thanks to the Brocade staff for their unparalleled support of this residency in terms of equipment and support in many areas:

Brian Steffler
Silviano Gaona
Marcus Thordal
Steven Tong
Jim Baldyga
Brocade Communications Systems

John McKibben
Cisco Systems

Now you can become a published author, too!

Here's an opportunity to spotlight your skills, grow your career, and become a published author—all at the same time! Join an ITSO residency project and help write a book in your area of expertise, while honing your experience using leading-edge technologies. Your efforts will help to increase product acceptance and customer satisfaction, as you expand your network of technical contacts and relationships. Residencies run from two to six weeks in length, and you can participate either in person or as a remote resident working from your home base.

Find out more about the residency program, browse the residency index, and apply online at: ibm.com/redbooks/residencies.html

Comments welcome

Your comments are important to us!

We want our books to be as helpful as possible. Send us your comments about this book or other IBM Redbooks publications in one of the following ways:

- ▶ Use the online **Contact us** review Redbooks form found at:

ibm.com/redbooks

- ▶ Send your comments in an email to:

redbooks@us.ibm.com

- ▶ Mail your comments to:

IBM Corporation, International Technical Support Organization
Dept. HYTD Mail Station P099
2455 South Road
Poughkeepsie, NY 12601-5400

Stay connected to IBM Redbooks

- ▶ Find us on Facebook:

<http://www.facebook.com/IBMRedbooks>

- ▶ Follow us on Twitter:

<http://twitter.com/ibmredbooks>

- ▶ Look for us on LinkedIn:

<http://www.linkedin.com/groups?home=&gid=2130806>

- ▶ Explore new Redbooks publications, residencies, and workshops with the IBM Redbooks weekly newsletter:

<https://www.redbooks.ibm.com/Redbooks.nsf/subscribe?OpenForm>

- ▶ Stay current on recent Redbooks publications with RSS Feeds:

<http://www.redbooks.ibm.com/rss.html>



Introduction

Computing is based on information. Information is the underlying resource on which all computing processes are based; it is a company asset. Information is stored on storage media and is accessed by applications that are running on a server. Often, the information is a unique company asset. Information is created and acquired every second of every day. Information is the currency of business.

To ensure that any business delivers the expected results, they must have access to accurate information, and without delay. The management and protection of business information is vital for the availability of business processes.

This chapter introduces the concept of a network, storage, and the storage area network (SAN), which is regarded as the ultimate response to all these needs.

1.1 What is a network

A computer network, often simply called a *network*, is a collection of computers and devices that are interconnected by communication channels. These channels allow for the efficient sharing of resources, services, and information among it.

Even though this definition is simple, understanding how to make a network work might be complicated for people who are not familiar with information technology (IT), or who are just starting out in the IT world. Because of this unfamiliarity, we explain the basic concepts that need to be understood to facilitate our understanding of the networking world.

1.1.1 The importance of communication

It is impossible to imagine the human world as stand-alone humans, with nobody that talks or does anything for each other. Much more importantly, it is hard trying to imagine how a human can work without using their senses. In our human world, we are sure you would agree with us that communication between individuals makes a significant difference in all aspects of life.

First of all, communication in any form is not easy, and we need a number of components. Factors consist of a common language, something to be communicated, a medium where the communication flows, and finally we need to be sure that whatever was communicated was received and understood. To do that in the human world, we use language as a communication protocol, and sounds and writing are the communication medium.

Similarly, a computer network needs almost the same components as our human example, but a difference is that all factors need to be governed in some way to ensure effective communications. This monitoring is achieved by the use of industry standards, and companies adhere to those standards to ensure that communication can take place.

There is a wealth of information that is devoted to networking history and its evolution, and we do not intend to give a history lesson in this book. This publication focuses on the prevalent interconnection models, storage, and networking concepts.

1.2 Interconnection models

An interconnection model is a standard that is used to connect sources and targets in a network, and there are some well-known models in the IT industry such as the open systems interconnection model (OSI), Department of Defense (DoD), TCP/IP protocol suite, and Fibre Channel. Each model has its advantages and disadvantages. Its model is applied where it has the maximum benefit in terms of performance, reliability, availability, cost benefits, and so on.

1.2.1 The open systems interconnection model

The open systems interconnection model (OSI model) was a product of the open systems interconnection effort at the International Organization for Standardization (ISO). It is a way of subdividing a communications system into smaller parts called *layers*. Similar communication functions are grouped into logical layers. A layer provides services to its upper layer while it receives services from the layer below. At each layer, an instance provides service to the instances at the layer above and requests service from the layer below.

For this book, we focus on the Physical, DataLink, Network, and Transport layers.

Layer 1: Physical Layer

The *Physical Layer* defines electrical and physical specifications for devices. In particular, it defines the relationship between a device and a transmission medium, such as a copper or optical cable. This relationship includes the layout of pins, voltages, cable specifications, and more.

Layer 2: DataLink Layer

The *DataLink Layer* provides the functional and procedural means to transfer data between network entities. This layer also detects and possibly corrects errors that might occur in the Physical Layer.

Layer 3: Network Layer

The *Network Layer* provides the functional and procedural means of transferring variable length data sequences from a source host on one network to a destination host on a different network. The Network layer provides this functionality while it maintains the quality of service requested by the Transport Layer (in contrast to the DataLink layer, which connects hosts within the same network). The Network Layer performs network routing functions. This layer might also perform fragmentation and reassembly, and report delivery errors. Routers operate at this layer by sending data throughout the extended network and making the Internet possible.

Layer 4: Transport Layer

The *Transport Layer* provides transparent transfer of data between users, providing reliable data transfer services to the upper layers. The Transport Layer controls the reliability of a specific link through flow control, segmentation and desegmentation, and error control. Some protocols are state- and connection-oriented. This means that the Transport Layer can track the segments and retransmit the ones that fail. The Transport layer also provides the acknowledgement of the successful data

Now that you know what an interconnection model is, what it does, and how important it is in a network, we can compare the OSI model with other models. Figure 1-1 shows a comparison table of various models.

OSI layer #	name	TCP/IP	Fibre Channel
5-7	application	telnet, ftp, SCSI-3 (iSCSI)	IP, SCSI-3 (FCP)
4	transport	TCP, UDP	FC-4
3	network	IP, ICMP, IGMP	FC-3
2	data link	Ethernet, Token Ring	FC-2, most of FC-1
1	physical	media	FC-0

Figure 1-1 OSI, TCP/IP, and FC models comparison table

The Fibre Channel model is covered later in this book.

1.2.2 Translating the OSI model to the physical world

To make a translation from theoretical models to reality, we introduce physical devices which perform certain tasks for each layer on each model.

Local area networks (LANs) are a good place to start. We define LANs as a small or large network that is limited within the same physical site. This site might be a traditional office or a corporate building.

In Figure 1-2, you see a basic network where computers and a printer are interconnected by using physical cables and interconnection devices.

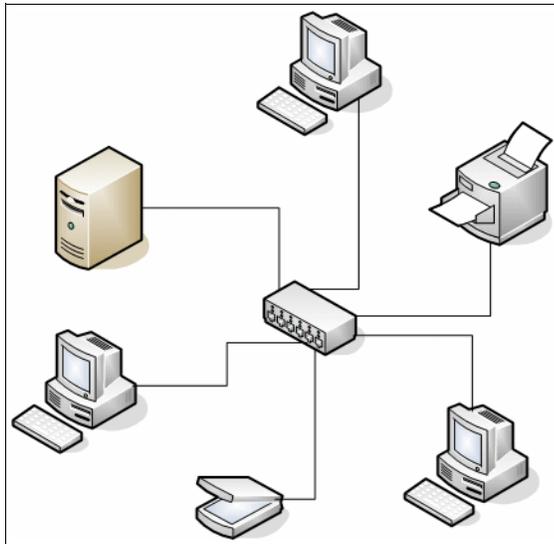


Figure 1-2 Basic network topology

We must keep in mind that any model we choose defines the devices, cables, connectors, and interface characteristics that we must implement to make it work. We must also support the protocols for each model layer.

All the network components are categorized into five groups:

- ▶ End devices: An end device is a computer system which has a final purpose like desktop computers, printers, storage, or servers.
- ▶ Network interface: It is an interface between the media and end devices which can interact with other network interfaces and understands an interconnection model.
- ▶ Connector: This is the physical element at the end of the media which allows a connection to the network interface.
- ▶ Media: This is the physical path that is used to transmit an electrical or optical signal. It might be wired or wireless, copper, or a fiber optic cable.
- ▶ Network devices: These are used to interconnect multiple end devices as a single point of interconnection, route communication through different networks, or for providing network security. Examples of network devices are switches, routers, firewalls, and directors.

Each network component executes a particular role within a network and all of them are required to reach the final goal of making communication possible.

1.3 What do we mean by storage?

To understand what storage is, and because understanding it is a key point for this book, we start from a basic *hard disk drive (HDD)*. We then progress through to storage systems that are high performance, fault tolerant, and highly available. During this explanation, instructional examples are used that might sometimes not reflect reality. However, the examples make it easier to understand for individuals that are just beginning to enter the world of storage systems.

1.3.1 Storing data

Data is stored on HDDs on which can be read and written. Depending on the methods that are used to run those tasks, and the HDD technology on which they were built, the read and write function can be faster or slower. The evolution of HDDs is incredible. We can now store hundreds of gigabytes on a single HDD, which allows us to keep all the data we can ever imagine. Even though this approach seems to bring us only advantages so far, one question might be what happens if for any reason we are unable to access the HDD?

The first solution might be to have a secondary HDD where we can manually copy our primary HDD to our secondary HDD. Immediately, we can see that now our data is safe. But, how often must we run those manual copies if we expect not to lose data and to keep it as up-to-date as possible? To keep it as current as possible, every time we change something we must make another copy. But, must we copy the entire amount of data from one HDD to the other, or must we copy only what changes?

Figure 1-3 shows a manual copy of data for redundancy.

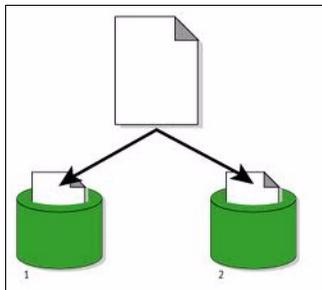


Figure 1-3 Manual copy of data

1.3.2 Redundant Array of Independent Disks

Fortunately, technology exists that can help us. That technology is the *Redundant Array of Independent Disks (RAID)* concept, which presents a possible solution to our problem. It is clear that data needs to be copied every time that it changes to provide us with a reliable fault tolerant system. It is also clear that it cannot be done in a manual way. A *RAID controller* can maintain disks in synchronization and can also manage all the writes and reads (input/output (I/O)) to and from the disks.

Now our RAID system looks like the diagram that is shown in Figure 1-4. The A, B, and C values in the diagram represent user data such as documents or pictures.

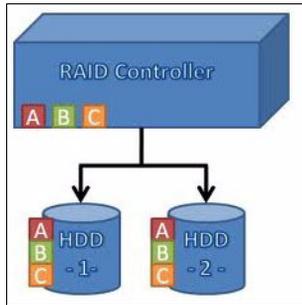


Figure 1-4 Typical RAID scenario

This type of RAID scenario is known as *RAID 1* or as a *mirrored disk*.

Stand-alone disks provide the following advantages:

- ▶ Provides redundancy to disk failure
- ▶ Faster when reading data because it can be taken from either disk

Stand-alone disks provide the following disadvantages:

- ▶ Slower when they are writing because data needs to be written twice
- ▶ Only half of the total capacity can be used

This naturally leads us to the next question. Is there any other RAID type that can improve things further, while it conserves the advantages and removes the disadvantages of RAID 1?

The answer is yes, and this type is known as *RAID 5*. This scenario consists of dividing the user data into $N-1$ parts (where N is the number of disks that is used to build the RAID) and then calculating a parity part. This part permits RAID to rebuild the user data if there is a disk failure.

RAID 5 uses *parity* or redundant information. If a block fails, enough parity information is available to recover the data. The parity information is spread across all the disks. If a disk fails, the RAID requires a rebuild and the parity information is used to re-create the data lost. This example is shown in Figure 1-5.

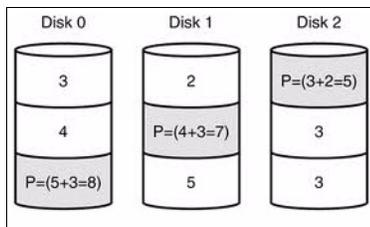


Figure 1-5 Example of RAID 5 with parity

RAID 5 requires a minimum of three disks and in theory there are no limitations to add disks. This RAID type combines data safety with efficient use of disk space. Disk failure does not result in a service interruption because data is read from parity blocks. RAID 5 is useful for people who need performance and constant access to their data.

In RAID 5+Spare, disk failure does not require immediate attention because the system rebuilds itself using the hot spare. However, the failed disk must be replaced as soon as

possible. A *spare disk* is an empty disk which is used by the RAID controller only when a disk fails. Figure 1-6 shows this example.

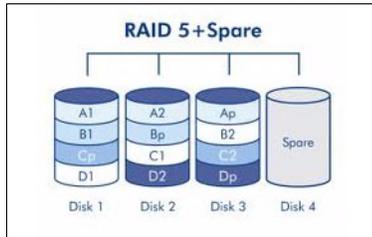


Figure 1-6 RAID 5 with hot spare

RAID 5 has better performance for I/O than RAID 1. And, depending on the number of disks that are used to build the RAID, the array disk space utilization is more than two-thirds. RAID 5 is also managed by a RAID controller which performs the same role as in RAID 1.

Figure 1-7 shows a brief comparison between the most common RAID levels.

RAID types: RAID 1 and 5 are the most common RAID levels. However, there are many other levels that are not covered in this book. Levels that are not described include: RAID 0, 3, 4, and 6; or nested (hybrid) types, such as RAID 0+1 or RAID 5+1. These hybrid types are used in environments where reliability and performance are key points to be covered from the storage perspective.

Features	RAID 0	RAID 1	RAID 5
Minimum # Drives	2	2	3
Data Protection	No Protection	Single-drive failure	Single-drive failure
Read Performance	High	High	High
Write Performance	High	Medium	Low
Read Performance (degraded)	N/A	Medium	Low
Write Performance (degraded)	N/A	High	Low
Capacity Utilization	100%	50%	67% - 94%
Typical Applications	High End Workstations, data logging, real-time rendering, very transitory data	Operating System, transaction databases	Data warehousing, web serving, archiving

Figure 1-7 RAID level comparison table

Our disk systems now seem to be ready to support failures, and they are also high performing. But what if our RAID controller fails? We might not lose data, but it is not accessible. Is there a solution to access this data?

It is almost the same scenario we initially faced when having only one disk as a storage system. This type of scenario is known as a *single point of failure (SPOF)*. We must add redundancy by introducing a secondary RAID controller to our storage system.

Now we are sure that no matter what happens, data is available to be used.

RAID controller role: The RAID controller role in some cases is performed by the software. This solution is less expensive than a hardware solution because it does not require controller hardware; however, it is a slower solution.

We now have a number of physical HDDs that are managed by two controllers.

Disk pools

When a logical storage volume needs to be provisioned to servers, first the storage RAID needs to be created. To create the RAID, select the available HDDs and group them together for a single purpose. The number of grouped HDDs depends on the RAID type that we choose and the space that is required for provisioning.

To understand what is meant, a basic example is shown that uses the following assumptions:

- ▶ There are 10 HDDs, named A, B, C, D, E, F, G, H, I, and J.
- ▶ There are two RAID controllers that support any RAID level, named RC1 and RC2.
- ▶ Each RAID controller can manage any HDD.
- ▶ Each RAID controller can act as a backup of the other, at any time.

The following tasks can now be performed:

- ▶ Select HDDs A, B, and F, and create a RAID 5 array that is managed by RC1. We call it G1.
- ▶ Select HDDs E, I, and J, and create a RAID 5 array that is managed by RC2. We call it G2.

Figure 1-8 shows these steps.

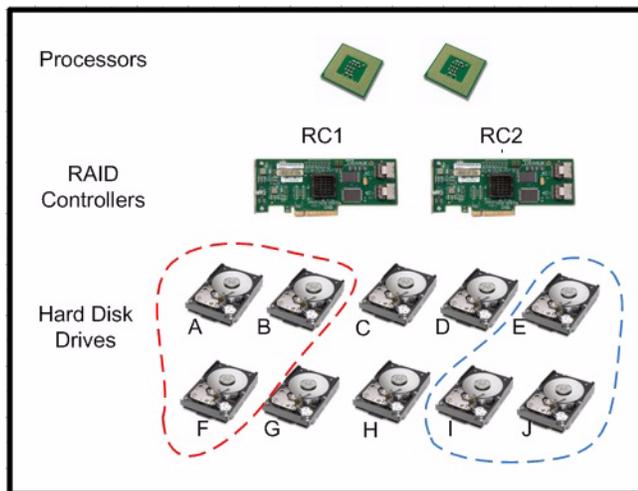


Figure 1-8 Disk pool creation

What we are doing by issuing these simple steps is creating disk pools. These pools basically consist of grouping disks together for a single purpose such as creating a RAID level, in our case, RAID 5.

In 1.3.2, “Redundant Array of Independent Disks” on page 5, we mentioned nested (hybrid) RAIDs such as 5+0. Solutions such as this are used when the amount of storage data is significant and is important for business continuity. RAID 50 consists of RAID 0 striping across lower-level RAID 5 arrays. The benefits of RAID 5 are gained while the spanned RAID 0 allows the incorporation of many more disks into a single logical drive. Up to one drive in each subarray might fail without loss of data. Also, rebuild times are substantially less than a single large RAID 5 array, as shown Figure 1-9 on page 9.

Nested (hybrid) RAIDs: Nested or hybrid RAIDs are a combination of existing RAID levels that create a RAID to reap the benefits of two different RAID levels.

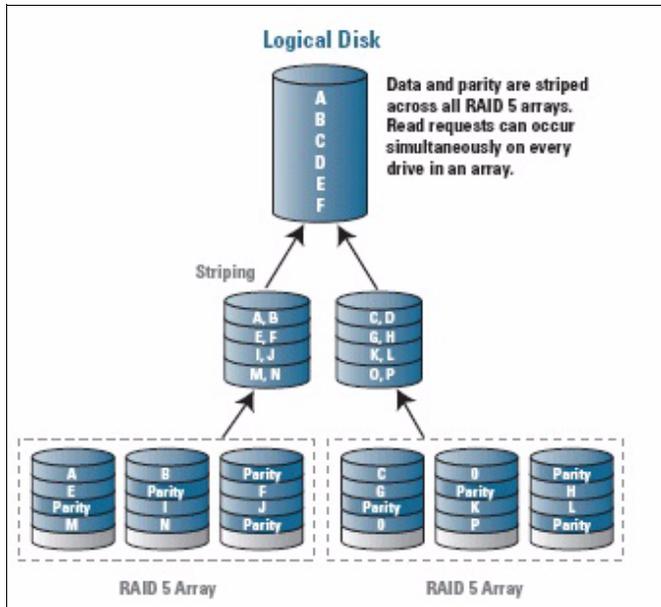


Figure 1-9 Nested (hybrid) RAID 5+0 or RAID 50

This nested RAID 50 can now be managed by RC1 or RC2 so we have full redundancy.

Storage systems

We are now not far away from building our basic storage system. However, to answer our previous questions, we need to add two new components and an enclosure.

One of those two components is a CPU which processes all the required instructions to allow data to flow. Adding one CPU creates a *single point of failure (SPOF)*, so we add two CPUs.

Furthermore, now that we almost have an independent system, and referring to the networking section, this system must be able to communicate with other systems in a network. Therefore, it requires a minimum of two network interfaces to be able to avoid a SPOF.

Now there is only one step that is left. The last step is to put all these hardware components into an enclosure. Now our storage system looks like what is shown in Figure 1-10.

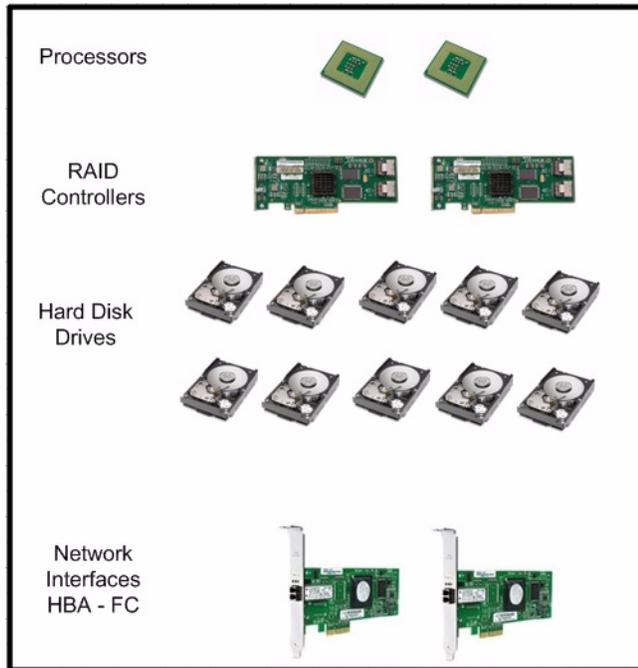


Figure 1-10 Basic storage system

This is only an example: This basic storage configuration is presented as an example. However, a configuration can have as many CPUs, RAID controllers, network interfaces, and HDDs as needed.

1.4 What is a storage area network?

The Storage Networking Industry Association (SNIA) defines the *storage area network (SAN)* as a network whose primary purpose is the transfer of data between computer systems and storage elements. A SAN consists of a communication infrastructure, which provides physical connections. It also includes a management layer, which organizes the connections, storage elements, and computer systems so that data transfer is secure and robust. The term SAN is usually (but not necessarily) identified with block I/O services rather than file access services.

In simple terms, a SAN is a specialized, high-speed network that attaches servers and storage devices. For this reason, it is sometimes referred to as the *network behind the servers*. A SAN allows an *any-to-any* connection across the network, by using interconnect elements such as switches and directors. It eliminates the traditional dedicated connection between a server and storage, and the concept that the server effectively *owns and manages* the storage devices. It also eliminates any restriction to the amount of data that a server can access, currently limited by the number of storage devices that are attached to the individual server. Instead, a SAN introduces the flexibility of networking to enable one server or many heterogeneous servers to share a common storage utility. A network might include many storage devices, including disk, tape, and optical storage. Additionally, the storage utility might be located far from the servers that it uses.

The SAN can be viewed as an extension to the storage *bus* concept. This concept enables storage devices and servers to be interconnected by using similar elements, such as LANs and wide area networks (WANs).

The diagram in Figure 1-11 shows a tiered overview of a SAN that connects multiple servers to multiple storage systems.

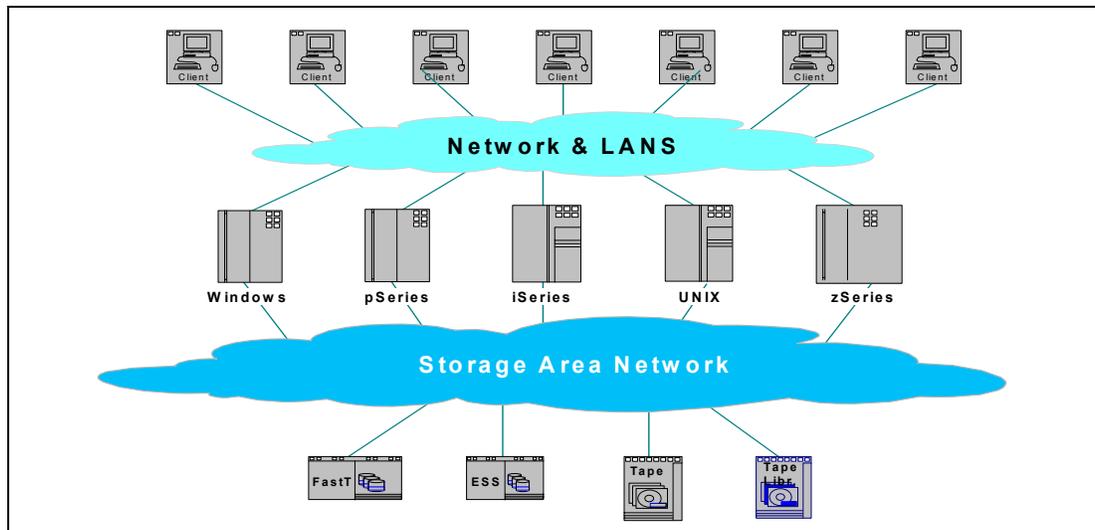


Figure 1-11 A SAN

SANs create new methods of attaching storage to servers. These new methods can enable great improvements in both availability and performance. The SANs of today are used to connect shared storage arrays and tape libraries to multiple servers, and are used by clustered servers for failover.

A SAN can be used to bypass traditional network bottlenecks. It facilitates direct, high-speed data transfers between servers and storage devices, potentially in any of the following three ways:

- ▶ Server to storage: This is the traditional model of interaction with storage devices. The advantage is that the same storage device might be accessed serially or concurrently by multiple servers.
- ▶ Server to server: A SAN might be used for high-speed, high-volume communications between servers.
- ▶ Storage to storage: This outboard data movement capability enables data to be moved without server intervention, therefore freeing up server processor cycles for other activities like application processing. Examples include a disk device that backs up its data to a tape device without server intervention, or a remote device mirroring across the SAN.

SANs allow applications that move data to perform better; for example, by having the data sent directly from the source to the target device with minimal server intervention. SANs also enable new network architectures where multiple hosts access multiple storage devices that are connected to the same network. Using a SAN can potentially offer the following benefits:

- ▶ Improvements to application availability: Storage is independent of applications and accessible through multiple data paths for better reliability, availability, and serviceability.
- ▶ Higher application performance: Storage processing is offloaded from servers and moved onto a separate network.
- ▶ Centralized and consolidated storage: Simpler management, scalability, flexibility, and availability.
- ▶ Data transfer and vaulting to remote sites: Remote copy of data that is enabled for disaster protection and against malicious attacks.
- ▶ Simplified centralized management: Single image of storage media simplifies management.

1.5 Storage area network components

Fibre Channel is the predominant architecture upon which most storage area network (SAN) implementations are built. IBM FICON® is the standard protocol for IBM z/OS® systems and Fibre Channel Protocol (FCP) is the standard protocol for open systems. The SAN components described in the following sections are Fibre Channel-based, and are shown in Figure 1-12 on page 13.

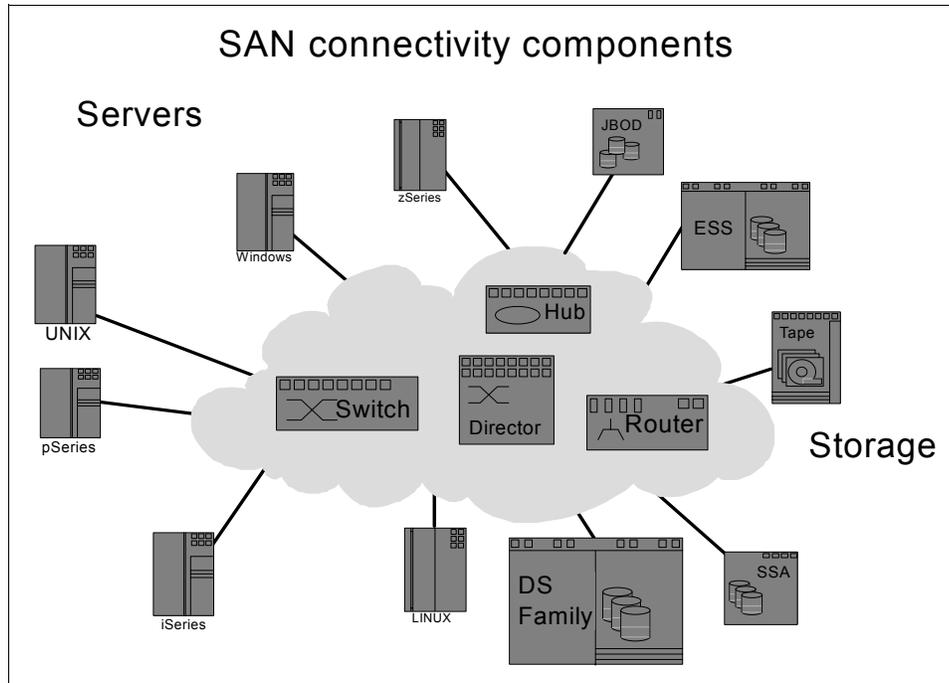


Figure 1-12 SAN components

1.5.1 Storage area network connectivity

The first element that must be considered in any storage area network (SAN) implementation is the connectivity of the storage and server components, which typically use Fibre Channel. The components that are listed in Figure 1-12 are typically used for LAN and WAN implementations. SANs, such as LANs, interconnect the storage interfaces together into many network configurations and across longer distances.

Much of the terminology that is used for SAN has its origins in Internet Protocol (IP) network terminology. In some cases, the industry and IBM use different terms that mean the same thing, and in some cases, mean different things.

1.5.2 Storage area network storage

The storage area network (SAN) liberates the storage device so it is not on a particular server bus, and attaches it directly to the network. In other words, storage is externalized and can be functionally distributed across the organization. The SAN also enables the centralization of storage devices and the clustering of servers. This has the potential to achieve easier and less expensive centralized administration that lowers the total cost of ownership (TCO).

The storage infrastructure is the foundation on which information relies, and therefore must support the business objectives and business model of a company. In this environment, simply deploying more and faster storage devices is not enough. A SAN infrastructure provides enhanced network availability, data accessibility, and system manageability. It is important to remember that a good SAN begins with a good design. This is not only a maxim, but must be a philosophy when we design or implement a SAN.

1.5.3 Storage area network servers

The server infrastructure is the underlying reason for all storage area network (SAN) solutions. This infrastructure includes a mix of server platforms such as Microsoft Windows, UNIX (and its various versions), and z/OS. With initiatives such as server consolidation and e-business, the need for SANs increase, making the importance of storage in the network greater.

1.6 The importance of standards or models

Why do we care about standards? Standards are the starting point for the potential interoperability of devices and software from different vendors in the SAN marketplace. SNIA, among others, defined and ratified the standards for the SANs of today, and will keep defining the standards for tomorrow. All of the players in the SAN industry are using these standards now, because these standards are the basis for the wide acceptance of SANs. Widely accepted standards potentially allow for the heterogeneous, cross-platform, and multivendor deployment of SAN solutions.

Because all vendors accepted these SAN standards, there *should* be no problem in connecting the different vendors into the same SAN network. However, nearly every vendor has an interoperability lab where it tests all kinds of combinations between their products and products of other vendors. Some of the most important aspects in these tests are the reliability, error recovery, and performance. If a combination passes the test, that vendor is going to certify or support this combination.

IBM participates in many industry standards organizations that work in the field of SANs. IBM believes that industry standards must be in place; and if necessary, redefined for SANs to be a major part of the IT business mainstream.

Probably the most important industry standards organization for SANs is the SNIA. IBM is a founding member and a board officer in SNIA. The SNIA, other standards organizations, and IBM are active participants in many of these organizations.



Why, and how, can we use a storage area network?

In Chapter 1, we introduced the basics by presenting a network and storage system introduction. We also worked on a standard *storage area network (SAN)* definition and brief description of the underlying technologies and concepts that are behind a SAN implementation.

In this chapter, we extend this discussion by presenting real-life SANs alongside well-known technologies and platforms that are used in SAN implementations. We also describe some of the trends that are driving the SAN evolution, and how they might affect the future of storage technology.

And although SAN technology is different, many of the concepts can also be applied in the Ethernet networking environment, which is covered in more depth later in this book.

2.1 Why use a storage area network?

This section describes the main motivators that drive storage area network (SAN) implementations, and present some of the key benefits that this technology might bring to data-dependent business.

2.1.1 The problem

Distributed clients and servers are frequently chosen to meet specific application needs. They might, therefore, run different operating systems (such as Windows Server, various UNIX offerings, IBM VMware vSphere, VMS). They might also run different database software (for example, IBM DB2®, Oracle, IBM Informix®, SQL Server). Therefore, they have different file systems and different data formats.

Managing this multi-platform, multivendor, networked environment is increasingly complex and costly. Software tools for multiple vendors and appropriately skilled human resources must be maintained to handle data and storage resource management on the many differing systems in the enterprise. Surveys that are published by industry analysts consistently show that management costs that are associated with distributed storage are much greater. The costs are shown to be up to 10 times more than the cost of managing consolidated or centralized storage. This comparison includes the costs of backup, recovery, space management, performance management, and disaster recovery planning.

Disk storage is often purchased from the processor vendor as an integral feature. It is difficult to establish if the price you pay per gigabyte (GB) is competitive, compared to the market price of disk storage. Disks and tape drives, directly attached to one client or server, cannot be used by other systems, leading to inefficient use of hardware resources. Organizations often find that they need to purchase more storage capacity, even though free capacity is available in other platforms.

Additionally, it is difficult to scale capacity and performance to meet rapidly changing requirements, such as the explosive growth in server, application, and desktop virtualization. There is also the need to manage information over its entire lifecycle, from conception to intentional destruction.

Information that is stored on one system cannot readily be made available to other users. One exception is to create duplicate copies and move the copy to the storage that is attached to another server. Movement of large files of data might result in significant degradation of performance of the LAN and WAN, causing conflicts with mission-critical applications. Multiple copies of the same data might lead to inconsistencies between one copy and another. Data that is spread on multiple small systems is difficult to coordinate and share for enterprise-wide applications. Some examples of this type of application include: E-business, enterprise resource planning (ERP), data warehouse, and business intelligence (BI).

Backup and recovery operations across a LAN might also cause serious disruption to normal application traffic. Even when using fast Gigabit Ethernet transport, the sustained throughput from a single server to tape is about 25 GB per hour. It would take approximately 12 hours to fully back up a relatively moderate departmental database of 300 GBs. This timeframe might exceed the available window of time in which the backup must be completed. And, it might not be a practical solution if business operations span multiple time zones. It is increasingly evident to IT managers that these characteristics of client/server computing are too costly and too inefficient. The islands of information that result from the distributed model of computing does not match the needs of the enterprise.

New ways must be found to control costs, improve efficiency, and simplify the storage infrastructure to meet the requirements of the modern business world.

2.1.2 The requirements

With this scenario in mind, there are a number of requirements that the storage infrastructures of today might consider. The following factors are some of the most important requirements to consider:

- ▶ **Unlimited and just-in-time scalability:** Businesses require the capability to flexibly adapt to the rapidly changing demands for storage resources without performance degradation.
- ▶ **System simplification:** Businesses require an easy-to-implement infrastructure with the minimum amount of management and maintenance. The more complex the enterprise environment, the more costs that are involved in terms of management. Simplifying the infrastructure can save costs and provide a greater return on investment (ROI).
- ▶ **Flexible and heterogeneous connectivity:** The storage resource must be able to support whatever platforms are within the IT environment. This resource is essentially an investment protection requirement that allows for the configuration of a storage resource for one set of systems. It later configures part of the capacity to other systems on an as-needed basis.
- ▶ **Security:** This requirement guarantees that data from one application or system does not become overlaid or corrupted by other applications or systems. Authorization also requires the ability to fence off the data of one system from other systems.
- ▶ **Encryption:** When sensitive data is stored, it must be read or written only from those authorized systems. If for any reason the storage system is stolen, data must never be available to be read from the system.
- ▶ **Hypervisors:** This requirement is for the support of the server, application, and desktop virtualization hypervisor features for cloud computing.
- ▶ **Speed:** Storage networks and devices must be able to manage the high number of gigabytes and intensive I/O that is required by each business industry.
- ▶ **Availability:** This is a requirement that implies both the protection against media failure and the ease of data migration between devices, without interrupting application processing. This requirement certainly implies improvements to backup and recovery processes. Attaching disk and tape devices to the same networked infrastructure allows for fast data movement between devices, which provides enhanced backup and recovery capabilities, such as:
 - **Serverless backup.** This is the ability to back up your data without using the computing processor of your servers.
 - **Synchronous copy.** This ensures that your data is at two or more places before your application goes to the next step.
 - **Asynchronous copy.** This ensures that your data is at two or more places within a short time. It is the disk subsystem that controls the data flow.

In the following section, we describe the use of SANs as a response to these business requirements.

2.2 How can we use a storage area network?

The key benefits that a storage area network (SAN) might bring to a highly data-dependent business infrastructure can be summarized into three concepts: Infrastructure simplification, information lifecycle management, and business continuity. They are an effective response to the requirements presented in the previous section, and are strong arguments for the adoption of SANs. These three concepts are briefly described, as follows.

2.2.1 Infrastructure simplification

There are four main methods by which infrastructure simplification can be achieved. An overview is provided for each of the main methods of infrastructure simplification:

- ▶ *Consolidation*

Concentrating the systems and resources into locations with fewer, but more powerful, servers and storage pools can help increase IT efficiency and simplify the infrastructure. Additionally, centralized storage management tools can help improve scalability, availability, and disaster tolerance.

- ▶ *Virtualization*

Storage virtualization helps in making complexity nearly transparent, and at the same time, can offer a composite view of storage assets. This feature might help reduce capital and administrative costs, and it provides users with better service and availability. Virtualization is designed to help make the IT infrastructure more responsive, scalable, and available.

- ▶ *Automation*

Choosing storage components with autonomic capabilities can improve availability and responsiveness, and can help protect data as storage needs grow. As soon as day-to-day tasks are automated, storage administrators might be able to spend more time on critical, higher-level tasks that are unique to the company's business mission.

- ▶ *Integration*

Integrated storage environments simplify system management tasks and improve security. When all servers have secure access to all data, your infrastructure might be better able to respond to the information needs of your users.

Figure 2-1 illustrates the consolidation movement from the distributed islands of information toward a single, and, most importantly, simplified infrastructure.

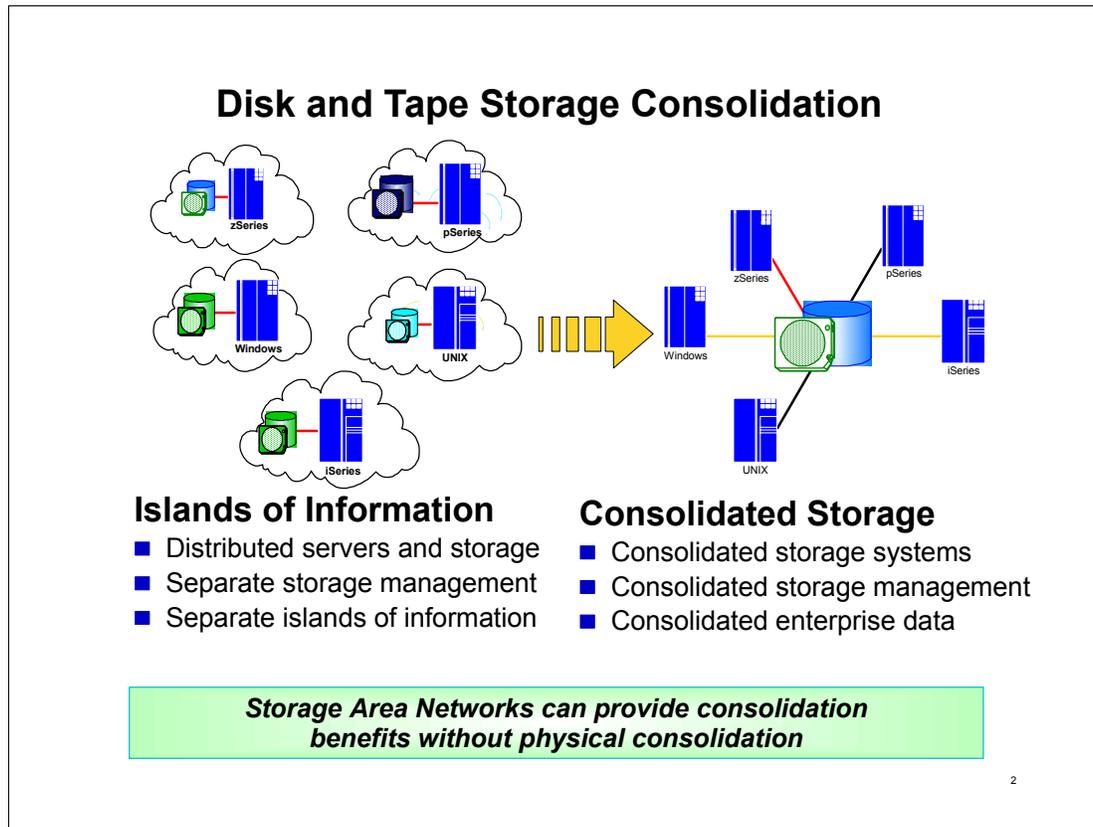


Figure 2-1 Disk and tape storage consolidation

Simplified storage environments have fewer elements to manage. This type of environment leads to increased resource utilization, simplifies storage management, and can provide economies of scale for owning disk storage servers. These environments can be more resilient and provide an infrastructure for virtualization and automation.

2.2.2 Information lifecycle management

Information is an increasingly valuable asset, but as the amount of information grows, it becomes increasingly costly and complex to store and manage it. Information lifecycle management (ILM) is a process for managing information through its lifecycle, from conception until intentional disposal. The ILM process manages this information in a manner that optimizes storage and maintains a high level of access at the lowest cost.

A SAN implementation makes it easier to manage the information lifecycle because it integrates applications and data into a single-view system, in which the information resides. This single-view location can be managed more efficiently.

IBM Tivoli Productivity Center For Data was specially designed to support ILM.

2.2.3 Business continuity

It goes without saying that the business climate in the on-demand era of today is highly competitive. Clients, employees, suppliers, and IBM Business Partners expect to be able to

tap into their information at any hour of the day, from any corner of the globe. Continuous business operations are no longer optional; they are a business imperative to becoming successful and maintaining a competitive advantage. Businesses must also be increasingly sensitive to issues of client privacy and data security so that vital information assets are not compromised. Also, factor in the legal and regulatory requirements, the inherent demands of participating in the global economy, and accountability. All of a sudden, the lot of an IT manager is not a happy one.

Currently, with natural disasters seemingly occurring with more frequency, a disaster recovery (DR) plan is essential. Implementing the correct SAN solution can help not only in real-time recovery techniques, but it also can reduce the recovery time objective (RTO) for your current DR plan.

There are many specific vendor solutions in the market which require a SAN running in the background like IBM VMware Site Recovery Manager (SRM) for business continuity.

It is little wonder that a sound and comprehensive business continuity strategy is now considered a business imperative, and SANs play a key role in this continuity. By deploying a consistent and safe infrastructure, SANs make it possible to meet any availability requirements.

2.3 Using the storage area network components

The foundation that a storage area network (SAN) is built on is the interconnection of storage devices and servers. This section further describes storage, interconnection components, and servers, and how the different types of servers and storage are used in a typical SAN environment.

2.3.1 Storage

This section briefly describes the main types of storage devices that can be found in the market.

Disk systems

By being contained within a single *box*, a disk system usually has a central control unit that manages all of the I/O. This configuration simplifies the integration of the system with other devices, such as other disk systems or servers.

We introduced you to what a storage system consists of in Chapter 1, “Introduction” on page 1. Depending on the specific functionality that is offered by a particular storage system, it is possible to make it behave as a small, mid-size, or enterprise solution. The decision of which type of disk system is more suitable for a SAN implementation strongly depends on the performance capacity and availability requirements for the particular SAN. We describe the product portfolio in Chapter 12, “The IBM product portfolio” on page 245.

Tape systems

Tape systems, in much the same way as disk systems, are devices that consist of all the necessary apparatus to manage the use of tapes for storage purposes. In this case, however, the serial nature of a tape makes it not possible for them to be treated in parallel. This treatment is because Redundant Array of Independent Disks (RAID) devices are leading to a simpler architecture to manage and use.

There are basically three types of tape systems: Drives, autoloaders, and libraries. An overview of each type of system is provided.

Tape drives

As with disk drives, tape drives are the means by which tapes can be connected to other devices. They provide the physical and logical structure for reading from, and writing to tapes.

Tape autoloaders

Tape autoloaders are autonomous tape drives that can manage tapes and perform automatic backup operations. They are usually connected to high-throughput devices that require constant data backup.

Tape libraries

Tape libraries are devices that can manage multiple tapes simultaneously and, as such, can be viewed as a set of independent tape drives or autoloaders. They are usually deployed in systems that require massive storage capacity, or that need some type of data separation that would result in multiple single-tape systems. Because a tape is not a random-access media, tape libraries cannot provide parallel access to multiple tapes as a way to improve performance. However, they can provide redundancy as a way to improve data availability and fault-tolerance.

The circumstances under which each of these systems, or even a disk system, might be used, strongly depend on the specific requirements of a particular SAN implementation. However, disk systems are usually used for online storage because of their superior performance. Whereas, tape systems are ideal for offline, high-throughput storage, because of the lower cost of storage per byte.

The next section describes the prevalent connectivity interfaces, protocols, and services for building a SAN.

2.3.2 Storage area network connectivity

Storage area network (SAN) connectivity consists of hardware and software components that make possible the interconnection of storage devices and servers. You are now introduced to the Fibre Channel (FC) model for SANs.

Standards and models for storage connectivity

Networking is governed by adherence to standards and models. Data transfer is also governed by standards. By far the most common is Small Computer System Interface (SCSI).

SCSI is an American National Standards Institute (ANSI) standard that is one of the leading I/O buses in the computer industry.

An industry effort was started to create a stricter standard allowing devices from different vendors to work together. This effort is recognized in the ANSI SCSI-1 standard. The SCSI-1 standard (circa 1985) is rapidly becoming obsolete. The current standard is SCSI-2. The SCSI-3 standard is in the production stage.

The SCSI bus is a parallel bus, which comes in a number of variants, as shown in Figure 2-2 on page 22.

Fibre Channel: For more information about parallel and serial data transfer, see Chapter 3, “Fibre Channel internals” on page 31.

SCSI Standard	Cable Length	Speed (MBps)	Devices Supported
SCSI-1	6	5	8
SCSI-2	6	5 to 10	8 or 16
Fast SCSI-2	3	10 to 20	8
Wide SCSI-2	3	20	16
Fast Wide SCSI-2	3	20	16
Ultra SCSI-3,8-bit	1.5	20	8
Ultra SCSI-3,16-bit	1.5	40	16
Ultra-2 SCSI	12	40	8
Wide Ultra-2 SCSI	12	80	16
Ultra-3 (Ultra160/m)	12	160	16

Figure 2-2 SCSI standards comparison table

In addition to a physical interconnection standard, SCSI defines a logical (command set) standard to which disk devices must adhere. This standard is called the Common Command Set (CCS) and was developed more or less in parallel with ANSI SCSI-1.

The SCSI bus not only has data lines, but also a number of control signals. An elaborate protocol is part of the standard to allow multiple devices to share the bus efficiently.

In SCSI-3, even faster bus types are introduced, along with serial SCSI buses that reduce the cabling overhead and allow a higher maximum bus length. It is at this point where the Fibre Channel model is introduced.

As always, the demands and needs of the market push for new technologies. In particular, there is always a push for faster communications without limitations on distance or on the number of connected devices.

Fibre Channel is a serial interface (primarily implemented with fiber-optic cable) and is the primary architecture for most SANs. To support this interface, there are many vendors in the marketplace that produce Fibre Channel adapters and other Fibre Channel devices. Fibre Channel brought these advantages by introducing a new protocol stack and by keeping the SCSI-3 CCS on top of it.

Figure 2-3 shows the evolution of Fibre Channel speeds. Fibre Channel is described in greater depth throughout this publication.

NAME	Throughput (Full duplex) (MBps)	Availability
1GFC	200	1997
2GFC	400	2001
4GFC	800	2005
8GFC	1600	2008
10GFC Serial	2550	2004
16GFC	3200	2011

Figure 2-3 Fibre Channel (FC) evolution

Figure 2-4 on page 23 shows an overview of the Fibre Channel model. The diagram shows the Fibre Channel, which is divided into four lower layers (FC-0, FC-1, FC-2, and FC-3) and

one upper layer (FC-4). FC-4 is where the upper level protocols are used, such as SCSI-3, Internet Protocol (IP), and Fibre Channel connection (FICON).

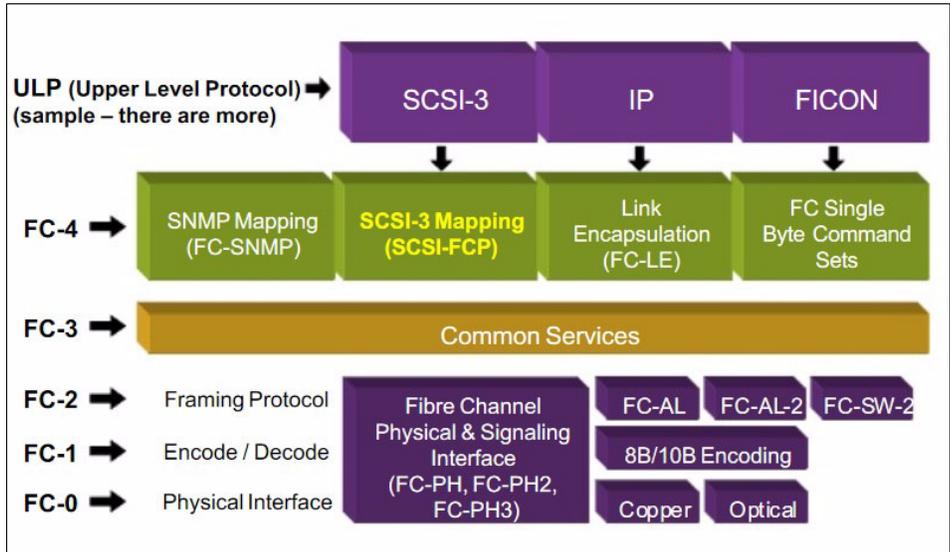


Figure 2-4 Fibre Channel (FC) model overview

Options for storage connectivity

In this section, we divided these components into three sections according to the abstraction level to which they belong: Lower-level layers, middle-level layers, and higher-level layers. Figure 2-5 provides an idea of each networking stack.

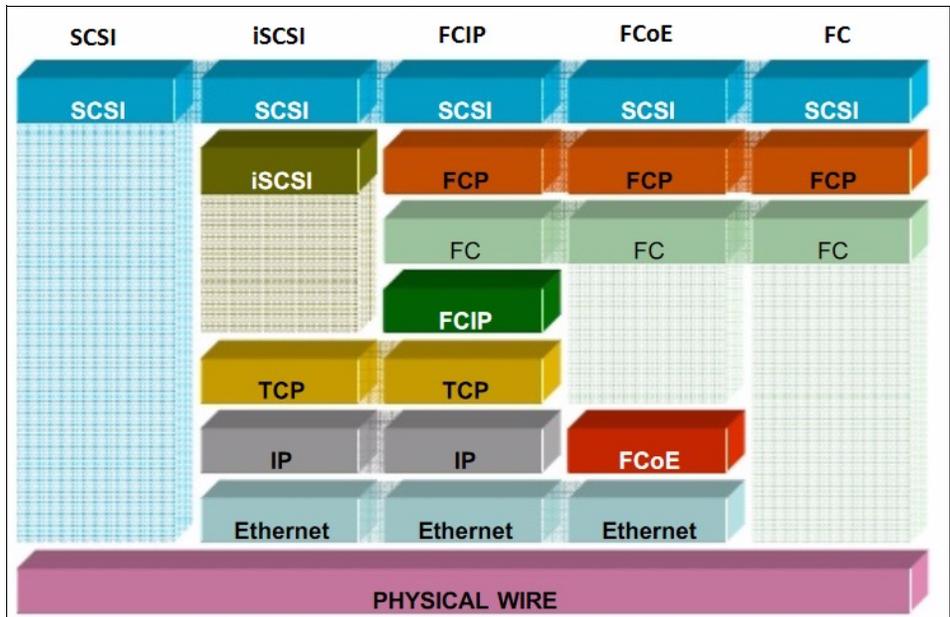


Figure 2-5 Networking stack comparison

Lower-level layers

As Figure 2-5 shows, there are only three stacks that can directly interact with the physical wire: Ethernet, SCSI, and Fibre Channel. Because of this configuration, these models are considered the *lower-level layers*. All of the other stacks are combinations of the layers, such

as Internet SCSI (iSCSI), Fibre Channel over IP (FCIP), and Fibre Channel over Ethernet (FCoE), also called the *middle-level layers*.

We are assuming a basic knowledge of Ethernet, which is typically used on conventional server-to-server or workstation-to-server network connections. The connections build up a common-bus topology by which every attached device can communicate with each other using this common bus. Ethernet speed is increasing as it becomes more pervasive in the data center. Key concepts of Ethernet are described later in this book.

Middle-level layers

This section consists of the transport protocol and session layers.

Fibre Channel over Ethernet (FCoE): FCoE is described later in this book. It is a vital model for the Converged Network Adapter (CNA).

Internet Small Computer System Interface

Internet Small Computer System Interface (iSCSI) is a transport protocol that carries SCSI commands from an initiator to a target. It is a data storage networking protocol that transports standard SCSI requests over the standard Transmission Control Protocol/Internet Protocol (TCP/IP) networking technology.

iSCSI enables the implementation of IP-based SANs, enabling clients to use the same networking technologies, for both storage and data networks. Because it uses TCP/IP, iSCSI is also suited to run over almost any physical network. By eliminating the need for a second network technology just for storage, iSCSI has the potential to lower the costs of deploying networked storage.

Fibre Channel Protocol

The *Fibre Channel Protocol (FCP)* is the interface protocol of SCSI on Fibre Channel (FC). It is a gigabit speed network technology that is primarily used for storage networking. Fibre Channel is standardized in the T11 Technical Committee of the InterNational Committee of Information Technology Standards (INCITS), an ANSI accredited standards committee. It started for use primarily in the supercomputer field, but is now the standard connection type for SANs in enterprise storage. Despite its name, Fibre Channel signaling can run on both twisted-pair copper wire and fiber optic cables.

Fibre Channel over IP

Fibre Channel over IP (FCIP) is also known as Fibre Channel tunneling or storage tunneling. It is a method to allow the transmission of Fibre Channel information to be tunneled through the IP network. Because most organizations already have an existing IP infrastructure, the attraction of being able to link geographically dispersed SANs, at a relatively low cost, is enormous.

FCIP encapsulates Fibre Channel block data and then transports it over a TCP socket. TCP/IP services are used to establish connectivity between remote SANs. Any congestion control and management, and data error and data loss recovery, is handled by TCP/IP services and does not affect Fibre Channel fabric services.

The major consideration with FCIP is that it does not replace Fibre Channel with IP; it allows deployments of Fibre Channel fabrics by using IP tunneling. The assumption that this might lead to is that the industry decided that Fibre Channel-based SANs are more than appropriate. Another possible assumption is that the only need for the IP connection is to facilitate any distance requirement that is beyond the current scope of an FCP SAN.

Fibre Channel connection

Fibre Channel connection (FICON) architecture is an enhancement of, rather than a replacement for, the traditional IBM Enterprise Systems Connection (ESCON®) architecture. A SAN is Fibre Channel-based (FC-based). Therefore, FICON is a prerequisite for IBM z/OS systems to fully participate in a heterogeneous SAN, where the SAN switch devices allow the mixture of open systems and mainframe traffic.

FICON is a protocol that uses Fibre Channel as its physical medium. FICON channels can achieve data rates up to 200 MBps full duplex and extend the channel distance (up to 100 km). FICON can also increase the number of control unit images per link and the number of device addresses per control unit link. The protocol can also retain the topology and switch management characteristics of ESCON.

Higher-level layers

This section consists of the presentation and application layers.

Server-attached storage

The earliest approach was to tightly couple the storage device with the server. This *server-attached storage* approach keeps performance overhead to a minimum. Storage is attached directly to the server bus by using an adapter, and the storage device is dedicated to a single server. The server itself controls the I/O to the device, issues the low-level device commands, and monitors device responses.

Initially, disk and tape storage devices had no onboard intelligence. They just ran the I/O requests of the server. Subsequent evolution led to the introduction of control units. These units are storage offload servers that contain a limited level of intelligence. They are able to perform functions, such as I/O request caching for performance improvements, or dual copy of data (RAID 1) for availability. Many advanced storage functions are developed and implemented inside the control unit.

Network-attached storage

Network-attached storage (NAS) is basically a LAN-attached file server that serves files by using a network protocol such as *Network File System (NFS)*. NAS is a term that is used to refer to storage elements that connect to a network and provide file access services to computer systems. A NAS storage element consists of an engine that implements the file services (by using access protocols such as NFS or Common Internet File System (CIFS)) and one or more devices, on which data is stored. NAS elements might be attached to any type of network. From a SAN perspective, a SAN-attached NAS engine is treated just like any other server. However, a NAS does not provide any of the activities that a server in a server-centric system typically provides, such as email, authentication, or file management.

NAS allows more hard disk storage space to be added to a network that already uses servers, without shutting them down for maintenance and upgrades. With a NAS device, storage is not a part of the server. Instead, in this storage-centric design, the server still handles all of the processing of data, but a NAS device delivers the data to the user. A NAS device does not need to be located within the server, but can exist anywhere in the LAN and can be made up of multiple networked NAS devices. These units communicate to a host by using Ethernet and file-based protocols. This method is in contrast to the disk units that are already described, which use Fibre Channel protocol and block-based protocols to communicate.

NAS storage provides acceptable performance and security, and it is often less expensive for servers to implement (for example, Ethernet adapters are less expensive than Fibre Channel adapters).

To bridge the two worlds and open up new configuration options for clients, some vendors, including IBM, sell NAS units that act as a gateway between IP-based users and SAN-attached storage. This configuration allows for the connection of the storage device and shares it between your high-performance database servers (attached directly through FC) and your users (attached through IP). These users do not have performance requirements nearly as strict.

NAS is an ideal solution for serving files that are stored on the SAN to users in cases where it would be impractical and expensive to equip users with Fibre Channel adapters. NAS allows those users to access your storage through the IP-based network that they already have.

2.3.3 Servers

Each of the different server platforms (IBM System z®, UNIX, IBM AIX®, HPUX, Sun Solaris, Linux, IBM i, and Microsoft Windows Server) implement SAN solutions by using various interconnects and storage technologies. The following sections review these solutions and the different implementations on each of the platforms.

Mainframe servers

In simple terms, a mainframe is a single, monolithic, and possibly multi-processor high-performance computer system. Apart from the fact that IT evolution is pointing toward a more distributed and loosely coupled infrastructure, mainframes still play an important role on businesses that depend on massive storage capabilities.

The IBM System z is a processor and operating system mainframe set. Historically, System z servers supported many different operating systems, such as z/OS, IBM OS/390®, VM, VSE, and TPF, which have been enhanced over the years. The processor to storage device interconnection also evolved from a bus and tag interface to ESCON channels, and now to FICON channels. Figure 2-6 shows the various processor-to-storage interface connections.

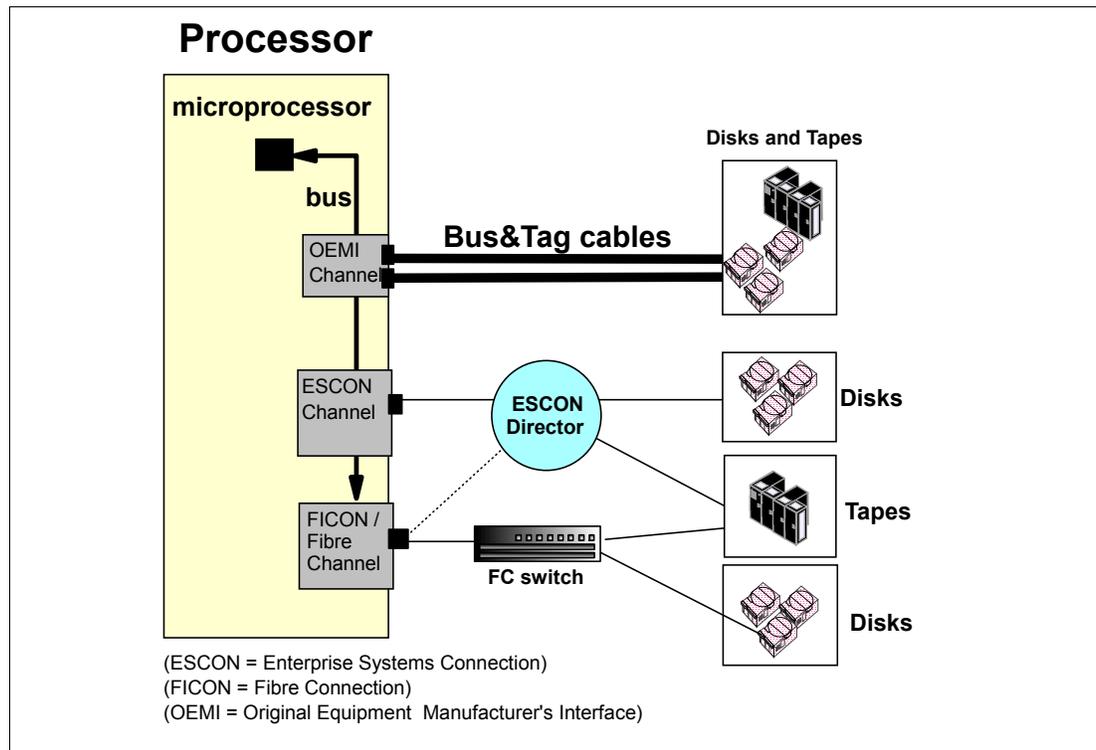


Figure 2-6 Processor-to-storage interface connections

Because of architectural differences, and strict data integrity and management requirements, the implementation of FICON is somewhat behind that of FCP on open systems. However, at the time of writing, FICON is caught up with FCP SANs, and they coexist amicably.

For the latest news on IBM zSeries® FICON connectivity, see this website:

<http://www-03.ibm.com/systems/z/hardware/connectivity/index.html>

In addition to FICON for traditional zSeries operating systems, IBM has standard Fibre Channel adapters for use with zSeries servers that can implement Linux.

UNIX based servers

Originally designed for high-performance computer systems, such as mainframes, the UNIX operating systems of today present on a great variety of hardware platforms, ranging from Linux based PCs to dedicated large-scale stations. Because of its popularity and maturity, it also plays an important role on both existing and earlier IT infrastructures.

The IBM System p® line of servers, running a UNIX operating system that is called *AIX*, offers various processor to storage interfaces, including SCSI, SAS (Serial Attached SCSI), and Fibre Channel. The Serial Storage Architecture (SSA) interconnection is primarily used for disk storage. Fibre Channel adapters are able to connect to tape and disk. Figure 2-7 shows the various processor-to-storage interconnect options for the System p family.

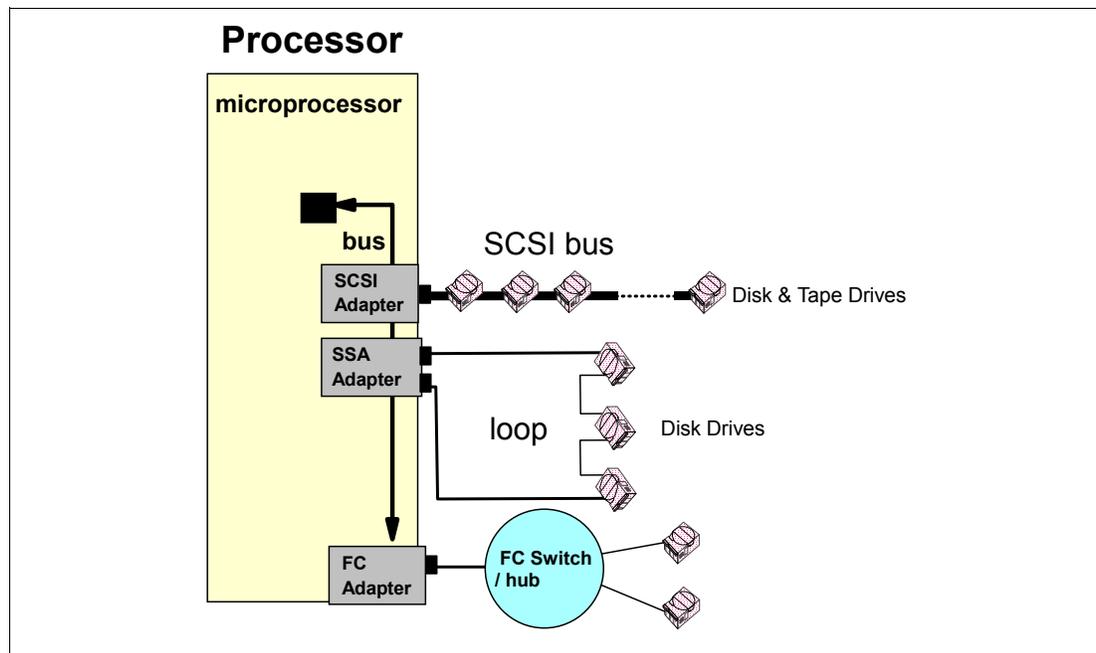


Figure 2-7 System p processor-to-storage interconnections

The various UNIX system vendors in the market deploy different variants of the UNIX operating system, each having some unique enhancements, and often supporting different file systems such as the journaled file system (JFS), enhanced journaled file system (JFS2), and the IBM Andrew File System (AFS®). The server-to-storage interconnect is similar to System p, as shown in Figure 2-7.

For the latest System p IBM Power Systems™ products, see this website:

<http://www.ibm.com/systems/storage/product/power.html>

Windows based servers

Based on the reports of various analysts regarding growth in the Windows server market (both in the number and size of Windows servers), Windows will become the largest market for SAN solution deployment. More and more Windows servers will host mission-critical applications that benefit from SAN solutions, such as disk and tape pooling, tape sharing, multipathing, and remote copy.

The processor-to-storage interfaces on IBM System x® servers (IBM Intel-based processors that support the Microsoft Windows Server operating system) are similar to the interfaces that are supported on UNIX servers, including SCSI and Fibre Channel.

For more information, see the IBM System x SAN website:

<http://www.ibm.com/systems/storage/product/x.html>

Single-level storage

Single-level storage (SLS) is probably the most significant differentiator in a SAN solution implementation on an IBM System i® server. This System i differentiator is a factor when compared to other systems such as z/OS, UNIX, and Windows. In IBM i, both the main storage (memory) and the auxiliary storage (disks) are treated as a large virtual address space that is known as SLS.

Figure 2-8 compares the IBM i SLS addressing with the way that Windows or UNIX systems work, by using the processor local storage. With 32-bit addressing, each process (job) has 4 GB of addressable memory. With 64-bit SLS addressing, over 18 million terabytes (18 exabytes) of addressable storage is possible. Because a single page table maps all virtual addresses to physical addresses, task switching is efficient. SLS further eliminates the need for address translation, thus speeding up data access.

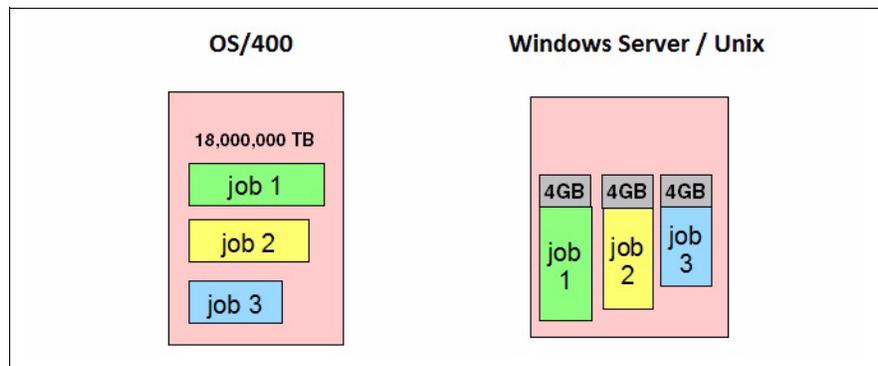


Figure 2-8 IBM i versus Windows Server 32 bits or UNIX storage addressing

The System i SAN support was rapidly expanded. System i servers now support attachment to switched fabrics and to most of the IBM SAN-attached storage products.

For more information, see the IBM System i SAN website:

<http://www.ibm.com/systems/i/hardware/storage/>

2.3.4 Putting the components together

After going through a myriad of technologies and platforms, we can easily understand why it is a challenge to implement true heterogeneous storage and data environments across different hardware and operating system platforms. Examples of such environments include: Disk and tape sharing across z/OS, IBM i, UNIX, and Windows Server.

One of the SAN principles, which is infrastructure simplification, cannot be easily achieved. Each platform, along with its operating system, treats data differently at various levels in the system architecture, thus creating some of these many challenges:

- ▶ Different attachment interfaces and protocols, such as SCSI, ESCON, and FICON.
- ▶ Different data formats, such as extended count key data (IBM ECKD™), blocks, clusters, and sectors.
- ▶ Different file systems, such as Virtual Storage Access Method (VSAM), JFS, JFS2, AFS, and Windows Server file system (NTFS).
- ▶ IBM i, with the concept of single-level storage.
- ▶ Different file system structures, such as catalogs and directories.
- ▶ Different file naming conventions, such as AAA.BBB.CCC and DIR/Xxx/Yyy.
- ▶ Different data encoding techniques, such as EBCDIC, ASCII, floating point, and little or big endian.

Figure 2-9 shows a brief summary of these differences for several systems.

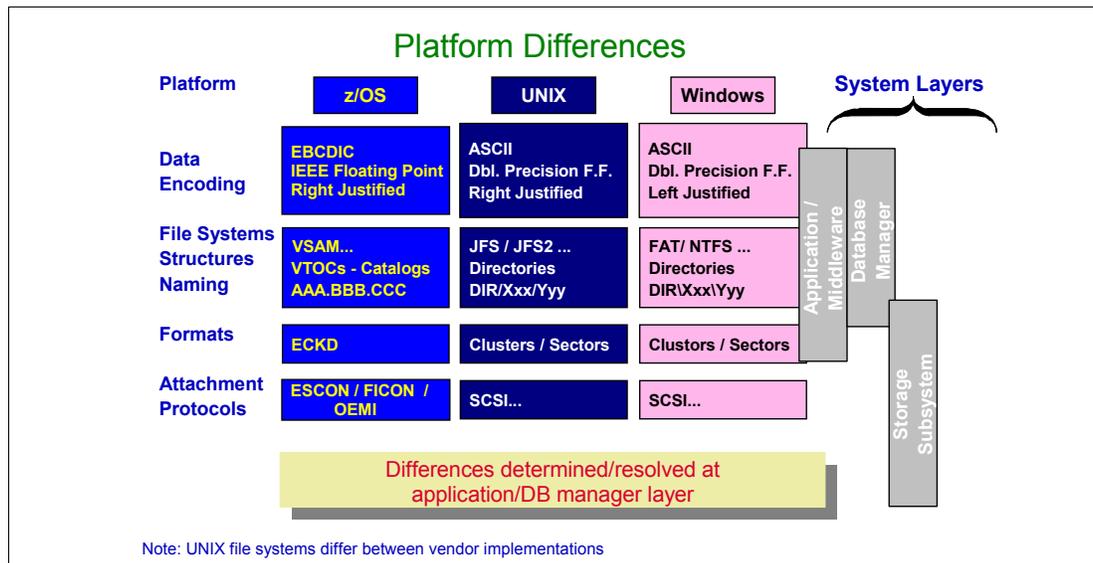


Figure 2-9 Hardware and operating systems differences



Fibre Channel internals

Fibre Channel (FC) is the predominant architecture upon which SAN implementations are built. Fibre Channel is a technology standard that allows data to be transferred at extremely high speeds. Current implementations support data transfers at up to 16 Gbps or even more. The Fibre Channel standard is accredited by many standards bodies, technical associations, vendors, and industry-wide consortiums. There are many products on the market that take advantage of the high-speed and high-availability characteristics of the Fibre Channel architecture.

Fibre Channel was developed through industry cooperation, unlike Small Computer System Interface (SCSI), which was developed by a vendor and submitted for standardization afterward.

Fibre or Fiber?: Fibre Channel was originally designed to support fiber optic cabling only. When copper support was added, the committee decided to keep the name in principle, but to use the UK English spelling (Fibre) when referring to the standard. The US English spelling (Fiber) is retained when referring generically to fiber optics and cabling.

Some people refer to Fibre Channel architecture as the fibre version of SCSI. Fibre Channel is an architecture that is used to carry intelligent peripheral interface (IPI) traffic, Internet Protocol (IP) traffic, Fibre Channel connection (FICON) traffic, Fibre Channel Protocol (FCP) (SCSI) traffic. Fibre Channel architecture might also carry traffic that uses other protocols, all on the standard Fibre Channel transport. An analogy might be Ethernet, where IP, Network Basic Input/Output System (NetBIOS), and Systems Network Architecture (SNA) are all used simultaneously over a single Ethernet adapter. This configuration is possible because these protocols all have mappings to Ethernet. Similarly, there are many protocols that are mapped onto Fibre Channel.

FICON is the standard protocol for z/OS, and will replace all Enterprise Systems Connection (ESCON) environments over time. FCP is the standard protocol for open systems, both using Fibre Channel architecture to carry the traffic.

3.1 First, why the Fibre Channel architecture?

Before we delve into the internals of Fibre Channel, we describe why Fibre Channel became the predominant SAN architecture.

3.1.1 The Small Computer Systems Interface legacy

The Small Computer System Interface (SCSI) is the conventional, server-centric method of connecting peripheral devices (disks, tapes, and printers) in the open client/server environment. As its name indicates, it was designed for the personal computer (PC) and small computer environment. It is a bus architecture, with dedicated, parallel cabling between the host and storage devices, such as disk arrays. This configuration is similar in implementation to the Original Equipment Manufacturer's Information (OEMI) bus and tag interface that was commonly used by mainframe computers until the early 1990s.

In addition to being a physical transport, SCSI is also a protocol. This protocol specifies commands and controls for sending blocks of data between the host and the attached devices. SCSI commands are issued by the host operating system, in response to user requests for data. Some operating systems, for example, Microsoft Windows NT, treat all attached peripherals as SCSI devices, and issue SCSI commands to deal with all read and write operations. SCSI was used in direct-attached storage (DAS) with internal and external devices that are connected via the SCSI channel in daisy chain fashion.

Typical SCSI device connectivity is shown in Figure 3-1.

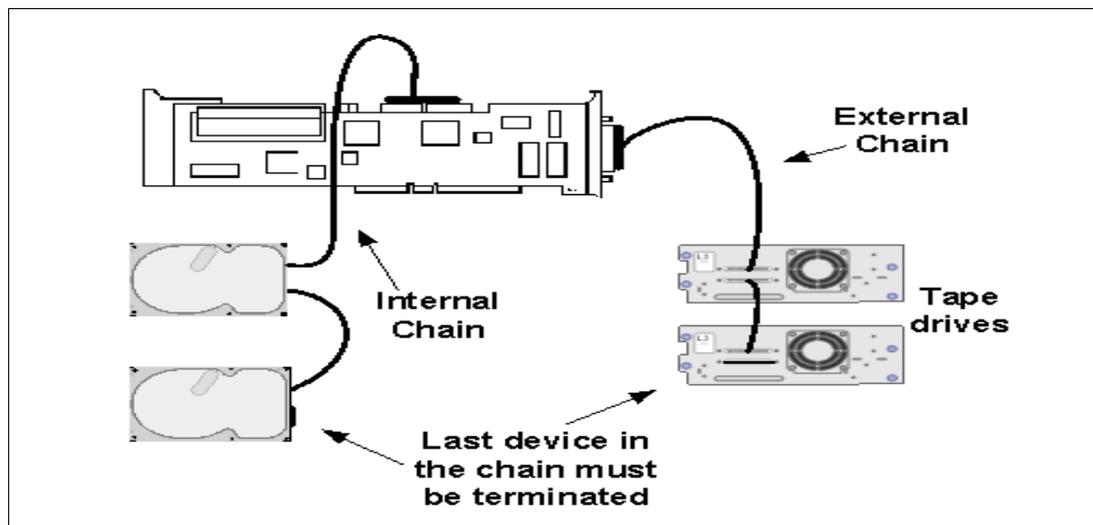


Figure 3-1 SCSI device connectivity

3.1.2 Limitations of the Small Computer System Interface

Some of the limitations of the Small Computer System Interface (SCSI) are described in the following topics.

Scalability limitations

The amount of data that is available to the server is determined by the number of devices which can attach to the bus. The amount is also determined by the number of buses that are attached to the server. Up to 15 devices can be attached to a server on a single SCSI bus. In

practice, because of performance limitations due to arbitration, it is common for no more than four or five devices to be attached in this way. This factor limits the scalability in terms of the number of devices that are able to be connected to the server.

Reliability and availability limitations

The SCSI shares aspects with bus and tag; for example, the cables and connectors are bulky, relatively expensive, and are prone to failure. Access to data is lost in the event of a failure of any of the SCSI connections to the disks. Data is also lost in the event of reconfiguration or servicing of a disk device that is attached to the SCSI bus. This loss is because all of the devices in the string must be taken offline. In today's environment, when many applications need to be available continuously, this downtime is unacceptable.

Speed and latency limitations

The data rate of the SCSI bus is determined by the number of bits transferred, and the bus cycle time (measured in megahertz (MHz)). Decreasing the cycle time increases the transfer rate. However, because of limitations inherent in the bus architecture, it might also reduce the distance over which the data can be successfully transferred. The physical transport was originally a parallel cable that consisted of eight data lines to transmit 8 bits in parallel, plus control lines. Later implementations widened the parallel data transfers to 16-bit paths (SCSI wide) to achieve higher bandwidths.

A SCSI propagation delay in sending data in parallel along multiple lines, leads to a phenomenon known as *skew*. Skew means that all bits might not arrive at the target device at the same time. Figure 3-2 shows this result.

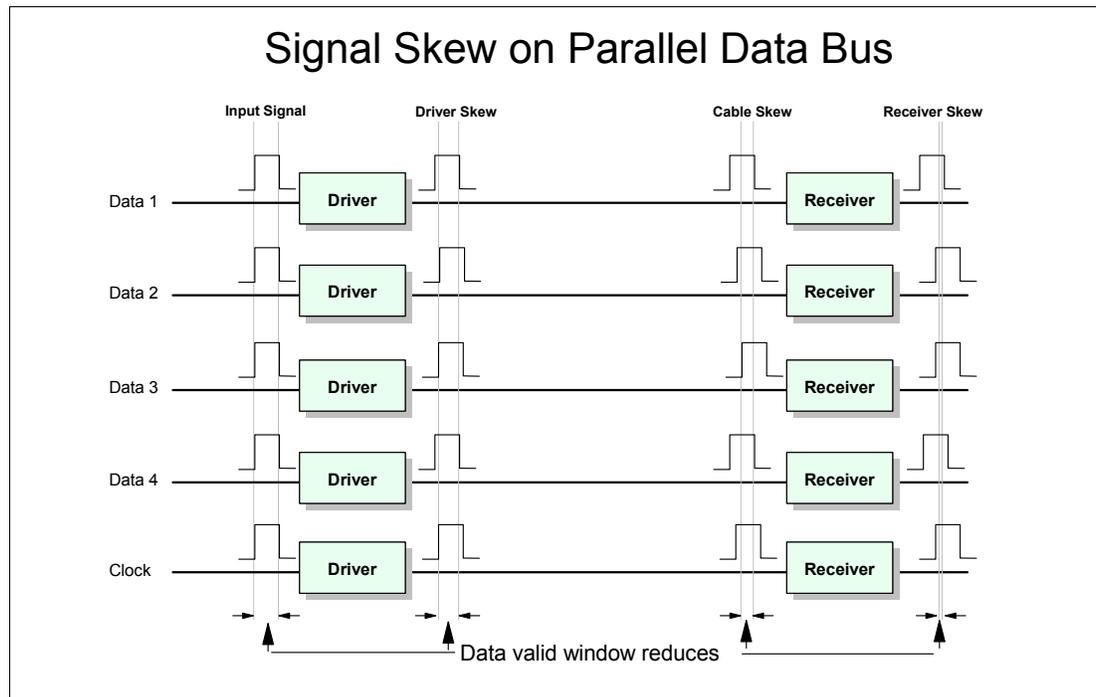


Figure 3-2 A SCSI propagation delay results in skew

Arrival occurs during a small window of time, depending on the transmission speed and the physical length of the SCSI bus. The need to minimize the skew limits the distance that devices can be positioned away from the initiating server to 2 - 25 meters. The distance depends on the cycle time. Faster speed means shorter distance.

Distance limitations

The distances refer to the maximum length of the SCSI bus, including all attached devices. The SCSI distance limitations are shown in Figure 3-3. These limitations might severely restrict the total GB capacity of the disk storage which can be attached to an individual server.

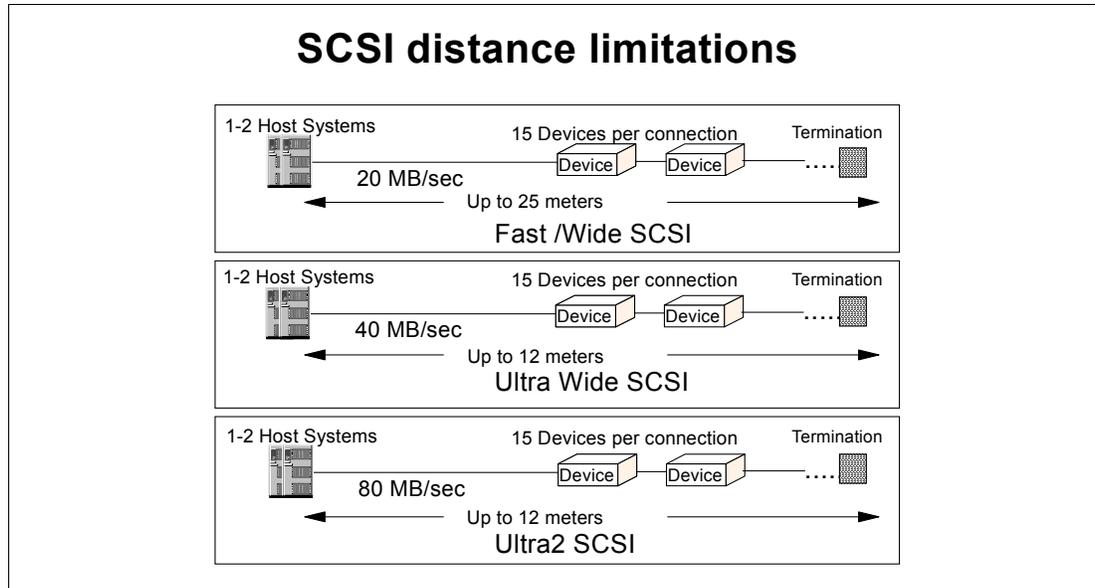


Figure 3-3 SCSI bus distance limitations

Device sharing

Many applications require the system to access several devices, or for several systems to share a single device. SCSI can enable this sharing by attaching multiple servers or devices to the same bus. This structure is known as a *multi-drop configuration*.

Figure 3-4 shows a multi-drop bus structure configuration.

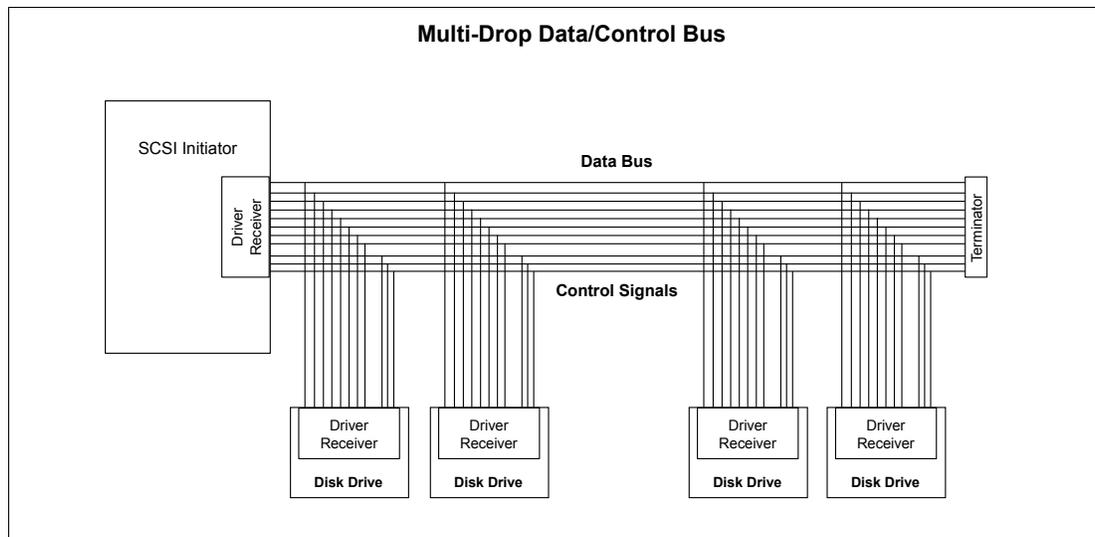


Figure 3-4 Multi-drop bus structure

To avoid signal interference, and therefore possible data corruption, all unused ports on a parallel SCSI bus must be properly terminated. Incorrect termination can result in transaction errors or failures.

Normally, only a single server can access data on a specific disk with a SCSI bus. In a shared bus environment, it is clear that all devices cannot transfer data at the same time. SCSI uses an arbitration protocol to determine which device can gain access to the bus. Arbitration occurs before and after every data transfer on the bus. While arbitration takes place, no data movement can occur. This loss of movement represents an additional performance overhead which reduces bandwidth utilization, substantially reducing the effective data rate achievable on the bus. Actual rates are typically less than 50% of the rated speed of the SCSI bus.

It is clear that the physical parallel SCSI bus architecture has a number of significant speed, distance, and availability limitations. These limits make it increasingly less suitable for many applications in the networked IT infrastructure of today. However, the SCSI protocol is deeply embedded in the way that commonly encountered operating systems handle user requests for data. Therefore, it would be a major inhibitor to progress if you had to move to new protocols.

3.1.3 Why Fibre Channel?

Fibre Channel is an open, technical standard for networking which incorporates the *channel transport* characteristics of an I/O bus, with the flexible connectivity and distance characteristics of a traditional network.

Because of its channel-like qualities, hosts and applications see storage devices that are attached to the SAN as though they are locally attached storage. Because of its network characteristics, it can support multiple protocols and a broad range of devices, and it can be managed as a network. Fibre Channel can use either optical fiber (for distance) or copper cable links (for short distance at low cost).

Fibre Channel is a multi-layered network, which is based on a series of American National Standards Institute (ANSI) standards that define characteristics and functions for moving data across the network. These standards include definitions of physical interfaces. Examples include: Cabling, distances, and signaling; data encoding and link controls; data delivery in terms of frames; flow control and classes of service; common services; and protocol interfaces.

Like other networks, information is sent in structured packets or frames, and data is serialized before transmission. But, unlike other networks, the Fibre Channel architecture includes a significant amount of hardware processing to deliver high performance.

Fibre Channel uses a serial data transport scheme, similar to other computer networks, which stream packets (frames) of bits, one behind the other, in a single data line to achieve high data rates.

Serial transfer by its very nature, does not suffer from the problem of skew, so speed and distance are not restricted like with parallel data transfers. Figure 3-5 shows the process of parallel versus serial data transfers.

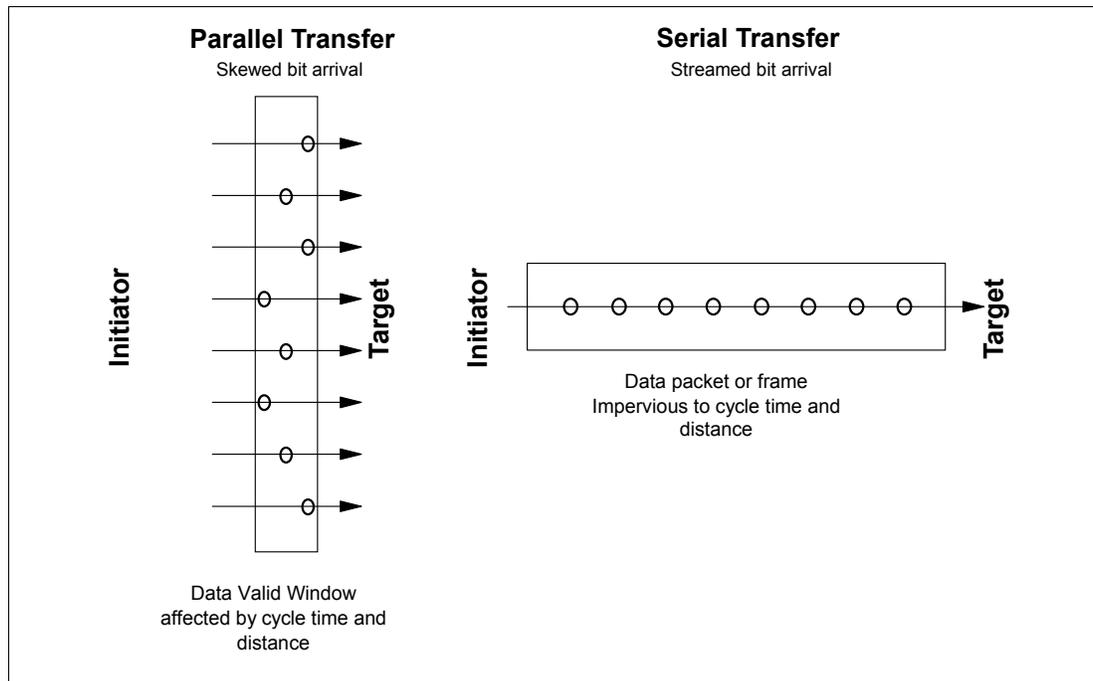


Figure 3-5 Parallel data transfers versus serial data transfers

Serial transfer enables simpler cabling and connectors, and also routing of information through switched networks. Fibre Channel can operate over longer distances, both natively and by implementing cascading, and longer with the introduction of repeaters. Just as LANs can be interlinked in WANs by using high speed gateways, so can campus SANs be interlinked to build enterprise-wide SANs.

Whatever the topology, information is sent between two nodes, which are the source (transmitter or initiator) and destination (receiver or target). A *node* is a device, such as a server (PC, workstation, or mainframe) or peripheral device, such as a disk or tape drive, or a video camera. Frames of information are passed between nodes, and the structure of the frame is defined by a protocol. Logically, a source and target node must use the same protocol, but each node might support several different protocols or data types.

Therefore, Fibre Channel architecture is extremely flexible in its potential application. Fibre Channel transport layers are protocol independent, enabling the transmission of multiple protocols.

Using a credit-based flow control approach, Fibre Channel is able to deliver data as fast as the destination device buffer is able to receive it. And low transmission overheads enable high sustained utilization rates without loss of data.

Therefore, Fibre Channel combines the best characteristics of traditional I/O channels with the characteristics of computer networks:

- ▶ High performance for large data transfers by using simple transport protocols and extensive hardware assists.
- ▶ Serial data transmission.
- ▶ A physical interface with a low error rate definition.

- ▶ Reliable transmission of data with the ability to guarantee or confirm error free delivery of the data.
- ▶ The ability to package data in packets (frames, in Fibre Channel terminology).
- ▶ Flexibility in terms of the types of information which can be transported in frames (such as data, video, and audio),
- ▶ Use of existing device-oriented command sets, such as SCSI and FCP.
- ▶ A vast expansion in the number of devices which can be addressed when compared to I/O interfaces: A theoretical maximum of more than 15 million ports.

There are several factors which make the Fibre Channel architecture ideal for the development of enterprise SANs. One example is the high degree of flexibility, availability, and scalability of the architecture. Other factors include the combination of multiple protocols at high speeds over long distances, and the broad acceptance of the Fibre Channel standards by vendors throughout the IT industry.

The topics that follow describe some of the key concepts that are mentioned in the previous pages and that are behind Fibre Channel SAN implementations. We also introduce more Fibre Channel SAN terminology and jargon that you can expect to encounter.

3.2 Layers

Fibre Channel (FC) is broken up into a series of five layers. The concept of *layers*, starting with the International Organization for Standardization/open systems interconnection (ISO/OSI) seven-layer model, allows the development of one layer to remain independent of the adjacent layers. Although a Fibre Channel contains five layers, those layers follow the general principles that are stated in the ISO/OSI model.

The series of five layers that make up a Fibre Channel, can be categorized into these two layers:

- ▶ Physical and signaling layer
- ▶ Upper layer

Fibre Channel is a layered protocol. Figure 3-6 on page 38 shows the upper and physical layers.

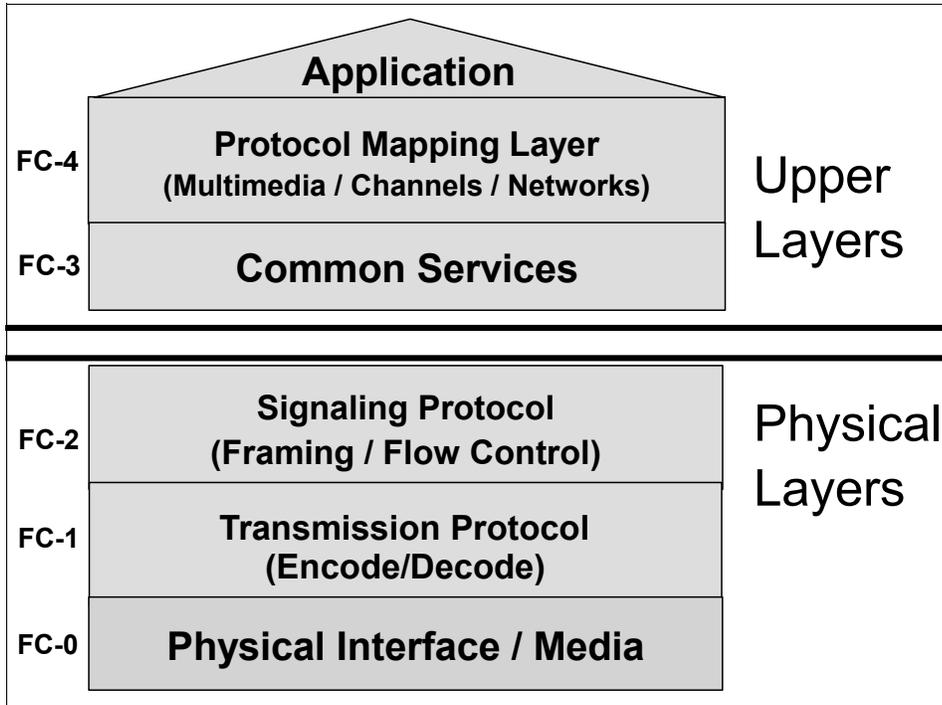


Figure 3-6 Fibre Channel upper and physical layers

The FC layers are briefly described in the following definitions:

Physical and signaling layers

The physical and signaling layers include the three lowest layers: FC-0, FC-1, and FC-2.

Physical interface and media: FC-0

The lowest layer, *FC-0*, defines the physical link in the system, including the cabling, connectors, and electrical parameters for the system at a wide range of data rates. This level is designed for maximum flexibility, and allows the use of many technologies to match the needs of the configuration.

A communication route between two nodes can be made up of links of different technologies. For example, in reaching its destination, a signal might start out on copper wire and become converted to single-mode fiber for longer distances. This flexibility allows for specialized configurations, depending on IT requirements.

Laser safety

Fibre Channel often uses lasers to transmit data, and can, therefore, present an optical health hazard. The FC-0 layer defines an open fiber control (OFC) system, and acts as a safety interlock for point-to-point fiber connections that use semiconductor laser diodes as the optical source. If the fiber connection is broken, the ports send a series of pulses until the physical connection is re-established and the necessary handshake procedures are followed.

Transmission protocol: FC-1

The second layer, *FC-1*, provides the methods for adaptive 8B/10B encoding to bind the maximum length of the code, maintain DC-balance, and provide word alignment. This layer is used to integrate the data with the clock information required by serial transmission technologies.

Framing and signaling protocol: FC-2

Reliable communications result from the *FC-2* framing and signaling protocol of the FC. FC-2 specifies a data transport mechanism that is independent of upper layer protocols. FC-2 is self-configuring and supports point-to-point, arbitrated loop, and switched environments.

FC-2, which is the third layer of the *Fibre Channel Physical and Signaling interface (FC-PH)*, provides the transport methods to determine the following factors:

- ▶ Topologies that are based on the presence or absence of a fabric
- ▶ Communication models
- ▶ Classes of service that are provided by the fabric and the nodes
- ▶ General fabric model
- ▶ Sequence and exchange identifiers
- ▶ Segmentation and reassembly

Data is transmitted in 4-byte ordered sets that contain data and control characters. Ordered sets provide the availability to obtain bit and word synchronization, which also establishes word boundary alignment.

Together, FC-0, FC-1, and FC-2 form the FC-PH.

Upper layers

The upper layer includes two layers: FC-3 and FC-4.

Common services: FC-3

FC-3 defines functions that span multiple ports on a single-node or fabric. Functions that are currently supported include the following features:

- ▶ Hunt groups

A *hunt group* is a set of associated N_Ports that is attached to a single node. This set is assigned an alias identifier that allows any frames that contain the alias to be routed to any available N_Port within the set. This process decreases the latency in waiting for an N_Port to become available.

- ▶ Striping

Striping is used to multiply bandwidth by using multiple N_Ports in parallel to transmit a single information unit across multiple links.

- ▶ Multicast

Multicast delivers a single transmission to multiple destination ports. This method includes the ability to broadcast to all nodes or a subset of nodes.

Upper layer protocol mapping: FC-4

The highest layer, *FC-4*, provides the application-specific protocols. Fibre Channel is equally adept at transporting both the network and channel information and allows both protocol types to be concurrently transported over the same physical interface.

Through mapping rules, a specific FC-4 describes how upper layer protocol (ULP) processes of the same FC-4 type interoperate.

A channel example is FCP. This protocol is used to transfer SCSI data over Fibre Channel. A networking example is sending IP packets between the nodes. FICON is another ULP in use today for mainframe systems. FICON is a contraction of *Fibre Connection* and refers to running ESCON traffic over Fibre Channel.

3.3 Optical cables

An optical fiber is a thin strand of silica glass and its geometry is quite like a human hair. In reality, it is a narrow, long glass cylinder with special characteristics. When light enters one end of the fiber, it travels (confined within the fiber) until it leaves the fiber at the other end. Two critical factors stand out:

- ▶ Little light is lost in its journey along the fiber.
- ▶ Fiber can bend around corners and the light stays within it and is guided around the corners.

An optical fiber consists of two parts: The core and the cladding. See Figure 3-7 on page 43. The core is a narrow cylindrical strand of glass and the cladding is a tubular jacket that surrounds it. The core has a (slightly) higher refractive index than the cladding. This means that the boundary (interface) between the core and the cladding acts as a perfect mirror. Light traveling along the core is confined by the mirror to stay within it, even when the fiber bends around a corner.

When light is transmitted on a fiber, the most important consideration is *what kind of light?* The electromagnetic radiation that we call *light* exists at many wavelengths. These wavelengths go from invisible infrared through all the colors of the visible spectrum to invisible ultraviolet. Because of the attenuation characteristics of fiber, we are only interested in infrared *light* for communication applications. This light is usually invisible, since the wavelengths used are usually longer than the visible limit of around 750 nanometers (nm).

If a short pulse of light from a source such as a laser or an LED is sent down a narrow fiber, it is changed (degraded) by its passage. It emerges (depending on the distance) much weaker, lengthened in time (*smearred out*), and distorted in other ways. The reasons for this transformation are described in the topics that follow.

3.3.1 Attenuation

The pulse is weaker because all glass absorbs light. More accurately, impurities in the glass can absorb light but the glass itself does not absorb light at the wavelengths of interest. In addition, variations in the uniformity of the glass cause scattering of the light. Both the rate of light absorption and the amount of scattering are dependent on the wavelength of the light and the characteristics of the particular glass. Most light loss in a modern fiber is caused by scattering.

3.3.2 Maximum power

There is a practical limit to the amount of power that can be sent on a fiber. This limit is about half a watt (in a standard single-mode fiber). This size is because of a number of non-linear effects that are caused by the intense electromagnetic field in the core when high power is present.

Polarization

Conventional communication optical fiber is cylindrically symmetric, but contains imperfections. Light traveling down such a fiber is changed in polarization (in current optical systems, this change does not matter, but in future systems it might become a critical issue).

Dispersion

Dispersion occurs when a pulse of light is spread out during transmission on the fiber. A short pulse becomes longer and ultimately joins with the pulse behind, making recovery of a reliable bit stream not possible. (In most communications systems, bits of information are sent as pulses of light: 1 = light, 0 = dark. But even in analog transmission systems where information is sent as a continuous series of changes in the signal, dispersion causes distortion.) There are many kinds of dispersion, each of which works in a different way. The three most important kinds of dispersion are described:

Material dispersion (chromatic dispersion)

Both lasers and LEDs produce a range of optical wavelengths (a band of light) rather than a single narrow wavelength. The fiber has different refractive index characteristics at different wavelengths; therefore, each wavelength travels at a different speed in the fiber. Thus, some wavelengths arrive before others and a signal pulse disperses (or smears out).

Modal dispersion

When you use a multimode fiber, the light is able to take many different paths or *modes* as it travels within the fiber. The distance traveled by light in each mode is different from the distance traveled in other modes. When a pulse is sent, parts of that pulse (rays or quanta) take many different modes (usually all available modes). Therefore, some components of the pulse arrive before others. The difference between the arrival time of light that takes the fastest mode, versus the slowest, obviously gets greater as the distance gets greater.

Waveguide dispersion

Waveguide dispersion is a complex effect and is caused by the shape and index profile of the fiber core. However, this effect can be controlled by careful design and, in fact, waveguide dispersion can be used to counteract material dispersion.

Noise

One of the great benefits of fiber optical communications is that the fiber does not pick up noise from outside the system. However, there are various kinds of noise that can come from components within the system itself. Mode partition noise can be a problem in single-mode fiber and modal noise is a phenomenon in multimode fiber.

It is not our intention to delve any deeper into optical than what is already described.

3.3.3 Fiber in the storage area network

Fibre Channel can be run over optical or copper media, but fiber-optic cables enjoy a major advantage in noise immunity. It is for this reason that fiber-optic cabling is preferred. However, copper is also used. It is likely that in the short term, a mixed environment needs to be tolerated and supported. Although, a mixed environment is less likely to be needed as storage area networks (SANs) mature.

In addition to the noise immunity, fiber-optic cabling provides a number of distinct advantages over copper transmission lines that make it an attractive medium for many applications. The following advantages are at the forefront:

- ▶ Greater distance capability than is generally possible with copper
- ▶ Insensitive to induced electromagnetic interference (EMI)
- ▶ No emitted electromagnetic radiation (RFI)
- ▶ No electrical connection between two ports
- ▶ Not susceptible to crosstalk
- ▶ Compact and lightweight cables and connectors

However, fiber-optic and optical links do have some drawbacks. Some of the drawbacks include the following considerations:

- ▶ Optical links tend to be more expensive than copper links over short distances.
- ▶ Optical connections do not lend themselves to backplane-printed circuit wiring.
- ▶ Optical connections might be affected by dirt and other contamination.

Overall, optical fibers provide a high-performance transmission medium, which was refined and proven over many years.

Mixing fiber-optical and copper components in the same environment is supported, although not all products provide that flexibility. Product flexibility needs to be considered when you plan a SAN. Copper cables tend to be used for short distances, up to 30 meters (98 feet), and can be identified by their DB-9, 9 pin connector.

Normally, fiber-optic cabling is referred to by mode or the frequencies of light waves that are carried by a particular cable type. Fiber cables come in two distinct types, as shown in Figure 3-7.

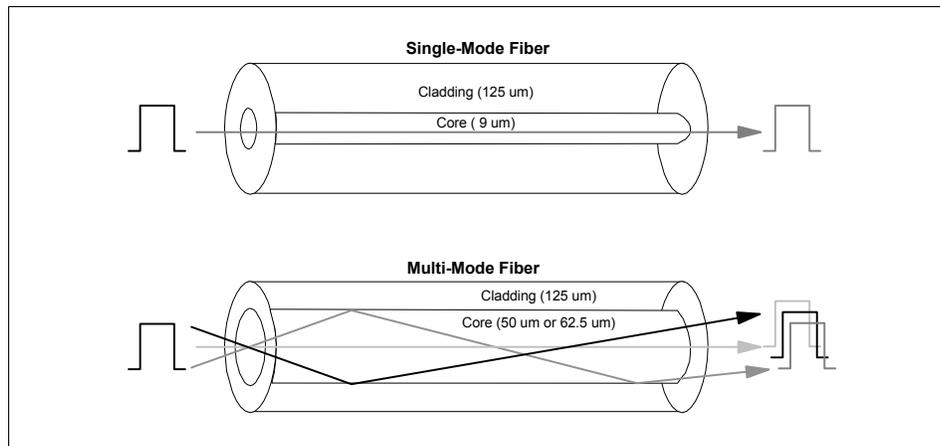


Figure 3-7 Cable types

► Multi-mode fiber for shorter distances

Multi-mode cabling is used with shortwave laser light and has either a 50-micron or a 62.5-micron core with a cladding of 125 micron. The 50-micron or 562.5-micron diameter is sufficiently large for injected light waves to be reflected off the core interior.

Multi-mode fiber (MMF) allows more than one mode of light. Common multi-mode core sizes are 50 micron and 62.5 micron. MMF fiber is better suited for shorter distance applications. Where costly electronics are heavily concentrated, the primary cost of the system does not lie with the cable. In such a case, MMF is more economical because it can be used with inexpensive connectors and laser devices, which reduces the total system cost.

► Single-mode fiber for longer distances

Single-mode fiber (SMF) allows only one pathway, or mode, of light to travel within the fiber. The core size is typically 8.3 micron. SMFs are used in applications where low signal loss and high data rates are required. An example of this type of application is on long spans between two system or network devices, where repeater and amplifier spacing needs to be maximized.

Fibre Channel architecture supports both short wave and long wave optical transmitter technologies in the following ways:

► Short wave laser

This technology uses a wavelength of 780 nanometers and is only compatible with MMF.

► Long wave laser

This technology uses a wavelength of 1300 nanometers. It is compatible with both SMF and MMF.

Different cable types along with their speed and distance, are listed in Table 3-1 on page 44.

Table 3-1 Fibre Channel modes, speeds, and distances

Fiber mode	Speed (Mbps)	Transmitter	Medium	Distance
Single-mode fiber	1600	1310 nm longwave light	1600-SM-LC-L	0.5 m - 10 km
		1490 nm longwave light	1600-SM-LZ-I	0.5 m - 2 km
	800	1310 nm longwave light	800-SM-LC-L	2 m - 10 km
			800-SM-LC-I	2 m - 1.4 km
	400	1310 nm longwave light	400-SM-LC-L	2 m - 10 km
			400-SM-LC-M	2 m - 4 km
			400-SM-LL-I	2 m - 2 km
	200	1550 nm longwave light	200-SM-LL-V	2 m - 50 km
		1310 nm longwave light	200-SM-LC-L	2 m - 10 km
			200-SM-LL-I	2 m - 2 km
	100	1550 nm longwave light	100-SM-LL-V	2 m - 50 km
		1310 nm longwave light	100-SM-LL-L	2 m - 10 km
			100-SM-LC-L	2 m - 10 km
			100-SM-LL-I	2 m - 2 km

Fiber mode	Speed (Mbps)	Transmitter	Medium	Distance
Multi-mode fiber ^a	1600	850 nm shortwave light	1600-M5F-SN-I	0.5 m - 125 m
			1600-M5E-SN-I	0.5 - 100 m
			1600-M5-SN-S	0.5 - 35 m
			1600-M6-SN-S	0.5 - 15 m
	800		800-M5F-SN-I	0.5 - 190 m
			800-M5E-SN-I	0.5 - 150 m
			800-M5-SN-S	0.5 - 50 m
			800-M6-SN-S	0.5 - 21 m
	400		400-M5F-SN-I	0.5 - 400 m
			400-M5E-SN-I	0.5 - 380 m
			400-M5-SN-I	0.5 - 150 m
			400-M6-SN-I	0.5 - 70 m
	200		200-M5E-SN-I	0.5 - 500 m
			200-M5-SN-I	0.5 - 300 m
			200-M6-SN-I	0.5 - 150 m
	100		100-M5E-SN-I	0.5 - 860 m
100-M5-SN-I		0.5 - 500 m		
100-M6-SN-I		0.5 - 300 m		
100-M5-SL-I		2 - 500 m		
100-M6-SL-I		2 - 175 m		

a. See Table 3-2 for multi-mode fiber (MMF) details

Table 3-2 shows MMF designations, optical multi-mode (OM) numbering, fiber-optic cable diameters, and FC media designation.

Table 3-2 Optical multimode designations

Multi-mode fiber	Fiber diameter (microns)	FC media designation
OM1	62.5 µm	M6
OM2	50 µm	M5
OM3	50 µm	M5E
OM4	50 µm	M5F

3.3.4 Dark fiber

To connect one optical device to another, some form of fiber-optic link is required. If the distance is short, then a standard fiber cable suffices. Over a slightly longer distance, for example from one building to the next, then a fiber link might need to be laid. This fiber might need to be laid underground or through a conduit, but it is not as simple as connecting two switches together in a single rack.

If the two units which need to be connected are in different cities, then the problem is much larger. Larger, in this case, is typically associated with more expensive. Because most businesses are not in the business of laying cable, they lease fiber-optic cables to meet their needs. When a company leases equipment, the fiber-optic cable that they lease is known as *dark fiber*.

Dark fiber generically refers to a long, dedicated fiber-optic link that can be used without the need for any additional equipment. It can be used while the particular technology supports the need.

Some forward-thinking services companies laid fiber-optic links alongside their pipes and cables. For example, a water company might be digging up a road to lay a mains pipe. Other examples include an electric company that might be taking a power cable across a mountain range by using pylons. Or, a cable TV company might be laying cable to all of the buildings in a city. While they carry out the work to support their core business, they might also lay fiber-optic links.

But these cables are just cables. They are not used in any way by the company who owns them. They remain dark until the user puts their own light down the fiber. Hence, the term *dark fiber*.

3.4 Classes of service

Applications might require different levels of service and guarantees regarding delivery, connectivity, and bandwidth. Some applications need to have bandwidth that is dedicated to them during the data exchange. An example of this type of application would be a tape backup. Other applications might be *bursty* in nature and not require a dedicated connection, but they might insist that an acknowledgement is sent for each successful transfer. The Fibre Channel standards provide different classes of service to accommodate different application needs. Table 3-3 provides brief details of the different classes of service.

Table 3-3 Fibre Channel classes of service

Class	Description	Requires acknowledge
1	Dedicated connection with full bandwidth.	Yes
2	Connectionless switch to switch communication for frame transfer and delivery.	Yes
3	Connectionless switch to switch communication for frame transfer and delivery.	No

Class	Description	Requires acknowledge
4	Dedicated connection with a fraction of bandwidth between ports by using virtual circuits.	Yes
6	Dedicated connection for multicast.	Yes
F	Switch to switch communication.	Yes

3.4.1 Class 1

In *class 1* service, a dedicated connection source and destination is established through the fabric during the transmission. It provides acknowledged service. This class of service ensures that the frames are received by the destination device in the same order in which they are sent. This class reserves full bandwidth for the connection between the two devices. It does not provide for a good utilization of the available bandwidth, since it is blocking another possible contender for the same device. Because of this blocking and the necessary dedicated connections, class 1 is rarely used.

3.4.2 Class 2

Class 2 is a connectionless, acknowledged service. Class 2 makes better use of available bandwidth since it allows the fabric to multiplex several messages on a frame-by-frame basis. As frames travel through the fabric they can take different routes, so class 2 service does not guarantee in-order delivery. Class 2 relies on upper layer protocols to take care of frame sequence. The use of acknowledgments reduces available bandwidth, which needs to be considered in large-scale busy networks.

3.4.3 Class 3

There is no dedicated connection in class 3 and the received frames are not acknowledged. *Class 3* is also called *datagram connectionless service*. It optimizes the use of fabric resources, but it is now up to the upper layer protocol to ensure that all frames are received in the correct order. The upper layer protocol also needs to request to the source device the retransmission of missing frames. Class 3 is a commonly used class of service in Fibre Channel networks.

3.4.4 Class 4

Class 4 is a connection-oriented service like class 1. The main difference is that class 4 allocates only a fraction of the available bandwidth of path through the fabric that connects two N_Ports. Virtual circuits (VCs) are established between two N_Ports with guaranteed quality of service (QoS), including bandwidth and latency. Like class 1, class 4 guarantees in-order delivery of frames and provides acknowledgment of delivered frames. However, now the fabric is responsible for multiplexing frames of different VCs. Class 4 service is intended for multimedia applications such as video and for applications that allocate an established bandwidth by department within the enterprise. Class 4 is included in the FC-PH-2 standard.

3.4.5 Class 5

Class 5 is called *isochronous service*, and is intended for applications that require immediate delivery of the data as it arrives, with no buffering. It is not clearly defined yet and is not included in the FC-PH documents.

3.4.6 Class 6

Class 6 is a variant of class 1, and is known as a *multicast class of service*. It provides dedicated connections for a reliable multicast. An N_Port might request a class 6 connection for one or more destinations. A multicast server in the fabric establishes the connections and gets acknowledgment from the destination ports, and sends it back to the originator. When a connection is established, it is retained and guaranteed by the fabric until the initiator ends the connection. Class 6 was designed for applications like audio and video that require multicast functionality. It is included in the FC-PH-3 standard.

3.4.7 Class F

Class F service is defined in the Fibre Channel Switched Fabric (FC-SW) and the FC-SW-2 standard for use by switches that communicate through inter-switch links (ISLs). It is a connectionless service with notification of non-delivery between E_Ports that are used for control, coordination, and configuration of the fabric. Class F is similar to class 2. The main difference is that class 2 deals with N_Ports that send data frames, while class F is used by E_Ports for control and management of the fabric.

3.5 Fibre Channel data movement

To move data bits with integrity over a physical medium, there must be a mechanism to check that this movement happened and that integrity is not compromised. This review is provided by a reference clock, which ensures that each bit is received as it was transmitted. In parallel topologies, this review can be accomplished by using a separate clock or strobe line. As data bits are transmitted in parallel from the source, the strobe line alternates between high or low to signal the receiving end that a full byte was sent. If there are 16-bit and 32-bit wide parallel cables, it would indicate that multiple bytes were sent.

The reflective differences in fiber-optic cabling mean that intermodal, or modal, dispersion (signal degradation) might occur.

This dispersion might result in frames that arrive at different times. This bit error rate (BER) is referred to as the *jitter budget*. No products are entirely jitter free. This jitter budget is an important consideration when you select the components of a SAN.

Because serial data transports have only two leads, transmit and receive, clocking is not possible by using a separate line. Serial data must carry the reference timing, which means that clocking is embedded in the bit stream.

Embedded clocking, though, can be accomplished by different means. Fibre Channel uses a byte-encoding scheme (covered in more detail in 3.5.1, “Byte encoding schemes” on page 49) and clock and data recovery (CDR) logic to recover the clock. From this recovery, it determines the data bits that comprise bytes and words.

Gigabit speeds mean that maintaining valid signaling, and ultimately valid data recovery, is essential for data integrity. Fibre Channel standards allow for a single bit error to occur only once in a million bits (1 in 10^{12}). In the real IT world, this rate equates to a maximum of one bit error every 16 minutes. However, actual occurrence is a lot less frequent.

3.5.1 Byte encoding schemes

To transfer data over a high-speed serial interface, the data is encoded before transmission and decoded upon reception. The encoding process ensures that sufficient clock information is present in the serial data stream. This information allows the receiver to synchronize to the embedded clock information and successfully recover the data at the required error rate. This 8b/10b encoding finds errors that a parity check cannot. A parity check does not find the even-numbered bit errors, only the odd numbers. The 8b/10b encoding logic finds almost all errors.

First developed by IBM, the 8b/10b encoding process converts each 8-bit byte into two possible 10-bit characters.

This scheme is called *8b/10b encoding* because it refers to the number of data bits input to the encoder and the number of bits output from the encoder.

The format of the 8b/10b character is of the format *Ann.m*:

- ▶ *A* represents D for data or K for a special character.
- ▶ *nn* is the decimal value of the lower 5 bits (EDCBA).
- ▶ “.” is a period.
- ▶ *m* is the decimal value of the upper 3 bits (HGF).

Figure 3-8 illustrates an encoding example.

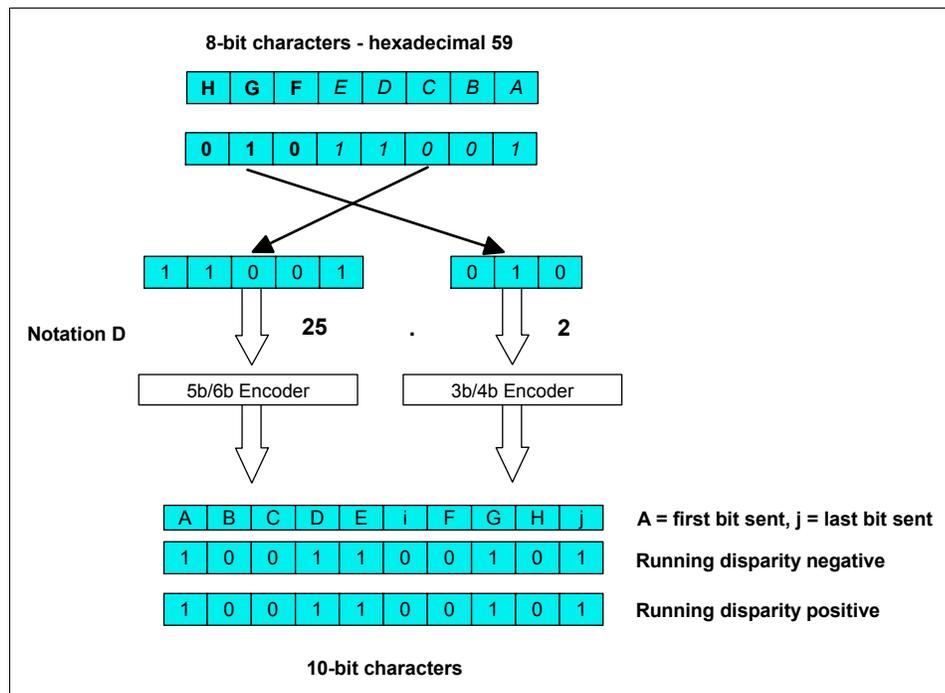


Figure 3-8 8b/10b encoding logic

The following steps occur in the encoding example, which is shown in Figure 3-8 on page 49:

1. Hexadecimal representation x'59' is converted to binary: 01011001.
2. The upper 3 bits are separated from the lower 5 bits: 010 11001.
3. The order is reversed and each group is converted to decimal: 25 2.
4. The letter notation D (for data) is assigned and becomes: D25.2.

Running disparity

As illustrated, the conversion of the 8-bit data bytes results in two 10-bit results. The encoder must choose one of these results to use. This decision is achieved by monitoring the running disparity of the previously processed character. For example, if the previous character had a positive disparity, then the next character that is issued might have an encoded value that represents negative disparity.

In the example that is used in Figure 3-8 on page 49, the encoded value, when the running disparity is either positive or negative, is the same. This outcome is legitimate. In some cases, the encoded value differs, and in others it is the same.

In Figure 3-8 on page 49, the encoded 10-bit byte has 5 bits that are set and 5 bits that are unset. The only possible results of the 8b/10b encoding are as follows:

- ▶ If 5 bits are set, then the byte is said to have neutral disparity.
- ▶ If 4 bits are set and 6 are unset, then the byte is said to have negative disparity.
- ▶ If 6 bits are set and four are unset, then the byte is said to have positive disparity.

The rules of Fibre Channel define that a byte that is sent cannot take the positive or negative disparity above one unit. Thus, if the current running disparity is negative, then the next byte that is sent must either have one of these properties:

- ▶ Neutral disparity
 - Keeping the current running disparity negative.
 - The subsequent byte needs to have either neutral or positive disparity.
- ▶ Positive disparity
 - Making the new current running disparity neutral.
 - The subsequent byte has either positive, negative, or neutral disparity.

Number of bits: At any time or at the end of any byte, the number of set bits and unset bits that pass over a Fibre Channel link, differ only by a maximum of two.

K28.5

In addition to the fact that many 8-bit numbers encode to *two* 10-bit numbers under the 8b/10b encoding scheme, there are some other key features.

Some 10-bit numbers cannot be generated from any 8-bit number. Thus, it is not possible to see these particular 10-bit numbers as part of a flow of data. This outcome is a useful fact because it means that these particular 10-bit numbers can be used by the protocol for signaling or control purposes.

These characters are referred to as *comma* characters, and rather than having the prefix D, they have the prefix K.

The only one character that gets used in Fibre Channel is the one that is known as *K28.5*, and it has a special property.

The two 10-bit encodings of K28.5 are shown in Table 3-4 on page 51.

Table 3-4 10-bit encoding of K28.5

Name of character	Encoding for current running disparity of	
	Negative	Positive
K28.5	001111 1010	110000 0101

All of the 10-bit bytes that are possible by using the 8b/10b encoding scheme have either 4, 5, or 6 bits that are set. The K28.5 character is special because it is the only character that is used in Fibre Channel that has five consecutive bits set or unset. All other characters have four or less consecutive bits of the same setting.

When you determine the significance of the bit settings, there are two things to consider:

- ▶ The first consideration is that the ones and zeros are actually representing light and dark on the fiber (assuming fiber optic medium). A 010 pattern would effectively be a light pulse between two periods of darkness. A 0110 would represent the same, except that the pulse of light would last for twice the length of time.

Because the two devices have their own clocking circuitry, the number of consecutive set bits, or consecutive unset bits, becomes important. For example, device 1 is sending to device 2 and the clock on device 2 is running 10% faster than that on device 1. If device 1 sent 20 clock cycles worth of set bits, device 2 would count 22 set bits. (This example is provided just to illustrate the point.) The worst possible case that you can have in Fibre Channel is five consecutive bits of the same setting within 1 byte: The K28.5.

- ▶ The other key consideration is that because K28.5 is the *only* character with five consecutive bits of the same setting, Fibre Channel hardware can look out for it specifically. Because K28.5 is used for control purposes, this setting is useful and allows the hardware to be designed for maximum efficiency.

64b/66b encoding

Communications of 10 and 16 Gbps use 64/66b encoding. Sixty-four bits of data are transmitted as a 66-bit entity. The 66-bit entity is made by prefixing one of two possible 2-bit *preambles* to the 64 bits to be transmitted. If the preamble is *01*, the 64 bits are entirely data.

If the preamble is *10*, an 8-bit type field follows, plus 56 bits of control information and data. The preambles *00* and *11* are not used, and generate an error if seen.

The use of the *01* and *10* preambles guarantees a bit transmission every 66 bits, which means that a continuous stream of zeros or ones cannot be valid data. It also allows easier clock and timer synchronization because a transmission must be seen every 66 bits.

The overhead of the 64B/66B encoding is considerably less than the more common 8b/10b encoding scheme.

3.6 Data transport

For Fibre Channel devices to be able to communicate with each other, there must be strict definitions regarding the way that data is sent and received. Because of this need, some data structures are defined. It is fundamental to understanding Fibre Channel that you have at least minimal knowledge of the way that data is moved around. You also need a basic understanding of the mechanisms that are used to accomplish this data movement.

3.6.1 Ordered set

Fibre Channel uses a command syntax, which is known as an *ordered set*, to move the data across the network. The ordered sets are 4-byte transmission words that contain data and special characters which have a special meaning. Ordered sets provide the availability to obtain bit and word synchronization, which also establishes word boundary alignment. An ordered set always begins with the special character K28.5. Three major types of ordered sets are defined by the signaling protocol.

The frame delimiters, the start-of-frame (SOF), and end-of-frame (EOF) ordered sets establish the boundaries of a frame. They immediately precede or follow the contents of a frame. There are 11 types of SOF and eight types of EOF delimiters that are defined for the fabric and N_Port sequence control.

The two primitive signals, idle and receiver ready (R_RDY), are ordered sets that are designated by the standard to have a special meaning. An *idle* is a primitive signal that is transmitted on the link to indicate that an operational port facility is ready for frame transmission and reception. The *R_RDY* primitive signal indicates that the interface buffer is available for receiving further frames.

A *primitive sequence* is an ordered set that is transmitted and repeated continuously to indicate specific conditions within a port. Or the set might indicate conditions that are encountered by the receiver logic of a port. When a primitive sequence is received and recognized, a corresponding primitive sequence or idle is transmitted in response. Recognition of a primitive sequence requires consecutive detection of three instances of the same ordered set. The primitive sequences that are supported by the standard include the following settings:

- ▶ Offline state (OLS)

The offline primitive sequence is transmitted by a port to indicate one of the following conditions: The port is beginning the link initialization protocol, the port received and recognized the NOS protocol, or the port is entering the offline status.

- ▶ Not operational (NOS)

The not operational primitive sequence is transmitted by a port in a point-to-point or fabric environment to indicate that the transmitting port detected a link failure. Or, the NOS might indicate an offline condition, waiting for the OLS sequence to be received.

- ▶ Link reset (LR)

The link reset primitive sequence is used to initiate a link reset.

- ▶ Link reset response (LRR)

Link reset response is transmitted by a port to indicate that it recognizes a link reset sequence and performed the appropriate link reset.

Data transfer

To send data over Fibre Channel, though, we need more than just the control mechanisms. Data is sent in frames. One or more related frames make up a sequence. One or more related sequences make up an exchange.

3.6.2 Frames

Fibre Channel places a restriction on the length of the data field of a frame, at 528 transmission words, which is 2112 bytes. See Table 3-5. Larger amounts of data must be transmitted in several frames. This larger unit that consists of multiple frames is called a *sequence*. An entire transaction between two ports is made up of sequences that are administered by an even larger unit that is called an *exchange*.

Frame arrival: Some classes of Fibre Channel communication guarantee that the frames arrive at the destination in the same order in which they were transmitted. Other classes do not. If the frames do arrive in the same order in which they were sent, it is considered to be an *in order* delivery of frames.

A frame consists of the following elements:

- ▶ SOF delimiter
- ▶ Frame header
- ▶ Optional headers and payload (data field)
- ▶ Cyclic redundancy check (CRC) field
- ▶ EOF delimiter

Figure 3-9 shows the layout of a Fibre Channel frame.

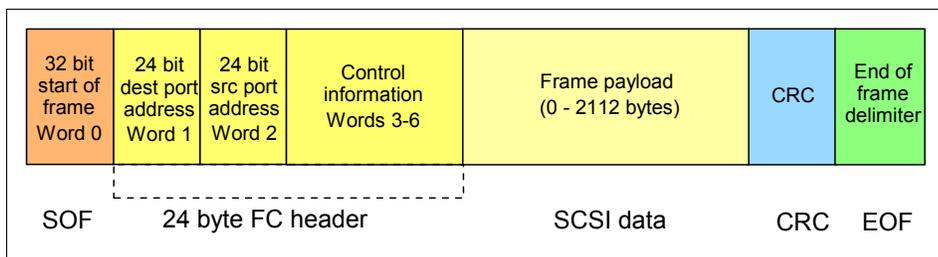


Figure 3-9 Fibre Channel frame structure

Framing rules

The following rules apply to the framing protocol:

- ▶ A frame is the smallest unit of information transfer.
- ▶ A sequence has at least one frame.
- ▶ An exchange has at least one sequence.

Transmission word

A *transmission word* is the smallest transmission unit that is defined in Fibre Channel. This unit consists of four transmission characters, 4 x 10 or 40 bits. When information that is transferred is not an even multiple of 4 bytes, the framing protocol adds fill bytes. The fill bytes are stripped at the destination.

Frames are the building blocks of Fibre Channel. A *frame* is a string of transmission words that are prefixed by a SOF delimiter and followed by an EOF delimiter. The way that transmission words make up a frame is shown in Table 3-5.

Table 3-5 Transmission words in a frame

SOF	Frame Header	Data Payload Transmission Words	CRC	EOF
1 TW	6 TW	0-528 TW	1 TW	1 TW

Frame header

Each frame includes a header that identifies the source and destination of the frame. The frame also includes control information that manages the frame and the sequences and exchanges that are associated with that frame. The structure of the frame header is shown in Table 3-6. The abbreviations are explained below the table.

Table 3-6 The frame header

	Byte 0	Byte 1	Byte 2	Byte 3
Word 0	R_CTL	Destination_ID (D_ID)		
Word 1	Reserved	Source_ID (S_ID)		
Word 2	Type	Frame Control (F_CTL)		
Word 3	SEQ_ID	DF_CTL	SequenceCount (SEQ_CNT)	
Word 4	Originator X_ID (OX_ID)		Responder X_ID (RX_ID)	
Word 5	Parameter			

Routing control (R_CTL)

This field identifies the type of information that is contained in the payload and where in the destination node it might be routed.

Destination ID

This field contains the address of the frame destination and is referred to as the D_ID.

Source ID

This field contains the address of where the frame is coming from and is referred to as the S_ID.

Type

The type field identifies the protocol of the frame content for data frames, such as SCSI, or a reason code for control frames.

F_CTL

This field contains control information that relates to the frame content.

SEQ_ID

The sequence ID is assigned by the sequence initiator and is unique for a specific D_ID and S_ID pair while the sequence is open.

DF_CTL

The data field control specifies whether there are optional headers that are present at the beginning of the data field.

SEQ_CNT

This count identifies the position of a frame within a sequence and is incremented by one for each subsequent frame that is transferred in the sequence.

OX_ID

This field identifies the exchange ID that is assigned by the originator.

RX_ID

This field identifies the exchange ID that is assigned to the responder.

Parameter

The parameter field specifies the relative offset for data frames, or information that is specific to link control frames.

3.6.3 Sequences

The information in a sequence moves in one direction, from a source N_Port to a destination N_Port. Various fields in the frame header are used to identify the beginning, middle, and end of a sequence. Other fields in the frame header are used to identify the order of frames, in case they arrive out of order at the destination.

3.6.4 Exchanges

Two other fields of the frame header identify the exchange ID. An exchange is responsible for managing a single operation that might span several sequences, possibly in opposite directions. The source and destination can have multiple exchanges active at a time.

Using SCSI as an example, a SCSI task is an exchange. The SCSI task is made up of one or more information units. The following information units (IUs) would be relevant for this SCSI task:

- ▶ Command IU
- ▶ Transfer ready IU
- ▶ Data IU
- ▶ Response IU

Each IU is one sequence of the exchange. Only one participant sends a sequence at a time. Figure 3-10 indicates the flow of the exchange, sequence, and frames.

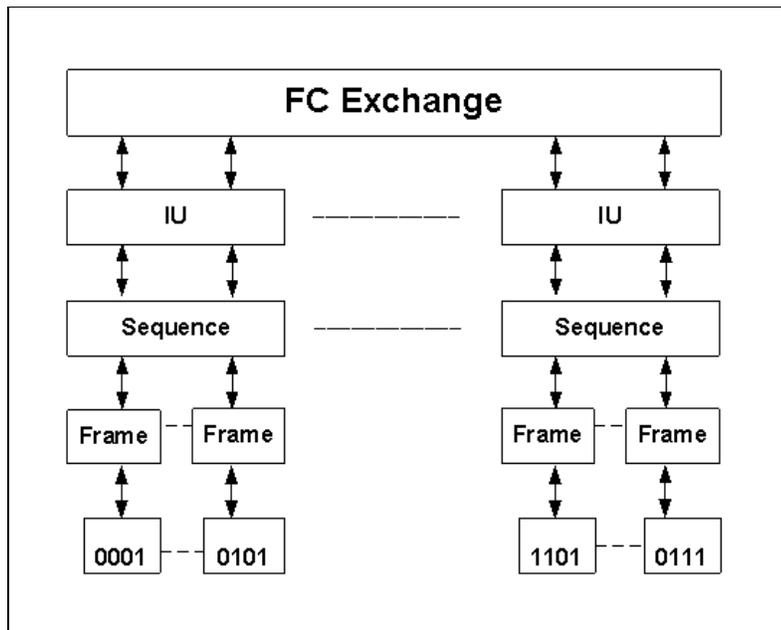


Figure 3-10 Fibre Channel (FC) exchange, sequence, and frame flow

3.6.5 In order and out of order

When data is transmitted over Fibre Channel, it is sent in frames. These frames carry a maximum of only 2112 bytes of data, often not enough to hold the entire set of information to be communicated. In this case, more than one frame is needed. Some classes of Fibre Channel communication guarantee that the frames arrive at the destination in the same order that they were transmitted. Other classes do not. If the frames do arrive in the same order that they were sent, this is considered to be an *in-order* delivery of frames.

In some cases, it is critical that the frames arrive in the correct order, and in others, it is not so important. In the latter case, which is considered *out of order*, the receiving port can reassemble the frames into the correct order before it passes the data out to the application. It is, however, common for switches and directors to guarantee in-order delivery, even if the particular class of communication allows for the frames to be delivered out of sequence.

3.6.6 Latency

The term *latency* means the delay between an action that is requested and the action taking place.

Latency occurs almost everywhere. A simple fact is that it takes time and energy to perform an action. The following areas highlight where you need to be particularly aware of the latency in a SAN:

- ▶ Ports
- ▶ Switches and directors
- ▶ Inter-Chassis Links in a DCX director
- ▶ Long-distance links
- ▶ Inter-Switch Links
- ▶ Application-specific integrated circuits (ASICs)

3.6.7 Open fiber control

When you deal with lasers, there are potential dangers to the eyes. Generally, the lasers in use in Fibre Channel are low-powered devices that are designed for quality of light and signaling rather than for maximum power. However, they can still be dangerous.

CAUTION: Never look into a laser light source. And never look into the end of a fiber optic cable unless you know exactly where the other end is. You also need to know that nobody can connect a light source to it.

To add a degree of safety, the concept of open fiber control (OFC) was developed. The following actions describe this concept:

- ▶ A device is turned on and it sends out low powered light.
- ▶ If it does not receive light back, then it assumes that there is no fiber connected. This feature is a fail-safe option.
- ▶ When it receives light, it assumes that there is a fiber that is connected and then switches the laser to full power.
- ▶ If one of the devices stops receiving light, then it reverts to the low-power mode.

When a device is transmitting at low power, it is not able to send data. The device is just waiting for a completed optical loop.

OFC ensures that the laser does not emit light which would exceed the class 1 laser limit when no fiber is connected. Non-OFC devices are guaranteed to be below class 1 limits at all times.

The key factor is that the devices at each end of a fiber link must either both be OFC or both be non-OFC.

All modern equipment uses non-OFC optics, but it is possible that some legacy (or existing) equipment might be using OFC optics.

3.7 Flow control

Now that you know that data is sent in frames, you also must understand that devices need to temporarily store the frames as they arrive. The data frames must be stored until they are assembled in sequence and then delivered to the upper layer protocol. The reason for this is that because of the potential high bandwidth of the Fibre Channel, it would be possible to inundate and overwhelm a target device with frames. There must be a mechanism to stop this from happening. The ability of a device to accept a frame is called its *credit*. This credit is usually referred to as the number of buffers (its buffer credit) that a node maintains for accepting incoming data.

3.7.1 Buffer to buffer

Buffer-to-buffer credits are the maximum number of frame transfers that a port can support. During login, N_Ports and F_Ports at both ends of a link establish its buffer-to-buffer credit (BB_Credit). Each port states the maximum BB_Credit that they can offer and the lower of the two is used.

3.7.2 End to end

At login, all N_Ports establish an end-to-end credit (EE_Credit) with each other. During data transmission, a port must not send more frames than the buffer of the receiving port can handle. It must first get an indication from the receiving port that it processed a previously sent frame.

3.7.3 Controlling the flow

Two counters are used to accomplish successful flow control: The BB_Credit_CNT and EE_Credit_CNT, and both are initialized to 0 during login. Each time that a port sends a frame, it increments the BB_Credit_CNT and EE_Credit_CNT by one. When it receives a receiver ready (R_RDY) indication from the adjacent port, it decrements the BB_Credit_CNT by one, and when it receives an acknowledgement (ACK) from the destination port, it decrements the EE_Credit_CNT by one. At certain times, the BB_Credit_CNT might become equal to the BB_Credit, or the EE_Credit_CNT might become equal to the EE_Credit of the receiving port. If this happens, the transmitting port must stop sending frames until the respective count is decremented.

The previous statements are true for class 2 service. Class 1 is a dedicated connection, so it does not need to care about the BB_Credit; only the EE_Credit is used (EE Flow Control). Class 3, however, is an unacknowledged service, so it uses only the BB_Credit (BB Flow Control), but the mechanism is the same on all cases.

3.7.4 Performance

Here, you can see the importance that the number of buffers has in overall performance. You need enough buffers to ensure that the transmitting port can continue sending frames without stopping in order to be able to use the full bandwidth. Having enough buffers is especially true with distance. At 1 Gbps, a frame occupies about 75 m - 4 km of fiber, which depends on the size of the data payload. In a 100 km link, many frames can be sent before the first one reaches its destination. You need an ACK back to start replenishing the EE_Credit or an R_RDY indication to replenish the BB_Credit.

For a moment, consider frames with 2 KB of data. These frames occupy approximately 4 km of fiber. You are able to send about 25 frames before the first frame arrives at the far end of the 100 km link. You are able to send another 25 frames before the first R_RDY or ACK indication is received. Therefore, you need at least 50 buffers to allow for non-stop transmission at a 100 km distance with frames of this size. If the frame size is reduced, more buffers would be required to allow non-stop transmission. In brief, the buffer credit management is critical in case of long-distance communication, hence the appropriate buffer credit allocation is important to obtain optimal performance. Inappropriate allocation of the buffer credit might result in a delay of transmission over the Fibre Channel link. As a preferred practice, always refer to the default buffer and maximum buffer credit values for each model of switch from each vendor.



Ethernet and system networking concepts

In this chapter, we introduce you to Ethernet and system networking concepts. We also describe the storage area network (SAN) Internet Protocol (IP) networking options and how we arrive at converged networks.

4.1 Ethernet

In Chapter 1, “Introduction” on page 1, we gave a brief introduction to what a network is and the importance of the models. The Ethernet standard fits into layer 2 of the open systems interconnection (OSI) model. The standard refers to the media access layer that devices are connected to (the cable) and compete for access by using the Carrier Sense Multiple Access with Collision Detection (CSMA/CD) protocol.

Ethernet is a standard communications protocol that is embedded in software and hardware devices, intended for building a local area network (LAN). Ethernet was designed by Bob Metcalfe in 1973, and through the efforts of Digital, Intel, and Xerox (for whom Metcalfe worked), *DIX Ethernet* became the standard model for LANs worldwide.

The formal designation for standardization of the Ethernet protocol is sometimes referred to as *IEEE 802.3*. The Institute of Electrical and Electronics Engineers (IEEE) proposed a working group in February 1980 to standardize network protocols. The third subcommittee worked on a flavor essentially identical to Ethernet, though there are insignificant variances. Consequently, generic use of the term Ethernet might refer to IEEE 802.3 or DIX Ethernet.

Ethernet was originally based on the idea of computers that communicate over a shared coaxial cable that acts as a broadcast transmission medium. The methods used were similar to those used in radio systems, with the common cable providing the communication channel likened to the luminiferous aether (light-bearing aether) in 19th century physics, and it was from this reference that the name *Ethernet* was derived.

4.1.1 Shared media

Since all communications happen on the same wire, any information sent by one computer is received by all, even if that information is intended for just one destination. The network interface card (NIC) interrupts the CPU only when applicable packets are received. The card ignores information that is not addressed to it. Use of a single cable also means that the bandwidth is shared, so that network traffic can be very slow when many stations are simultaneously active.

Collisions reduce throughput by their very nature. In the worst case, when there are numerous hosts with long cables that attempt to transmit many short frames, excessive collisions can reduce throughput dramatically.

Ethernet networks are composed of broadcast domains and there is no clock signal on the wire, as serial connections often have. Instead, Ethernet systems must determine if the wire is in use, and if not, the system must send enough data to enable the remote station to allow it to synchronize properly. This synchronization mechanism that is combined with the ability to detect other computers that are attempting to access the wire, is a formalized protocol that is called *Carrier Sense Multiple Access with Collision Detection (CSMA/CD)*.

4.1.2 Ethernet frame

Figure 4-1 shows an Ethernet frame.

IEEE 803.2 / 802.2								
7 bytes	1 byte	2 or 6 bytes	2 or 6 bytes	2 bytes	4-1500 bytes		4 bytes	
Preamble	Start Frame Delimiter	Dest. MAC address	Source MAC address	Length	(Data / Pad)			FCS
					DSAP	SSAP	CTRL	

Figure 4-1 Ethernet frame

The Ethernet frame that is shown in Figure 4-1, contains the following components.

Preamble

A *preamble* is a stream of bits that are used to allow the transmitter and receiver to synchronize their communication. The preamble is an alternating pattern of 56 binary ones and zeros. The preamble is immediately followed by the Start Frame Delimiter.

Start Frame Delimiter

The *Start Frame Delimiter* is always 10101011 and is used to indicate the beginning of the frame information.

Destination Media Access Control

The *destination Media Access Control (MAC)* is the address of the system that is receiving data. When a NIC is listening to the wire, it is checking this field for its own MAC address.

Source Media Access Control

The *source Media Access Control (MAC)* is the MAC address of the system that is transmitting data.

Length

This is the *length* of the entire Ethernet frame, in bytes. Although this field can hold any value 0 - 65,534, it is rarely larger than 1500. This smaller value is because it is usually the maximum transmission frame size for most serial connections. Ethernet networks tend to use serial devices to access the Internet.

Data/pad (also known as payload)

The data is inserted in the *data/pad* (also known as *payload*). This location is where the IP header and data are placed if you are running IP over Ethernet. This field contains Internetwork Packet Exchange (IPX) information if you are running IPX/Sequenced Packet Exchange (SPX) protocol (Novell). Contained within the data/pad section of an IEEE 803.2 frame are the following specific fields:

- ▶ DSAP: destination service access point
- ▶ SSAP: source service access point
- ▶ CTRL: control bits for Ethernet communication
- ▶ NLI: network layer interface
- ▶ FCS: frame check sequence

Frame check sequence

This field contains the frame check sequence (FCS) which is calculated by using a cyclic redundancy check (CRC). The FCS allows Ethernet to detect errors in the Ethernet frame and reject the frame if it appears damaged.

4.1.3 How Ethernet works

When a device that is connected to an Ethernet network wants to send data, it first checks to ensure that it has a carrier on which to send its data (usually a piece of copper cable that is connected to a hub or another machine). This step is known as *Carrier Sense*.

All machines on the network are free to use the network whenever they choose if no one else is transmitting. This setup is known as *Multiple Access*.

There also needs to be a means of ensuring that when two machines start to transmit data simultaneously, that the resultant corrupted data is discarded. Also, retransmissions must be generated at differing time intervals. This assurance is known as *Collision Detection*.

Figure 4-2 shows a bus Ethernet network.

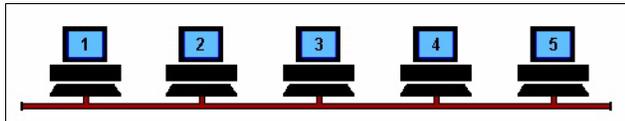


Figure 4-2 Bus Ethernet network

In Figure 4-2, assume that machine 2 wants to send a message to machine 4, but first it “listens” to make sure that no one else is using the network.

If the path is all clear, machine 2 starts to transmit its data on to the network. Each packet of data contains the destination address, the sender address, and the data to be transmitted.

The signal moves down the cable and is received by every machine on the network. But, because the signal it is only addressed to machine 4, the other machines ignore it. Machine 4 then sends a message back to machine 2 to acknowledge receipt of the data.

But what happens when two machines try to transmit at the same time? A collision occurs, and each machine has to “back off” for a random period before trying to transmit again. Figure 4-3 illustrates what happens when two machines are transmitting simultaneously.

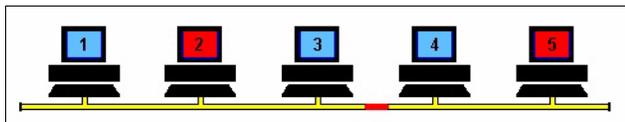


Figure 4-3 Machine 2 and machine 5 are both trying to transmit simultaneously

The resulting collision in Figure 4-3 destroys both signals and each machine knows this has happened because they do not “hear” their own transmission within a given period. This time period is the propagation delay and is equivalent to the time that it takes for a signal to travel to the furthest part of the network and back again.

Both of the machines then wait for a random period before trying to transmit again. On small networks, this all happens so quickly that it is virtually unnoticeable. However, as more machines are added to a network, the number of collisions rise dramatically and eventually

result in slow network response. The exact number of machines that a single Ethernet segment can handle depends upon the applications that are being used, but it is generally considered that 40 - 70 users are the limit before network speed is compromised.

Figure 4-4 shows two different scenarios: hub and switch. The *hub* is where all the machines are interconnected so that only one machine at a time can use the media. In the *switch* network, more than one machine can be using the media at the time.

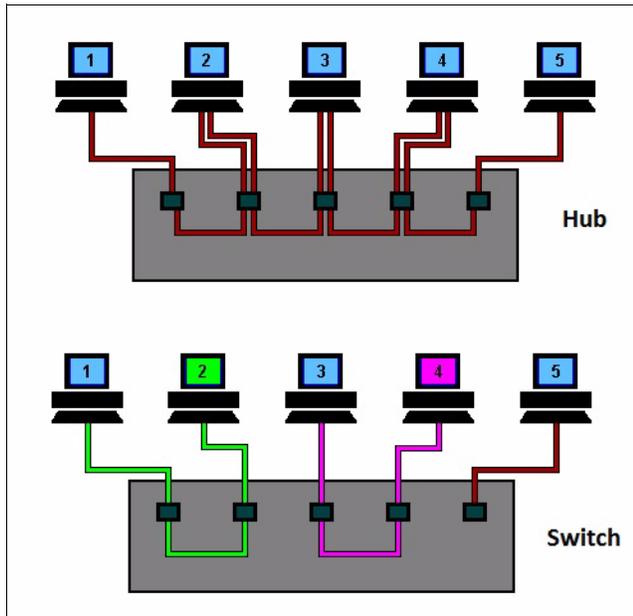


Figure 4-4 Hub and switch scenarios

An Ethernet hub changes the topology from a “bus” to a “star wired bus”. As an example, assume again that machine 1 is transmitting data to machine 4. But this time, the signal travels in and out of the hub to each of the other machines.

As you can see, it is still possible for collisions to occur but hubs have the advantage of centralized wiring, and they can automatically bypass any ports that are disconnected or have a cabling fault. This makes the network much more fault tolerant than a coax-based system where disconnecting a single connection brings the whole network down.

With a switch, machines can transmit simultaneously. Each switch reads the destination addresses and *switches* the signals directly to the recipients without broadcasting to all of the machines on the network.

This *point-to-point switching* alleviates the problems that are associated with collisions and considerably improves the network speed.

4.1.4 Speed and bandwidth

By convention, network data rates are denoted either in bits (bits per second or bps) or bytes (bytes per second or Bps). In general, parallel interfaces are quoted in bytes and serial in bits.

The numbers below are simplex data rates, which might conflict with the duplex rates that vendors sometimes use in promotional materials. Where two values are listed, the first value is the downstream rate and the second value is the upstream rate.

All quoted figures are in metric decimal units:

- 1 Byte = 8 bits
- 1 Kbps = 1,000 bits per second
- 1 Mbps = 1,000,000 bits per second
- 1 Gbps = 1,000,000,000 bits per second
- 1 KBps = 1,000 bytes per second
- 1 MBps = 1,000,000 bytes per second
- 1 GBps = 1,000,000,000 bytes per second
- 1 TBps = 1,000,000,000,000 bytes per second

These figures go against the traditional use of binary prefixes for memory size. These decimal prefixes are established in data communications.

Table 4-1 shows the technology rates and the medium.

Table 4-1 Technology rates and medium

Technology	Rate (Bit/s)	Rate (Byte/s)	Media
Fast Ethernet (100BASE-X)	100 Mb/s	12.5 MB/s	UTP Cat 5
Gigabit Ethernet (1000BASE-X)	1000 Mb/s	125 MB/s	UTP Cat 5e / 6
10 Gigabit Ethernet (10GBASE-X)	10000 Mb/s	1250 MB/s	UTP Cat 7 - fiber

4.1.5 10 GbE

From its origin more than 25 years ago, Ethernet has evolved to meet the increasing demands of packet-based networks. Ethernet provides the benefits of proven low implementation cost, reliability, and relative simplicity of installation and maintenance. Because of these benefits, the popularity of Ethernet has grown to the point that nearly all of the traffic on the Internet originates or terminates with an Ethernet connection. Furthermore, as the demand for ever-faster network speeds has increased, Ethernet has adapted to handle these higher speeds and the surges in volume demand that accompany them.

The IEEE 802.3ae 2002 (the 10 Gigabit Ethernet (10 GbE) standard) is different in some respects from earlier Ethernet standards in that it operates only in full-duplex mode (collision-detection protocols are unnecessary).

Ethernet can now progress to 10 gigabits per second while retaining its critical Ethernet properties, such as the packet format. The current capabilities are easily transferable to the new standard.

The 10 Gigabit Ethernet technology continues the evolution of Ethernet in terms of speed and distance, while it retains the same Ethernet architecture used in other Ethernet specifications. However, there is one key exception. Since 10 Gigabit Ethernet is a full-duplex only technology, it does not need the CSMA/CD protocol that is used in other Ethernet technologies. In every other respect, 10 Gigabit Ethernet matches the original Ethernet model.

4.1.6 10 GbE copper versus fiber

Once the decision is made to implement 10 Gigabit Ethernet (10 GbE) functionality, organizations must consider the data carrying techniques that facilitate such bandwidth. Copper and fiber cabling are the preeminent technologies for data transmission and provide their own unique benefits and drawbacks.

Copper is the default standard for transmitting data between devices because of its low cost, easy installation, and flexibility. It also possesses distinct shortcomings. Copper is best when used in short lengths, typically 100 meters or less. When employed over long distances, electromagnetic signal characteristics hinder performance. In addition, bundling copper cabling can cause interference, making it difficult to employ as a comprehensive backbone. For these reasons, copper cabling is the principal data carrying technique for communication among PCs and LANs, but not campus or long-distance transmission.

Conversely, fiber cabling is typically used for remote campus connectivity, crowded telecommunications closets, long-distance communications, and environments that need protection from interference. An example of such an environment is a manufacturing area. Since it is reliable and less susceptible to attenuation, it is optimum for sending data beyond 100 meters.

However, fiber is also more costly than copper and its use is typically limited to those applications that demand it.

As a result, most organizations use a combination of copper and fiber cabling. As these companies transition to 10 GbE functionality, they must have a solid understanding of the various cabling technologies. Companies must also have a sound migration strategy to ensure that their cabling infrastructure will support their network infrastructure both today and tomorrow.

The IEEE 802.3 Higher Speed Study Group formed in 1998, and the development of *10GigE* began the following year. By 2002, the 10GigE standard was first published as *IEEE Std 802.3ae-2002*. This standard defines a normal data rate of 10 Gigabits, making it 10 times faster than the Gigabit Ethernet.

Subsequent standard updates ensued in relation to the first 10GigE version published in 2002. The IEEE 802.3ae-2002 fiber, followed by 802.3ak-2004 in 2004, were later consolidated into *IEEE 802.3-2005* in the year 2005. In 2006, 802.3an-2006, which is a 10 Gigabit Base-T copper twisted pair, and an enhanced version with fiber-LRM PMD followed, known as *802.3aq-2006*. Finally, in 2007, the *802.3ap-2007* with copper backplane evolved.

As a result of these standards, there are two main types of 10 Gigabit Ethernet cabling: fiber and copper.

The following standards pertain to the 10 Gigabit Ethernet fiber cabling:

- ▶ 10GBASE-LX4: It supports ranges of 240 - 300 meters (790 - 980 ft) over traditional multi-mode cabling. This range is achieved by using four separate laser sources that operate at 3.125 Gbit/s in the range of 1300 nm on unique wavelengths. The 10GBASE-LX4 standard also supports 10 kilometers (6.2 mi) over System Management Facilities (SMF).
- ▶ 10GBASE-SR: Over obsolete 62.5 micron multi-mode fiber cabling (OM1), it has a maximum range of 26 - 82 meters (85 - 269 ft), depending on the cable type. Over standard 50 μ m 2000 MHz-km OM3 multi-mode fiber (MMF), it has a maximum range of 300 meters (980 ft).

- ▶ 10GBASE-LR: Has a specified reach of 10 kilometers (6.2 mi), but 10GBASE-LR optical modules can often manage distances of up to 25 kilometers (16 mi) with no data loss.
- ▶ 10GBASE-LRM: Supports distances up to 220 meters (720 ft) on FDDI-grade 62.5 μ m MMF. This fiber was originally installed in the early 1990s for Fiber Distributed Data Interface (FDDI), 100BaseFX networks, and for 260 meters (850 ft) on OM3. The reach of 10GBASE-LRM is not as far as the older 10GBASE-LX4 standard.
- ▶ 10GBASE-ER: This extended range has a reach of 40 kilometers (25 mi).
- ▶ 10GBASE-ZR: Several manufacturers introduced 80 km (50 mi) range ER pluggable interfaces under the name *10GBASE-ZR*. This 80 km PHY is not specified within the IEEE 802.3ae standard. Manufacturers created their own specifications that are based upon the 80 km PHY described in the OC-192/STM-64 SDH/SONET specifications.

A 10G Ethernet connection can also run over twin-ax cabling and twisted pair cabling. The following standards pertain to the 10 Gigabit Ethernet copper cabling:

- ▶ 10GBASE-CX4: This was the first 10G copper standard that was published by 802.3 (as *802.3ak-2004*). It is specified to work up to a distance of 15 m (49 ft). Each lane carries 3.125 gigabaud (Gbaud) of signaling bandwidth.
- ▶ 10GBASE-T or IEEE 802.3an-2006: This standard was released in 2006 to provide 10 Gbit/s connections over unshielded or shielded twisted pair cables, over distances up to 100 meters (330 ft).

Cables needed to carry 10GBASE-T: Category 6A, or better, of balanced twisted-pair cables that are specified in ISO 11801 amendment 2 or ANSI/TIA-568-C.2, are needed to carry 10GBASE-T up to distances of 100 m. Category 6 cables can carry 10GBASE-T for shorter distances when qualified, according to the guidelines in ISO TR 24750 or TIA-155-A.

The following standards pertain to the 10 Gigabit Ethernet copper backplane cabling:

- ▶ 10GBASE-X
- ▶ 10GBASE-KX4
- ▶ 10GBASE-KR

4.1.7 Virtual local area network

A virtual local area network (VLAN) is a networking concept in which a network is logically divided into smaller virtual LANs. The Layer 2 traffic in one VLAN is logically isolated from other VLANs, as illustrated in Figure 4-5.

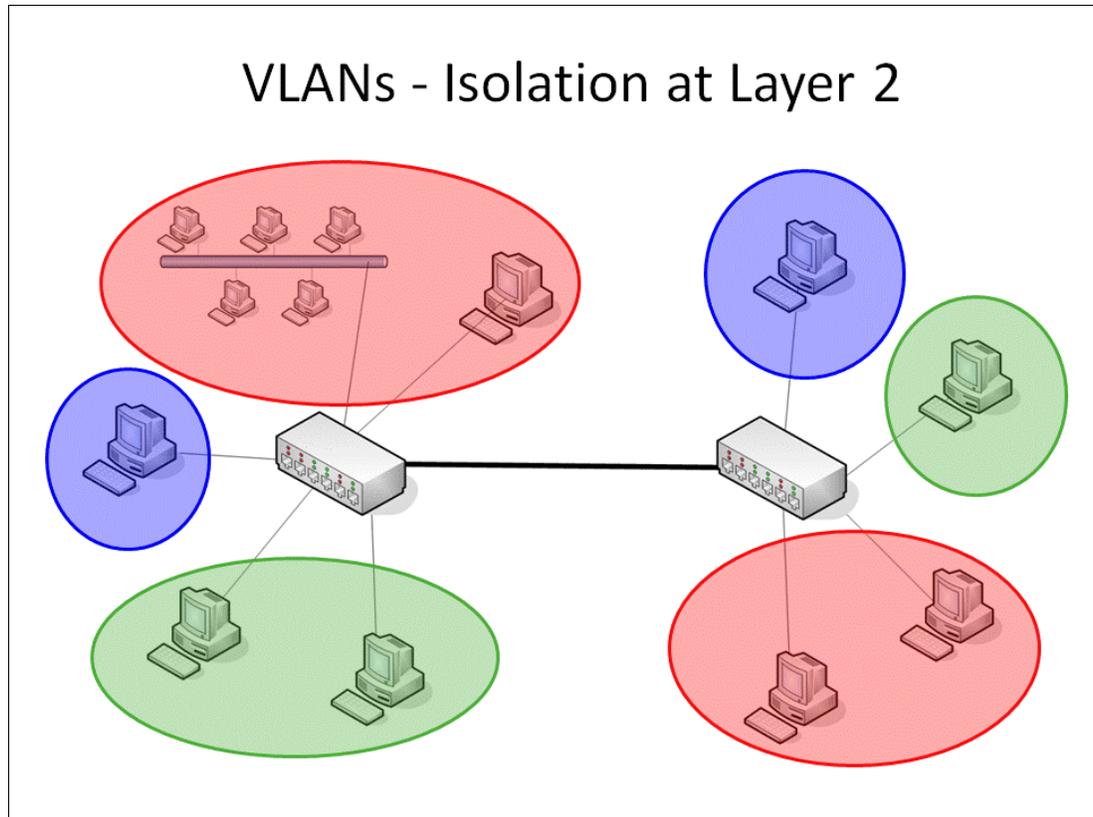


Figure 4-5 Isolation at Layer 2

Figure 4-6 shows two methods for maintaining isolation of VLAN traffic between switches.

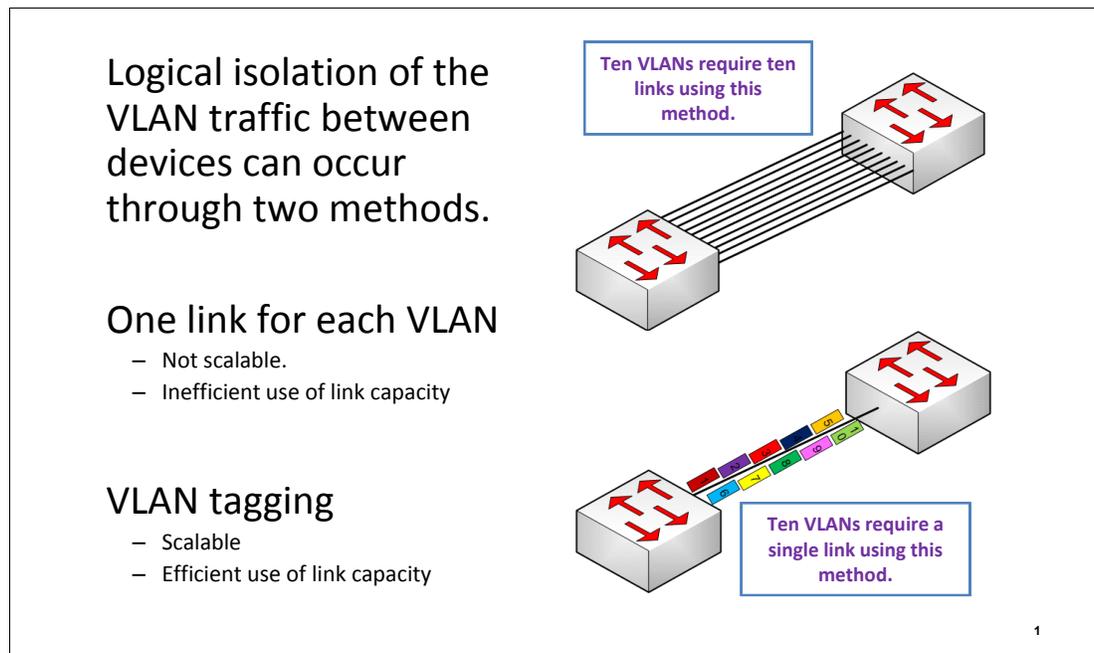


Figure 4-6 VLAN tagging

The first method uses a single link for each VLAN. This method does not scale well because it uses many ports in networks that have multiple VLANs and multiple switches. Also, this method does not use link capacity efficiently when traffic in the VLANs is not uniform.

The second method is VLAN tagging over a single link in which each frame is tagged with its VLAN ID. This method is highly scalable because only a single link is required to provide connectivity to many VLANs. This configuration provides for better utilization of the link capacity when VLAN traffic is not uniform.

The protocol for VLAN tagging of frames in a LAN environment is defined by the *IEEE 802.1p/q* standard (priority tagging and VLAN identifier tagging).

Inter-switch link (ISL): ISL is another protocol for providing the VLAN tagging function in a network. This protocol is not compatible with the IEEE 802.1p/q standard.

Tagged frames

The IEEE 802.1p/q standard provides a methodology for added information such as VLAN membership and priority to the frame, as shown in Figure 4-7.

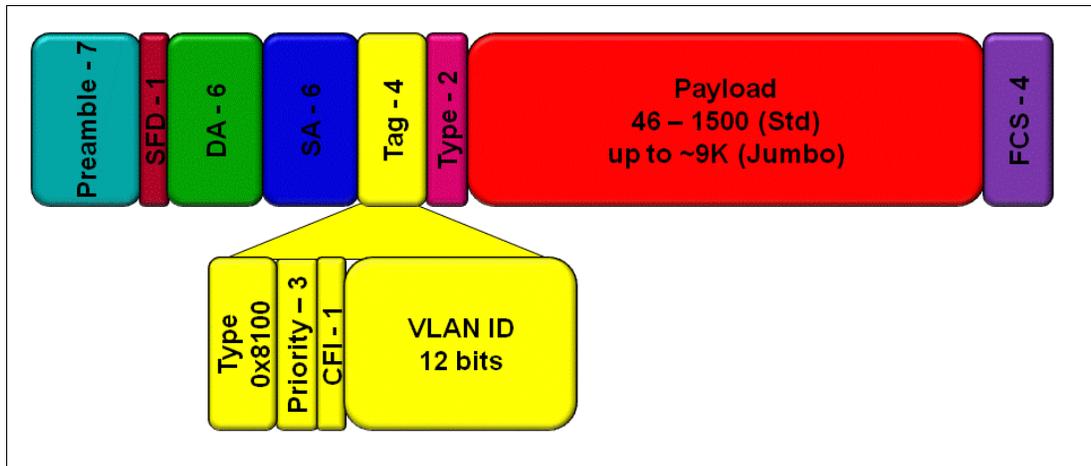


Figure 4-7 IEEE 802.1p/q tagged Ethernet frame

The standard provides an additional 4 bytes of information to be added to each Ethernet frame. A frame that includes this extra information is known as a *tagged frame*.

The 4-byte tag has four component fields:

- ▶ The *type field* is 2 bytes long and has the hexadecimal value of x8100 to identify the frame as an 802.1p/q tagged frame.
- ▶ The *priority field* is 3 bits long and allows a priority value of eight different values to be included in the tag. This field has the “P” portion of the 802.1p/q standard.
- ▶ The *Canonical Format Indicator field* is 1 bit long and identifies when the contents of the payload field are in canonical format.
- ▶ The *VLAN ID field* is 12 bits long and identifies which VLAN that the frame is a member of, with 4096 different VLANs possible.

4.1.8 Interface virtual local area network operation modes

Interfaces on a switch can operate in two virtual local area network (VLAN) modes: Single VLAN mode or multiple VLAN mode.

Single virtual local area network mode

The *single VLAN mode* operation is also referred to as *access mode*. A port that is operating in this mode is associated with a single VLAN. Incoming traffic does not have any VLAN identification. When the untagged frames enter the port, the VLAN identification for the VLAN that is configured for the port is added to the inbound frames.

Switch ports: Some vendors use terms other than *access mode* for ports that are operating in the single VLAN mode. The *switch ports* of those vendors might be configured to operate in the single VLAN mode. This step can be done by configuring a Port VLAN ID (PVID) and adding the port as a member of the VLAN.

Multiple virtual local area network mode

The *multiple VLAN mode* operation is also referred to as *trunk mode*. A port that is operating in this mode can receive frames that have VLAN tags. The port is also configured with VLANs to which the port is allowed to send and receive frames.

With the *IEEE 802.1Q* specification, untagged traffic on a multi-VLAN port can be associated with a single VLAN, which is referred to as the *native VLAN* for the port (Figure 4-8). By using this provision, traffic with no VLAN tag can be received and associated with the VLAN that is configured as the PVID or native VLAN. Outbound traffic for this VLAN on a port that is configured in this manner is transmitted with no tag. This method allows the receiving device to receive the frame in an untagged format.

This method provides compatibility with existing devices. Compatibility is also provided for devices that are configured in the single VLAN mode and are attached to a port that is configured as a multi-VLAN port.

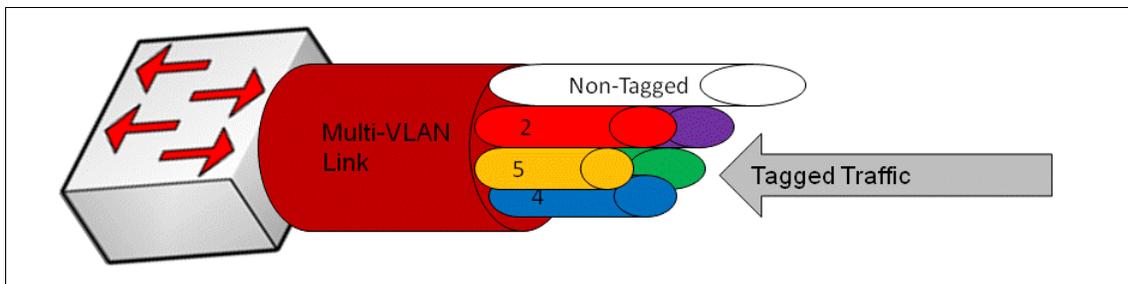


Figure 4-8 Multiple VLAN mode link

Variations in the meaning of trunk: The term *trunk* is used to express different ideas in the networking industry. When you use this term, keep in mind that others might use the term in a different manner. *Trunk* can mean that a port is operating in multiple VLAN mode or it can mean a link aggregated port.

4.1.9 Link aggregation

Link aggregation combines multiple physical links to operate as a single larger logical link. The *member links* no longer function as independent physical connections, but as members of the larger logical link (Figure 4-9).

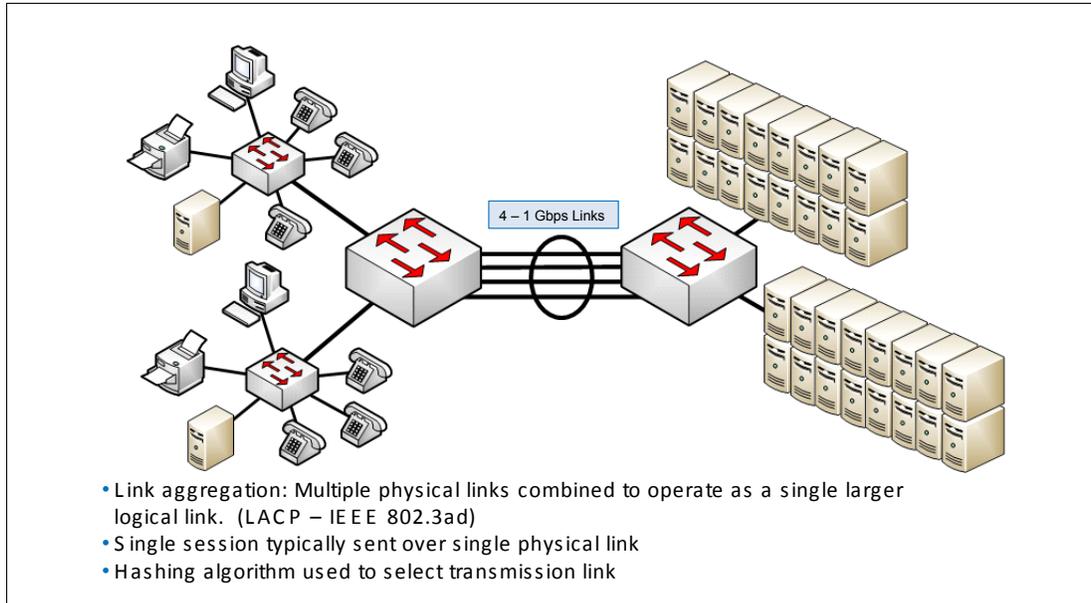


Figure 4-9 Link aggregation

Link aggregation provides greater bandwidth between the devices at each end of the aggregated link. Another advantage of link aggregation is increased availability because the aggregated link is composed of multiple member links. If one member link fails, the aggregated link continues to carry traffic over the remaining member links.

Each of the devices that is interconnected by the aggregated link uses a hashing algorithm to determine on which of the member links that frames will be transmitted. The hashing algorithm might use varying information in the frame to make the decision. This algorithm might include a source MAC, destination MAC, source IP, destination IP, and more. It might also include a combination of these values.

4.1.10 Spanning Tree Protocol

Spanning Tree Protocol (STP) provides Layer 2 loop prevention and is commonly in different forms, such as existing STP, Rapid STP (RSTP), Multiple STP (MSTP), and VLAN STP (VSTP). RSTP is a common default STP. This form provides faster convergence times than STP. However, some existing networks require the slower convergence times that basic STP provides.

The operation of Spanning Tree Protocol

STP uses Bridge Protocol Data Unit (BPDU) packets to exchange information with other switches. BPDUs send out hello packets at regular intervals to exchange information across bridges and detect loops in a network topology.

Two types of BPDUs are available:

- ▶ **Configuration BPDUs**

These BPDUs contain configuration information about the transmitting switch and its ports, including switch and port MAC addresses, switch priority, port priority, and port cost.

- ▶ **Topology Change Notification (TCN) BPDUs**

When a bridge must signal a topology change, it starts to send TCNs on its root port. The designated bridge receives the TCN, acknowledges it, and generates another one for its own root port. The process continues until the TCN reaches the root bridge.

STP uses the information that is provided by the BPDUs in several ways: To elect a root bridge, identify root ports for each switch, identify designated ports for each physical LAN segment, and prune specific redundant links to create a loop-free tree topology. All leaf devices calculate the best path to the root device. The devices place their ports in blocking or forwarding states that are based on the best path to the root. The resulting tree topology provides a single active Layer 2 data path between any two end stations.

Rapid Spanning Tree Protocol

RSTP provides better reconvergence time than the original STP. RSTP identifies certain links as *point to point*. When a point-to-point link fails, the alternate link can make the transition to the forwarding state.

An RSTP domain has the following components:

Root port	The “best path” to the root device.
Designated port	Indicates that the switch is the designated bridge for the other switch that is connecting to this port.
Alternate port	Provides an alternate root port.
Backup port	Provides an alternate designated port.

RSTP was originally defined in the IEEE 802.1w draft specification and later incorporated into the IEEE 802.1D-2004 specification.

Multiple Spanning Tree Protocol

Although RSTP provides faster convergence time than STP, it still does not solve a problem inherent in STP. This inherent issue is that all VLANs within a LAN must share the same spanning tree. To solve this problem, we use MSTP to create a loop-free topology in networks with multiple spanning-tree regions.

In an MSTP region, a group of bridges can be modeled as a single bridge. An MSTP region contains multiple spanning-tree instances (MSTIs). MSTIs provide different paths for different VLANs. This functionality facilitates better load sharing across redundant links.

An MSTP region can support up to 64 MSTIs, and each instance can support 1 - 4094 VLANs.

MSTP was originally defined in the IEEE 802.1s draft specification and later incorporated into the IEEE 802.1Q-2003 specification.

VLAN Spanning Tree Protocol

With VSTP, switches can run one or more STP or RSTP instances for each VLAN on which VSTP is enabled. For networks with multiple VLANs, VSTP enables more intelligent tree spanning. This level of tree spanning is possible because each VLAN can have interfaces that are enabled or disabled depending on the paths that are available to that specific VLAN.

By default, VSTP runs RSTP, but you cannot have both stand-alone RSTP and VSTP running simultaneously on a switch. VSTP can be enabled for up to 253 VLANs.

Bridge Protocol Data Unit (BPDU) protection

BPDU protection can help prevent STP misconfigurations that can lead to network outages. Receipt of BPDUs on certain interfaces in an STP, RSTP, VSTP, or MSTP topology, can lead to network outages.

BPDU protection is enabled on switch interfaces that are connected to user devices or on interfaces on which no BPDUs are expected, such as edge ports. If BPDUs are received on a protected interface, the interface is disabled and stops forwarding the frames.

Loop protection

Loop protection increases the efficiency of STP, RSTP, VSTP, and MSTP by preventing ports from moving into a forwarding state that might result in a loop opening in the network.

A blocking interface can transition to a forwarding state in error if the interface stops receiving BPDUs from its designated port on the segment. Such a transition error can occur when there is a hardware error on the switch or software configuration error between the switch and its neighbor.

When loop protection is enabled, the spanning tree topology detects root ports and blocked ports and ensures that both keep receiving BPDUs. If a loop protection-enabled interface stops receiving BPDUs from its designated port, it reacts as it might react to a problem with the physical connection on this interface. It does not transition the interface to a forwarding state, but instead transitions it to a loop-inconsistent state. The interface recovers and then transitions back to the spanning-tree blocking state as soon as it receives a BPDU.

You must enable loop protection on all switch interfaces that have a chance of becoming root or designated ports. Loop protection is the most effective when it is enabled in the entire switched network. When you enable loop protection, you must configure at least one action (**alarm**, **block**, or both).

An interface can be configured for either loop protection or root protection, but not for both.

Root protection

Root protection increases the stability and security of STP, RSTP, VSTP, and MSTP by limiting the ports that can be elected as root ports. A root port that is elected through the regular process has the possibility of being wrongly elected. A user bridge application that is running on a PC can also generate BPDUs and interfere with root port election. With root protection, network administrators can manually enforce the root bridge placement in the network.

Root protection is enabled on interfaces that should not receive superior BPDUs from the root bridge and should not be elected as the root port. These interfaces become designated ports and are typically on an administrative boundary. If the bridge receives superior STP BPDUs on a port that has root protection enabled, that port transitions to a root-prevented STP state (inconsistency state), and the interface is blocked. This blocking prevents a bridge that should not be the root bridge from being elected the root bridge. After the bridge stops receiving superior STP BPDUs on the interface with root protection, the interface returns to a listening state. This state is followed by a learning state and ultimately back to a forwarding state. Recovery back to the forwarding state is automatic.

When root protection is enabled on an interface, it is enabled for all of the STP instances on that interface. The interface is blocked only for instances for which it receives superior

BPDUs. Otherwise, it participates in the spanning tree topology. An interface can be configured for either root protection or loop protection, but not for both.

4.1.11 Link Layer Discovery Protocol

Link Layer Discovery Protocol (LLDP) is a vendor-independent protocol for network devices to advertise information about their identity and capabilities. It is referred to as *Station and Media Access Control Connectivity Discovery*, which is specified in the 802.1ab standard. With LLDP and Link Layer Discovery Protocol–Media Endpoint Discovery (LLDP-MED), network devices can learn and distribute device information about network links. With this information, the switch can quickly identify various devices, resulting in a LAN that interoperates smoothly and efficiently.

LLDP-capable devices transmit information in Type Length Value (TLV) messages to neighbor devices. Device information can include specifics such as chassis and port identification, and system name and system capabilities.

LLDP-MED goes one step further, exchanging IP-telephony messages between the switch and the IP telephone. These TLV messages provide detailed information about the Power over Ethernet (PoE) policy. With the PoE Management TLVs, the switch ports can advertise the power level and power priority that is needed. For example, the switch can compare the power that is needed by an IP telephone that is running on a PoE interface with available resources. If the switch cannot meet the resources that are required by the IP telephone, the switch can negotiate with the telephone until a compromise on power is reached.

The switch also uses these protocols to ensure that voice traffic gets tagged and prioritized with the correct values at the source itself. For example, the 802.1p class of service (COS) and 802.1Q tag information can be sent to the IP telephone.

4.1.12 Link Layer Discovery Protocol Type Length Values (LLDP TLVs)

The basic TLVs include the following information:

Chassis identifier	The MAC address that is associated with the local system.
Port identifier	The port identification for the specified port in the local system.
Port description	The user-configured port description. This description can be a maximum of 256 characters.
System name	The user-configured name of the local system. The system name can be a maximum of 256 characters.
System description	The system description contains information about the software and current image that is running on the system. This information is not configurable, but taken from the software.
System capabilities	The primary function that is performed by the system. The capabilities that the system supports; for example, bridge or router. This information is not configurable, but is based on the model of the product.
Management address	The IP management address of the local system.

Additional 802.3 TLVs include the following details:

Power via MDI A TLV that advertises MDI power support, a Power Sourcing Equipment (PSE) power pair, and power class information.

MAC/PHY configuration status

A TLV that advertises information about the physical interface, such as auto-negotiation status, support, and MAU type. The information is not configurable, but based on the physical interface structure.

Link aggregation A TLV that advertises if the port is aggregated and its aggregated port ID.

Maximum frame size A TLV that advertises the maximum transmission unit (MTU) of the interface that is sending LLDP frames.

Port VLAN A TLV that advertises the VLAN name that is configured on the interface.

LLDP-MED provides the following TLVs:

LLDP MED capabilities

A TLV that advertises the primary function of the port. The capability values range 0 - 15. The device class values range 0 - 255:

- 0 Capabilities
- 1 Network policy
- 2 Location identification
- 3 Extended power via MDI-PSE
- 4 Inventory
- 5 through 15** . . . Reserved

LLDP-MED device class values

- 0 Class not defined
- 1 Class 1 device
- 2 Class 2 device
- 3 Class 3 device
- 4 Network connectivity device
- 5 through 255** . . Reserved

Network policy A TLV that advertises the port VLAN configuration and associated Layer 2 and Layer 3 attributes. Attributes include the policy identifier, application types such as voice or streaming video, 802.1Q VLAN tagging, 802.1p priority bits, and Diffserv code points.

Endpoint location A TLV that advertises the physical location of the endpoint.

Extended power via MDI

A TLV that advertises the power type, power source, power priority, and power value of the port. It is the responsibility of the PSE device (network connectivity device) to advertise the power priority on a port.

4.2 Storage area network IP networking

Now that we introduced the protocols at a high level, what are the strategic differences between them all? Do I need them all, any, or none? What are some of the benefits that these technologies can bring me? The following list provides just some of the benefits you can gain:

- ▶ Departmental isolation and resource sharing alleviation
- ▶ Technology migration and integration
- ▶ Remote replication of disk systems
- ▶ Remote access to disk and tape systems
- ▶ Low-cost connection to storage area networks (SANs)
- ▶ Inter-fabric routing
- ▶ Overcoming distance limitations

People do not want to make a large financial investment without any guarantees that there will be some form of return. However, the appeal of these protocols is that they immediately bring benefits. Because these are standards-based protocols, they allow the use of both the existing TCP/IP and FCP infrastructure, they support existing Fibre Channel devices, and enable simplification of the infrastructure by removing any SAN islands.

4.2.1 The multiprotocol environment

As with any technology, it comes with its unique jargon and terminology. Typically it is borrowed from the networking world, but might have a different meaning. It is not our intent to cover every unique description. However, we do make some distinctions that we believe are important for a basic introduction to routing in an IP SAN.

4.2.2 Fibre Channel switching

A Fibre Channel *switch* filters and forwards packets between Fibre Channel connections on the *same* fabric, but it cannot transmit packets between fabrics. As soon as you join two switches together, you merge the two fabrics into a single fabric with one set of fabric services.

4.2.3 Fibre Channel routing

A *router* forwards data packets *between* two or more fabrics. Routers use headers and forwarding tables to determine the best path for forwarding the packets.

Separate fabrics each have their own addressing schemes. When they are joined by a router, there must be a way to translate the addresses between the two fabrics. This mechanism is called *network address translation* (NAT) and is inherent in all the IBM System Storage multiprotocol switch/router products. NAT is sometimes referred to as *FC-NAT* to differentiate it from a similar mechanism which exists in IP routers.

4.2.4 Tunneling

Tunneling is a technique that allows one network to send its data via the connection of another network. Tunneling works by encapsulating a network protocol within packets that are carried by the second network. For example, in a Fibre Channel over Internet Protocol (FCIP) solution, Fibre Channel packets can be encapsulated inside IP packets. Tunneling raises issues of packet size, compression, out-of-order packet delivery, and congestion control.

4.2.5 Routers and gateways

When a Fibre Channel router needs to provide protocol conversion or tunneling services, it is a *gateway* rather than a router. However, it is common usage to broaden the term *router* to include these functions. FCIP is an example of tunneling, while Internet Small Computer System Interface (iSCSI) and Internet Fibre Channel Protocol (iFCP) are examples of protocol conversion.

4.2.6 Internet Storage Name Service

The Internet Storage Name Service (iSNS) protocol facilitates automated discovery, management, and configuration of iSCSI and Fibre Channel devices that exist on an IP network. iSNS provides storage discovery and management services comparable to those that are found in Fibre Channel networks. What this means is that the IP network seems to operate in a similar capacity as a SAN. Coupling this capability with its ability to emulate Fibre Channel fabric services, iSNS allows for a transparent integration of IP and Fibre Channel networks. This integration is possible because it can manage both iSCSI and Fibre Channel devices.

4.3 Delving deeper into the protocols

We introduced all of the protocols at a high level. Now, in greater depth, we show the method of what they do with the Fibre Channel traffic.

4.3.1 Fibre Channel over Internet Protocol (FCIP)

FCIP is a method for tunneling Fibre Channel packets through an IP network. FCIP encapsulates Fibre Channel block data and transports it over a TCP socket, or tunnel. TCP/IP services are used to establish connectivity between remote devices. The Fibre Channel packets are not altered in any way. They are simply encapsulated in IP and transmitted.

Figure 4-10 shows FCIP tunneling, assuming that the Fibre Channel packet is small enough to fit inside a single IP packet.

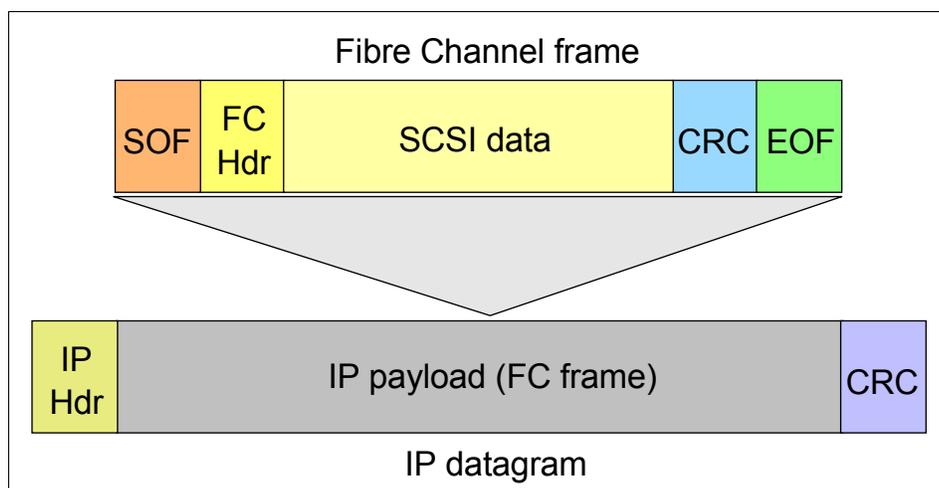


Figure 4-10 FCIP encapsulates the Fibre Channel frame into IP packets

The main advantage of FCIP is that it overcomes the distance limitations of basic Fibre Channel. It also enables geographically distributed devices to be linked by using the existing IP infrastructure, while it keeps the fabric services intact.

The architecture of FCIP is outlined in the Internet Engineering Task Force (IETF) Request for Comments (RFC) 3821, "Fibre Channel over TCP/IP (FCIP)", available at this website:

<http://ietf.org/rfc/rfc3821.txt>

Because FCIP simply tunnels Fibre Channel, creating an FCIP link is like creating an ISL. And, the two fabrics at either end are merged into a single fabric. This merge creates issues in situations where you do not want to merge the two fabrics for business reasons, or where the link connection is prone to occasional fluctuations.

Many corporate IP links are robust, but it can be difficult to be sure because traditional IP-based applications tend to be retry-tolerant. Fibre Channel fabric services are not as retry-tolerant. Each time the link disappears or reappears, the switches renegotiate and the fabric is reconfigured.

By combining FCIP with FC-FC routing, the two fabrics can be left "unmerged", each with its own separate Fibre Channel services.

4.3.2 Internet Fibre Channel Protocol (iFCP)

iFCP is a gateway-to-gateway protocol. It provides Fibre Channel fabric services to Fibre Channel devices over an IP network. iFCP uses TCP to provide congestion control, error detection, and recovery. The primary purpose of iFCP is to allow interconnection and networking of existing Fibre Channel devices at wire speeds over an IP network.

Under iFCP, IP components and technology replace the Fibre Channel switching and routing infrastructure. iFCP was originally developed by Nishan Systems, acquired by McDATA in September 2003. McDATA was then acquired by Brocade.

To learn more about the architecture and specification of iFCP, see the document at this website:

<http://tools.ietf.org/html/draft-ietf-ips-ifcp-14>

There is a myth that iFCP does not use encapsulation. In fact, iFCP encapsulates the Fibre Channel packet in much the same way that FCIP does. In addition, it maps the Fibre Channel header to the IP header and a TCP session, as shown in Figure 4-11.

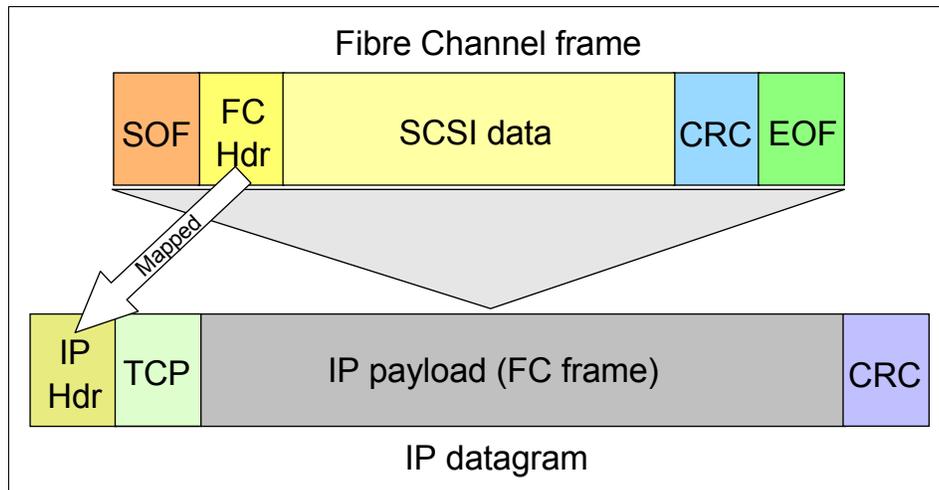


Figure 4-11 iFCP encapsulation and header mapping

iFCP uses the same iSNS mechanism that is used by iSCSI.

iFCP also allows data to fall across IP packets and share IP packets. Some FCIP implementations can achieve a similar result when they run software compression, but not otherwise. FCIP typically breaks each large Fibre Channel packet into two dedicated IP packets. iFCP compression is payload compression only. Headers are not compressed to simplify diagnostics.

iFCP uses one TCP connection per *fabric login (FLOGI)*, while FCIP typically uses one connection per router link (although more are possible). A FLOGI is the process by which an N_PORT registers its presence on the fabric, obtains fabric parameters such as classes of service that are supported, and receives its N_PORT address. Because under iFCP there is a separate TCP connection for each N_PORT to N_PORT couple, each connection can be managed to have its own quality of service (QoS) identity. A single incidence of congestion does not have to drop the sending rate for all connections on the link.

While all iFCP traffic between a given remote and local N_PORT pair must use the same iFCP session, that iFCP session can be shared across multiple gateways or routers.

4.3.3 Internet Small Computer System Interface (iSCSI)

The Small Computer System Interface (SCSI) protocol has a client/server architecture. Clients (called *initiators*) issue SCSI commands to request services from logical units on a server that is known as a *target*. A SCSI *transport* maps the protocol to a specific interconnect.

The SCSI protocol is mapped over various transports, including Parallel SCSI, Intelligent Peripheral Interface (IPI), IEEE-1394 (firewire), and Fibre Channel. All of these transports are ways to pass SCSI commands. Each transport is I/O specific and has limited distance capabilities.

The iSCSI protocol is a means of transporting SCSI packets over TCP/IP to take advantage of the existing Internet infrastructure.

A session between a iSCSI initiator and an iSCSI target is defined by a session ID. This ID is a combination of an initiator part (ISID) and a target part (Target Portal Group Tag).

The iSCSI transfer direction is defined with respect to the initiator. Outbound or outgoing transfers are transfers from an initiator to a target. Inbound or incoming transfers are transfers from a target to an initiator.

For performance reasons, iSCSI allows a *phase-collapse*. A command and its associated data might be shipped together from initiator to target, and data and responses might be shipped together from targets.

An iSCSI name specifies a logical initiator or target. It is not tied to a port or hardware adapter. When multiple NICs are used, they generally all present the same iSCSI initiator name to the targets because they are paths to the same SCSI layer. In most operating systems, the named entity is the operating system image.

The architecture of iSCSI is outlined in IETF RFC 3720, "Internet Small Computer Systems Interface (iSCSI)", at this website: <http://www.ietf.org/rfc/rfc3720.txt>

Figure 4-12 shows the format of the iSCSI packet.

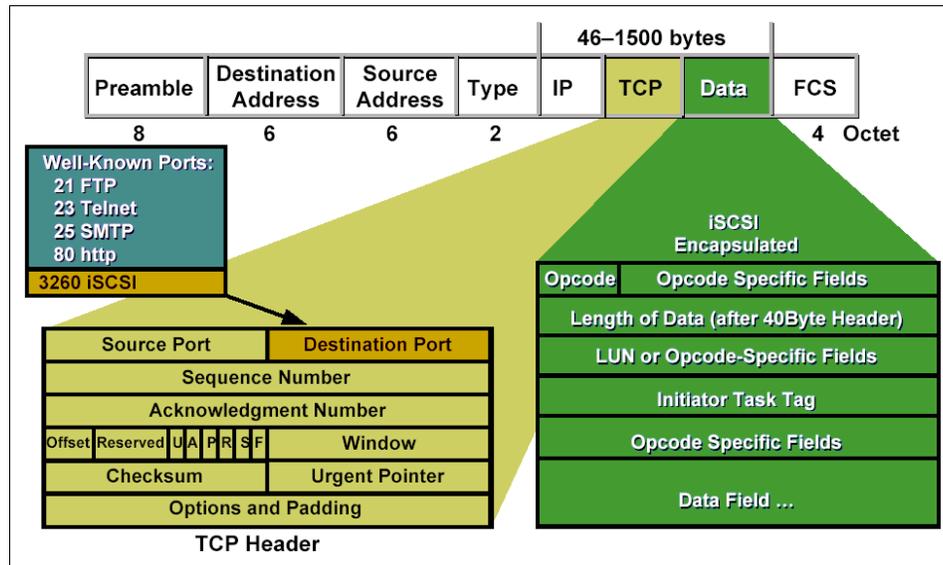


Figure 4-12 iSCSI packet format

Testing on iSCSI latency shows a difference of up to 1 ms of additional latency for each disk I/O as compared to Fibre Channel. This delay does not include such factors as trying to do iSCSI I/O over a shared, congested, or long-distance IP network, all of which might be tempting for some clients. iSCSI generally uses a shared 1 Gbps network.

iSCSI naming and discovery

Although we do not propose to go into an in-depth analysis of iSCSI in this book, there are ways for an iSCSI initiator to understand which devices are in the network:

- ▶ In small networks, you can use the **sendtargets** command.
- ▶ In larger networks, you can use the Service Location Protocol (SLP, multicast discovery).
- ▶ In large networks, we suggest that you use iSNS.

iSNS: At time of writing, not all vendors have delivered iSNS.

You can find a range of drafts that cover iSCSI naming, discovery, and booting at this website:

<http://www.ietf.org/proceedings/02mar/220.htm>

4.3.4 Routing considerations

As you might expect with any technology, there are going to be a unique set of characteristics that must be given consideration. The topics that follow briefly describe some of the issues, or items, that are considerations in a multiprotocol Fibre Channel environment.

4.3.5 Packet size

The standard size of a Fibre Channel packet is 2148 bytes, and the standard IP packet size is 1500 bytes (with a 1460-byte payload). It is evident that one packet is larger than the other and must be accommodated somehow.

When you transport Fibre Channel over IP, you can use jumbo IP packets to accommodate larger Fibre Channel packets. Keep in mind that jumbo IP packets must be turned on for the whole data path. In addition, a jumbo IP packet is not compatible with any devices in the network that do not have jumbo IP packets enabled.

Alternatively, you can introduce various schemes to split Fibre Channel packets across two IP packets. Some compression algorithms can allow multiple small Fibre Channel packets or packet segments to share a single IP packet.

Each technology and each vendor might implement this procedure differently. But the key point is that they all try to avoid sending small inefficient packets.

4.3.6 TCP congestion control

Sometimes standard TCP congestion mechanisms might not be suitable for tunneling storage. Standard TCP congestion control is designed to react quickly and severely to network congestion, but to recover slowly. This is well suited to traditional IP networks that are variable and unreliable. But for storage applications, this approach is not always appropriate and might cause disruption to latency-sensitive applications.

When three duplicate unanswered packets are sent on a traditional TCP network, the sending rate backs-off by 50%. When packets are successfully sent, it does a slow-start linear ramp-up again.

Some vendors tweak the back-off and recovery algorithms. For example, the tweak causes the send rate to drop by 12.5% each time that congestion is encountered. And the algorithm is tweaked so that the network can recover rapidly to the full sending rate by doubling each time until the full rate is regained.

Other vendors take a simpler approach to achieve a similar outcome.

If you are sharing your IP link between storage and other IP applications, then using either of these storage-friendly congestion controls might affect your other applications.

For more information about the specification for TCP congestion control, see this website:

<http://www.ietf.org/rfc/rfc2581.txt>

4.3.7 Round-trip delay

Round-trip link latency is the time that it takes for a packet to make a round trip across the link. The term *propagation delay* is also sometimes used. Round-trip delay generally includes both inherent latency and delays because of congestion.

Fibre Channel cable has an inherent latency of approximately 5 microseconds per kilometer each way. Typical Fibre Channel devices, like switches and routers, have inherent latencies of around 5 microseconds each way. IP routers might vary 5 - 100 microseconds in theory, but when tested with filters applied, the results are more likely to be measured in milliseconds.

This type of measurement is the essential problem with tunneling Fibre Channel over IP. Fibre Channel applications are generally designed for networks that have round-trip delays that are measured in microseconds. IP networks generally deliver round-trip delays that are measured in milliseconds or tens of milliseconds. Internet connections often have round-trip delays that are measured in hundreds of milliseconds.

Any round-trip delay that is caused by more routers and firewalls along the network connection also has to be added to the total delay. The total round-trip delay varies considerably depending on the models of routers or firewalls that are used, and the traffic congestion on the link.

So how does this latency affect you? If you are purchasing the routers or firewalls yourself, we recommend that you include the latency of any particular product in the criteria that you use to choose the products. If you are provisioning the link from a service provider, we recommend that you include at least the maximum total round-trip latency of the link in the service level agreement (SLA).

Time of frame in transit

The time of frame in transit is the actual time that it takes for a specific frame to pass through the slowest point of the link. Therefore, it depends on both the frame size and link speed.

The maximum size of the payload in a Fibre Channel frame is 2112 bytes. The Fibre Channel headers add 36 bytes to this measurement, for a total Fibre Channel frame size of 2148 bytes. When you transfer data, Fibre Channel frames at or near the full size are usually used.

If we assume that we are using jumbo frames in the Ethernet, the complete Fibre Channel frame can be sent within one Ethernet packet. The TCP and IP headers and the Ethernet medium access control (MAC) add a minimum of 54 bytes to the size of the frame. This addition gives a total Ethernet packet size of 2202 bytes, or 17616 bits.

For smaller frames, such as the Fibre Channel acknowledgement frames, the time in transit is much shorter. The minimum possible Fibre Channel frame is one with no payload. With FCIP encapsulation, the minimum size of a packet with only the headers is 90 bytes, or 720 bits.

Table 4-2 details the transmission times of this FCIP packet over various common wide area network (WAN) link speeds.

Table 4-2 FCIP packet transmission times over different WAN links

Link type	Link speed	Large packet	Small packet
Gigabit Ethernet	1250 Mbps	14 μ s	0.6 μ s
OC-12	622.08 Mbps	28 μ s	1.2 μ s
OC-3	155.52 Mbps	113 μ s	4.7 μ s

Link type	Link speed	Large packet	Small packet
T3	44.736 Mbps	394 μ s	16.5 μ s
E1	2.048 Mbps	8600 μ s	359 μ s
T1	1.544 Mbps	11 400 μ s	477 μ s

If we cannot use jumbo frames, each large Fibre Channel frame must be divided into two Ethernet packets. This requirement doubles the amount of TCP, IP, and Ethernet MAC overhead for the data transfer.

Normally, each Fibre Channel operation transfers data in only one direction. The frames that move in the other direction are close to the minimum size.

4.4 Multiprotocol solution briefs

The solution briefs in the following sections show how you can use multiprotocol routers.

4.4.1 Dividing a fabric into subfabrics

Assume that you have eight switches in your data center, and they are grouped into two fabrics of four switches each. Two of the switches are used to connect the development and test environment, two are used to connect a joint-venture subsidiary company, and four are used to connect the main production environment.

The development and test environment does not follow the same change control disciplines as the production environment. Also, systems and switches can be upgraded, downgraded, or rebooted on occasions, usually unscheduled and without any form of warning.

The joint-venture subsidiary company is up for sale. The mandate is to provide as much separation and security as possible between it and the main company, and the subsidiary. The backup and restore environment is shared between the three environments.

In summary, we have a requirement to provide a degree of isolation, and a degree of sharing. In the past, this requirement was accommodated through zoning. Some fabric vendors might still recommend that approach as the simplest and most cost-effective. However, as the complexity of the environment grows, zoning can become complex. Any mistakes in setup can disrupt the entire fabric. Adding FC-FC routing to the network allows each of the three environments to run separate fabric services and provides the capability to share the tape backup environment.

In larger fabrics with many switches and separate business units, for example, in a shared services-hosting environment, separation and routing are valuable. These features are beneficial in creating a larger number of simple fabrics, rather than fewer more complex fabrics.

4.4.2 Connecting a remote site over IP

Suppose that you want to replicate your disk system to a remote site, perhaps 50 km away synchronously, or 500 km away asynchronously. Using FCIP tunneling or iFCP conversion, you can transmit your data to the remote disk system over a standard IP network. The router includes Fibre Channel ports to connect network devices, or switches and IP ports to connect to a standard IP WAN router. Standard IP networks are generally much lower in cost to

provision than traditional high-quality dedicated *dense wavelength division multiplexing (DWDM)* networks. Standard IP networks also often have the advantage of being understood by internal operational staff.

Similarly, you might want to provision storage volumes from your disk system to a remote site by using FCIP or iFCP.

Low-cost connections: FCIP and iFCP can provide a low-cost way to connect remote sites using familiar IP network disciplines.

4.4.3 Connecting hosts using Internet Small Computer System Interface

Many hosts do not require high bandwidth, low latency access to storage. For such hosts, Internet Small Computer System Interface (iSCSI) might be a more cost-effective connection method. iSCSI can be thought of as an IP SAN. There is no requirement to provide a Fibre Channel switch port for every server. Nor is there a need to purchase Fibre Channel host bus adapters (HBAs), or to lay Fibre Channel cable between storage and servers.

The iSCSI router has both Fibre Channel and Ethernet ports to connect to servers located either locally on the Ethernet, or remotely, over a standard IP WAN connection.

The iSCSI connection delivers block I/O access to the server so that it is application independent. That is, an application cannot really tell the difference between direct SCSI, iSCSI, or Fibre Channel, since all three are delivery SCSI block I/Os.

Different router vendors quote different limits on the number of iSCSI connections that are supported on a single IP port.

iSCSI places a significant packetizing and depacketizing workload on the server CPU. This workload can be mitigated by using the TCP/IP offload engine (TOE) Ethernet cards. However, since these cards can be expensive, they somewhat undermine the low-cost advantage of iSCSI.

iSCSI provides low-cost connections: iSCSI can be used to provide low-cost connections to the SAN for servers that are not performance critical.



Topologies and other fabric services

In this chapter, we introduce Fibre Channel (FC) topologies and other fabric services that are commonly encountered in a storage area network (SAN).

We also provide an insight into the emerging converged topology and the option to merge FC to Fibre Channel over Ethernet (FCoE).

5.1 Fibre Channel topologies

Fibre Channel-based networks support three types of base topologies: point-to-point, arbitrated loop, and switched fabric. A *switched fabric* is the most commonly encountered topology today and it has subclassifications of topology. Figure 5-1 depicts the various classifications of SAN topology.

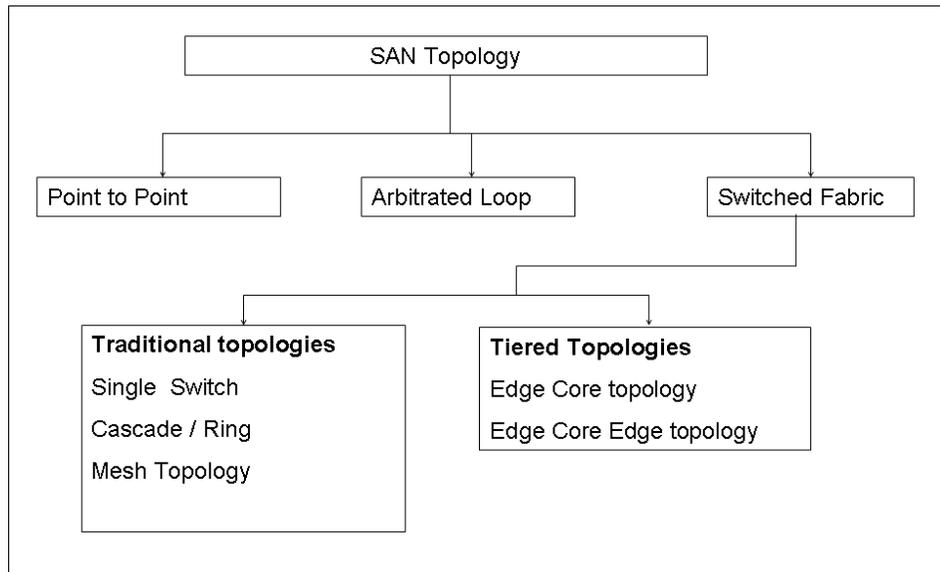


Figure 5-1 SAN topologies

5.1.1 Point-to-point topology

A *point-to-point connection* is the simplest topology. It is used when there are exactly two nodes, and future expansion is not predicted. There is no sharing of the media, which allows the devices to use the total bandwidth of the link. A simple link initialization is needed before communications can begin.

Fibre Channel is a *full duplex protocol*, which means both paths transmit data simultaneously. As an example, Fibre Channel connections that are based on the 1 Gbps standard are able to transmit at 100 MBps and receive at 100 MBps simultaneously. As another example, for Fibre Channel connections that are based on the 2 Gbps standard, they can transmit at 200 MBps and receive at 200 MBps simultaneously. This speed extends to 4 Gbps, 8 Gbps, and 16 GBPS technologies as well.

Figure 5-2 shows an illustration of a simple point-to-point connection.

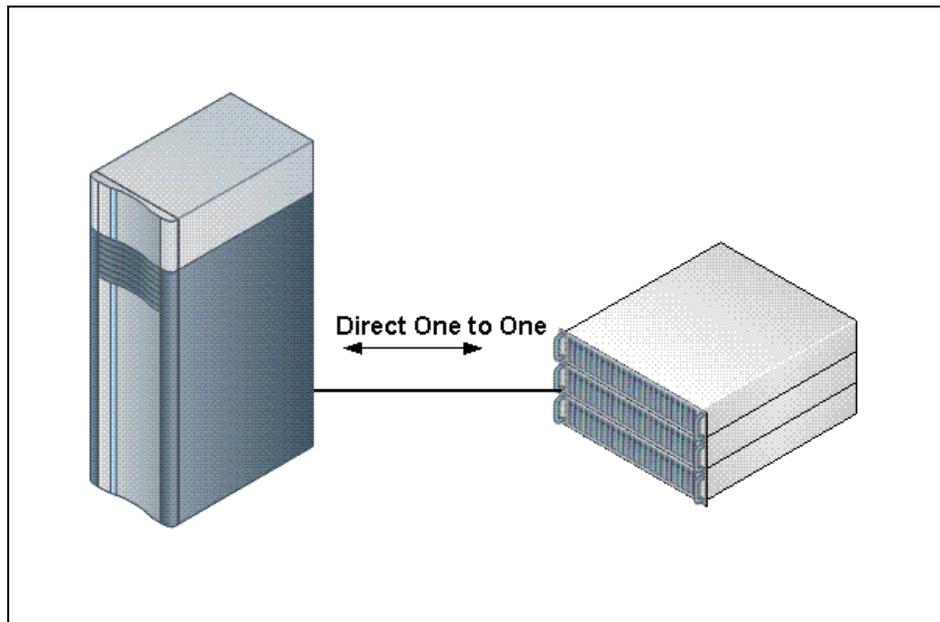


Figure 5-2 Point-to-point connection

5.1.2 Arbitrated loop topology

Arbitrated loop topology: Although this topology is rarely encountered anymore, and is considered as a legacy topology, we include it for historical reasons only.

Our second topology is *Fibre Channel Arbitrated Loop (FC-AL)*. FC-AL is more useful for storage applications. It is a loop of up to 126 nodes (NL_Ports) that is managed as a shared bus. Traffic flows in one direction, carrying data frames and primitives around the loop with a total bandwidth of 400 MBps (or 200 MBps for a loop-based topology on 2 Gbps technology).

Using arbitration protocol, a single connection is established between a sender and a receiver, and a data frame is transferred around the loop. When the communication comes to an end between the two connected ports, the loop becomes available for arbitration and a new connection might be established. Loops can be configured with hubs to make connection management easier. A distance of up to 10 km is supported by the Fibre Channel standard for both of these configurations. However, latency on the arbitrated loop configuration is affected by the loop size.

A simple loop, which is configured by using a hub, is shown in Figure 5-3.

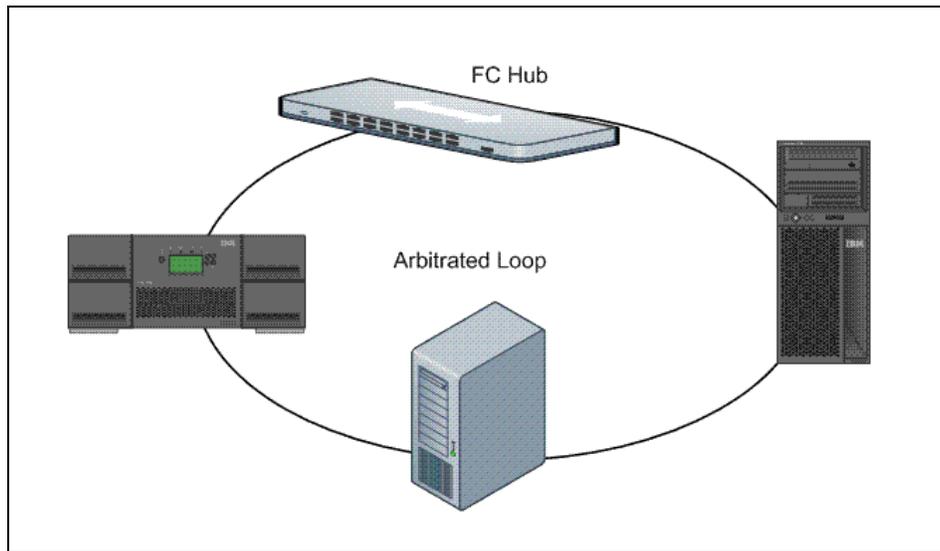


Figure 5-3 Arbitrated loop

We describe FC-AL in more depth in 5.4, “Fibre Channel Arbitrated Loop protocols” on page 108.

5.1.3 Switched fabric topology

Our third, and the most useful topology that is used in SAN implementations, is *Fibre Channel Switched Fabric (FC-SW)*. It applies to switches and directors that support the FC-SW standard; that is, it is not limited to switches as its name suggests. A Fibre Channel fabric is one or more fabric switches in a single, sometimes extended, configuration. Switched fabrics provide full bandwidth for each port that is compared to the shared bandwidth for each port in arbitrated loop implementations.

One of the key differentiators is that if you add a device into the arbitrated loop, you further divide the shared bandwidth. However, in a switched fabric, adding a device or a new connection between existing ones actually increases the bandwidth. For example, an eight-port switch (assume that it is based on 2 Gbps technology) with three initiators and three targets, can support three concurrent 200 MBps conversations or a total of 600 MBps throughput. This equates to 1,200 MBps, if full-duplex applications were available.

A switched fabric configuration is shown in Figure 5-4.

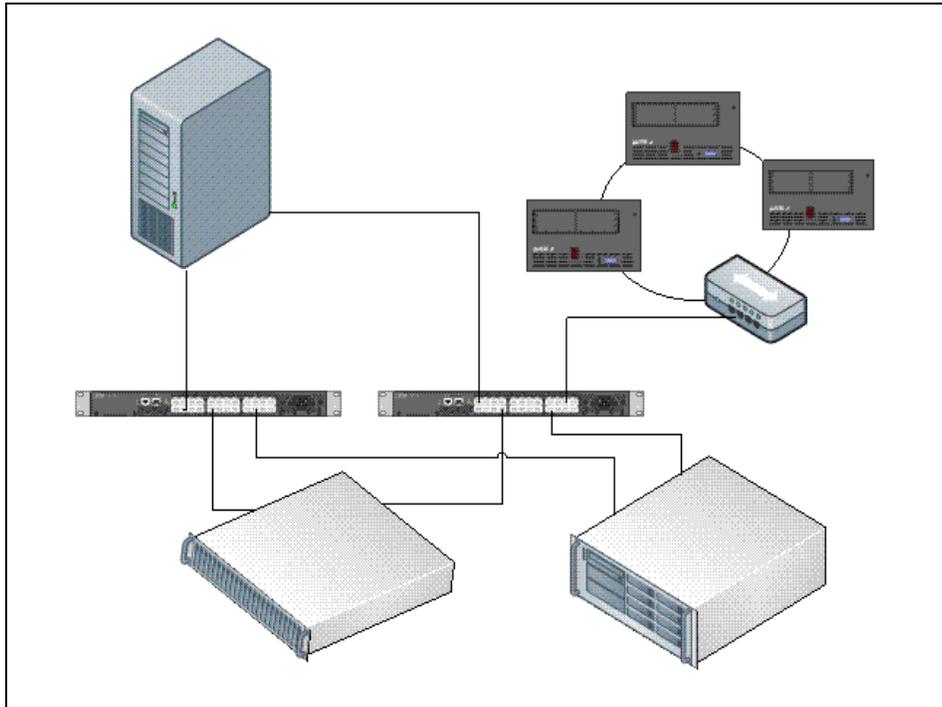


Figure 5-4 Sample switched fabric topology

This configuration is one of the major reasons why arbitrated loop is considered a legacy SAN topology. A *switched fabric* is usually referred to as a *fabric*.

In terms of switch interconnections, the switched SAN topologies can be classified as the following types:

- ▶ Single switch topology
- ▶ Cascaded and ring topology
- ▶ Mesh topology

5.1.4 Single switch topology

The *single switch topology* has only one switch and has no *inter-switch links (ISLs)*. It is the simplest design for infrastructures which do not need any redundancy. Because of the issues of it introducing a single point of failure, this topology is rarely used.

Figure 5-5 indicates a single switch topology with all devices connected to same switch.

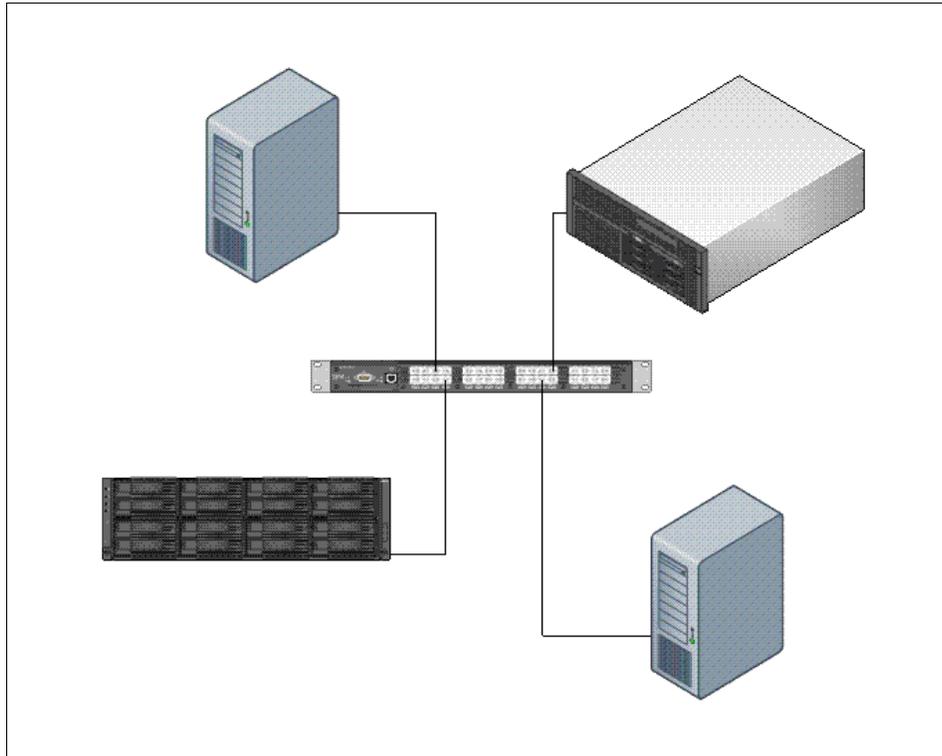


Figure 5-5 Single switch topology

5.1.5 Cascaded and ring topology

In a *cascaded topology*, switches are connected in a *queue fashion*, as shown in Figure 5-6.

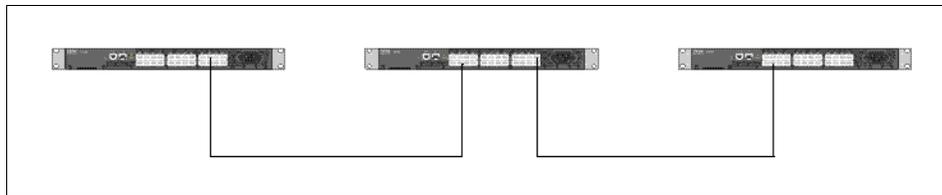


Figure 5-6 Cascade topology

Even in a *ring topology*, the switches are connected in a queue fashion, but it forms a closed ring with an additional inter-switch link (ISL), as shown in Figure 5-7.

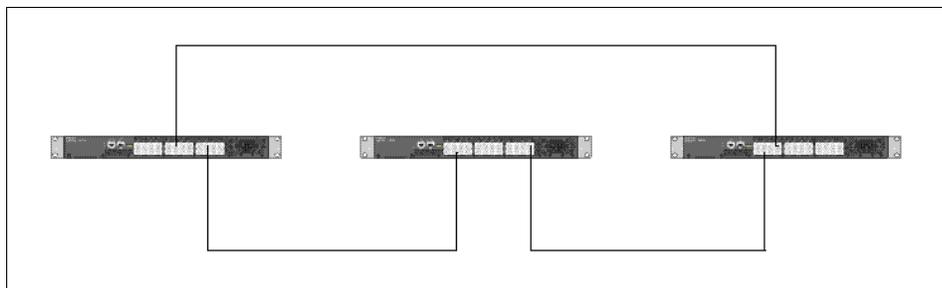


Figure 5-7 Ring topology

5.1.6 Mesh topology

In a full *mesh topology*, each switch is connected to every other switch in the fabric, as shown in Figure 5-8.

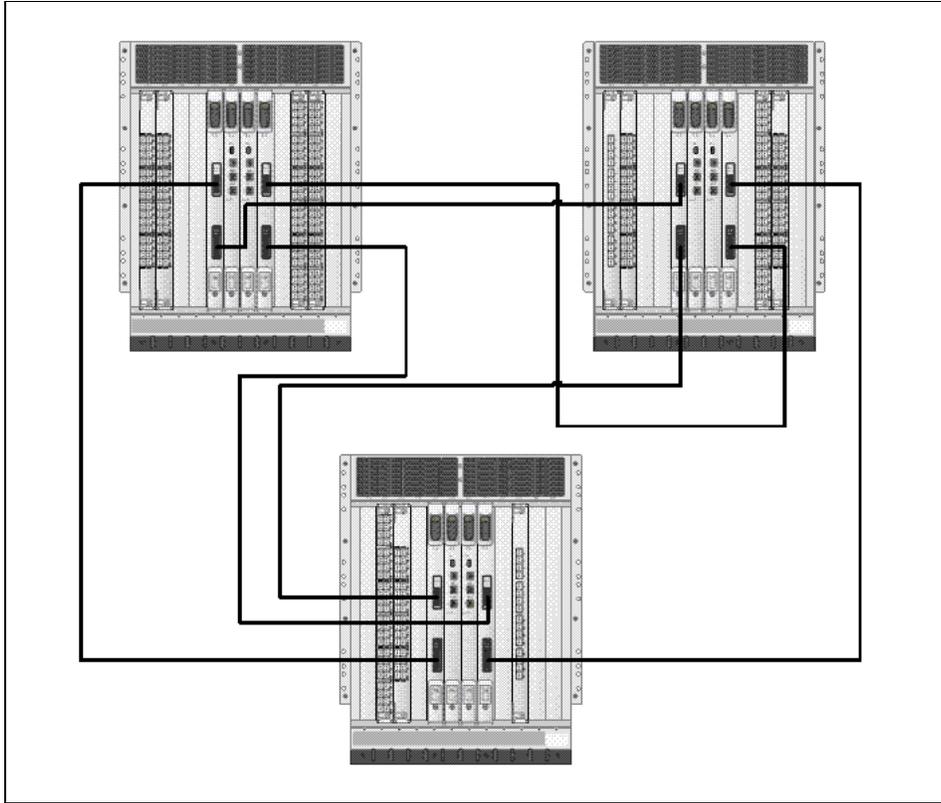


Figure 5-8 IBM SAN768B connected to form a mesh topology

In terms of a tiered approach, the switched fabric can be further classified with the following topology:

- ▶ Core-edge topology
- ▶ Edge-core-edge topology

5.1.7 Core-edge topology

In *core-edge topology*, the servers are connected to the edge fabric and the storage is connected to core switches. Figure 5-9 depicts the core-edge topology.

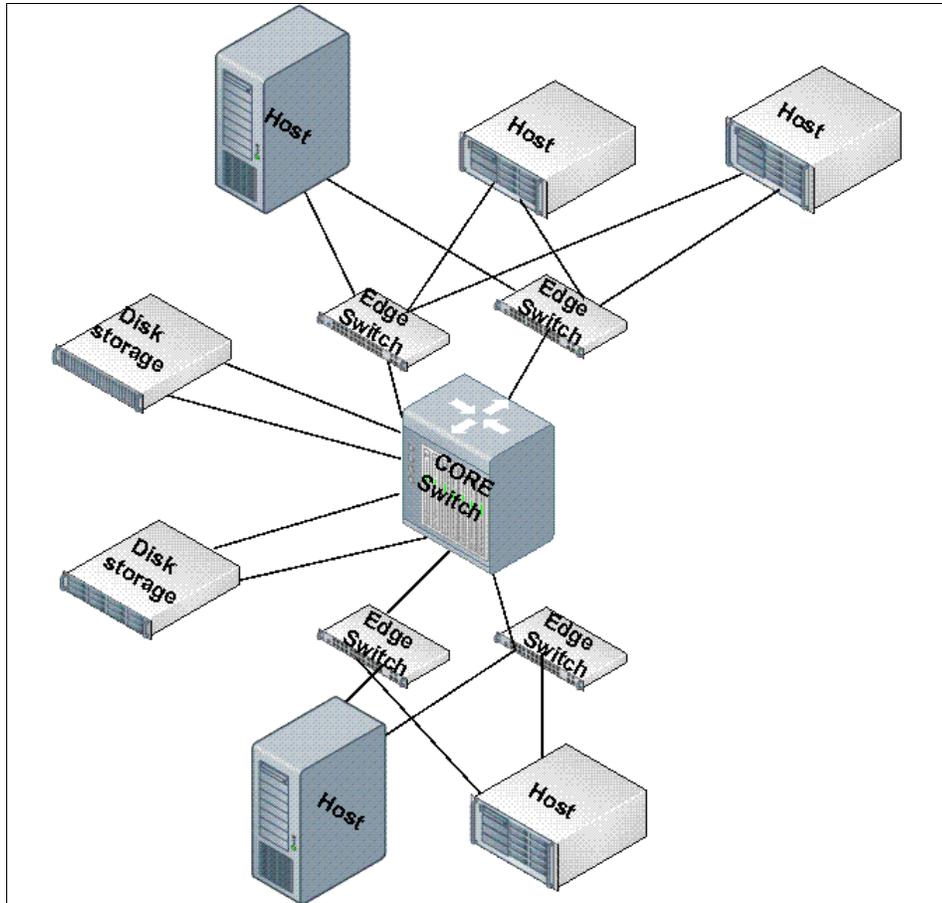


Figure 5-9 Core-edge topology

5.1.8 Edge-core-edge topology

In this topology, the server and storage are connected to the edge fabric and the core switch connectivity is used only for scalability in terms of connecting to edge switches. This configuration expands the SAN traffic flow to long distance by dense wavelength division multiplexing (DWDM), connecting to virtualization appliances, and encryption switches. Also, the servers might be isolated to one edge and storage can be at the other edge which helps with management.

Figure 5-10 shows the edge-core-edge topology.

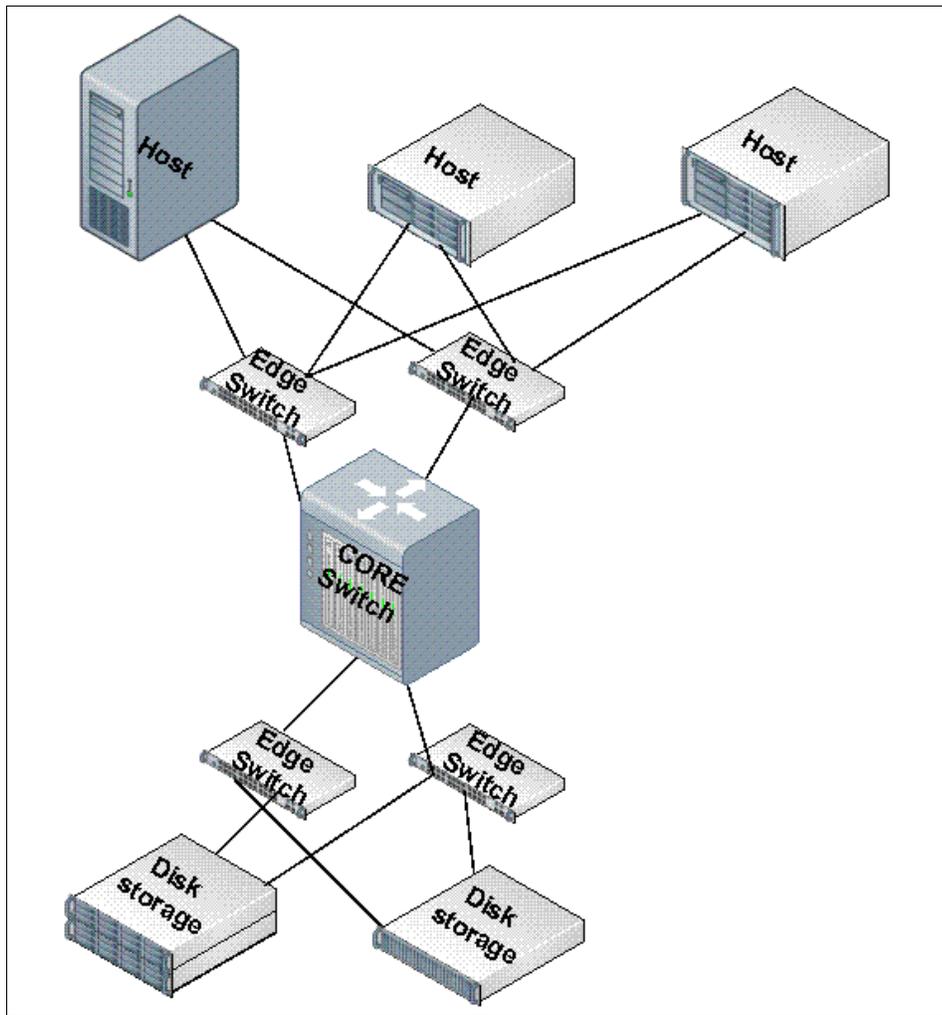


Figure 5-10 Edge-core-edge topology

5.2 Port types

The basic building block of the Fibre Channel is the *port*. There are various kinds of Fibre Channel port types.

5.2.1 Common port types

The following list provides the various kinds of Fibre Channel port types and their purpose in switches, servers, and storage:

- ▶ F_Port. This is a fabric port that is connected to an N_Port point-point to a switch.
- ▶ FL_Port. This is a fabric port that is counted to a loop device. It is used to connect an NL_Port to the switch in a public loop configuration.
- ▶ TL_port. A Cisco specific port type. It is a Translative loop port that is connected with non-fabric aware, private loop devices.

- ▶ **G_Port.** This is a generic port that can operate as either an E_Port or an F_Port. A port is defined as a G_Port after it is connected but has not received a response to *loop initialization* or has not yet completed the link initialization procedure with the adjacent Fibre Channel device.
- ▶ **L_Port.** This is a loop-capable node or switch port.
- ▶ **U_Port.** This is a universal port: a more generic switch port than a G_Port. It can operate as either an E_Port, F_Port, or FL_Port. A port is defined as a U_Port when it is not connected or has not yet assumed a specific function in the fabric.
- ▶ **N_Port.** This is a node port that is not loop capable. It is a host end port that is used to connect to the fabric switch.
- ▶ **NL_Port.** This is a node port that is loop capable. It is used to connect an equipment port to the fabric in a loop configuration through an L_Port or FL_Port.

Figure 5-11 depicts different common port types of switch and nodes.

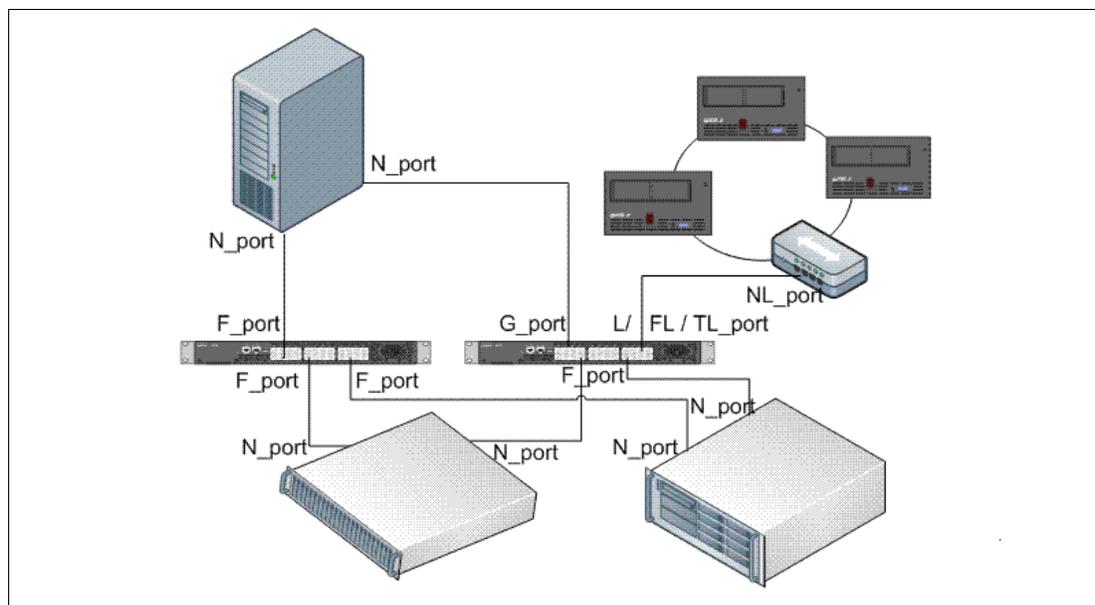


Figure 5-11 Common port types

5.2.2 Expansion port types

The following ports are found in a multi-switch fabric where switches are interconnected via an FC link:

- ▶ **E_Port.** This is an expansion port. A port is designated an E_Port when it is used as an ISL to connect to the E_Port of another switch to enlarge the switch fabric.
- ▶ **Ex_port.** The type of E_Port used to connect a Multiprotocol Router to an edge fabric. An EX_Port follows standard E_Port protocols, and supports FC-NAT, but does not allow fabric merging across EX_Ports.
- ▶ **VE_port.** A virtual E port is a port that emulates an E_Port over an FCIP link. VE port connectivity is supported over point-to-point links.

- ▶ VEX_port. VEX_Ports are routed VE_Ports, just as Ex_Ports are routed E_Ports. VE_Ports and VEX_Ports have the same behavior and functionality.
- ▶ TE_port. The TE_port provides not only standard E_port functions, but allows for routing of multiple VSANs (virtual SANs). This capability is accomplished by modifying the standard Fibre Channel frame (VSAN tagging) upon ingress and egress of the VSAN environment. It is also known as a *Trunking E_port*.

Figure 5-12 shows a fabric with expansion ports.

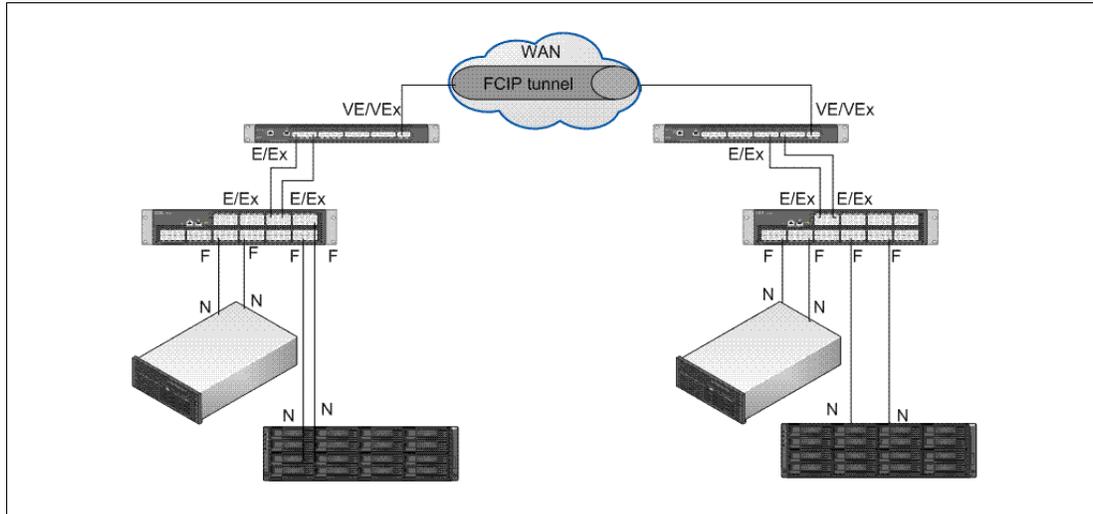


Figure 5-12 Fabric with expansion ports

5.2.3 Diagnostic port types

- ▶ D_port is a diagnostic port type which can be enabled only on the 16 Gbps b-type switches with Fabric Operating System 7.0. This system uses the *Spinfab* test and performs electrical loop back, optical loop back, measures link distance, and also stress tests with a link saturation test.

Figure 5-13 describes the different test options. Long-distance cable checks also can be done with D_Port diagnostic capabilities.

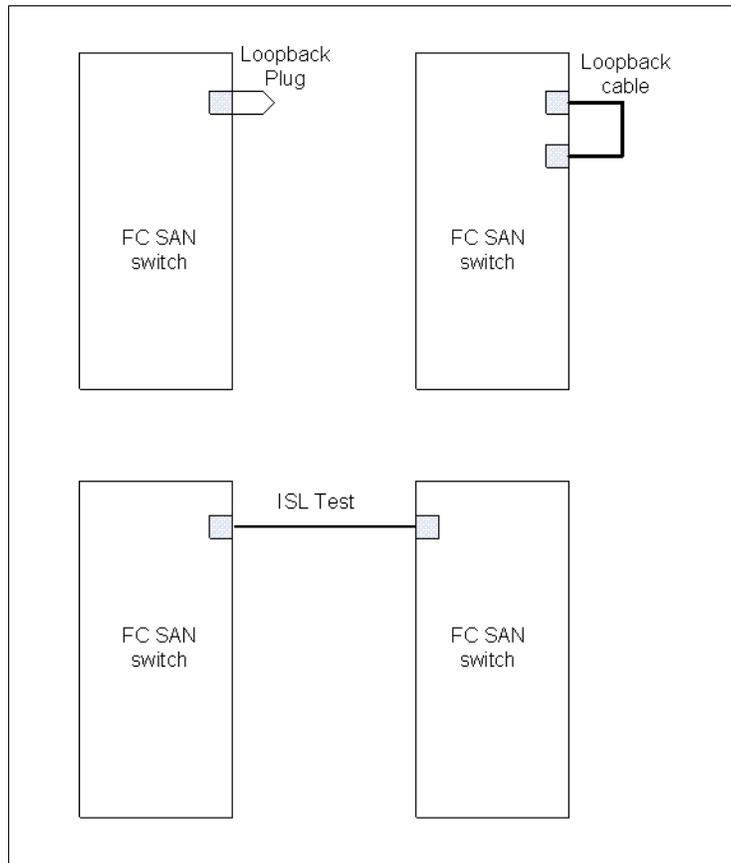


Figure 5-13 D_Port type diagnostics

- ▶ *MTx_Port* is a CNT port that is used as a mirror for viewing the transmit stream of the port to be diagnosed.
- ▶ *MRx_Port* is a CNT port that is used as a mirror for viewing the receive stream of the port to be diagnosed.
- ▶ *SD_Port* is a Cisco SPAN diagnostic port that is used for diagnostic capture with a connection to SPAN- switch port analyzer.
- ▶ *ST_Port* is the Cisco port type for Remote Strategic Position Analysis (RSPAN) monitoring in a source switch. This switch is an undedicated port that is used for RSPAN analysis, and is not connected to any other device.

Figure 5-14 represents the Cisco specific Fibre Channel port types.

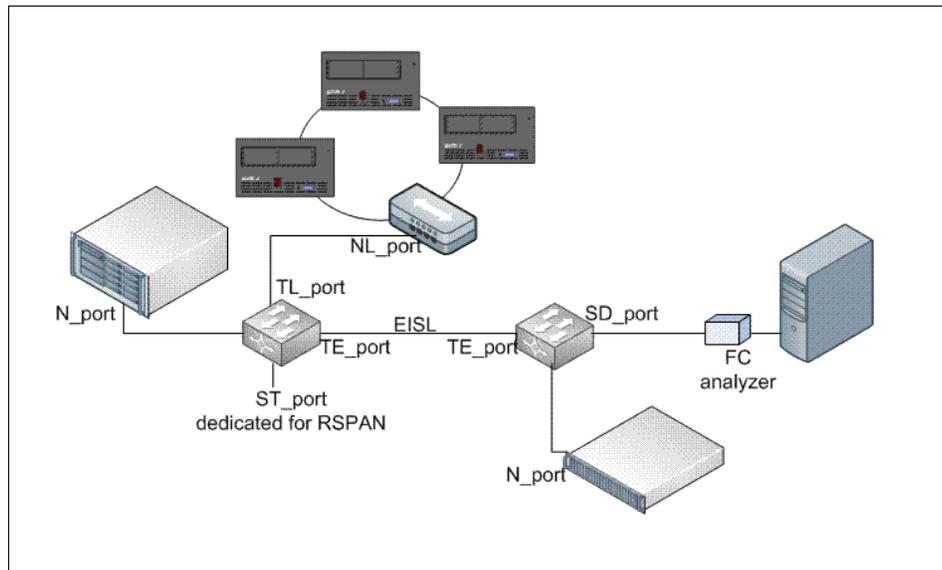


Figure 5-14 Cisco specific Fibre Channel ports

5.3 Addressing

All devices in a Fibre Channel environment have an identity. The way that the identity is assigned and used depends on the format of the Fibre Channel fabric. For example, there is a difference between the way that addressing is done in an arbitrated loop and a fabric.

5.3.1 Worldwide name

All Fibre Channel devices have a unique identity that is called a *worldwide name (WWN)*. This identification is similar to the way that all Ethernet cards have a unique *Media Access Control (MAC)* address.

Each N_Port has its own WWN, but it is also possible for a device with more than one Fibre Channel adapter to have its own WWN as well. Thus, for example, a storage server can have its own WWN and incorporate the WWNs of the adapter within it. This means that a soft zone can be created by using the entire array, or individual zones can be created by using particular adapters. In the future, this ability will be the case for the servers as well.

This WWN is a 64-bit address, and if two WWN addresses are put into the frame header, this leaves 16 bytes of data just for identifying the destination and source address. So 64-bit addresses can affect routing performance.

Each device in the SAN is identified by a unique WWN. The WWN contains a vendor identifier field, which is defined and maintained by the Institute of Electrical and Electronics Engineers (IEEE), and a vendor-specific information field.

Currently, there are two formats of the WWN as defined by the IEEE. The original format contains either a hex 10 or hex 20 in the first 2 bytes of the address. This address is then followed by the vendor-specific information.

Both the old and new WWN formats are shown in Figure 5-15 on page 98.

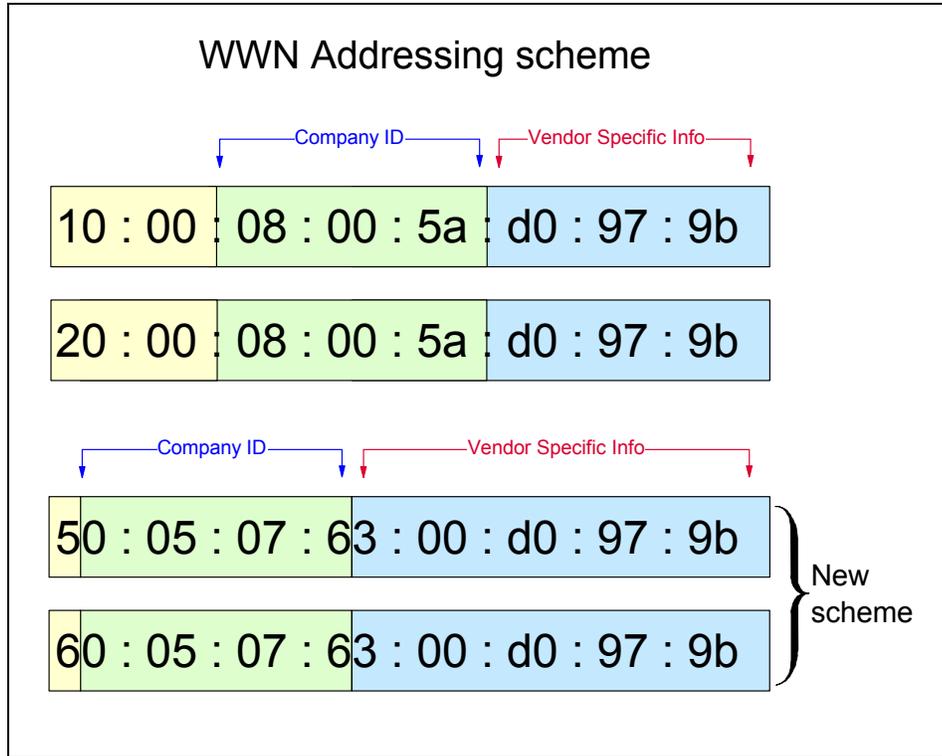


Figure 5-15 Worldwide name (WWN) addressing scheme

The new addressing scheme starts with a hex 5 or 6 in the first half-byte followed by the vendor identifier in the next 3 bytes. The vendor-specific information is then contained in the following fields. Both of these formats are currently in use and depend on the hardware manufacturer standards to follow either of the formats. However, the Vendor ID and company ID are assigned unique by the IEEE standards and each vendor and their identifier can be found in the following text file:

<http://standards.ieee.org/develop/regauth/oui/oui.txt>

A *worldwide node name (WWNN)* is a globally unique 64-bit identifier that is assigned to each Fibre Channel *node* or *device*. For servers and hosts, the WWNN is unique for each *host bus adapter (HBA)*, and in a case of a server with two HBAs, they have two WWNNs. For a SAN switch, the WWNN is common for the chassis. For storage, the WWNN is common for each controller unit of midrange storage. And, in a case of high-end enterprise storage, the WWNN is unique for the entire array.

A *worldwide port number (WWPN)* is a unique identifier for each FC port of any Fibre Channel device. For a server, we have a WWPN for each port of the HBA. For a switch, the WWPN is available for each port in the chassis; and for storage, each host port has an individual WWPN.

Server WWNN and WWPN

Figure 5-16 indicates that there is a WWNN for each HBA. And, every port in the HBA has an individual WWPN.

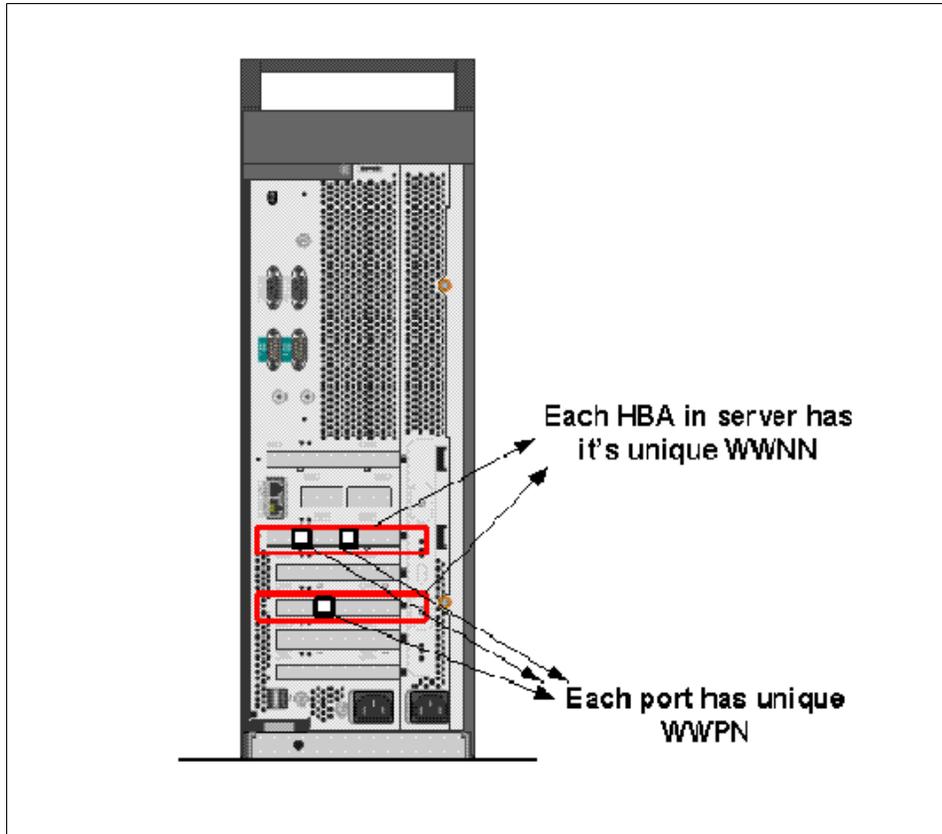


Figure 5-16 Server WWNN and WWPN

SAN WWNN and WWPN

Figure 5-17 on page 100 indicates that the WWNN is for the entire SAN switch chassis and the WWPN is for each FC port in the SAN switch chassis.

Fabric Assigned PWWNs: The new 16 Gbps b-type switches with FOS 7.0 can also have a virtual WWPN defined by switches called Fabric Assigned PWWNs (FAPWWNs). These FAPWWNs can be used for pre-configuring zoning before the physical servers are connected. This feature helps to simplify and accelerate server deployment and improve operational efficiency by avoiding the wait time for physical connectivity to be done. This feature also requires servers to be using Brocade HBAs/Adapters with an HBA driver version 3.0.0.0 or higher, which can be configured to use FAPWWN.

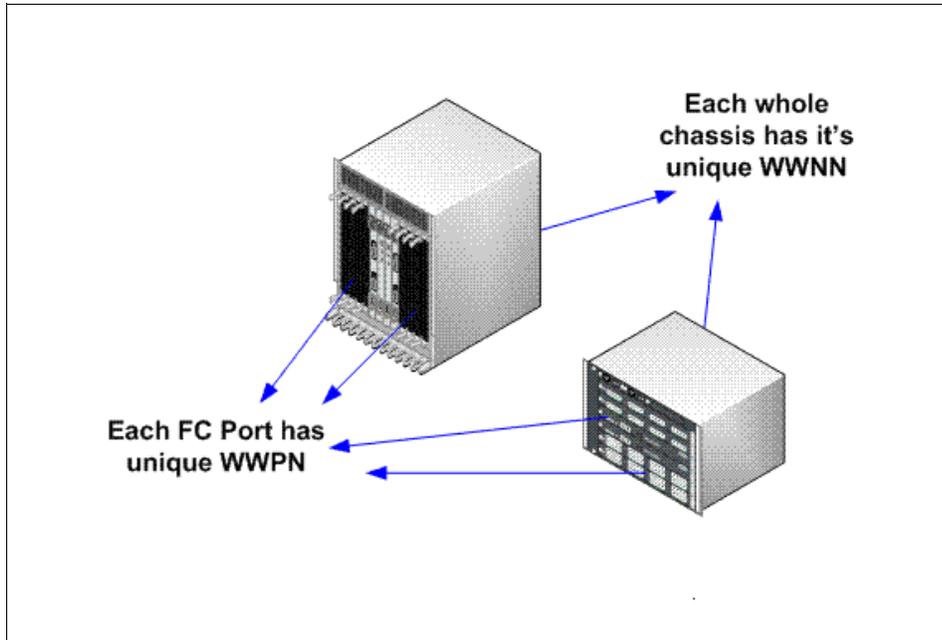


Figure 5-17 SAN switch WWNN and WWPN

Storage WWNN and WWPN

Disk storage has an individual WWNN for the entire storage system and the individual FC host ports have a unique WWPN, as indicated in Figure 5-18. The diagram shows a dual controller module.

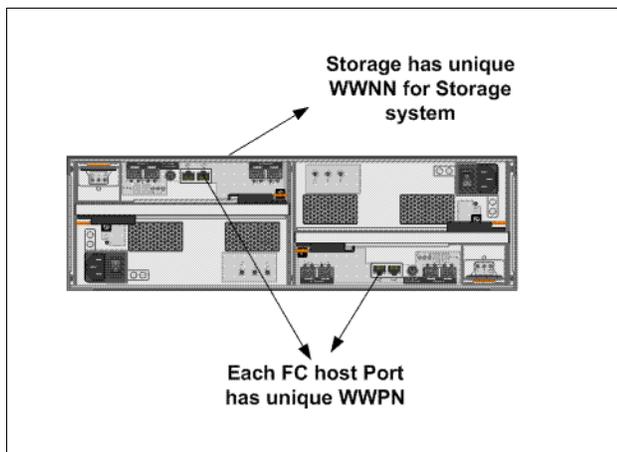


Figure 5-18 Storage WWNN and WWPN

Worldwide node name (WWNN): The IBM virtualization storage systems have a different WWNN usage. For example, each node in a SAN Volume Controller or the IBM Storwize® V7000, has an individual and unique WWNN.

For the IBM DS8000®, each Storage Facility Image has a unique individual WWNN.

5.3.2 Tape Device WWNN and WWPN

For tape devices, each drive inside the tape library has an individual WWPN and WWNN. Figure 5-19 indicates that multiple drive libraries have an individual WWNN and WWPN for each drive.

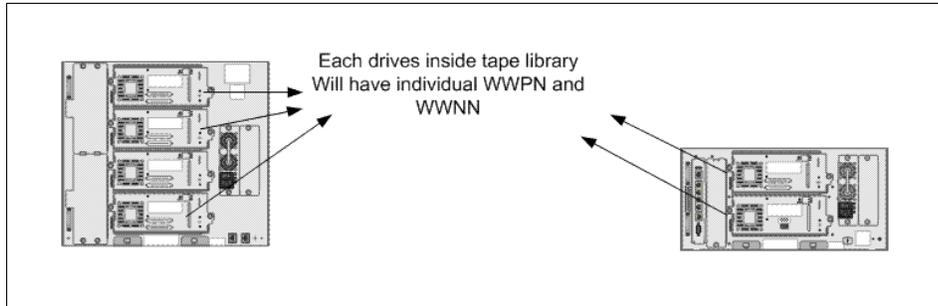


Figure 5-19 Tape device WWNN and WWPN

5.3.3 Port address

Because of the potential affect on routing performance by using 64-bit addressing, there is another addressing scheme that is used in Fibre Channel networks. This scheme is used to address ports in the switched fabric. Each port in the switched fabric has its own unique 24-bit address. With this 24-bit address scheme, there is a smaller frame header, and this configuration can speed up the routing process. With this frame header and routing logic, the Fibre Channel is optimized for high-speed switching of frames.

With a 24-bit addressing scheme, this configuration allows for up to 16 million addresses, which are an address space larger than any practical SAN design in existence in today's world. There has to be some relationship between this 24-bit address and the 64-bit address that is associated with WWNs. We explain this relationship in the section that follows.

5.3.4 The 24-bit port address

The 24-bit address scheme removes the performance overhead of the manual administration of addresses by allowing the topology itself to assign addresses. This configuration is *not* like WWN addressing where the addresses are assigned to manufacturers by the IEEE standards committee and are built into the device at the time of manufacture. If the topology itself is assigning the 24-bit addresses, then something must be responsible for maintaining the addressing scheme from WWN addressing to port addressing.

In the switched fabric environment, the switch itself is responsible for assigning and maintaining the port addresses. When a device with a WWN logs in to the switch on a specific port, the switch assigns the port address to that port. The switch also maintains the correlation between the port address and the WWN address of the device of that port. This function of the switch is implemented by using the name server.

The *name server* is a component of the fabric operating system, which runs inside the switch. It is essentially a database of objects in which a fabric-attached device registers its values.

Dynamic addressing also removes the partial element of human error in addressing maintenance and provides more flexibility in additions, moves, and changes in the SAN.

A 24-bit port address consists of three parts:

- ▶ Domain (bits 23 - 16)
- ▶ Area (bits 15 - 08)
- ▶ Port or Arbitrated Loop physical address: AL_PA (bits 07 - 00)

We show how the address is built in Figure 5-20.

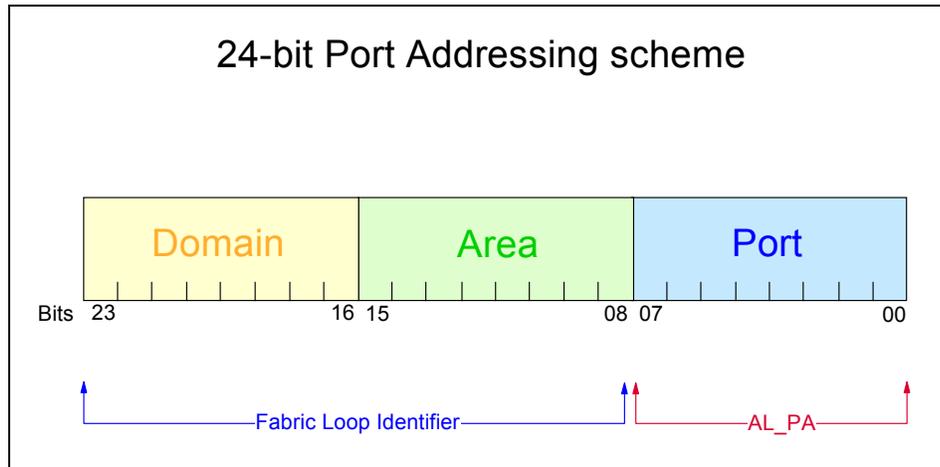


Figure 5-20 Fabric port address

The following functions provide the significance of some of the bits that make up the port address:

▶ Domain

The most significant byte of the port address is the *domain*. This byte is the address of the switch itself. A *domain ID* is a unique number that identifies the switch or director to a fabric. It can be either *static* or *dynamic*. *Static* (insistent) domain IDs are a requirement for Fibre Channel connection (FICON). Each manufacturer has a range of numbers and a maximum number of domain IDs that can be used in a fabric.

One byte allows up to 256 possible addresses. Because some of these addresses are reserved, as for the one for broadcast, there are only 239 addresses available. This number means that you can theoretically have as many as 239 switches in your SAN environment. The domain number allows each switch to have a unique identifier if you have multiple interconnected switches in your environment.

▶ Area

The *area* field provides 256 addresses. This part of the address is used to identify the individual ports. Hence, to have more than 256 ports in one switch in a director class of switches, you must follow the shared area addressing.

▶ Port

The final part of the address provides 256 addresses for identifying attached N_Ports and NL_Ports.

To arrive at the number of available addresses is a simple calculation, which is based on:

$$\text{Domain} \times \text{Area} \times \text{Ports}$$

This calculation means that there are $239 \times 256 \times 256 = 15,663,104$ addresses that are available.

Depending on the fabric topology, the fabric addressing format of device differs.

In a fabric topology, devices have an addressing format type of *DDAA00*. For example, the address 020300 indicates that the device belongs to the switch with domain id 02. This switch is connected to port 03 and the ALPA address is 00, indicating that this device is not a loop fabric device, That is, it is a switched fabric device. For any switched fabric device, the ALPA ID is always *00*.

5.3.5 Loop address

An *NL_Port*, like an *N_Port*, has a 24-bit port address. If no switch connection exists, the two upper bytes of this port address are zeros (x'00 00') and referred to as a *private loop*. The devices on the loop have no connection with the outside world. If the loop is attached to a fabric and an *NL_Port* supports a fabric login, the upper 2 bytes are assigned a positive value by the switch. We call this mode a *public loop*.

Because fabric-capable *NL_Ports* are members of both a local loop and the greater fabric community, a 24-bit address is needed as an identifier in the network. In this case of public loop assignment, the value of the upper 2 bytes represents the loop identifier, and this ID is common to all *NL_Ports* on the same loop that logged in to the fabric.

In both public and private arbitrated loops, the last byte of the 24-bit port address refers to the *arbitrated loop physical address (AL_PA)*. The *AL_PA* is acquired during initialization of the loop and might, in the case of a fabric-capable loop device, be modified by the switch during login.

The total number of the *AL_PAs* available for arbitrated loop addressing is 127. This number is based on the requirements of 8b/10b running disparity between frames.

5.3.6 The b-type addressing modes

IBM b-type (the IBM OEM agreement with Brocade is referred to as *b-type*) has three different addressing modes: native mode, core PID mode, and shared area addressing mode.

Native mode is used in traditional switches which support a maximum of 16 ports. This number is used because in native mode, the fabric addressing format that is used is *DDIA00*. The area part of the fabric address always has a prefix of 1 and hence, it supports a port count from hexadecimal 10 to 1F (maximum of 16 ports).

Core PID mode is used to support a maximum of 256 ports per domain/switch. This number is used because in core PID mode, the area part of the fabric address supports addresses from hexadecimal 00 to FF (maximum of 256 ports). The fabric addressing format that is used for this mode is *DDAA00*.

Figure 5-21 explains the native and core PID modes with the example FC address of two devices.

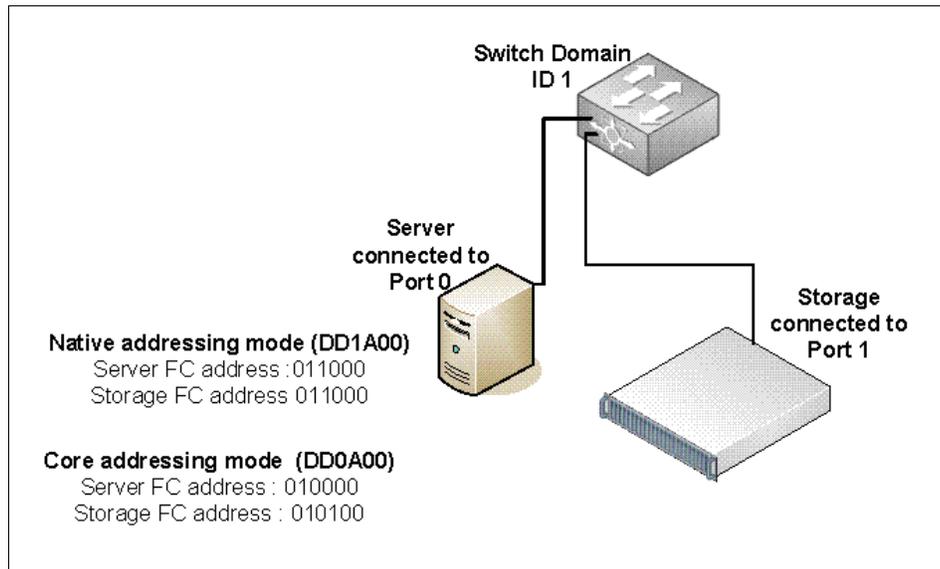


Figure 5-21 Native versus core addressing mode

Shared addressing mode is used when more than 256 ports are used in the same domain/switch. This mode is used in directors with high port density. The port addressing in these directors use the same area numbers for two ports by having the third byte of the FC address (node addresses) as 80 for higher port numbers. By having the Area ID used more than once, this mode enables more than 256 ports to exist in a single domain.

Figure 5-22 shows port 24. Port 25 shares the Area ID with port 32, and 33 of FC4-48 port.

Index	Slot	Port	Address	Media	Speed	State
168	3	24	01a800	--	N8	No_Module
169	3	25	01a900	--	N8	No_Module
<truncated output>						
288	3	32	01a880	--	N8	No_Module
289	3	33	01a980	--	N8	No_Module

Figure 5-22 Shared addressing mode

5.3.7 FICON address

FICON generates the 24-bit FC port address field in yet another way. When communication is required from the FICON channel port to the FICON CU port, the FICON channel (by using FC-SB-2 and FC-FS protocol information) provides the address of its port, the source port address identifier (S_ID), and the address of the CU port. This CU port address is the destination port address identifier (D_ID) when the communication is from the channel N_Port to the CU N_Port.

The Fibre Channel architecture does not specify how a server N_Port determines the destination port address of the storage device N_Port with which it requires communication.

This determination is node and N_Port implementation dependent. Basically, there are two ways that a server can determine the address of the N_Port with which it wants to communicate:

- ▶ The *discovery method*. The address is determined by knowing the WWN of the target Node N_Port, and then requesting a WWN for the N_Port port address from a Fibre Channel Fabric Service. This service is called the *fabric name server*.
- ▶ The *defined method*. The address is determined by the server (processor channel) N_Port having a known predefined port address of the storage device (CU) N_Port, with which it requires communication. This later approach is referred to as *port address definition*. It is the approach that is implemented for the FICON channel in the FICON native (FC) mode. This method is done by using either the z/OS hardware configuration definition (HCD) function or an input/output configuration program (IOCP). These functions are used to define a 1-byte switch port, which is a 1-byte FC area field of the 3-byte Fibre Channel N_Port port address.

The Fibre Channel architecture (*FC-FS*) uses a 24-bit FC port address (3 bytes) for each port in an FC switch. The switch port addresses in a FICON native (FC) mode are always assigned by the switch fabric.

For the FICON channel in FICON native (FC) mode, the *Accept (ACC ELS)* response to the Fabric Login (FLOGI) in a switched point-to-point topology provides the channel with the 24-bit N_Port address to which the channel is connected. This N_Port address is in the ACC destination address field (D_ID) of the FC-2 header.

The FICON CU port also performs a fabric login to obtain its 24-bit FC port address. Figure 5-23 shows how the FC-FS 24-bit FC port address identifier is divided into three fields.

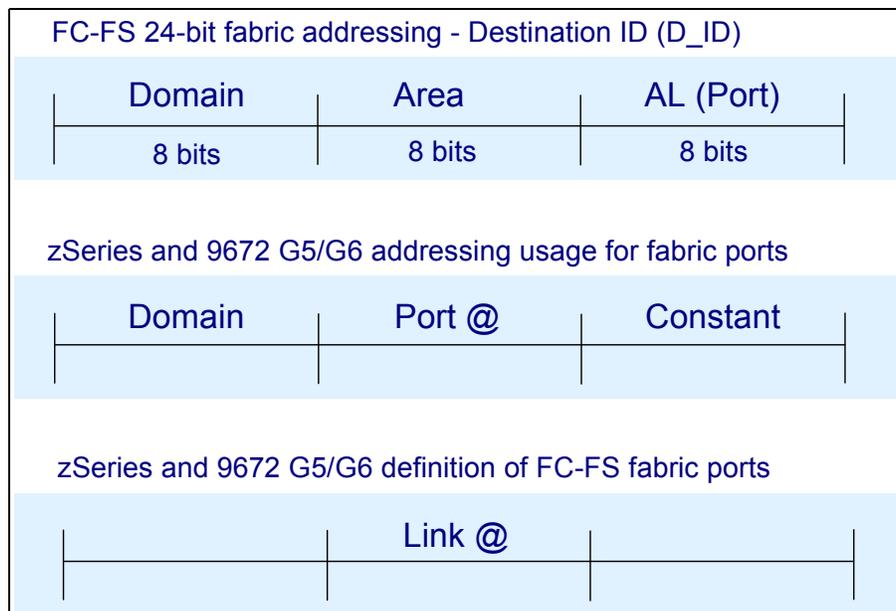


Figure 5-23 FICON port addressing

Figure 5-23 shows the FC-FS 24-bit port address and the definition of usage of that 24-bit address in a zSeries and 9672 G5/G6 environment. Only the 8 bits that make up the FC port address are defined for the zSeries and 9672 G5/G6 to access a FICON CU. The FICON channel in FICON native (FC) mode that works with a switched point-to-point FC topology (single switch), provides the other 2 bytes that make up the 3-byte FC port address of the CU to be accessed.

The zSeries and 9672 G5/G6 processors, when working with a switched point-to-point topology, require that the *Domain* and the *AL_Port* (*Arbitrated Loop*) field values be the same for all of the FC F_Ports in the switch. Only the area field value is different for each switch F_Port.

For the zSeries and 9672 G5/G6, the *area* field is referred to as the *port address field* of the F_Port. This field is just a 1-byte value, and when defining access to a CU that is attached to this port, by using the zSeries HCD or IOCP, the port address is referred to as the *link address*.

As shown in Figure 5-24, the 8 bits for the domain address and the 8-bit constant field are provided from the *Fabric Login* initialization result. Although, the 8 bits, 1 byte for the port address (1-byte link address) are provided from the zSeries or 9672 G5/G6 CU link definition (by using HCD and IOCP).

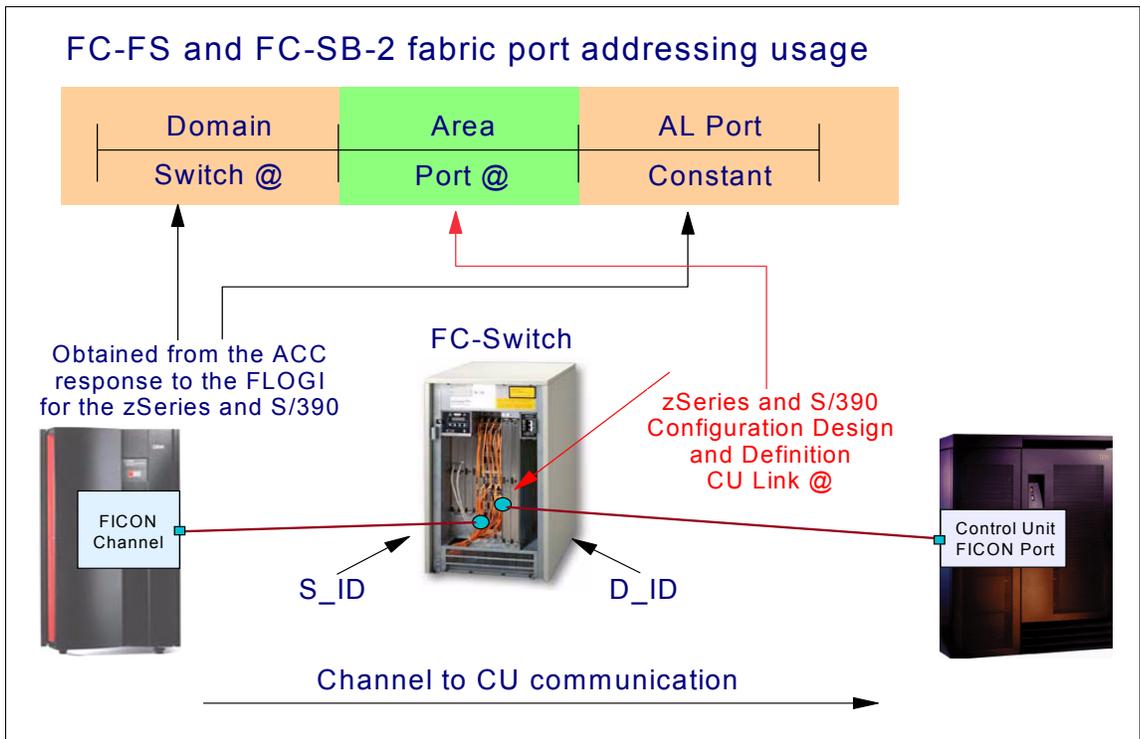


Figure 5-24 FICON single switch: Switched point-to-point link address

FICON address support for cascaded switches

The Fibre Channel architecture (*FC-FS*) uses a 24-bit FC port address of 3 bytes for each port in an FC switch. The switch port addresses in a FICON native (FC) mode are always assigned by the switch fabric.

For the FICON channel in FICON native (FC) mode, the Accept (ACC ELS) response to the Fabric Login (FLOGI) in a two-switch cascaded topology, provides the channel with the 24-bit N_Port address to which the channel is connected. This N_Port address is in the ACC destination address field (D_ID) of the FC-2 header.

The FICON CU port also performs a fabric login to obtain its 24-bit FC port address.

Figure 5-25 shows that the FC-FS 24-bit FC port address identifier is divided into three fields:

- ▶ Domain
- ▶ Area
- ▶ AL Port

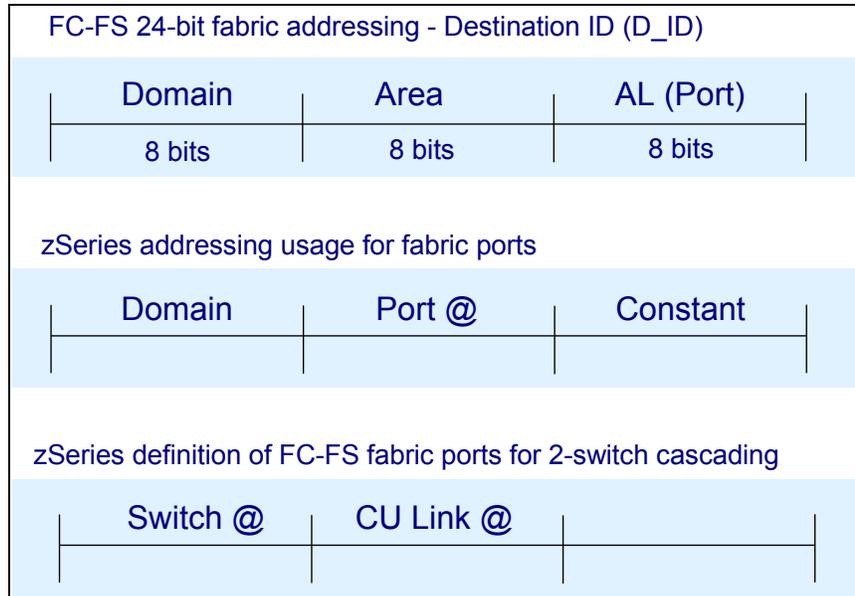


Figure 5-25 FICON addressing for cascaded directors

Figure 5-25 shows the FC-FS 24-bit port address and the definition usage of that 24-bit address in a zSeries environment. Here, the 16 bits that make up the FC port address must be defined for the zSeries to access a FICON CU in a cascaded environment. The FICON channel in the FICON native (FC) mode that works with a cascaded FC topology, two-switch, provides the remaining byte that makes up the full 3-byte FC port address of the CU to be accessed.

It is required that the Domain, switch @, AL_Port, and the Arbitrated Loop field values, be the same for all the FC F_Ports in the switch. Only the area field value is different for each switch F_Port.

The zSeries domain and area fields are referred to as the *port address field* of the F_Port. This field is a 2-byte value, and when defining access to a CU that is attached to this port (by using the zSeries HCD or IOCP), the port address is referred to as the *link address*.

As shown in Figure 5-26, the 8 bits for the constant field are provided from the Fabric Login initialization result. Although, the 16 bits for the port address and 2-byte link address are provided from the zSeries CU link definition by using HCD and IOCP.

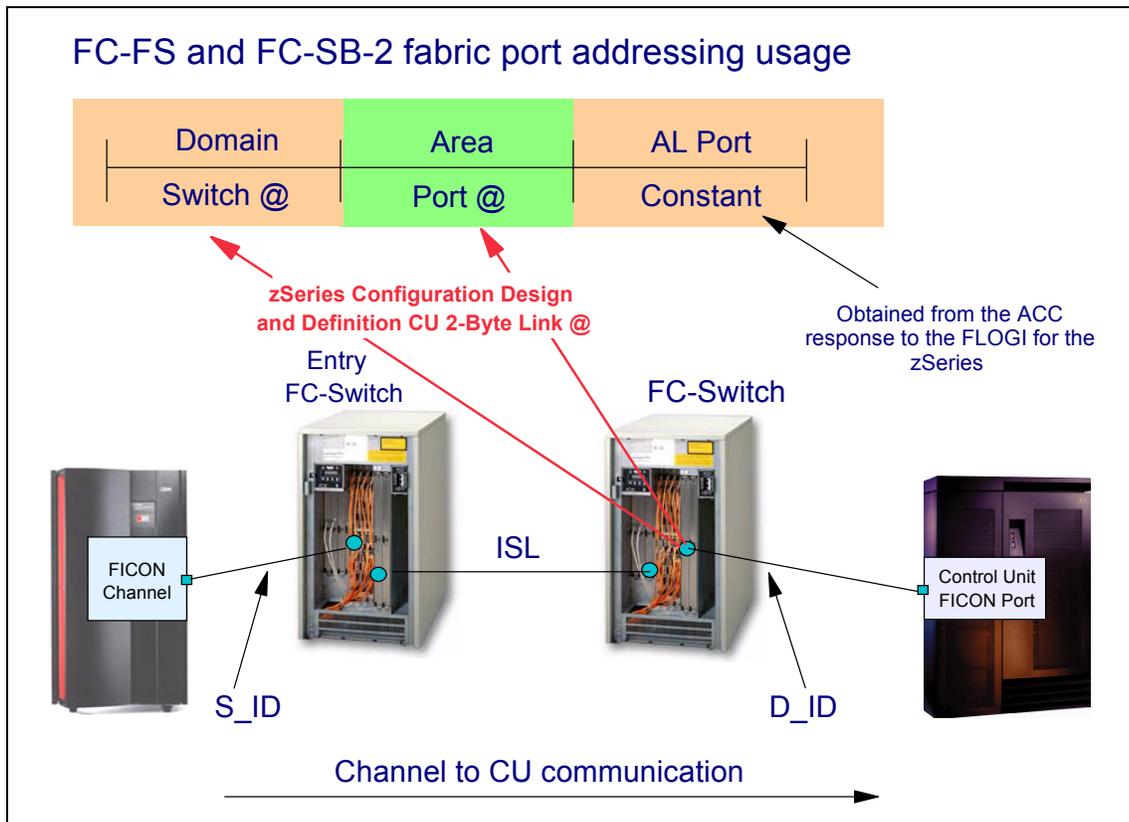


Figure 5-26 Two cascaded director FICON addressing

As a footnote, FCP connectivity is device-centric and is defined in the fabric by using the WWPN of the devices that are allowed to communicate. When an FCP device attaches to the fabric, it queries the name server for the list of devices that it is allowed to form connections with (that is, the *zoning information*). FICON devices do not query the name server for accessible devices because the allowable port and device relationships are defined in the host. Therefore, the zoning and name server information does not need to be retrieved.

5.4 Fibre Channel Arbitrated Loop protocols

To support the shared behavior of *Fibre Channel Arbitrated Loop (FC-AL)*, a number of loop-specific protocols are used. These protocols are used in the following ways:

- ▶ Initialize the loop and assign addresses.
- ▶ Arbitrate for access to the loop.
- ▶ Open a loop circuit with another port in the loop.
- ▶ Close a loop circuit when two ports complete their current use of the loop.
- ▶ Implement the access fairness mechanism to ensure that each port has an opportunity to access the loop.

We describe some of these topics in the sections that follow.

5.4.1 Fairness algorithm

The way that the *fairness algorithm* works is based around the IDLE ordered set, and the way that arbitration is carried out. To determine that the loop is not in use, an NL_Port waits until it sees an IDLE go by and it can arbitrate for the loop by sending an *arbitrate primitive signal (ARB)* ordered set. If a higher priority device arbitrates before the first NL_Port sees its own ARB come by, then it loses the arbitration. But, if it sees that its own ARB has gone all the way around the loop, then it has won arbitration. It can then open a communication to another NL_Port. When it has finished, it can close the connection and either arbitrate for the loop or send one or more IDLEs. If it complies with the fairness algorithm, it takes the option of sending IDLEs. That forces lower priority NL_Ports to successfully arbitrate for sending IDLEs, and that allows lower priority NL_Ports to successfully arbitrate for the loop. However, there is no rule that forces any device to operate the fairness algorithm.

5.4.2 Loop addressing

An *NL_Port*, like an *N_Port*, has a 24-bit port address. If no switch connection exists, the two upper bytes of this port address are zeros (x'00 00') and referred to as a *private loop*. The devices on the loop have no connection with the outside world. If the loop is attached to a fabric and the NL_Port supports a fabric login, the upper 2 bytes are assigned a positive value by the switch. We call this mode a *public loop*.

Because fabric-capable NL_Ports are members of both a local loop and a greater fabric community, a 24-bit address is needed as an identifier in the network. If there is a public loop assignment, the value of the upper 2 bytes represents the loop identifier. This identifier is common to all NL_Ports on the same loop that logged in to the fabric.

In both public and private Arbitrated Loops, the last byte of the 24-bit port address refers to the *Arbitrated Loop physical address (AL_PA)*. The AL_PA is acquired during initialization of the loop and might, in the case of fabric-capable loop devices, be modified by the switch during login.

The total number of the AL_PAs available for Arbitrated Loop addressing is 127, which is based on the requirements of 8b/10b running disparity between frames.

As a frame terminates with an *end-of-frame (EOF)* character, this forces the current running disparity to be *negative*. In the Fibre Channel standard, each transmission word between the end of one frame and the beginning of another frame also leaves the running disparity negative. If all 256 possible 8-bit bytes are sent to the 8b/10b encoder, 134 emerge with neutral disparity characters. Of these 134, seven are reserved for use by Fibre Channel. The 127 neutral disparity characters that left are assigned as AL_PAs. Stated another way, the 127 AL_PA limit is the maximum number, minus the reserved values, of neutral disparity addresses that can be assigned for use by the loop. This number does not imply that we recommend this amount, or load, but only that it is possible.

Arbitrated Loop assigns priority to AL_PAs, based on numeric value. The lower the numeric value, the higher the priority is.

It is the Arbitrated Loop initialization that ensures that each attached device is assigned a unique AL_PA. The possibility for address conflicts arises only when two separated loops are joined without initialization.

IBM System z: IBM System z9® and zSeries servers do not support the arbitrated loop topology.

5.5 Fibre Channel port initialization and fabric services

You learned that there are different port types. At a high level, port initialization starts with *port type detection*. Then, there is *speed and active state detection* where the speed is negotiated according to the device that is connected, and then the *port initializes* to an active state. In this active state, every F port/FL port that has an N port/NL port that is connected, as well as the Extended Link Service (ELS) and the Fibre Channel Common Transport (FCCT) protocol, are used for further *switch-port to node-port* communication. Only after this initialization completes, can data flow happen. We now review the services that are responsible for the port initialization in a fabric switch.

The following three login types are available for fabric devices:

- ▶ Fabric login (FLOGI)
- ▶ Port login (PLOGI)
- ▶ Process login (PRLI)

Aside from these login types, we also describe the roles of other fabric services such as the fabric controller, management server, and time server.

5.5.1 Fabric login (FLOGI)

After the fabric-capable Fibre Channel device is attached to a fabric switch, it carries out a *fabric login (FLOGI)*.

Similar to port login, FLOGI is an extended link service command that sets up a session between two participants. A session is created between an N_Port or NL_Port and the switch. An N_Port sends a FLOGI frame that contains its Node Name, its N_Port Name, and service parameters to a well-known address of *0xFFFFFE*.

The switch accepts the login and returns an *accept (ACC) frame* to the sender. If some of the service parameters that are requested by the N_Port or NL_Port are not supported, the switch sets the appropriate bits in the ACC frame to indicate this status.

NL_Ports derive their AL_PA during the *loop initialization process (LIP)*. The switch then decides if it accepts this AL_PA, if it does not conflict with any previously assigned AL_PA on the loop. If not, a new AL_PA is assigned to the NL_Port, which then causes the start of another LIP.

Figure 5-27 shows nodes that are performing FLOGI.

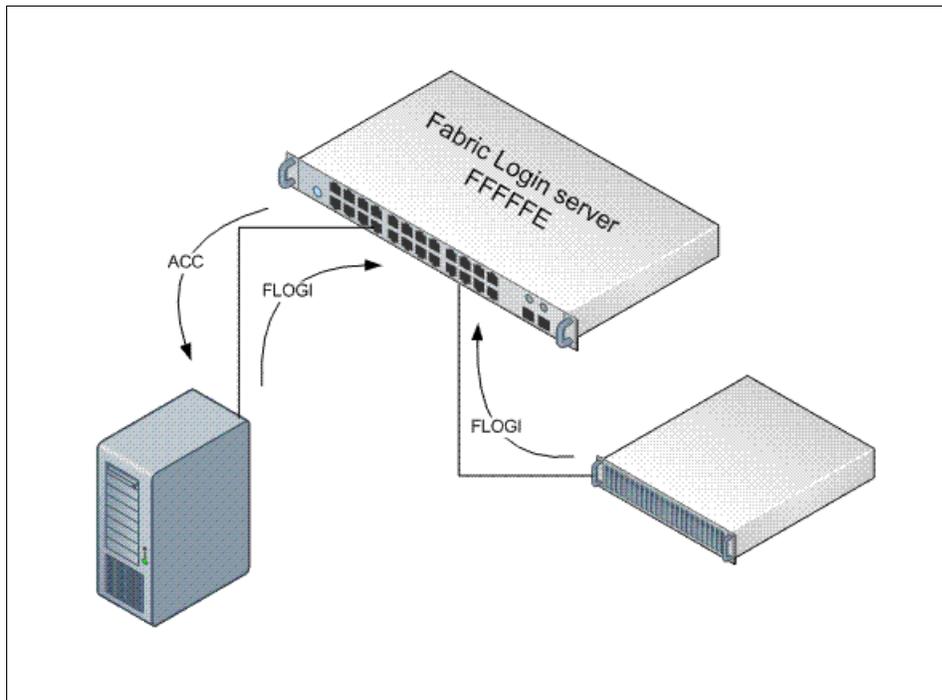


Figure 5-27 FLOGI of nodes

5.5.2 Port login (PLOGI)

Port login (PLOGI), is used to establish a session between two N_Ports, and is necessary before any upper level commands or operations can be performed. During port login, two N_Ports (devices) swap service parameters and make themselves known to each other by performing a port login to the well-known address of $0xFFFFFC$. The device might register values for all or some of its objects, but the most useful include the following objects:

- ▶ 24-bit port address
- ▶ 64-bit port name
- ▶ 64-bit node name
- ▶ Buffer-to-buffer credit capability
- ▶ Maximum frame size
- ▶ Class-of-service parameters
- ▶ FC-4 protocols that are supported
- ▶ Port type

When the communication parameters and identities of other devices are discovered, they are able to establish logical sessions between devices (initiator and targets).

Figure 5-28 shows the PLOGI of a host.

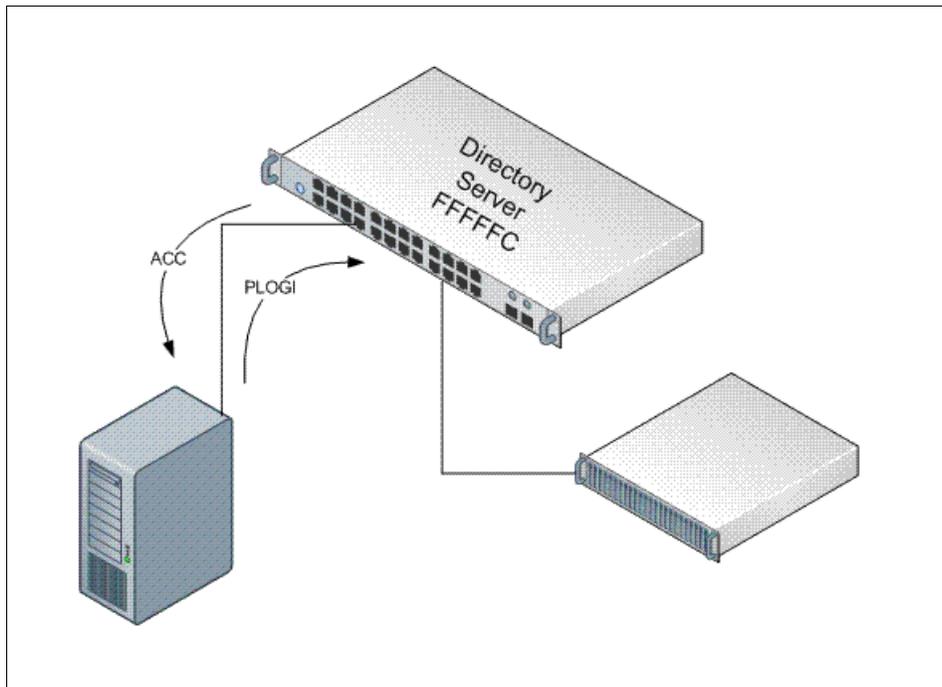


Figure 5-28 Node PLOGI to probe other nodes in the fabric

5.5.3 Process login (PRLI)

Process login (PRLI) is used to set up the environment between related processes on an originating N_Port and a responding N_Port. A group of related processes is collectively known as an *image pair*. The processes that are involved can be system processes and system images, such as mainframe logical partitions, control unit images, and FC-4 processes. Use of process login is optional from the perspective of the Fibre Channel FC-2 layer. However, it might be required by a specific upper-level protocol, as in the case of SCSI-FCP mapping.

Figure 5-29 indicates the PRLI from server to storage.

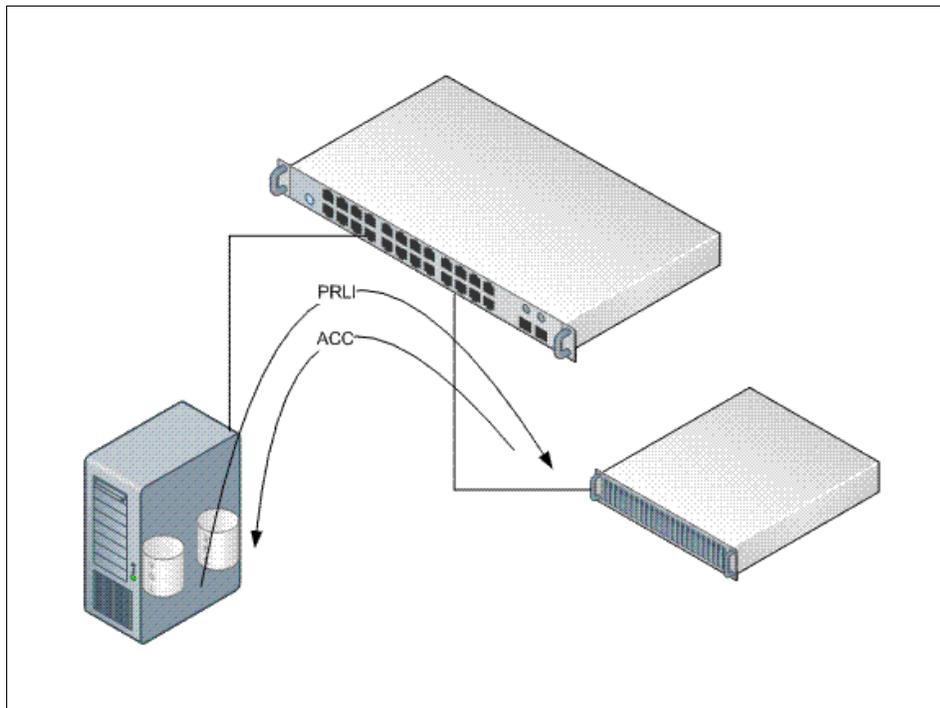


Figure 5-29 PRLI request from the initiator to the target

5.6 Fabric services

There is a set of services that are available to all devices that participate in a Fibre Channel fabric. They are known as *fabric services*, and include the following functions:

- ▶ Management services
- ▶ Time services
- ▶ Simple name server
- ▶ Login services
- ▶ Registered State Change Notification (RSCN)

These services are implemented by switches and directors that participate in the SAN. Generally speaking, the services are distributed across all of the devices, and a node can use whichever switching device to which it is connected.

All these services are addressed by FC-2 frames and are accessed by so called *well-known addresses*.

5.6.1 Management server

Management server is an in-band fabric service that allows data to be passed from the device to management platforms. This data includes such information as the topology of the SAN. A critical feature of this service is that it allows management software access to the *simple name server (SNS)*, bypassing any potential block that is caused by zoning. This means that a management suite can have a view of the entire SAN. The well-known port that is used for the management server is `0xFFFFFA`.

5.6.2 Time server

The *time service* or *time server* is provided to serve time information that is sufficient for managing expiration time. This service is provided at the well-known address identifier, *0xFFFFFB*.

The functional model of the time server consists of primarily two entities:

- ▶ Time Service Application. This is the entity that represents a user that is accessing the time service.
- ▶ Time Server. This is the entity that provides the time information through the time service.

There might be more than one distributed time server instance within the Fibre Channel network. However, from a perspective of the user, the time service seems to come from the entity that is accessible at the time service well-known address identifier. If the time service is distributed, it is transparent to the application.

5.6.3 Simple name server

Fabric switches implement a concept that is known as the *simple name server (SNS)*. All switches in the fabric keep the SNS updated, and are therefore aware of all devices in the SNS. After a node successfully logs in to the fabric, it performs a PLOGI into the well-known address of *0xFFFFFC*. This action allows it to register itself and pass on critical information such as class-of-service parameters, its WWN and address, and the upper layer protocols that it can support.

5.6.4 Fabric login server

To do a fabric login, a node communicates with the fabric login server at the well-known address *0xFFFFFE*.

5.6.5 Registered state change notification service

The service, *registered state change notification (RSCN)*, is critical because it propagates information about a change in the state of one node to all other nodes in the fabric. This communication means that in the event of, for example, a node being shut down, that the other nodes on the SAN are informed and can take the necessary steps to stop communicating with it. This notification prevents the other nodes from trying to communicate with the node that is shut down, timing out, and trying again.

The nodes register to the *fabric controller* with a *state change registration (SCR)* frame. The fabric controller, which maintains the fabric state with all of the registered device details, alerts registered devices with an RSCN. This alert is sent whenever there is any device that is added or removed, there is a zone change, a switch IP, or name change, and so on. The fabric controller has a well-known address of *0xFFFFFD*.

5.7 Routing mechanisms

A complex fabric can be made of interconnected switches and directors, even spanning a LAN or WAN connection. The challenge is to route the traffic with a minimum of performance overhead and latency, and to prevent an out-of-order delivery of frames, all while maintaining reliability. The following subsections describe some of the mechanisms.

5.7.1 Spanning tree

If there is a failure, it is important to consider having an alternative path that is available between the source and destination. This allows data to still reach its destination. However, having different paths available might lead to the delivery of frames that are out of order. This order might happen because of a frame taking a different path and arriving earlier than one of its predecessors.

A solution, which can be incorporated into the meshed fabric, is called a *spanning tree* and is an IEEE 802.1 standard. This concept means that switches stay on certain paths because the spanning tree protocol blocks certain paths to produce a simply connected active topology. Then, the shortest path in terms of hops is used to deliver the frames, and only one path is active at a time. This means that all associated frames go over the same path to the destination. The paths that are blocked can be held in reserve and used only if, for example, a primary path fails.

The most commonly used path selection protocol is *fabric shortest path first (FSPF)*. This type of path selection is usually performed at the time of booting, and no configuration is needed. All paths are established at the start time, and reconfiguration takes place only if no ISLs are broken or added.

5.7.2 Fabric shortest path first

According to the FC-SW-2 standard, *fabric shortest path first (FSPF)* is a link state path selection protocol. The concepts that are used in FSPF were first proposed by Brocade, and are incorporated into the FC-SW-2 standard. Since then, it has been adopted by most, if not all, manufacturers.

What fabric shortest path first is

FSPF tracks the links on all switches in the fabric and associates a cost with each link. The cost is always calculated as being directly proportional to the number of hops. The protocol computes paths from a switch to all other switches in the fabric by adding the cost of all links that are traversed by the path, and choosing the path that minimizes the cost.

How fabric shortest path first works

The collection of link states (including cost) of all switches in a fabric constitutes the topology database (or link state database). The topology database is kept in all switches in the fabric, and they are maintained and synchronized to each other. There is an initial database synchronization, and an update mechanism. The initial database synchronization is used when a switch is initialized, or when an ISL comes up. The update mechanism is used when there is a link state change. This mechanism ensures consistency among all switches in the fabric.

How fabric shortest path first helps

In the situation where there are multiple routes, FSPF ensures that the route that is used is the one with the lowest number of hops. If all of the hops have the same latency, operate at the same speed, and have no congestion, then FSPF ensures that the frames get to their destinations by the fastest route.

5.8 Zoning

Zoning allows for finer segmentation of the switched fabric. Zoning can be used to instigate a barrier between different environments. Only the members of the same zone can communicate within that zone; all other attempts from outside are rejected.

For example, it might be desirable to separate a Microsoft Windows NT environment from a UNIX environment. This is useful because of the manner in which Windows attempts to claim all available storage for itself. Because not all storage devices can protect their resources from any host that is seeking available resources, it makes good business sense to protect the environment in another manner. We show an example of zoning in Figure 5-30 where we separate AIX from NT and created Zone 1 and Zone 2. This diagram also shows how a device can be in more than one zone.

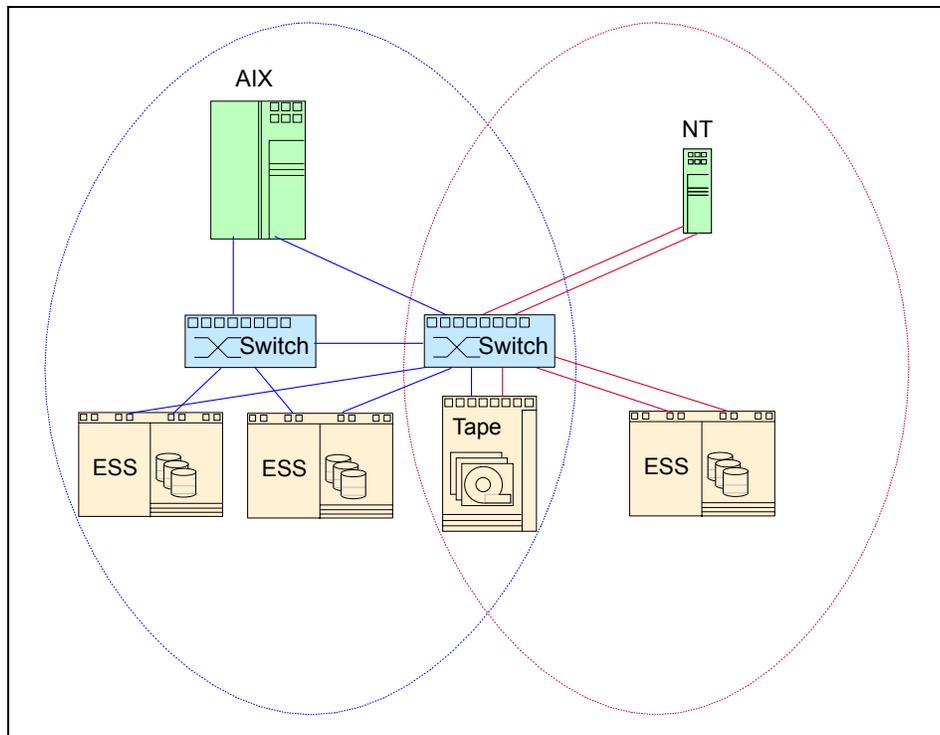


Figure 5-30 Zoning

Looking at zoning in this way, it can also be considered as a security feature, and not just for separating environments. Zoning can also be used for test and maintenance purposes. For example, not many enterprises mix their test and maintenance environments with their production environment. Within a fabric, you can easily separate your test environment from your production bandwidth allocation on the same fabric by using zoning.

An example of zoning is shown in Figure 5-31:

- ▶ Server A and Storage A can communicate with each other.
- ▶ Server B and Storage B can communicate with each other.
- ▶ Server A cannot communicate with Storage B.
- ▶ Server B cannot communicate with Storage A.
- ▶ Both servers and both storage devices can communicate with the tape.

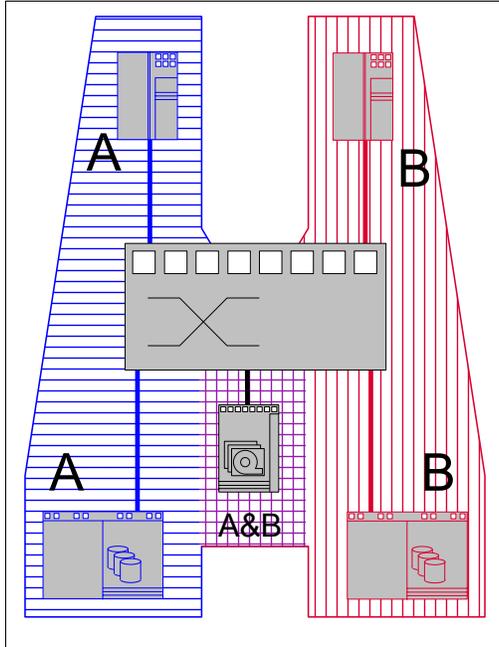


Figure 5-31 An example of zoning

Zoning also introduces the flexibility to manage a switched fabric to meet different user groups objectives.

Zoning can be implemented in two ways:

- ▶ Hardware zoning
- ▶ Software zoning

These forms of zoning are different, but are not necessarily mutually exclusive. Depending upon the particular manufacturer of the SAN hardware, it is possible for hardware zones and software zones to overlap. While this ability adds to the flexibility, it can make the solution complicated, increasing the need for good management software and documentation of the SAN.

5.8.1 Hardware zoning

Hardware zoning is based on the physical fabric port number. The members of a zone are physical ports on the fabric switch. It can be implemented in the following configurations:

- ▶ One-to-one
- ▶ One-to-many
- ▶ Many-to-many

Figure 5-32 on page 118 shows an example of zoning that is based on the switch port numbers.

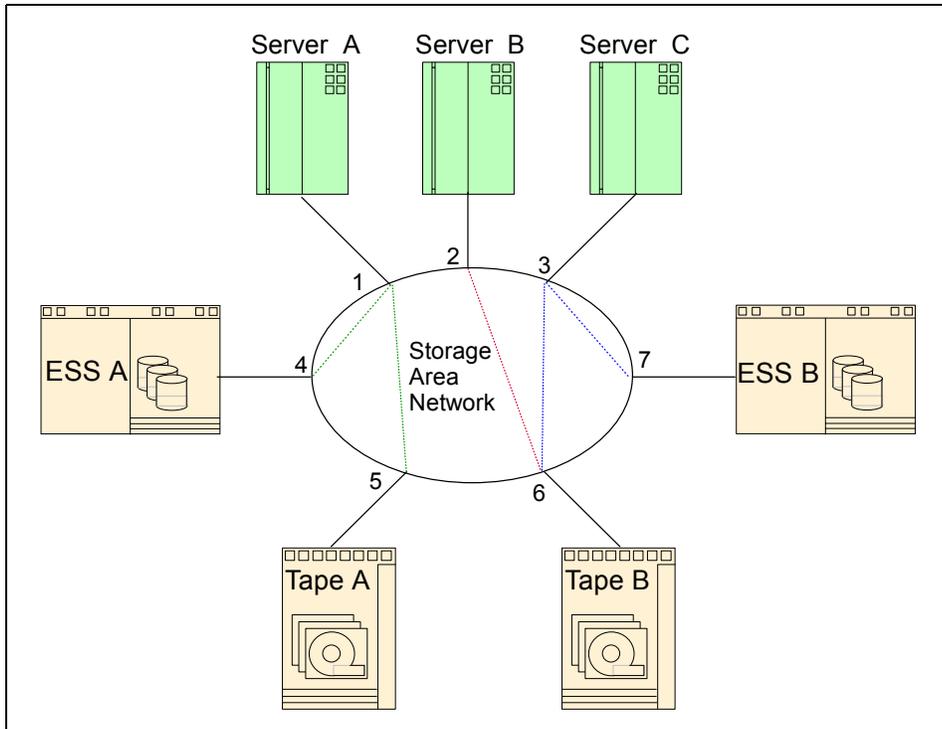


Figure 5-32 Zoning that is based on the switch port number

In Figure 5-32, zoning is based on the switch port number:

- ▶ Server A is restricted to see only storage devices that are zoned to port 1: ports 4 and 5.
- ▶ Server B is also zoned so that it can see only from port 2 to port 6.
- ▶ Server C is zoned so that it can see both ports 6 and 7, even though port 6 is also a member of another zone.
- ▶ A single port can also belong to multiple zones.

We show an example of hardware zoning in Figure 5-33 on page 119. This example illustrates another way of considering the hardware zoning as an array of connections.

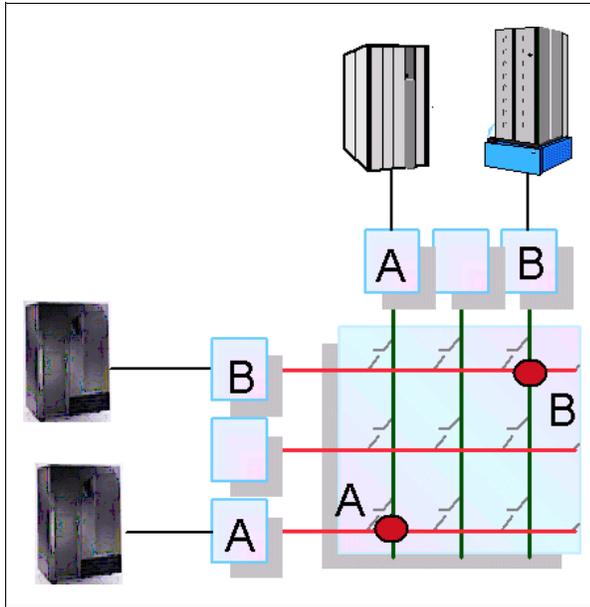


Figure 5-33 Hardware zoning

In Figure 5-33, device A can access only storage device A through connection A. Device B can access only storage device B through connection B.

In a hardware-enforced zone, *switch hardware*, usually at the *application-specific integrated circuit (ASIC)* level, ensures that there is no data that is transferred between unauthorized zone members. However, devices can transfer data between ports within the same zone. Consequently, hard zoning provides the highest level of security. The availability of hardware-enforced zoning and the methods to create hardware-enforced zones depends on the switch hardware.

One of the disadvantages of hardware zoning is that devices must be connected to a specific port, and the whole zoning configuration can become unusable when the device is connected to a different port. In cases where the device connections are not permanent, the use of software zoning is likely to simplify your configuration.

The advantage of hardware zoning is that it can be implemented into a routing engine by filtering. As a result, this type of zoning has a low affect on the performance of the routing process.

If possible, the designer can include some unused ports in a hardware zone. Therefore, in the event of a particular port failing, maybe caused by a gigabit interface converter (GBIC) or transceiver problem, the cable can be moved to a different port in the same zone. This means that the zone would not need to be reconfigured.

5.8.2 Software zoning

Software zoning is implemented by the fabric operating systems within the fabric switches. They are almost always implemented by a combination of the name server and the Fibre Channel Protocol. When a port contacts the name server, the name server replies only with information about the ports in the same zone as the requesting port. A *soft zone*, or *software zone*, is not enforced by hardware. What this means is that if a frame is incorrectly delivered (addressed) to a port that it was not intended, then it is delivered to that port. This type of zoning is in contrast to hard zones.

When using software zoning, the members of the zone can be defined by using their WWNs:

- ▶ Node WWN
- ▶ Port WWN

Usually, zoning software also allows you to create symbolic names for the zone members and for the zones themselves. Dealing with the symbolic name or aliases for a device is often easier than trying to use the WWN address.

The number of members possible in a zone is limited only by the amount of memory in the fabric switch. A member can belong to multiple zones. You can define multiple sets of zones for the fabric, but only one set can be active at any time. You can activate another zone set any time that you want, without needing to power down the switch.

With software zoning, there is no need to worry about the physical connections to the switch. If you use WWNs for the zone members, even when a device is connected to another physical port, it remains in the same zoning definition. This scenario is because the WWN of the device remains the same. The zone follows the WWN.

Important: As stated, there is no need to worry about your physical connections to the switch when you are using software zoning. However, this statement does not mean that if you unplug a device, such as a disk subsystem, and plug it into another switch port, that your host is still able to communicate with your disks. That is, you cannot assume that your host is still able to communicate until you either reboot or unload, and load your operating system device definitions. This is true even if the device remains a member of that particular zone. The connection depends on the components that you use in your environment, like the operating system and multipath software.

Figure 5-34 on page 121 shows an example of WWN-based zoning. In this example, symbolic names are defined for each WWN in the SAN to implement the same zoning requirements, as shown in Figure 5-32 on page 118 for port zoning:

- ▶ Zone_1 contains the aliases alex, ben, and sam, and is restricted to only these devices.
- ▶ Zone_2 contains the aliases robyn and ellen, and is restricted to only these devices.
- ▶ Zone_3 contains the aliases matthew, max, and ellen, and is restricted to only these devices.

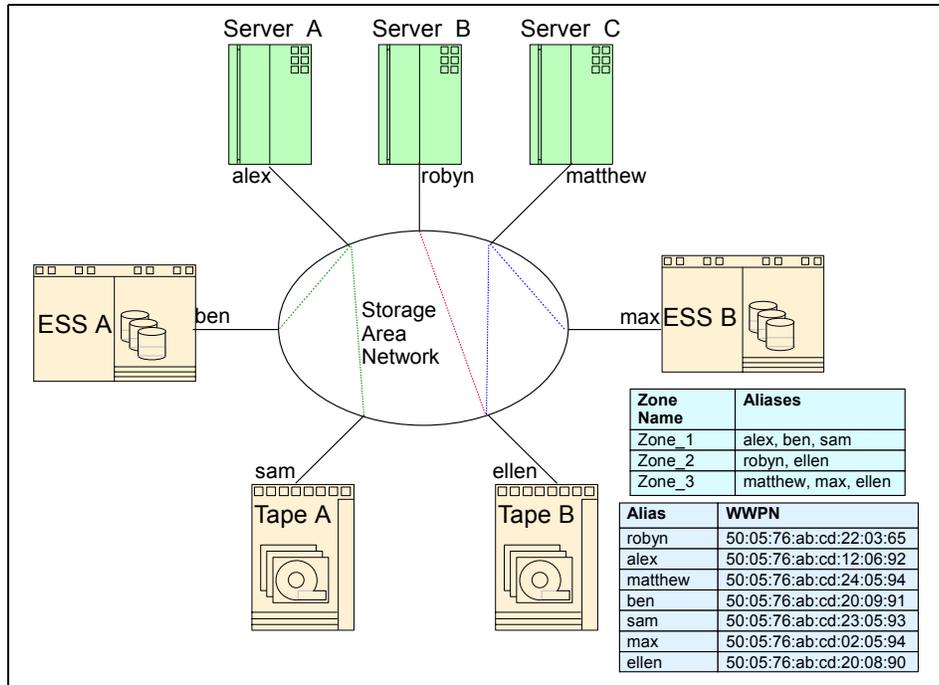


Figure 5-34 Zoning that is based on the WWNs of the devices

There are some potential security leaks with software zoning:

- ▶ When a specific host logs in to the fabric and asks for available storage devices, the simple name server (SNS) looks in the software zoning table to see which devices are allowable. The host sees only the storage devices that are defined in the software zoning table. But, the host can also make a direct connection to the storage device, by using device discovery, without asking SNS for the information.
- ▶ It is possible for a device to define the WWN that it uses, rather than using the one designated by the manufacturer of the HBA. This concept is known as *WWN spoofing*. An unknown server can masquerade as a trusted server and thus gain access to data on a particular storage device. Some fabric operating systems allow the fabric administrator to prevent this risk by allowing the WWN to be tied to a particular port.
- ▶ Any device that does any form of probing for WWNs is able to discover devices and talk to them. A simple analogy is that of an unlisted telephone number. Although the telephone number is not publicly available, there is nothing to stop a person from dialing that number, whether by design or accident. The same holds true for WWN. There are devices that randomly probe for WWNs to see if they can start a conversation with them.

A number of switch vendors offer hardware-enforced WWN zoning, which can prevent this security exposure. *Hardware-enforced zoning* uses hardware mechanisms to restrict access rather than relying on the servers to follow the Fibre Channel protocols.

Software zoning: When a device logs in to a software-enforced zone, it queries the name server for devices within the fabric. If zoning is in effect, only the devices in the same zone or zones are returned. Other devices are hidden from the name server query reply. When you use software-enforced zones, the switch does not control data transfer and there is no guarantee of data being transferred from unauthorized zone members. Use software zoning where flexibility and security are ensured by the cooperating hosts.

Frame filtering

Zoning is a fabric management service that can be used to create logical subsets of devices within a SAN. This service can also enable partitioning of resources for management and access control purposes. *Frame filtering* is another feature that enables devices to provide zoning functions with finer granularity. Frame filtering can be used to set up port-level zoning, WWN zoning, device-level zoning, protocol-level zoning, and LUN-level zoning. Frame filtering is commonly performed by an *application-specific integrated circuit (ASIC)*. This configuration has the result that, after the filter is set up, the complicated function of zoning and filtering can be achieved at wire speed.

5.8.3 Logical unit number masking

The term *logical unit number (LUN)* was originally used to represent the entity within a SCSI target which runs I/Os. A single SCSI device usually has only a single LUN, but some devices, such as tape libraries, might have more than one LUN.

With storage arrays, the array makes virtual disks available to the servers. These virtual disks are identified by LUNs.

It is possible for more than one host to see the same storage device or LUN. This is potentially a problem, both from a practical and a security perspective. Another approach to securing storage devices from hosts that want to take over already assigned resources is *LUN masking*. Every storage device offers its resources to the hosts with LUNs.

For example, each partition in the storage server has its own LUN. If the host server wants to access the storage, it must request access to the LUN in the storage device. The purpose of LUN masking is to control access to the LUNs. The storage device itself accepts or rejects access requests from different hosts.

The user defines which hosts can access which LUN with the storage device control program. Whenever the host accesses a particular LUN, the storage device checks its access list for that LUN, and it allows or disallows access to the LUN.



Storage area network as a service for cloud computing

While information can be your greatest asset, it can also be your greatest challenge as you struggle to keep up with explosive data growth. More data means more storage and more pressure to install another rack into the data center.

Cloud computing offers a new way of solution provisioning with significant cost savings and high reliability.

6.1 What is a cloud?

Cloud computing is a model for enabling ubiquitous, convenient, on-demand network access to a shared pool of configurable computing resources (for example: networks, servers, storage, applications, and services). These resources can be rapidly provisioned and released with minimal management effort or service provider interaction. Figure 6-1 shows an overview of cloud computing.

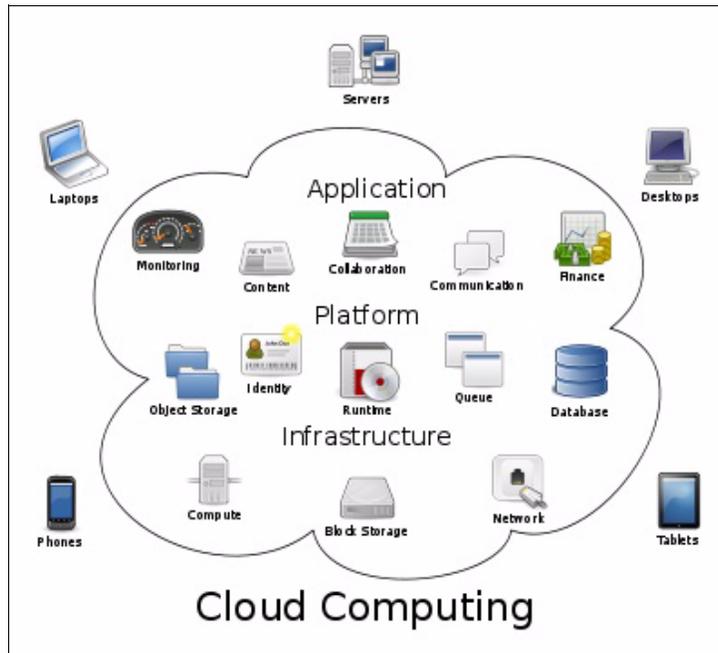


Figure 6-1 Cloud computing overview

Cloud computing provides computation, software, data access, and storage services that do not require user knowledge of the physical location and configuration of the system that delivers the services. Parallels to this concept can be drawn with the electricity grid, wherein users use power without needing to understand the component devices or infrastructure that is required to provide the service.

Cloud computing describes a new consumption and delivery model for IT services, and it typically involves provisioning of dynamically scalable and virtualized resources. The cloud introduces three key concepts: cost savings, service reliability, and infrastructure flexibility.

To cater to the increasing, on-demand needs of business, IT services and infrastructures are moving rapidly towards a flexible utility and consumer model by adopting new technologies.

One of these technologies is *virtualization*. Cloud computing is an example of a virtual, flexible delivery model. Inspired by consumer Internet services, cloud computing puts the user in the “driver’s seat”; that is, users can use Internet offerings and services by using this self-service, on-demand model.

Cloud computing has the potential to make an enormous affect to your business by providing the following benefits:

- ▶ Reducing IT labor costs for configuration, operations, management, and monitoring
- ▶ Improving capital utilization and significantly reducing license costs
- ▶ Reducing provisioning cycle times from weeks to minutes

- ▶ Improving quality and eliminating many software defects
- ▶ Reducing user IT support costs

From a technical perspective, cloud computing enables these capabilities, among others:

- ▶ Abstraction of resources
- ▶ Dynamic right-sizing
- ▶ Rapid provisioning

6.1.1 Private and public cloud

A cloud can be private or public. A *public cloud* sells services to anyone on the Internet. A private cloud is a proprietary network or a data center that supplies hosted services to a limited number of people. When a service provider uses public cloud resources to create their private cloud, the result is called a *virtual private cloud*. Whether private or public, the goal of cloud computing is to provide easy, scalable access to computing resources and IT services. A cloud has four basic components, as shown in Figure 6-2.

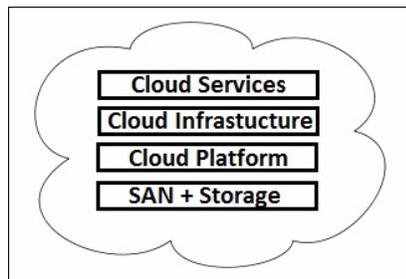


Figure 6-2 Cloud computing components

6.1.2 Cloud computing components

We describe the cloud computing components, or layers, in our model.

Cloud Services

This layer is the service that is delivered to the client, it can be an application, a desktop, a server, or disk storage space. The client does not need to know where or how their service is running, they just use it.

Cloud Infrastructure

This layer can be difficult to visualize depending on the final delivered service. If the final service is a chat application, the cloud infrastructure is the servers on which the chat application is running. In the other case, if the final service is a virtualized server, the cloud infrastructure is all the other servers that are required to provide “a server” as a service to the client. Examples of these types of servers include: domain name server (DNS), security services, and management.

Cloud Platform

This layer consists of the selected platform to build the cloud. There are many vendors like IBM Smart Business Storage Cloud, VMware vSphere, Microsoft Hyper V, and Citrix Xen Server, which are well known cloud solutions in the market.

SAN + Storage

This layer is where information flows and lives. Without it, nothing can happen. Depending on the cloud design, the storage can be any of the previously presented solutions. Examples include: Direct-attached storage (DAS), network-attached storage (NAS), Internet Small Computer System Interface (iSCSI), storage area network (SAN), or Fibre Channel over Ethernet (FCoE). For the purpose of this book, we describe Fibre Channel or FCoE for networking and compatible storage devices.

6.1.3 Cloud computing models

While cloud computing is still a relatively new technology, there are generally three cloud service models, each with a unique focus. The American National Institute of Standards and Technology (NIST) defined the following cloud service models:

- ▶ *Software as a service (SaaS)*: This capability that is provided to the consumer is to use the applications that a provider runs on a cloud infrastructure. The applications are accessible from various client devices through a thin client interface, such as a web browser (for example, web-based email). The consumer does not manage or control the underlying cloud infrastructure, including the network, servers, operating systems, storage, or even individual application capabilities. One possible exception is for the consumer to continue the control of limited user-specific application configuration settings.
- ▶ *Platform as a service (PaaS)*: This capability that is provided to the consumer is to deploy consumer-created or acquired applications onto the cloud infrastructure. Examples of these types of applications include those that are created by using programming languages and tools that are supported by the provider. The consumer does not manage or control the underlying cloud infrastructure, including the network, servers, operating systems, or storage. But, the consumer has control over the deployed applications and possibly application-hosting environment configurations.
- ▶ *Infrastructure as a service (IaaS)*: This capability that is provided to the consumer is to provision processing, storage, networks, and other fundamental computing resources where the consumer is able to deploy and run arbitrary software. These resources can include operating systems and applications. The consumer does not manage or control the underlying cloud infrastructure, but has control over operating systems, storage, and deployed applications. The consumer might also have limited control of select networking components (for example, hosts).

Figure 6-3 shows these cloud models.

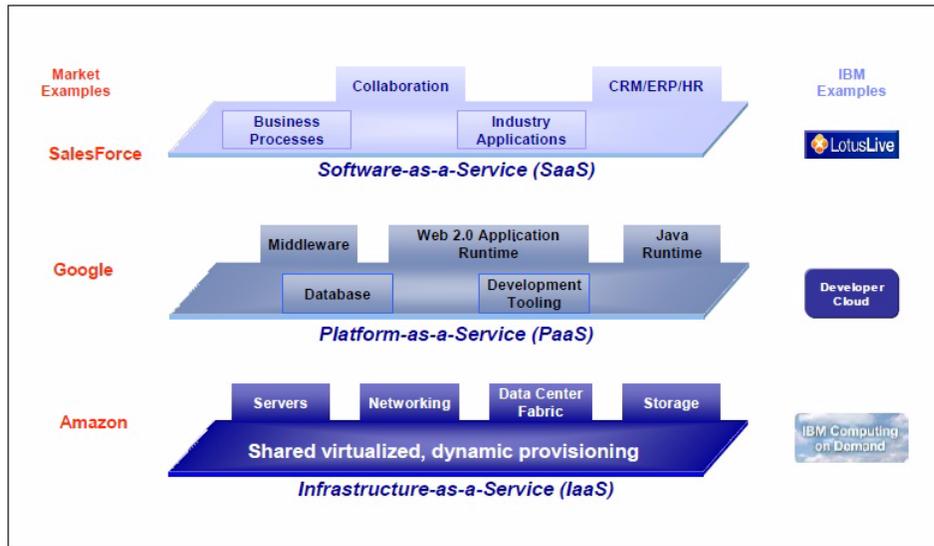


Figure 6-3 Examples of SaaS, PaaS, and IaaS services

In addition, NIST also defined the following models for deploying cloud services:

- ▶ **Private cloud:** The cloud infrastructure is owned or leased by a single organization and is operated solely for that organization.
- ▶ **Community cloud:** The cloud infrastructure is shared by several organizations and supports a specific community that shares (for example, mission, security requirements, policy, and compliance considerations).
- ▶ **Public cloud:** The cloud infrastructure is owned by an organization that sells cloud services to the general public or to a large industry group.
- ▶ **Hybrid cloud:** The cloud infrastructure is a composition of two or more clouds (internal, community, or public) that remain unique entities. However, these entities are bound together by standardized or proprietary technology that enables data and application portability (for example, cloud bursting).

Figure 6-4 shows cloud computing deployment models.

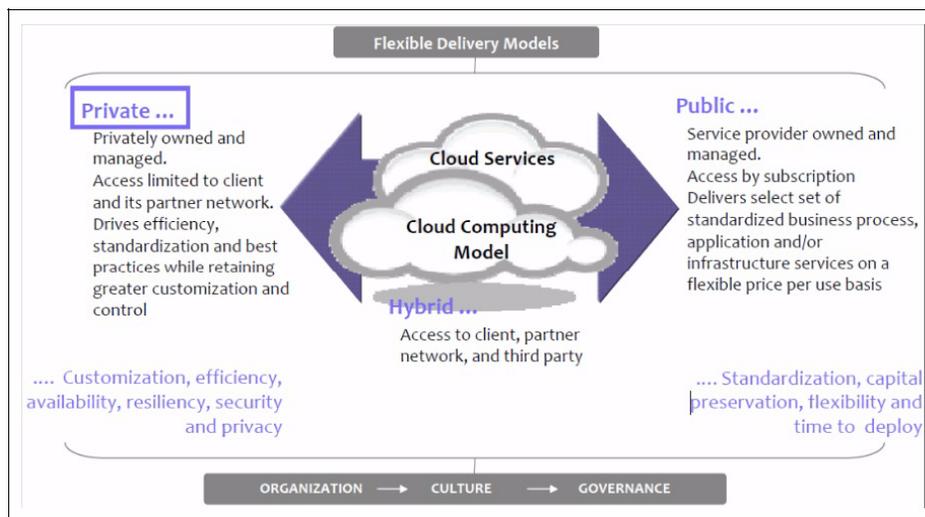


Figure 6-4 Cloud computing deployment models

From a storage perspective, IBM clients, based on their business requirements, can choose to adopt either a public or private storage cloud. The following definitions describe these types of storage clouds:

- ▶ *Public storage cloud:* This is designed for clients who do not want to own, manage, or maintain the storage environment, thus reducing their capital and operational expenditures for storage. IBM dictates the choice of technology and cloud location, shared infrastructure with variable monthly charges, dynamic physical capacity at the client level, and security measures to isolate client data. The public storage cloud allows for variable billing options and shared tenancy of the storage cloud, giving clients the flexibility to manage the use and growth of their storage needs. This type is the industry-standard view of a storage cloud offering and is comparable to storage cloud offerings by other vendors.
- ▶ *Private storage cloud:* With a private storage cloud, clients have the choice of technology and location on a dedicated infrastructure with fixed monthly charges and a physical capacity that is manageable by the client. Each application can use dynamic capacity by sharing the cloud storage among multiple applications.

Private storage cloud solution technology and services from IBM address multiple areas of functionality. For more information, see this website:

<http://www.ibm.com/cloud-computing/us/en/>

6.2 Virtualization and the cloud

When people talk about virtualization, they are usually referring to *server virtualization*, which means partitioning one physical server into several virtual servers, or machines. Each virtual machine can interact independently with other devices, applications, data, and users as though it were a separate physical resource.

Different virtual machines can run different operating systems and multiple applications while they are sharing the resources of a single physical computer. And, because each virtual machine is isolated from other virtualized machines, if one crashes, it does not affect the others.

Hypervisor software is the secret sauce that makes virtualization possible. This software sits between the hardware and the operating system, and de-couples the operating system and applications from the hardware. The hypervisor assigns the amount of access that the operating systems and applications have with the processor and other hardware resources, such as memory and disk input/output.

In addition to using virtualization technology to partition one machine into several virtual machines, you can also use virtualization solutions to combine multiple physical resources into a single virtual resource. A good example of this solution is storage virtualization. This type of virtualization is where multiple network storage resources are pooled into what is displayed as a single storage device for easier and more efficient management of these resources. Other types of virtualization you might hear about include the following:

- ▶ *Network virtualization* splits available bandwidth in a network into independent channels that can be assigned to specific servers or devices.
- ▶ *Application virtualization* separates applications from the hardware and the operating system, putting them in a container that can be relocated without disrupting other systems.
- ▶ *Desktop virtualization* enables a centralized server to deliver and manage individualized desktops remotely. This type of virtualization gives users a full client experience, but allows IT staff to provision, manage, upgrade, and patch them virtually, instead of physically.

Virtualization was first introduced in the 1960s by IBM. It was designed to boost utilization of large, expensive mainframe systems by partitioning them into logical, separate virtual machines that could run multiple applications and processes at the same time. In the 1980s and 1990s, this centrally shared mainframe model gave way to a distributed, client/server computing model, in which many low-cost x86 servers and desktops independently run specific applications.

6.2.1 Cloud infrastructure virtualization

This type consists of virtualizing three key parts: servers, desktops, or applications. The virtualization concept that is used for servers and desktops is almost the same, but for applications, the concept is different.

Virtualizing servers and desktops basically takes physical computers and makes them virtual. To make virtualization possible, a cloud platform is required. We show the traditional physical environment in Figure 6-5 on page 130. This model shows where one application maps to one operating system (OS), and one OS to one physical server, and one physical server to one storage.

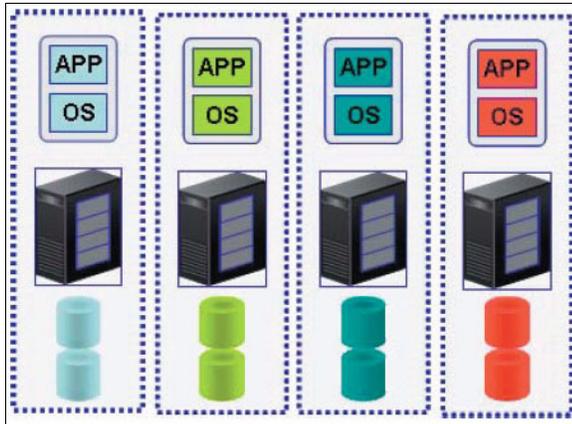


Figure 6-5 Traditional physical environment model

6.2.2 Cloud platforms

There must be a platform that can handle putting multiple virtual servers into a single physical computer. This platform is called the *hypervisor*. This platform is a layer in the computer stack between the virtual and physical components.

There are four core concepts in virtualization: encapsulation, isolation, partitioning, and hardware independence:

- ▶ *Encapsulation.* The entire machine becomes a set of files, and these files contain the operating system and application files plus the virtual machine configuration. The virtual machine files can be managed the same way that you manage other files.
- ▶ *Isolation.* Virtual machines (VMs) that run on a hardware platform cannot see or affect each other, so multiple applications can be run securely on a single server.
- ▶ *Partitioning.* VMware, for example, divides and actively manages the physical resources in the server to maintain optimum allocation.
- ▶ *Hardware independence.* The hypervisor provides a layer between the operating systems and hardware. This layer allows hardware from multiple vendors to run on the same physical resource, if the server is on *Hardware Compatibility List*.

Figure 6-6 shows the virtualized environment.



Figure 6-6 Virtualized environment model

Server virtualization

There are three popular approaches to server virtualization: the virtual machine model, the paravirtual machine model, and virtualization at the operating system layer.

Virtual machines (VMs) are based on the host/guest paradigm. Each guest runs on a virtual implementation of the hardware layer. This approach allows the guest operating system to run without modifications. It also allows the administrator to create guests that use different operating systems. The guest has no knowledge of the host operating system because it is not aware that it is not running on real hardware. It does, however, require real computing resources from the host so it uses a hypervisor to coordinate instructions to the CPU.

The paravirtual machine (PVM) model is also based on the host/guest paradigm and it uses a virtual machine monitor (VMM). In the paravirtual machine model, however, the VMM actually modifies the code of the guest operating system. This modification is called *porting*. Porting supports the VMM so it can use privileged systems calls sparingly. Like virtual machines, paravirtual machines can run multiple operating systems. Xen and UML both use the paravirtual machine model.

Virtualization at the OS level works a little differently. It is not based on the host/guest paradigm. In the OS level model, the host runs a single OS kernel as its core and exports the operating system functionality to each of the guests. Guests must use the same operating system as the host, although different distributions of the same system are allowed. This distributed architecture eliminates system calls between layers, which reduce CPU usage overhead. It also requires that each partition remains strictly isolated from its neighbors so that a failure or security breach in one partition is not able to affect any of the other partitions. In this model, common binary files and libraries on the same physical machine can be shared, allowing an OS-level virtual server to host thousands of guests at the same time. IBM AIX VIO and Solaris Zones both use OS-level virtualization.

Desktop Virtualization

This is sometimes referred to as *client virtualization*, and is defined as a virtualization technology that is used to separate a computer desktop environment from the physical

computer. Desktop virtualization is considered a type of client/server computing model because the virtualized desktop is stored on a centralized, or remote, server and not the physical machine that is being virtualized.

Desktop virtualization virtualizes desktop computers and these virtual desktop environments are “served” to users on the network. Users interact with a virtual desktop in the same way that a physical desktop is accessed and used. Another benefit of desktop virtualization is that it allows you to remotely log in to access your desktop from any location.

One of the most popular uses of desktop virtualization is in the data center, where personalized desktop images for each user are hosted on a data center server.

There are also options for using hosted virtual desktops, where the desktop virtualization services are provided to a business through a third party. The service provider provides the managed desktop configuration, security, and SAN.

Application Virtualization

Application virtualization is just like desktop virtualization, where individual desktop sessions (OS and applications) are virtualized and run from a centralized server. However, *Application virtualization* virtualizes the applications so that it can either be run from a centralized server or it can be streamed from a central server and run in an isolated environment in the desktop itself.

In the first type of application virtualization, the application image is loaded on to a central server and when a user requests the application, it is streamed to an isolated environment on the user’s computer for execution. The application starts running shortly after it gets sufficient data to start running, and since the application is isolated from other applications, there might not be any conflicts. The applications that can be downloaded can be restricted based on the user ID which is established by logging in to corporate directories such as Active Directory (AD) or Lightweight Directory Access Protocol (LDAP).

In the second type of application virtualization, the applications are loaded as an image in remote servers and they are run (executed) in the servers itself. Only the on-screen information that is required to be seen by the user is sent over the LAN. This is closer to desktop virtualization, but here only the application is virtualized instead of both the application and the operating system. The biggest advantage of this type of application virtualization is that it does not matter what the underlying OS is in the user’s computer because the applications are processed in the server. Another advantage is the effectiveness of mobile devices (mobile phones, tablet computers, and so on) that have lesser processing power while running processor hungry applications. This is because these applications are processed in the powerful processors of the servers.

6.2.3 Storage virtualization

Storage virtualization refers to the abstraction of storage systems from applications or computers. It is a foundation for the implementation of other technologies, such as thin provisioning, tiering, and data protection, which are transparent to the server.

These are some of the advantages of storage virtualization:

- ▶ Improved physical resource utilization: By consolidating and virtualizing storage systems, we can make more efficient use of previously wasted white spaces.
- ▶ Improved responsiveness and flexibility: De-coupling physical storage to virtual storage provides the ability to reallocate resources dynamically, as required by the applications or storage subsystems.

- ▶ Lower total cost of ownership: Virtualized storage allows more to be done with the same or less storage.

Several types of storage virtualization are available.

Block level storage virtualization

Block level storage virtualization refers to provisioning storage to your operating systems or applications in the form of virtual disks. Fibre Channel (FC) and Internet Small Computer System Interface (iSCSI) are examples of protocols that are used by this type of storage virtualization.

There are two types of block level virtualization:

- ▶ *Disk level virtualization*. This is an abstraction process from a physical disk to a logical unit number (LUN) that is presented as if it were a physical device.
- ▶ *Storage level virtualization*. Unlike disk level virtualization, storage level virtualization hides the physical layer of Redundant Array of Independent Disks (RAID) controllers and disks, and hides and virtualizes the entire storage system.

File level storage virtualization

File level storage virtualization refers to provisioning storage volumes to operating systems or applications in the form of files and directories. Access to storage is by network protocols, such as Common Internet File Systems (CIFS) and Network File Systems (NFS). It is a file presentation in a single global namespace, regardless of the physical file location.

Tape virtualization

Tape virtualization refers to the virtualization of tapes and tape drives that use specialized hardware and software. This type of virtualization can enhance backup and restore flexibility and performance because disk devices are used in the virtualization process, rather than tape media.

6.3 SAN virtualization

For SAN virtualization, we describe the virtualization features available in the IBM System Networking portfolio. These features enable the SAN infrastructure to support the requirements of scalability and consolidation, and combine this with a lower TCO and a higher ROI:

- ▶ IBM b-type Virtual Fabrics
- ▶ CISCO Virtual SAN (VSAN)
- ▶ N_Port ID Virtualization (NPIV) support for virtual nodes

6.3.1 IBM b-type Virtual Fabrics

The Virtual Fabric of the IBM b-type switches is a licensed feature which enables the logical partitioning of SAN switches. When Virtual Fabric is enabled, a default logical switch that is using all the ports is formed and this default logical switch can be then divided into multiple logical switches by grouping them together at a port level.

Figure 6-7 indicates the flow of Virtual Fabric creation.

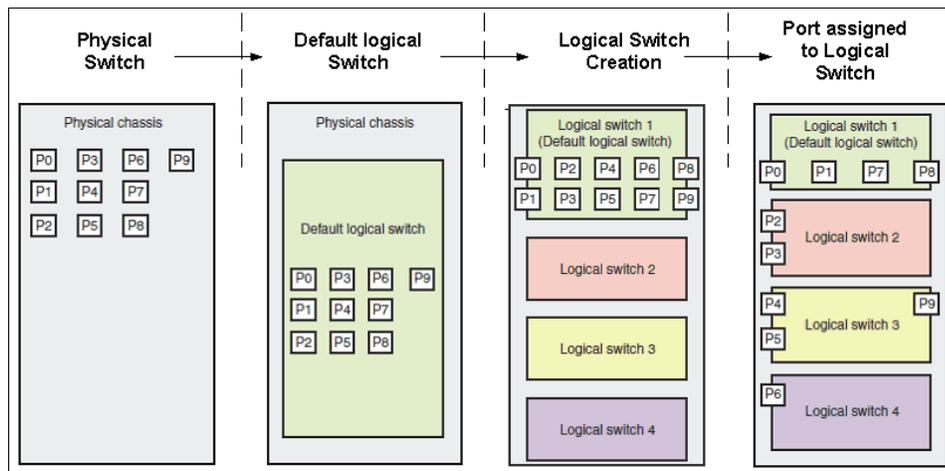


Figure 6-7 Virtual fabric creation

Logical fabric

When the fabric is formed with at least one logical switch, it is called a *logical fabric*. The logical fabric has two different ways of fabric connectivity.

- ▶ A logical fabric is connected with a dedicated ISL to another switch or a logical switch. Figure 6-8 shows a logical fabric that is formed between logical switches through a dedicated ISL for logical switches.

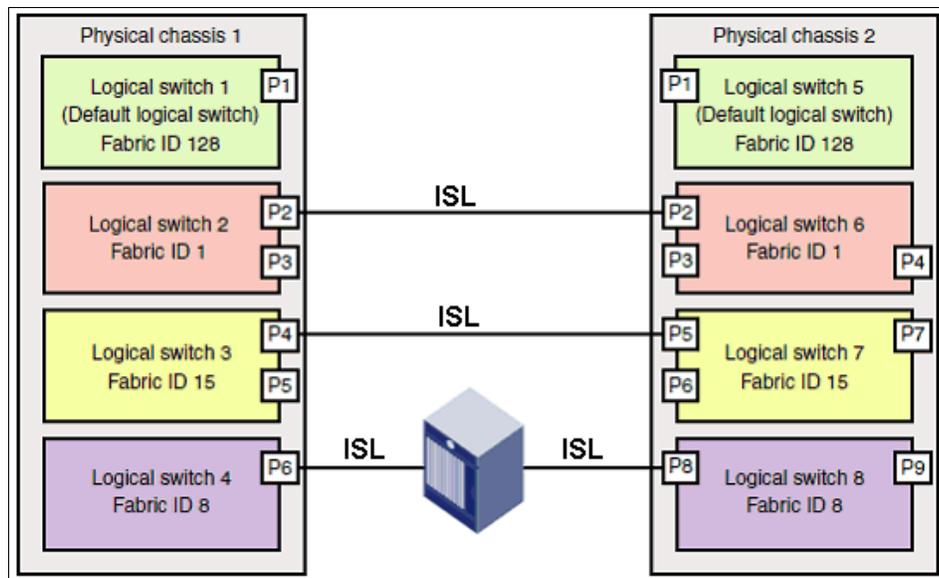


Figure 6-8 Logical fabrics with dedicated ISL

- ▶ Logical fabrics are connected by using a shared ISL (called extended ISL (XISL)) from a base logical switch. In this case, a separate logical switch is configured to be a base switch and is used only for XISL connectivity and not for device connectivity. Figure 6-9 on page 135 indicates a logical fabric that is formed through the XISL in the base switch.

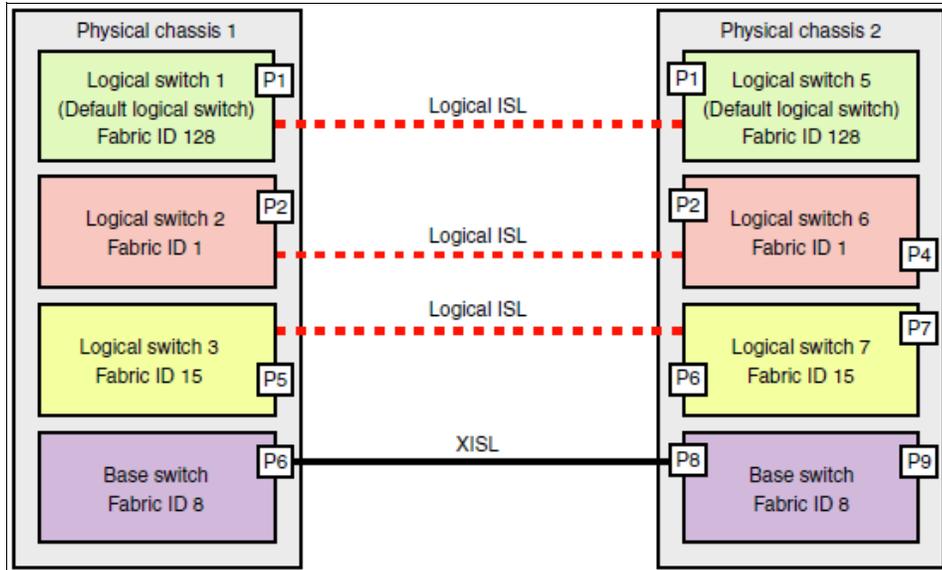


Figure 6-9 Logical ISL through XISL in the base switch

6.3.2 Cisco virtual storage area network

Cisco *virtual storage area network (VSAN)* is a feature which enables the logical partition of SAN switches. A VSAN provides the flexibility to partition, for example, a dedicated VSAN for disk and tape. Or, the VSAN provides the flexibility to have production and test devices in separate VSANs on the same chassis. Also, the VSAN can scale across the chassis, which allows it to overcome the fixed port numbers on the chassis.

Virtual storage area network in a single SAN switch

VSAN brings the ability to consolidate small fabrics into the same chassis. This consolidation can also enable more security by logical separation of the chassis into two individual VSANs. Figure 6-10 shows a single chassis that is divided into two logical VSANs.

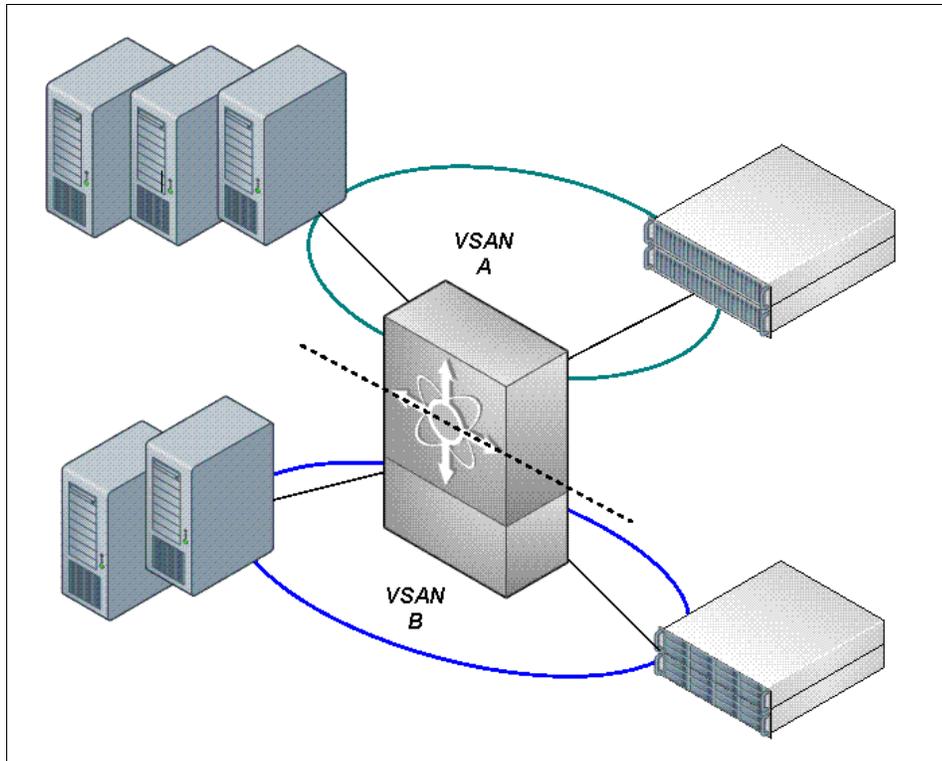


Figure 6-10 Two VSANs in a single chassis

Virtual storage area network across multiple chassis

In multiple chassis, the virtual storage area network (VSAN) can be formed with devices in one chassis to devices in another switch chassis through the *extended inter-switch link (XISL)*.

Figure 6-11 shows the VSAN across chassis with an *enhanced inter-switch link (EISL)* for VSAN communication.

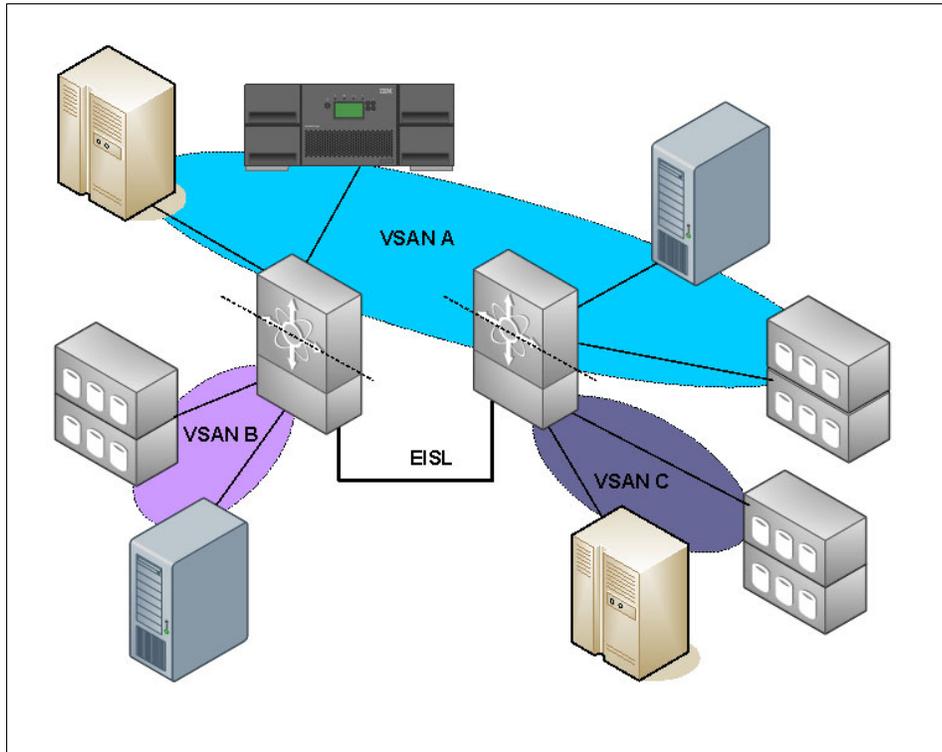


Figure 6-11 VSAN across multiple chassis

6.3.3 N-Port ID Virtualization

Server virtualization with blade servers provides enhanced scalability of servers and this scalability is supported equally in the SAN with something called *N_Port ID Virtualization (NPIV)*. NPIV allows SAN switches to have one port that is shared by many virtual nodes, which in turn supports a single HBA having many virtual nodes.

Figure 6-12 shows the sharing of a single HBA by multiple virtual nodes. In this case, the same HBA is defined with multiple virtual WWN and WWPNS.

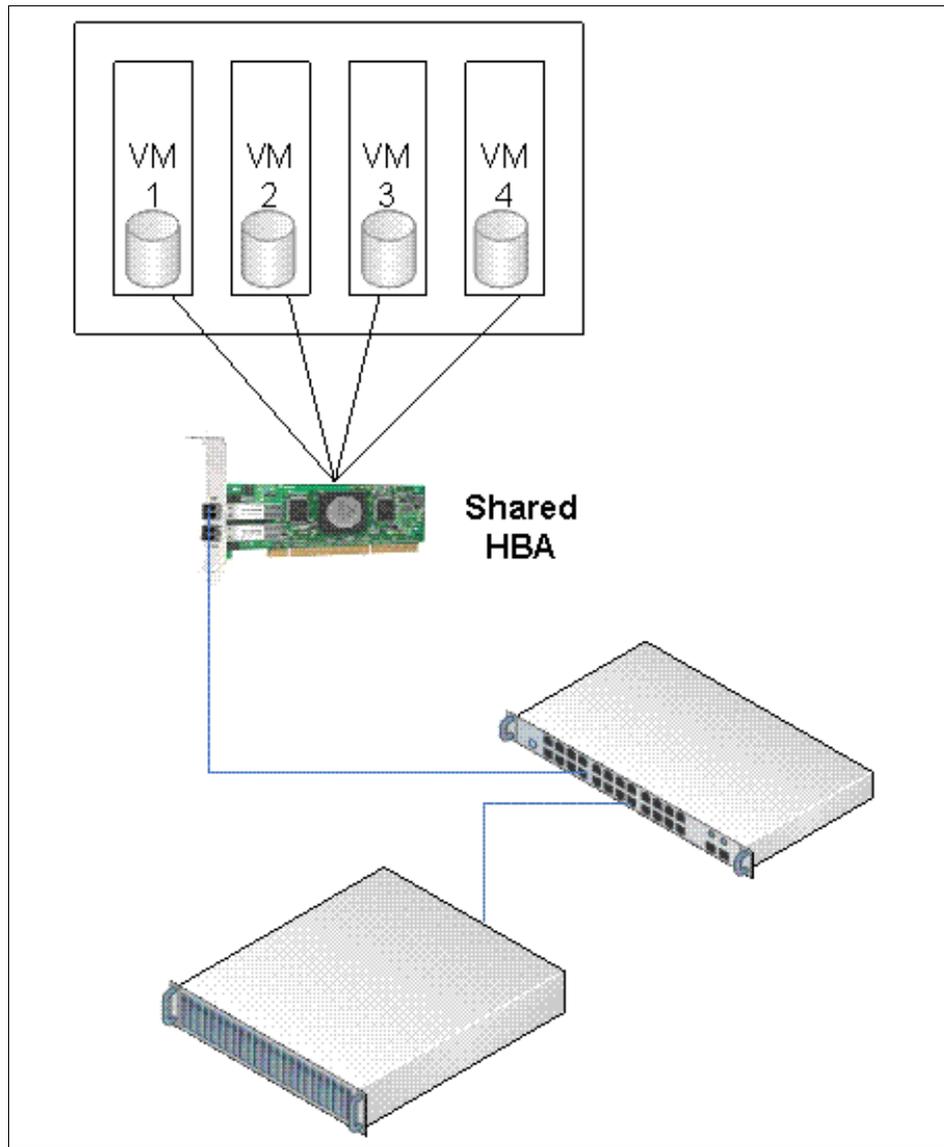


Figure 6-12 Single HBA with multiple virtual nodes

NPIV mode of blade server switch modules

On Blade servers, when enabled with the NPIV mode, the FC switch modules that are connected to an external SAN switch for access to storage, act as an HBA N_Port in this case (instead of a switch E_Port). The back-end ports are F_Ports which are connected to server blade modules.

Figure 6-13 shows the switch module in the NPIV mode.

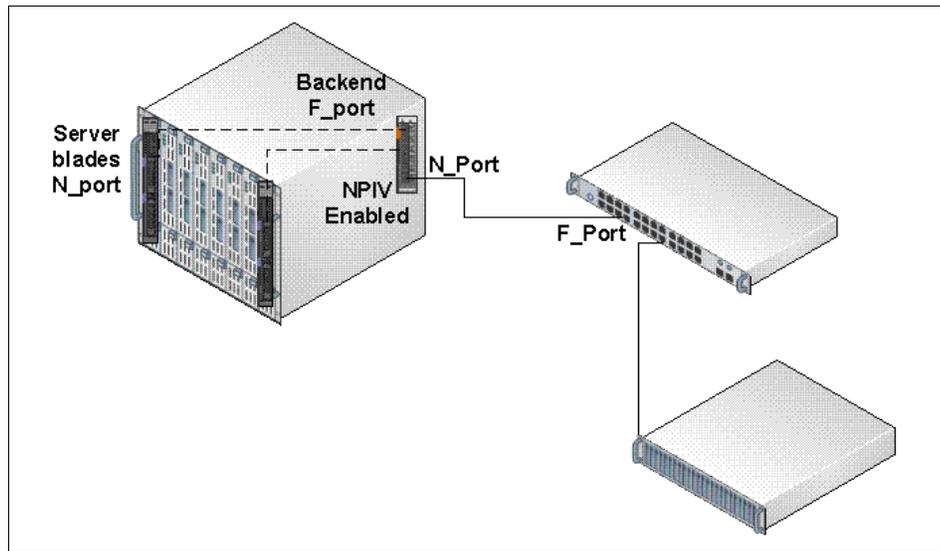


Figure 6-13 Blade server with FC switch module in the NPIV mode

With the NPIV mode, we can overcome the interoperability issues of merging external switches to the blade server switch module which might be from different vendors. Also, we have the benefit of easy management because the blade switch module becomes a node in the fabric. And, we can overcome the scalability limitations of having many switch domains for a switch module in blade servers.

6.4 Building a smarter cloud

Storage-as-a-Service is a business model in which a large company rents space in their storage infrastructure to a smaller company or individual. It is generally seen as a good alternative when a small business lacks the capital budget or technical personnel to implement and maintain their own storage infrastructure. In some circumstances, it is being promoted as a way for all businesses to mitigate risks in disaster recovery, provide long-term retention for records, and enhance both business continuity and availability.

6.4.1 Automated tiering

In modern and complex application environments, the increasing and often unpredictable demands for storage capacity and performance leads to relevant issues in terms of planning and optimization of storage resources.

Most of these issues can be managed by having spare resources available and by moving data, by using data mobility tools, or by using operating systems features (such as host level mirroring). However, all these corrective actions are expensive in terms of hardware resources, labor, and service availability. Relocating data among the physical storage resources dynamically, that is, transparently to hosts, is becoming increasingly important.

IBM Storage Solutions offer two types of automated tiering.

Automated tiering to optimize performance

The IBM Easy Tier® feature, available with the DS8000, SAN Volume Controller, and Storwize V7000, provides performance optimization. Easy Tier is a built-in, dynamic data relocation feature that provides optimal performance at the lowest cost. Easy Tier is designed to determine the appropriate tier of storage to use, based on data access patterns. Then, it automatically and non-disruptively moves data, at the sub-LUN or subvolume level, to the appropriate disk tier.

Automated tiering to optimize space management

The ability to optimize space management is an information lifecycle management (ILM) function that is available, for instance, with Scale Out Network Attached Storage and with Hierarchical Storage Management. Examples include functions that are provided by Tivoli Storage Manager and the IBM Data Facility Storage Management System (DFSMS) Hierarchical Storage Management (DFSMSHsm).

Policy-based automation is used to migrate less active data to lower-cost storage.

6.4.2 Thin provisioning

Traditional storage provisioning pre-allocates and dedicates physical storage space for use by the application or host. However, often not all space that is allocated to applications is needed resulting in wasted “white space”.

Thin provisioning allows a server to see logical volume sizes that are larger than the physical capacity that is dedicated to the volumes on the storage system. From the server or application perspective, thinly provisioned volumes are displayed and function just the same as fully provisioned volumes. However, physical disk drive capacity is allocated only as needed (on demand) for write activity to the volumes. Deallocated physical capacity is available for use as needed by all volumes in a storage pool or even across an entire storage system.

Some of the advantages of thin provisioning are:

- ▶ It allows higher storage systems utilization which in turn leads to a reduction in the amount of storage you need, lowering your direct capital expenditure (CAPEX).
- ▶ It lowers operational expenditure (OPEX) because your storage occupies less data center space and requires less electricity and cooling.
- ▶ It postpones the need to buy more storage, and as storage prices continue to drop over time, when more capacity is required, it will likely cost less.
- ▶ Capacity planning is simplified because you are able to manage a single pool of free storage. Multiple applications or users can allocate storage from the same free pool, avoiding the situation in which some volumes are capacity constrained while others have capacity to spare.
- ▶ Your storage environment becomes more agile and it becomes easier to react to change.

Thin provisioning increases utilization ratios

Thin provisioning increases storage efficiency by increasing storage utilization ratios. Real physical capacity is provided only as it is needed for writing data. This results in large potential savings in both storage acquisition and operational costs, including infrastructure costs such as power, space and cooling.

Storage utilization is measured by comparing the amount of physical capacity that is used for data with the total amount of physical capacity that is allocated to a server. Historically,

utilization ratios have been well under 50%, indicating a large amount of allocated but unused physical storage capacity. Often, neither the users or storage administrators are certain how much capacity is needed, but they must ensure that they do not run out of space, and they also must allow for growth. As a result, users might request more than they need and storage administrators might allocate more than is requested, resulting in significant over-allocation of storage capacity.

Thin provisioning increases storage utilization ratios by reducing the need to over-allocate physical storage capacity to prevent out of space conditions. Large logical or virtual volume sizes might be created and presented to applications without dedicating an equivalent amount of physical capacity. Physical capacity can be allocated on demand as needed for writing data. Deallocated physical capacity is available for multiple volumes in a storage pool or across the entire storage system.

Thin provisioning also increases storage efficiency by reducing the need to resize volumes or add volumes and restripe data as capacity requirements grow. Without thin provisioning, if an application requires capacity beyond what is provided by its current set of volumes, there are two options:

- ▶ Existing volumes might be increased in size
- ▶ Additional volumes might be provisioned

In many environments, these options are undesirable because of the steps and potential disruption that is required to make the larger or additional volumes visible and optimized for the application.

With thin provisioning, large virtual or logical volumes might be created and presented to applications while the associated physical capacity grows only as needed, transparent to the application.

Without thin provisioning, physical capacity was dedicated at the time of volume creation, and storage systems typically did not display or report how much of the dedicated physical capacity was used for data. As storage systems implemented thin provisioning, physical allocation and usage became visible. Thin provisioning increases storage efficiency by making it easy to see the amount of physical capacity that is needed and used because physical space is not allocated until it is needed for data.

6.4.3 Deduplication

Data deduplication emerged as a key technology to dramatically reduce the amount and the cost that is associated with storing large amounts of data. Deduplication is the art of intelligently reducing storage needs in order of magnitude. This method is better than common data compression techniques. Deduplication works through the elimination of redundant data so that only one instance of a data set is stored. IBM has the broadest portfolio of deduplication solutions in the industry, which gives IBM the freedom to solve client issues with the most effective technology. Whether its source or target, inline or post, hardware or software, disk or tape, IBM has a solution with the technology that best solves the problem:

- ▶ IBM ProtecTIER® Gateway and Appliance
- ▶ IBM System Storage N series Deduplication
- ▶ IBM Tivoli Storage Manager

Data deduplication is a technology that reduces the amount of space that is required to store data on disk. It achieves this space reduction by storing a single copy of data that is backed up repetitively.

Data deduplication products read data while they look for duplicate data. Data deduplication products break up data into elements, using their respective technique to create a signature or identifier for each data element. Then, they compare the data element signature to identify duplicate data. After they identify duplicate data, they retain one copy of each element. They create pointers for the duplicate items, and discard the duplicate items.

The effectiveness of data deduplication is dependent upon many variables, including the rate of data change, the number of backups, and the data retention period. For example, if you back up the exact same incompressible data once a week for six months, you save the first copy and do not save the next 24. This method would provide a 25 to 1 data deduplication ratio. If you back up an incompressible file on week one, then back up the exact same file again on week two and never back it up again, you have a 2 to 1 deduplication ratio. A more likely scenario is that some portion of your data changes from backup to backup so that your data deduplication ratio will change over time. With data deduplication, you can minimize your storage requirements.

Data deduplication can provide greater data reduction and storage space savings than other existing technologies.

Figure 6-14 shows the basic concept of data deduplication.

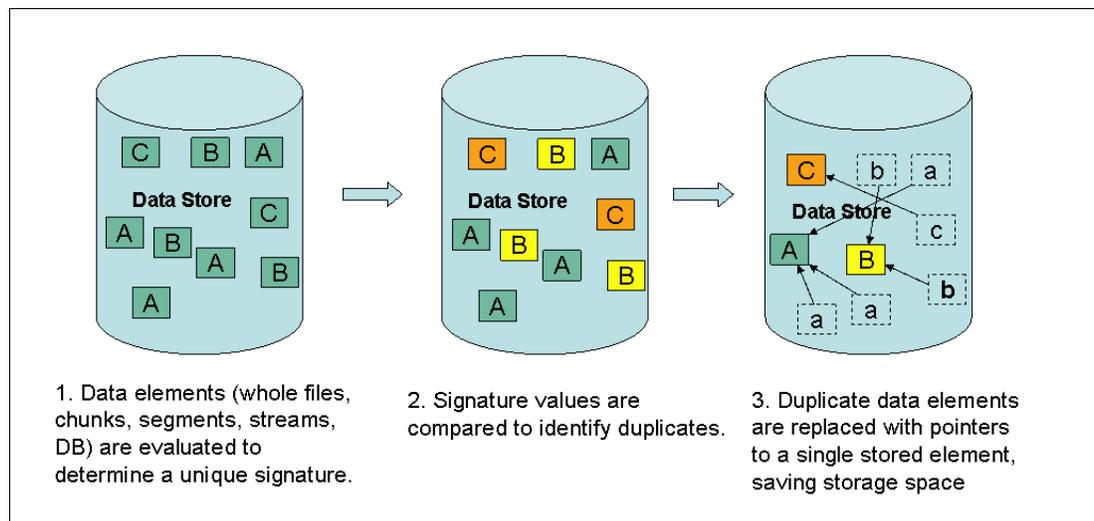


Figure 6-14 The basic concept of data deduplication

Data deduplication can reduce your storage requirements but the benefit you derive is determined by your data and your backup policies. Workloads with a high database content generally have the highest deduplication ratios. However, product functions like Tivoli Storage Manager Progressive Incremental, Oracle RMAN, or Light Speed, can reduce the deduplication ratio. Compressed, encrypted, or otherwise scrambled workloads typically do not benefit from deduplication. Good candidates for deduplication are typically text files, log files, uncompressed and non-encrypted database files, email files (PST, DBX, IBM Domino®), and Snapshots (Filer Snaps, BCVs, VMware images).

Types of data deduplication and HyperFactor

Many vendors offer products that perform deduplication. Various methods are used for deduplicating data. The following three methods are frequently used for data deduplication:

- **Hash based** deduplication uses a hashing algorithm to identify chunks of data. Commonly used process is Secure Hash Algorithm 1 (SHA-1) or Message-Digest Algorithm 5 (MDA-5). The details of each technique are beyond the intended scope of this publication.

- ▶ **Content aware** deduplication methods are aware of the structure of common patterns of data that is used by applications. It assumes the best candidate to de-duplicate against is an object with the same properties, such as a file name. When a file match is found, a bit by bit comparison is performed to determine if data has changed and saves the changed data.
- ▶ **IBM HyperFactor®** is a patented technology which is used in IBM System Storage ProtecTIER Enterprise Edition higher software. HyperFactor takes an approach that reduces the phenomenon of missed factoring opportunities, providing a more efficient process. With this approach, HyperFactor is able to surpass the reduction ratios attainable by any other data reduction method. HyperFactor can reduce any duplicate data, regardless of its location or how recently it was stored. HyperFactor data deduplication uses a 4 GB Memory Resident Index to track similarities for up to 1 petabyte (PB) of physical disk in a single repository.

HyperFactor technology uses a pattern algorithm that can reduce the amount of space that is required for storage by up to a factor of 25, based on evidence from existing implementations. The capacity expansion that results from data deduplication is often expressed as a ratio, essentially the ratio of nominal data to the physical storage used.

Figure 6-15 shows the HyperFactor technology.

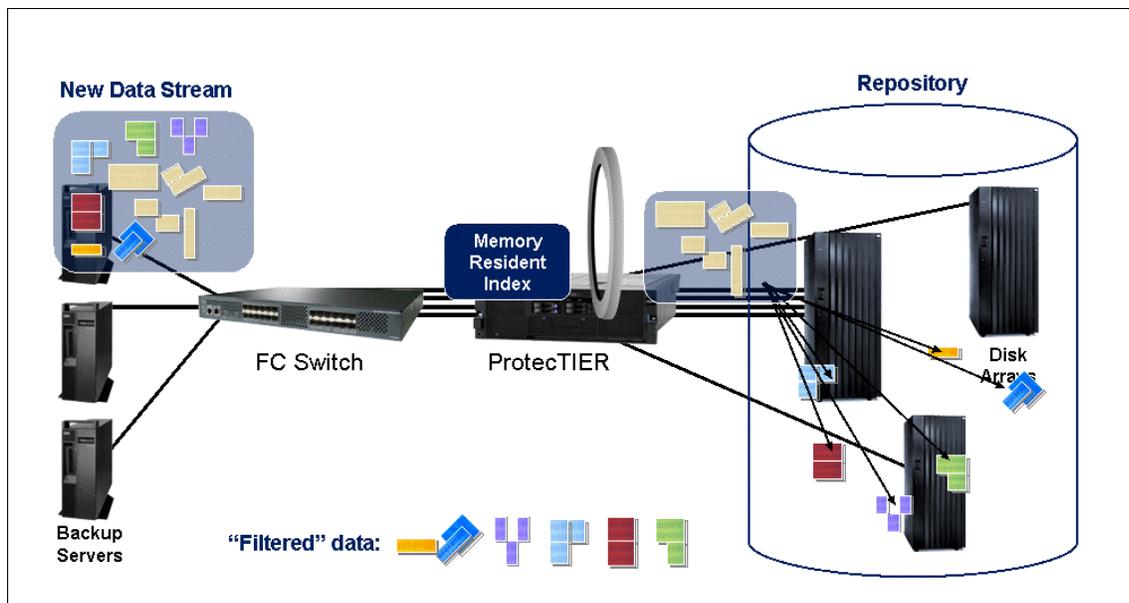


Figure 6-15 IBM HyperFactor technology

Data deduplication processing

Data deduplication can either occur while the data is being backed up to the storage media (real-time or inline) or after the data is written to the storage media (post-processing). Each method certainly brings positive and negative aspects, these considerations must be evaluated by the engineer or technical specialist that is responsible for the concrete solution architecture and deployment. IBM decided to use inline deduplication processing as it offers larger target storage space without any need of temporary disk cache pool for post processed deduplication data.

Bit comparison techniques such as the one used by ProtecTIER were designed to provide 100% data integrity by avoiding the risk of hash collisions.

6.4.4 New generation management tools

It is paramount that this new virtualized infrastructure be managed by new generation management tools because older tools generally lack the required features.

When used properly, these tools can make the adoption of virtualization technology easier and more cost effective. The tools have the following additional benefits:

- ▶ Enable line of business insight into storage utilization and allocation, enabling easier departmental charge back.
- ▶ Allow for more intelligent business decisions about storage efficiency, which enable you to respond faster to changing business demands, yet reduce costs.
- ▶ Provide better understanding of application storage performance and utilization patterns that enable better information lifecycle management (ILM) of application data.
- ▶ Allow an organization to perform infrastructure management pro-actively, through proper capacity management, rather than reactively.
- ▶ Improve operational efficiency that leads to cost savings for the organization.

6.4.5 Business continuity and disaster recovery

The importance of business continuity and disaster recovery remains at the forefront of thought for many executives and IT technical professionals. The most important factor to consider is how the choice of the technology affects the recovery time objective (RTO). For SANs, there are many possible solutions that are available. The cloud design drives the selections that are capable of meeting the requirements. A smart cloud must be capable of guaranteeing business continuity and any disaster recovery plans.

Disaster recovery: For more information about DR lessons learned and solutions, refer to the IBM Storage Infrastructure for Business Continuity.

<http://www.redbooks.ibm.com/abstracts/redp4605.html?Open>

6.4.6 Storage on demand

Scalable, pay-per-use cloud storage can help to manage massive data growth and storage budget. This leads to a costly procurement cycle with setup costs and implementation delays every time that more storage is needed. Cloud storage allows the ability to expand storage capacity on the spot, and later, shrink storage consumption if needed.

Cloud storage provides a ready-made data storage solution that helps in these areas:

- ▶ Reduce up front capital expenses
- ▶ Meet demands without expensive over-provisioning
- ▶ Supplement other storage systems more cost-effectively
- ▶ Align data storage costs with business activity
- ▶ Scale dynamically



Fibre Channel products and technology

In this chapter, we describe some of the most common Fibre Channel storage area network (SAN) products and technology that are encountered. For a description of the IBM products that are in the IBM System Storage and TotalStorage portfolio, see Chapter 12, “The IBM product portfolio” on page 245.

7.1 The environment

The *Storage Networking Industry Association (SNIA)* defines the meaning of SAN, Fibre Channel, and Storage:

► *Storage area network (SAN)*

A *network* whose primary purpose is the transfer of data between computer systems and storage elements and among storage elements.

A SAN consists of a communication infrastructure which provides physical connections and a management layer. This layer organizes the connections, storage elements, and computer systems so that data transfer is secure and robust. The term SAN is usually (but not necessarily) identified with block I/O services rather than file access services.

► *Fibre Channel*

A serial I/O interconnect that is capable of supporting multiple protocols, including access to open system storage (FCP), access to mainframe storage (FICON), and networking (TCP/IP). Fibre Channel supports point-to-point, arbitrated loop, and switched topologies with various copper and optical links that are running at speeds from 1 Gbps to 10 Gbps. The committee that is standardizing Fibre Channel is the INCITS Fibre Channel (T11) Technical Committee.

► *Storage system*

A storage system that consists of storage elements, storage devices, computer systems, and appliances, plus all control software, which communicates over a network.

Storage subsystems, storage devices, and server systems can be attached to a Fibre Channel SAN. Depending on the implementation, several different components can be used to build a SAN. It is as the name suggests, a *network*, so any combination of devices that is able to interoperate are likely to be used.

Given this definition, a Fibre Channel network might be composed of many different types of interconnect entities, including directors, switches, hubs, routers, gateways, and bridges.

It is the deployment of these different types of interconnect entities that allow Fibre Channel networks of varying scale to be built. In smaller SAN environments you can employ hubs for Fibre Channel arbitrated loop topologies, or switches and directors for Fibre Channel switched fabric topologies. As SANs increase in size and complexity, Fibre Channel directors can be introduced to facilitate a more flexible and fault-tolerant configuration. Each of the components that composes a Fibre Channel SAN, provides an individual management capability and participate in an often complex end-to-end management environment.

Figure 7-1 shows a generic SAN connection.

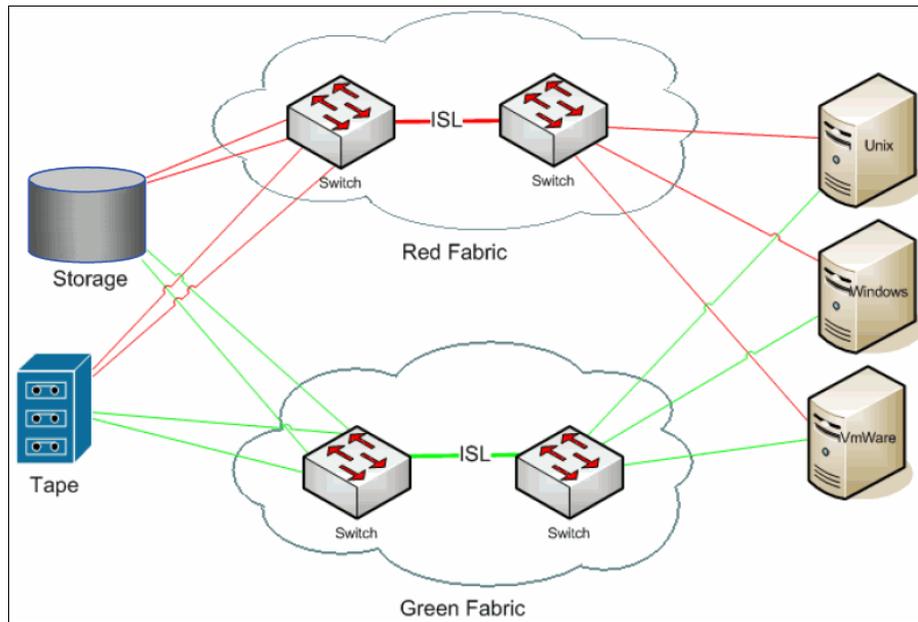


Figure 7-1 Generic SAN

7.2 Storage area network (SAN) devices

A Fibre Channel SAN employs a *fabric* to connect devices, or *end points*. A fabric can be as simple as a single cable that connects two devices, akin to server-attached storage. However, the term is most often used to describe a more complex network to connect servers and storage by using switches, directors, and gateways.

Independent from the size of the fabric, a good SAN environment starts with good planning, and always includes an up-to-date map of the SAN.

Some of the items to consider include the following questions:

- ▶ How many ports do I need now?
- ▶ How fast will I grow in two years?
- ▶ Are my servers and storage in the same building?
- ▶ Do I need long-distance solutions?
- ▶ Do I need redundancy for every server or storage?
- ▶ How high are my availability needs and expectations?
- ▶ Will I connect multiple platforms to the same fabric?
- ▶ What technology do I want to use, FC - FCoE - iScsi?

7.2.1 Fibre Channel bridges

Fibre Channel bridges allow the integration of traditional SCSI devices in a Fibre Channel network. Fibre Channel bridges provide the capability for Fibre Channel and SCSI interfaces to support both SCSI and Fibre Channel devices seamlessly. Therefore, they are often referred to as *FC-SCSI routers*.

Data Center Bridging: Fibre Channel bridges are not to be confused with *Data Center Bridging (DCB)*, though fundamentally they serve the same purpose, which is to interconnect different protocols.

A *bridge* is a device that converts signals and data from one form to another. You can imagine these devices in a similar way as the bridges that we use to cross rivers. They act as a translator (a bridge) between two different protocols. These protocols can include the following types:

- ▶ Fibre Channel
- ▶ Internet Small Computer System Interface (iSCSI)
- ▶ Serial Storage Architecture (SSA)
- ▶ Fibre Channel over IP (FCIP)

We do not see many of these devices today and they are considered legacy devices.

7.2.2 Arbitrated loop hubs and switched hubs

Arbitrated loop, also known as *FC-AL*, is a Fibre Channel topology in which devices are connected in a one-way loop fashion in a ring topology. This topology is also described in Chapter 5, “Topologies and other fabric services” on page 85.

Figure 7-2 shows an FC-AL topology.

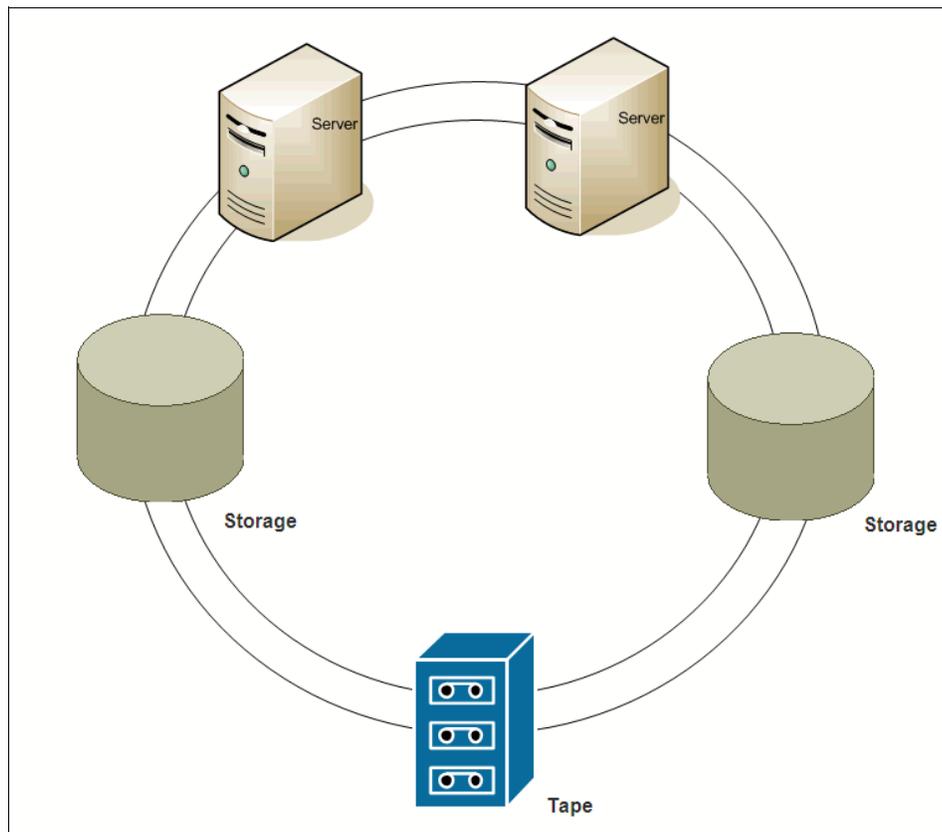


Figure 7-2 Arbitrated loop

In FC-AL, all devices on the loop share the bandwidth. The total number of devices that might participate in the loop is 126, without using any hubs or fabric. For practical reasons, however, the number tends to be limited to no more than 10 and 15.

Hubs are typically used in a SAN to attach devices or servers that do not support switched fabric-only FC-AL. They might be unmanaged hubs, managed hubs, or switched hubs.

Unmanaged hubs serve as cable concentrators and as a means to configure the Arbitrated Loop that is based on the connections it detects. When any of the interfaces on the hub, usually a *gigabit interface converter (GBIC)*, senses that no cable is connected, that interface shuts down. The hub port is then bypassed as part of the Arbitrated Loop configuration.

Managed hubs offer all the benefits of unmanaged hubs, but in addition, offer the ability to manage them remotely, using *Simple Network Management Protocol (SNMP)*.

By using FC-AL, you can connect many servers and storage devices without using costly Fibre Channel switches. FC-AL is not used much today because switched fabrics now lead in the Fibre Channel market.

Switched hubs

Switched hubs allow devices to be connected in their own Arbitrated Loop. These loops are then internally connected by a switched fabric.

A switched hub is useful to connect several FC-AL devices together, but to allow them to communicate at full Fibre Channel bandwidth rather than them all sharing the bandwidth.

Switched hubs are usually managed hubs.

FC-AL: In its early days, FC-AL was described as “SCSI on steroids”. Although FC-AL has the bandwidth advantage over SCSI, it does not come anywhere close to the speeds that can be achieved and sustained on a per port basis in a switched fabric. For this reason, FC-AL implementations are, by some observers, considered as legacy SANs.

7.2.3 Switches and directors

Switches and directors allow Fibre Channel devices to be connected (cascaded) together, implementing a switched fabric topology between them. The switch intelligently routes frames from the initiator to the responder and operates at full Fibre Channel bandwidth.

It is possible to connect switches together in cascades and meshes by using *inter-switch links (E_Ports)*. Keep in mind that devices from different manufacturers might not interoperate fully.

The switch also provides various fabric services and features. The following list provides some examples:

- ▶ Name service
- ▶ Fabric control
- ▶ Time service
- ▶ Automatic discovery and registration of host and storage devices
- ▶ Rerouting of frames, if possible, in the event of a port problem
- ▶ Storage services (virtualization, replication, and extended distances)

It is common to refer to switches as either core switches or edge switches, depending on where they are in the SAN. If the switch forms, or is part of the SAN backbone, then it is the *core switch*. If it is mainly used to connect to hosts or storage, then it is called an *edge switch*. Directors are also sometimes referred to as switches because they are switches, in essence. Directors are large switches with higher redundancy than most normal switches. Whether this analogy is appropriate or not, is a matter for debate that is not covered in this book.

7.2.4 Multiprotocol routing

There are also devices that are multiprotocol routers and devices. These provide improved scalability, security, and manageability by enabling devices in separate SAN fabrics to communicate *without* merging fabrics into a single, large Meta-SAN fabric. Depending on the manufacturer, they support a number of protocols and have their own features, such as zoning. As their name suggests, the following list provides the protocols that are supported:

- ▶ Fibre Channel Protocol (FCP)
- ▶ Fibre Channel over IP (FCIP)
- ▶ Internet Fibre Channel Protocol (iFCP)
- ▶ Internet Small Computer System Interface (iSCSI)
- ▶ Internet Protocol (IP)

7.2.5 Service modules

Increasingly, with the demand for the intermix of protocols and the introduction to the marketplace of new technologies, SAN vendors are starting to adopt a modular system approach to their components. What this means is that service modules can be plugged into a slot on the switch or director to provide functions and features such as virtualization, the combining of protocols, and storage services.

7.2.6 Multiplexers

Multiplexing is the process of simultaneously transmitting multiple signals over the same physical connection. There are common types of multiplexing used for fiber optic connections that are based on either time or wavelength:

- ▶ Time-division multiplexing (TDM)
- ▶ Wavelength division multiplexing (WDM)
- ▶ Dense wavelength division multiplexing (DWDM)

When you use multiplexers in a SAN environment, more parameters in the SAN switch configuration might be needed to ensure correct load balancing. Therefore, check with your SAN switch vendor for preferred practices.

Multiplexers: Usually multiplexers are transparent to the SAN fabric. If you are troubleshooting an ISL link that covers some distance, keep in mind that the multiplexer, if installed, plays an important role in that path.

7.3 Componentry

There are a number of components that must come together to make a SAN, a SAN. We identify some of the components that are likely to be encountered.

7.3.1 Application-specific integrated circuit

The fabric electronics use a personalized *application-specific integrated circuit (ASIC)* and its predefined set of elements. Examples of these types of elements include: logic functions, I/O circuits, memory arrays, and backplanes to create specialized fabric interface components.

An ASIC provides services to Fibre Channel ports. The circuit might be used to connect to external N_Ports (such as an F_Port or FL_Port), external loop devices (such as an FL_Port), or to other switches (such as an E_Port). The ASIC contains the Fibre Channel interface logic, message and buffer queuing logic, and receives buffer memory for the on-chip ports, as well as other support logic.

Frame filtering

Frame filtering is a feature that enables devices to provide zoning functions with finer granularity. Frame filtering can be used to set up port-level zoning, worldwide name (WWN) zoning, device level zoning, protocol level zoning, and logical unit number (LUN) level zoning. Frame filtering is commonly carried out by an ASIC. This solution has the result that, after the filter is set up, the complicated function of zoning and filtering can be achieved at wire speed.

7.3.2 Fibre Channel transmission rates

Fibre Channel transmission rates are sometimes referred to as *feeds and speeds*. The current set of vendor offerings for switches, host bus adapters (HBAs), and storage devices is constantly increasing. Currently, the 16 Gb FC port has the fastest line rate that is supported for an IBM SAN. The 16 Gb FC port uses a 14.025 Gbps transfer rate and uses 64b/66b encoding that provides approximately 1600 MBps in throughput. The 8 Gb FC port has a line rate of 8.5 Gbps that uses 8b/10b encoding, which results in approximately 800 MBps. When you compare feeds and speeds, the FC ports are sometimes referred to as *full duplex*. The transceiver and receive parts of the FC port are then added, therefore “doubling” the MBps.

Encoding: By introducing the 64b/66b encoding to Fibre Channel, the encoding overhead is reduced from approximately 20% by using 8b/10b encoding, down to approximately 3% with the 64b/66b encoding.

The new 16 Gb FC port is approved by the Fibre Channel Industry Association (FCIA). This approval ensures that each port speed is able to communicate with at least two previous approved port speeds; for example, 16 Gb is able to communicate with 8 Gb and 4 Gb.

The FCIA also created a roadmap for future feeds and speeds. For more information, see this website:

<http://www.fibrechannel.org/>

7.3.3 SerDes

The communication over a fiber, whether optical or copper, is serial. Computer busses, however, use parallel busses. This means that Fibre Channel devices must be able to convert between the two. For this conversion, the devices use a serializer/deserializer, which is commonly referred to as a *SerDes*.

7.3.4 Backplane and blades

Rather than having a single printed circuit assembly that contains all the components in a device, sometimes the design that is used is that of a backplane and blades. For example, directors and large core switches usually implement this technology.

The *backplane* is a circuit board with multiple connectors into which other cards can be plugged. These other cards are usually referred as *blades* or *modules*, but other terms can be used.

If the backplane is in the center of the unit with blades being plugged in at the back and front, then it would usually be referred to as a *midplane*.

7.4 Gigabit transport technology

In Fibre Channel technology, frames are moved from source to destination by using *gigabit transport*, which is a requirement to achieve fast transfer rates. To communicate with gigabit transport, both sides must support this type of communication. This support can be accomplished by installing this feature into the device or by using specially designed interfaces that can convert other communication transport into gigabit transport. The *bit error rate (BER)* allows for only a single bit error to occur once in every 1,000,000,000,000 bits in the Fibre Channel standard. Gigabit transport can be used in copper or fiber optic infrastructure.

Layer 1 of the *open systems interconnection (OSI) model* is the layer at which the physical transmission of data occurs. The unit of transmission at Layer 1 is a *bit*. This section explains some of the common concepts that are at the Layer 1 level.

7.4.1 Fibre Channel cabling

Fibre Channel cabling is one of two forms: *fiber optic cabling* or *copper cabling*. Fiber optic cabling is the usual cabling type, but the introduction of *Fibre Channel over Ethernet (FCoE)* copper cabling is introduced.

Fiber optic cabling is more expensive than copper cabling. The optical components for devices and switches and the cost of any client cabling is typically more expensive to install. However, the higher costs are often easily justified by the benefits of fiber optic cabling.

Fiber optic cabling provides for longer distance and is resistant to the signaling being distorted by electromagnetic interference.

Fiber optic cabling

In copper cabling, electric signals are used to transmit data through the network. The copper cabling is the medium for that electrical transmission. In fiber optic cabling, light is used to transmit the data. Fiber optic cabling is the medium for channeling the light signals between devices in the network.

Two modes of fiber optic signaling are explained in this chapter: *single-mode* and *multimode*. The difference between the modes is the wavelength of the light that is used for the transmission, as illustrated in Figure 7-3.

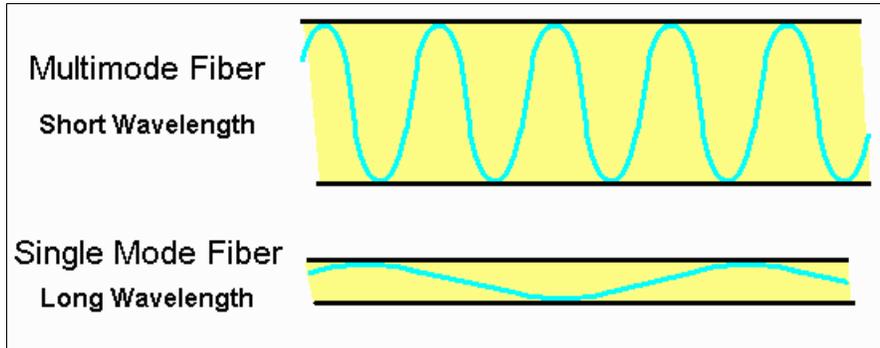


Figure 7-3 Multimode versus single-mode optic signaling

Single-mode fiber

Single-mode fiber (SMF) uses long wavelength light to transmit data and requires a cable with a small core for transmission (Figure 7-3). The core diameter for single-mode cabling is 9 microns in diameter (Figure 7-4).

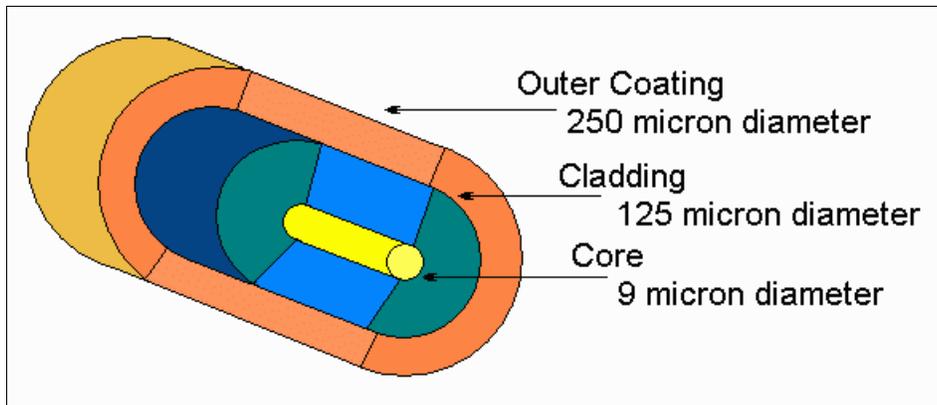


Figure 7-4 Single-mode fiber cable

Multimode fiber

Multimode fiber (MMF) uses short wavelength light to transmit data and requires a cable with a larger core for transmission (Figure 7-3 on page 153). The core diameter for multimode cabling can be 50 or 62.5 microns in diameter, as illustrated in Figure 7-5.

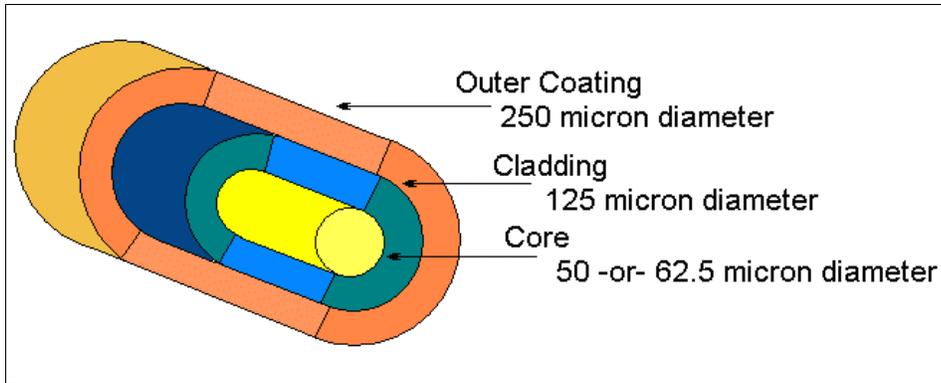


Figure 7-5 Multimode fiber cable

The color of the outer coating is sometimes used to identify if a cable is a multimode or single-mode fiber cable, but the color is not a reliable method. The *Telecommunications Industry Association-598C (TIA-598C)* standard suggests the outer coating to be yellow for single mode fiber and orange for multimode fiber for civilian applications. This guideline is not always implemented, as illustrated in Figure 7-6, which shows a blue cable. The reliable method is to look at the specifications of the cable that are printed on the outer coating of the cabling. See also Figure 7-7 and Figure 7-8.



Figure 7-6 Blue 62.5-micron MMF cable



Figure 7-7 Yellow SMF cable



Figure 7-8 Orange 50-micron MMF cable

Copper cabling

When we refer to *copper cabling*, we mean that the material that is used to transfer the signals are made of copper. The most common copper wire is the twisted-pair cable that is used for normal Ethernet. This type of cabling is explained in more depth in the following section.

Twisted-pair cabling

Twisted-pair copper cabling is a common media for Ethernet networking installations. Twisted-pair cabling is available as *unshielded twisted pair (UTP)* or *shielded twisted pair (STP)*. This shielding helps prevent electromagnetic interference.

Several different categories of twisted-pair cabling are available, as listed in Table 7-1. These categories indicate the signaling capabilities of the cabling.

Table 7-1 TIA/ Electronic Industries Alliance (EIA) cabling categories

TIA/EIA cabling category	Maximum network speeds supported
Cat 1	Telephone or ISDN
Cat 2	4 Mb Token Ring
Cat 3	10 Mb Ethernet
Cat 4	16 Mb Token Ring
Cat 5	100 Mb Ethernet
Cat 5e	1 Gb Ethernet
Cat 6	10 Gb Ethernet Short Distance - 55 m (180 ft.)
Cat 6a	10 Gb Ethernet

The connector that is used for Ethernet twisted-pair cabling is likely the one that most people recognize and associate with networking, the *RJ45 connector*. This connector is shown in Figure 7-9.



Figure 7-9 RJ45 copper connector

Twisted-pair cabling contains four pairs of wire inside the cable, as illustrated in Figure 7-10.

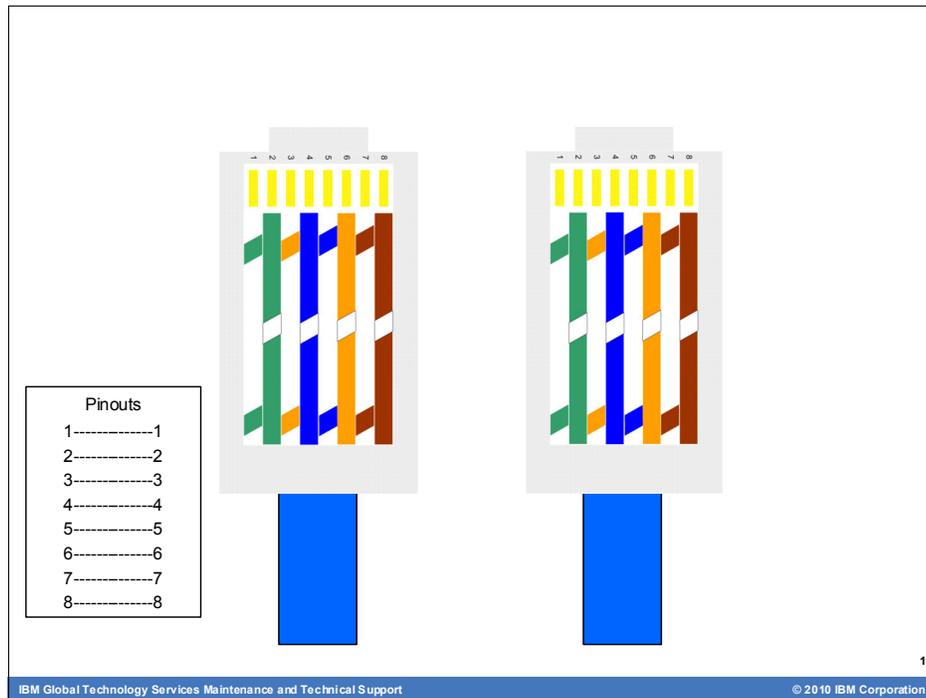


Figure 7-10 Straight through Ethernet cable

An Ethernet that is operating in 10/100 Mb mode, uses only two pairs: pairs 1-2 and 3-6. An Ethernet that is operating in 1 Gb mode, uses all four pairs: pairs 1-2, 3-6, 4-5, and 7-8. Distances up to 100 meters are supported.

Damaged twisted pair: If a twisted-pair cable is damaged so that pair 4-5 or pair 7-8 is unable to communicate, the link is unable to communicate in 1 Gbps mode. If the devices are set to auto negotiate speed, the devices successfully operate in 100 Mbps mode.

Supported maximum distances of cabling segment: The actual maximum distances of a cabling segment that are supported, vary on multiple factors. Examples of these factors include: vendor support, cabling type, electromagnetic interference, and the number of physical connections in the segment.

Twinax cabling

Twinax cables have been used by IBM for many years, but they have recently been reintroduced to the market as a transport media for 10 Gb Ethernet. One of the biggest benefits of a twinax cable is its low power consumption. Also, this cable costs less than standard fiber cables. The downside is the limited capability to connect over long distance.

Connector types

The most common connector type for fiber optic media that is used in networking today, is the *LC connector*, which is shown in Figure 7-11.

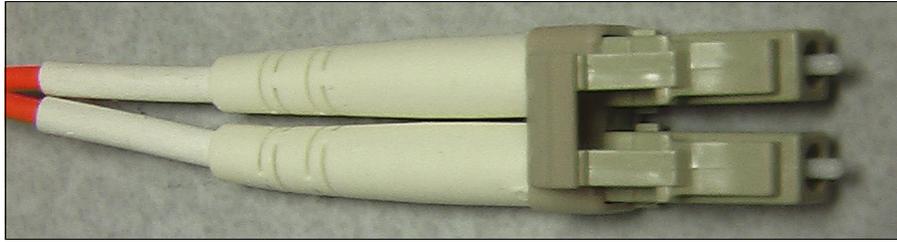


Figure 7-11 LC fiber connector

Other types of connectors are the SC connector (Figure 7-12) and the ST connector (not shown).

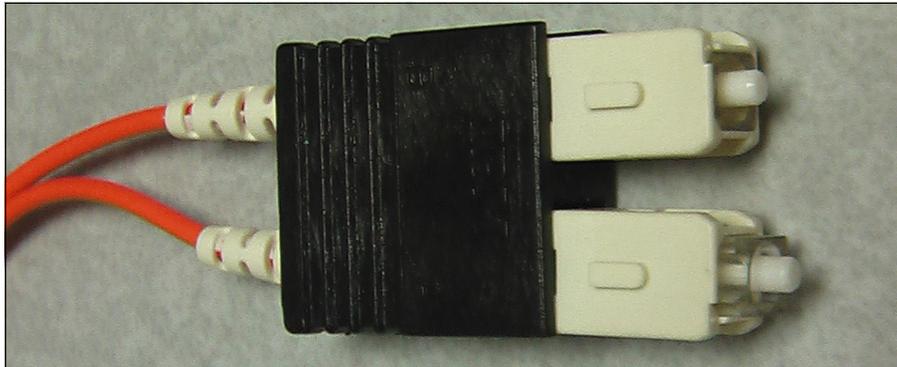


Figure 7-12 SC fiber connector

7.4.2 Transceivers

A *transceiver* or *transmitter/receiver* is the fiber optic port of a device. It is where the fiber optic cables connect. Occasionally, a device might have an integrated transceiver, which limits the flexibility in the type of cabling that can be used. Most devices provide a slot for a modular transceiver to be inserted, providing flexibility for single or multimode implementations to be selected.

Some equipment might use a larger transceiver that is known as a *Gigabit Interface Converter* (GBIC), which is shown in Figure 7-13 on page 158. As technology advances, smaller transceivers are introduced, which provide much higher port density, such as small form-factor pluggables (SFPs), 10 Gigabit SFP+, 10 Gigabit SFP-XFP, and Quad SFP (QSFP).



Figure 7-13 Gigabit Interface Converter (GBIC)

Figure 7-14 compares the different transceivers.



Figure 7-14 From left to right: SFP-MMF, SFP-SMF, SFP+-MMF, XFP-MMF, and XFP-SMF

Figure 7-15 shows a QSFP and cable.



Figure 7-15 QSFP and cable

7.4.3 Host bus adapters

The device that acts as an interface between the fabric of a SAN and either a host or a storage device, is a *host bus adapter (HBA)*.

Fibre Channel host bus adapter

The HBA connects to the bus of the host or storage system. Some devices offer more than one Fibre Channel connection and even have a built-in SFP that can be replaced. The function of the HBA is to convert the parallel electrical signals from the bus into a serial signal to pass to the SAN.

An HBA is shown in Figure 7-16.

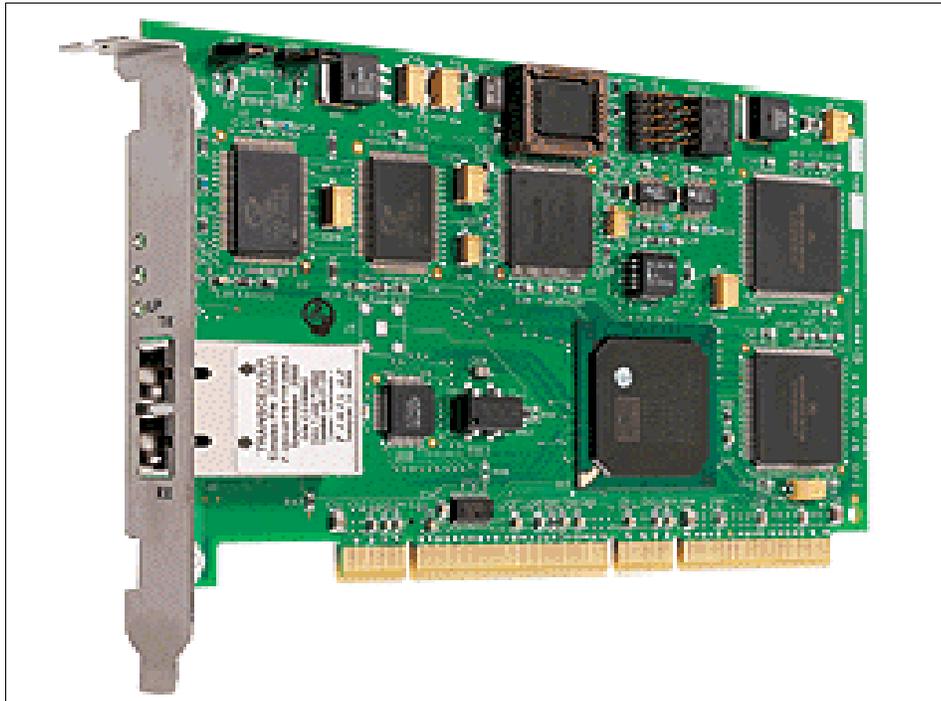


Figure 7-16 Host bus adapter (HBA)

There are several manufacturers of HBAs, and an important consideration when you plan a SAN, is the choice of HBAs. HBAs can have more than one port, they can be supported by some equipment and not others, and they can have parameters that can be used to tune the system. And, they can have many other features. HBAs also have a certain amount of buffer-to-buffer credits. If you are considering using a certain HBA with multiple virtual machines behind it, you need to be aware that the choice of an HBA is a critical one.

Converged Network Adapter host bus adapter

Converged Network Adapters (CNAs) can run both *Converged Enhanced Ethernet (CEE)* and Fibre Channel traffic at the same time. These CNAs combine the functions of an HBA and a Network Interface Card (NIC) on one card. CNAs fully support FCoE protocols and allow Fibre Channel traffic to converge onto 10 Gbps CEE networks. These adapters play a critical role in the FCoE implementation.

When you implement CNAs, they are significant enablers in reducing data center costs by converging data and storage networking. Standard TCP/IP and Fibre Channel traffic can both run on the same high-speed 10 Gbps Ethernet wire, which results in cost savings through reduced requirements for: adapters, switches, cabling, power, cooling, and management. CNAs gained rapid market traction because they deliver excellent performance, help reduce data center TCO, and protect the current data center investment. The cutting-edge 10 Gbps bandwidth can eliminate performance bottlenecks in the I/O path with a 10X data rate improvement versus existing 1 Gbps Ethernet solutions. Additionally, full hardware offload for FCoE protocol processing, reduces system processor utilization for I/O operations. This reduction leads to faster application performance and higher levels of consolidation in virtualized systems.

Figure 7-17 shows a dual port CNA adapter.

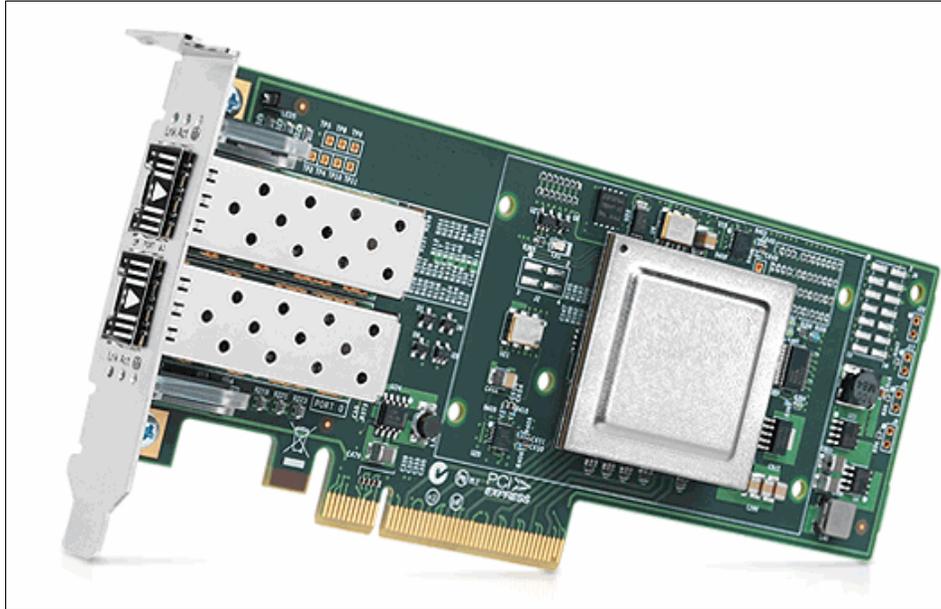


Figure 7-17 Dual port Converged Network Adapter (CNA)

7.5 Inter-switch links

A link that joins a port on one switch to a port on another switch (referred to as *E_Ports*) is called an *inter-switch link (ISL)*.

ISLs carry frames that originate from the node ports and frames that are generated within the fabric. The frames that are generated within the fabric, serve as control, management, and support for the fabric.

Before an ISL can carry frames that originate from the node ports, the joining switches need to go through a synchronization process on which the operating parameters are interchanged. If the operating parameters are not compatible, the switches cannot join, and the ISL becomes *segmented*. Segmented ISLs cannot carry traffic that originates on node ports, but they can still carry management and control frames.

There is also the possibility to connect an *E_Port* to a Fibre Channel router or a switch with embedded routing capabilities. Then the port becomes an *EX-port* on the router side. Brocade calls these ports an *inter-fabric link (IFL)*; however, Cisco uses *TE_Port (trunked E_Port)*, also known as an *EISL*, that allows traffic (from multiple VSANs) to be routed through that link.

7.5.1 Cascading

Expanding the fabric is called *switch cascading*, or *cascading*. Cascading is basically interconnecting Fibre Channel switches and directors by using ISLs. The cascading of switches provides the following benefits to a SAN environment:

- ▶ The fabric can be seamlessly extended. Additional switches can be added to the fabric, without powering down existing fabric.
- ▶ You can easily increase the distance between various SAN participants.

- ▶ By adding more switches to the fabric, you increase connectivity by providing more available ports.
- ▶ Cascading provides high resilience in the fabric.
- ▶ With ISLs, you can increase the bandwidth. The frames between the switches are delivered over all available data paths. Therefore, the more ISLs you create, the faster the frame delivery is. However, careful consideration must be employed to ensure that a bottleneck is not introduced.
- ▶ When the fabric grows, the name server is fully distributed across all the switches in the fabric.
- ▶ With cascading, you also provide greater fault tolerance within the fabric.

7.5.2 Hops

When Fibre Channel traffic traverses an ISL, this process is known as a *hop*. Or, to state it another way, traffic that is going from one E_Port over an ISL to another E_Port, is one hop. As stated in 7.5, “Inter-switch links” on page 161, ISLs are made from connecting an E_Port to an E_Port. Figure 7-18 on page 163 shows an illustration of the hop count from server to storage.

There is a hop count limit. This limit is set by the fabric operating system and is used to derive a frame hold time value for each switch. This value is the maximum amount of time that a frame can be held in a switch before it is dropped, or the fabric indicates that it is too busy. The hop count limits must be investigated and considered in any SAN design work because it has a major effect on the proposal.

7.5.3 Fabric shortest path first

Although this next topic is not a physical component, it is wise to introduce the concept of *fabric shortest path first (FSPF)* at this stage. According to the FC-SW-2 standard, FSPF is a link state path selection protocol. FSPF tracks the links on all switches in the fabric (in routing tables) and associates a cost with each link. The protocol computes paths from a switch to all the other switches in the fabric. This process is done by adding the cost of all the links that are traversed by the path, and choosing the path that minimizes the cost; that is, the shortest link.

For example, as shown in Figure 7-18 on page 163, if a server needs to connect to its storage through multiple switches, FSPF routes all traffic from this server to its storage through switch A directly to switch C. This path is taken because it has a lower cost than traveling through more hops via switch B.

Figure 7-18 shows hops in a fabric.

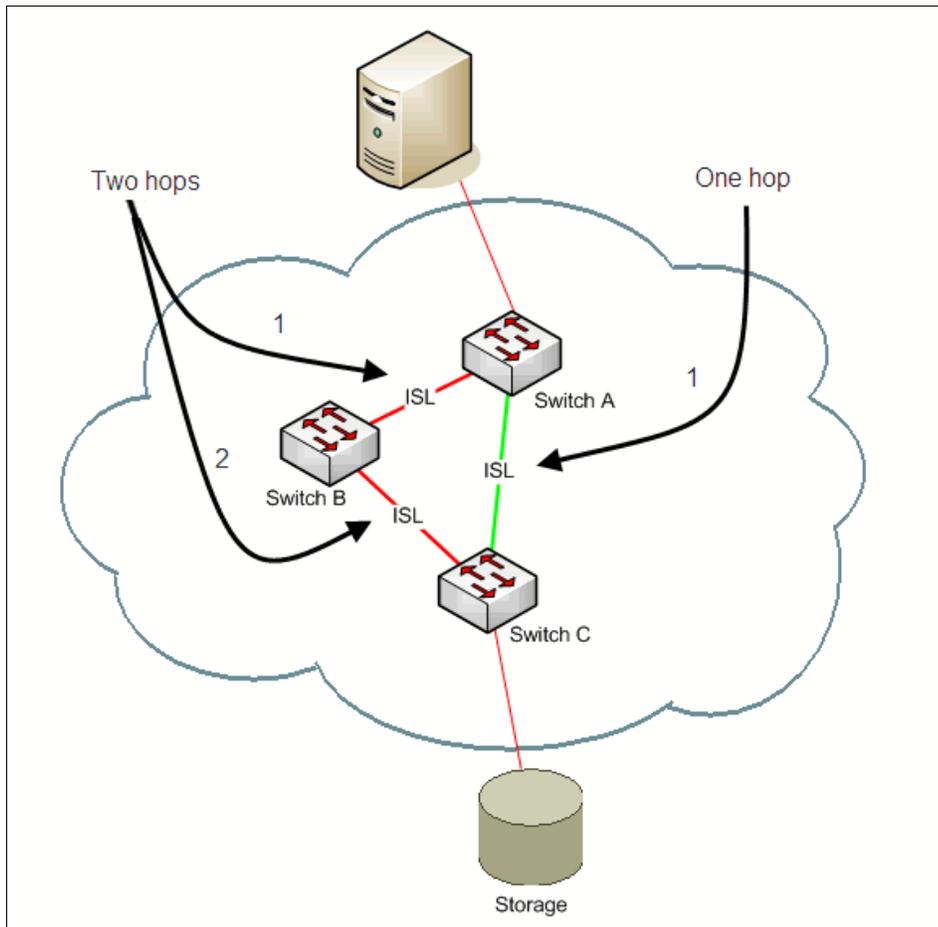


Figure 7-18 Hops explained

FSPF is based on the hop count cost.

The collection of link states, including the cost, of all the switches in a fabric constitutes the *topology database*, or *link state database*. The topology database is kept in all switches in the fabric, and they are maintained and synchronized to each other. There is an *initial database synchronization*, and an *update mechanism*. The initial database synchronization is used when a switch is initialized, or when an ISL comes up. The update mechanism is used when there is a link state change; for example, when an ISL is going down or coming up, and on a periodic basis. This mechanism ensures consistency among all switches in the fabric.

7.5.4 Non-blocking architecture

To support highly performing fabrics, the fabric components, switches, or directors must be able to move around data. These components must move the data without any affect to other ports, targets, or initiators that are on the same fabric. If the internal structure of a switch or director cannot do so without an effect, we end up with blocking.

Blocking

Blocking means that the data does not get to the destination. Blocking is not the same as *congestion* because with congestion, data is still being delivered, but with a delay. Currently, almost all Fibre Channel switches are created by using non-blocking architecture.

Non-blocking

A *non-blocking* architecture is now most commonly used by switch vendors. Non-blocking switches enable multiple connections that are traveling through the switch at the same time. Figure 7-19 illustrates this concept.

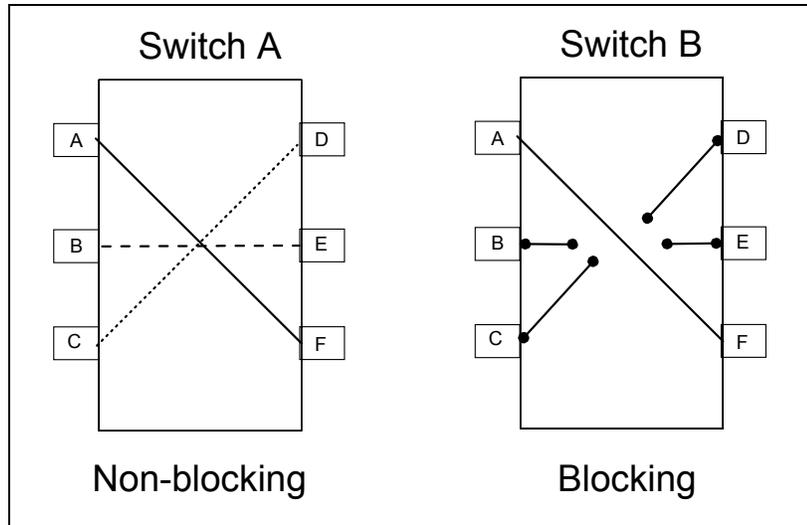


Figure 7-19 Non-blocking and blocking switching

As Figure 7-19 shows, port A in the non-blocking Switch A, speaks to port F; port B speaks to E, and C speaks to D, without any form of suspension of communication or delay. That is, the communication is not blocked. In the blocking switch, Switch B, while port A is speaking to F, is stopped or blocked from all other communication and does not continue until A finishes talking to F.

7.5.5 Latency

Typically, in the SAN world, *latency* is the time that it takes for a Fibre Channel frame to traverse the fabric. When we describe SAN, the latency in a SAN is rarely considered because it is in the low microsecond range. This concept is sometimes confused with *disk latency*, which is the measure on how fast or slow a storage target completes a read or write request that is sent from the server. However, when we describe very long distances, all latency, both storage and SAN, plays a significant role.

The more ISLs there are, the more the latency there is because the Fibre Channel frame must traverse the fabric by using ISLs. By fabric, we mean the Fibre Channel components and any latency discussion that is related to the SAN. Usually the time that is taken is expressed in microseconds, which gives an indication as to the performance characteristics of the SAN fabric. It is often given at a switch level, and sometimes at a fabric level.

7.5.6 Oversubscription

Another aspect of data flow is *fan-in ratio* (also called the *oversubscription ratio* and frequently the *fan-out ratio*, from the storage device perspective), both in terms of host ports to target ports, and device to ISL. This ratio is the number of device ports that need to share a single port.

For example, two servers each equipped with a 4 Gb port ($4+4=8$ Gb) are both communicating with a storage device through a single 4 Gb port, giving a 2:1 ratio. In other

words, the total theoretical input is higher than what that port can provide. Figure 7-20 on page 166 shows a typical oversubscription through an ISL.

Oversubscription can occur on storage device ports and ISLs. When you design a SAN, it is important to consider the possible traffic patterns to determine the possibility of oversubscription. An oversubscription might result in degraded performance. Oversubscription of an ISL can be overcome by adding an ISL between the switches to increase the bandwidth. Oversubscription to a storage device might be overcome by adding more ports from the storage device to the fabric.

Oversubscription: There is a difference on how vendors practice utilization on their stated overall bandwidth per chassis, though both use storage port and ISL oversubscription. Verify oversubscription preferred practices with your switch vendor.

7.5.7 Congestion

When oversubscription occurs, it leads to a condition called *congestion*. When a node is unable to use as much bandwidth as it wants, because of contention with another node, then there is congestion. A port, link, or fabric can be congested. This condition normally has a direct effect on the application in forms of poor performance.

Congestion can be difficult to detect because it can also be directly related to buffer-to-buffer credit starvation in the switch port. Therefore, when you look at the data throughput from the switch, it seems like normal or less traffic is flowing through the ports. However, the server I/O is unable to perform because the data cannot be transported because of a lack of buffer-to-buffer credits.

7.5.8 Trunking or port-channeling

One means of delivering high availability at the network level is aggregation of multiple physical ISLs into a single logical interface. This aggregation allows you to provide link redundancy, greater aggregated bandwidth, and load balancing. Cisco calls this technology *port channeling*, others call it, *trunking*.

We illustrate the concepts of trunking in Figure 7-20.

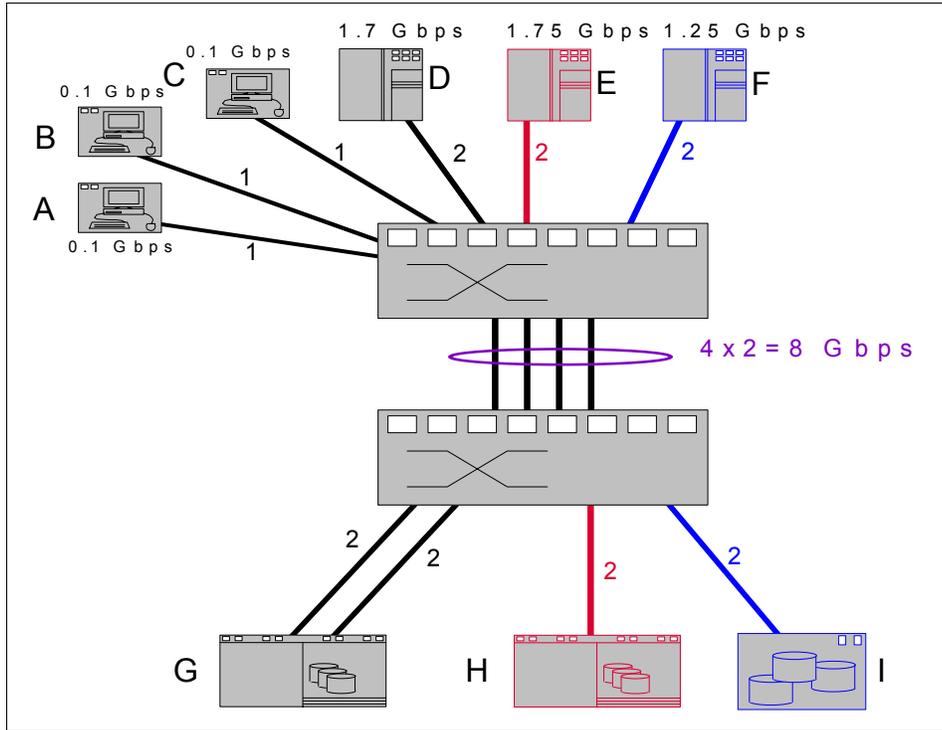


Figure 7-20 Trunking

In Figure 7-20, there are six computers that are accessing three storage devices. Computers A, B, C, and D are communicating with Storage G. Server E is communicating with storage H, and server F uses disks in storage device I.

The speeds of the links are shown in Gbps, and the target throughput for each computer is shown. If we allow FSPF alone to decide the routing, we might have a situation where servers D and E are both using the same ISL. This leads to oversubscription, and hence congestion because 1.7 added to 1.75 is greater than 2.

If all of the ISLs are gathered together into a trunk, then effectively they can be seen as a single, large ISL. In effect, they appear to be an 8 Gbps ISL. This bandwidth is greater than the total requirement of all of the servers. In fact, the nodes require an aggregate bandwidth of 5 Gbps. Therefore, you might even suffer a failure of one of the ISLs and still have enough bandwidth to satisfy their needs.

When the nodes come up, FSPF simply sees one route, and they are all assigned a route over the same trunk. The fabric operating systems in the switches share the load over the actual ISLs, which combine to make up the trunk. This process is done by distributing frames over the physical links, and then reassembling them at the destination switch so that an in-order delivery can be assured, if necessary. And to FSPF, a trunk is displayed as a single, low-cost ISL.



Management

Management is one of the key issues behind the concept of *infrastructure simplification*. The ability to manage heterogeneous systems at different levels as though they were a fully integrated infrastructure, is a goal that many vendors and developers are striving to achieve. Another goal is to offer the system administrator a unified view of the whole storage area network (SAN).

In this chapter, we look at some of the initiatives that have been developed, and are in the process of being developed in the field of SAN management. These solutions will incrementally smooth the way towards infrastructure simplification.

8.1 Management principles

SAN management systems typically comprise a set of multiple-level software components that provide tools for monitoring, configuring, controlling (performing actions), diagnosing, and troubleshooting a SAN. In this section, we briefly describe the different types and levels of management that can be found in a typical SAN implementation. We also describe the efforts that are being made towards the establishment of open and general-purpose standards for building interoperable, manageable components.

In this section, it is also shown that despite these efforts, the reality of a “one pill cures all” solution is a long way off. Typically, each vendor and each device has its own form of software and hardware management techniques. These techniques are usually independent of each other. To pretend that there is one SAN management solution that will provide a single point of control, capable of performing every possible action, would be premature at this stage.

This book does not aim at fully describing each vendor’s own standards, but at presenting the reader with an overview of the myriad of possibilities that they might find in the IT environment. That stated, the high-level features of any SAN management solution are likely to include most, if not all, of the following functions:

- ▶ Capacity management
- ▶ Device management
- ▶ Fabric management
- ▶ Pro-active monitoring
- ▶ Fault isolation and troubleshooting
- ▶ Centralized management
- ▶ Remote management
- ▶ Performance management
- ▶ Security and standard compliant

8.1.1 Management types

There are essentially two philosophies that are used for building management mechanisms: *in-band management* and *out-of-band management*. They can be defined as:

In-band management

This means that the management data, such as status information, action requests, and events flows through the same path as the storage data itself.

Out-of-band management

This means that the management data flows through a dedicated path, therefore not sharing the same physical path that is used by the storage data.

In-band and out-of-band models are illustrated as shown in Figure 8-1. These models are not mutually exclusive. In many environments, a combination of both models might be wanted.

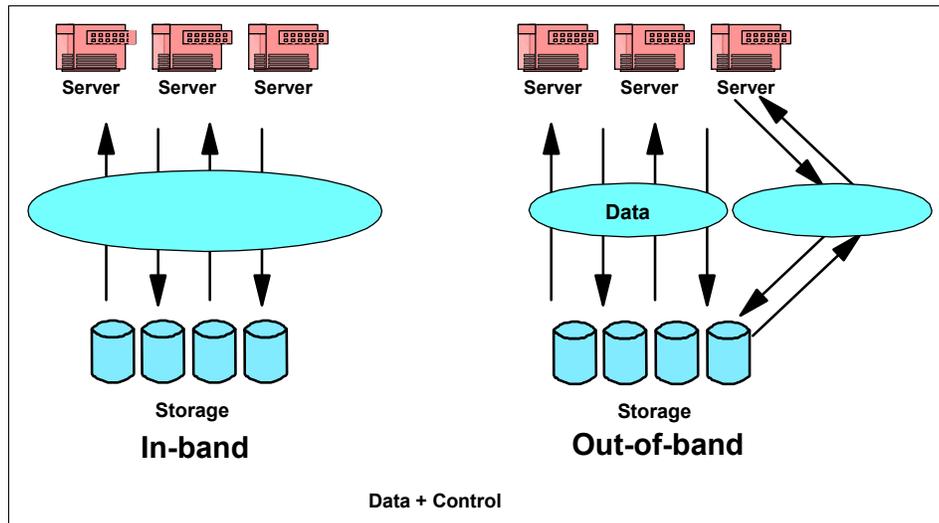


Figure 8-1 In-band and out-of-band models

The in-band approach is simple to implement, requires no dedicated channels (other than LAN connections), and has inherent advantages, such as the ability for a switch to initiate a SAN topology map with queries to other fabric components. However, in the event of a failure of the Fibre Channel transport itself, the management information cannot be transmitted. Therefore, access to devices is lost, as is the ability to detect, isolate, and recover from network problems. This problem can be minimized by a provision of redundant paths between devices in the fabric.

In-band management allows attribute inquiries on storage devices and configuration changes for all elements of the SAN. Since in-band management is performed over the SAN itself, administrators are not required to manage any additional connections.

Conversely, out-of-band management does not rely on the storage network; its main advantage is that management commands and messages can be sent even if a loop or fabric link fails. Integrated SAN management facilities are more easily implemented. However, unlike in-band management, it cannot automatically provide SAN topology mapping.

In summary, we can say that in-band management has these main advantages:

- Device installation, configuration, and monitoring
- Inventory of resources on the SAN
- Automated component and fabric topology discovery
- Management of the fabric configuration, including zoning configurations
- Health and performance monitoring

Out-of-band management has these main advantages:

- ▶ It keeps management traffic out of the Fibre Channel, so it does not affect the business-critical data flow on the storage network.
- ▶ It makes management possible, even if a device is down.
- ▶ It is accessible from anywhere in the routed network.

8.1.2 Connecting to storage area network management tools

A usual way of connecting to a SAN device (by SAN device, we mean Fibre Channel switches and storage devices that are connected to a SAN) is by connecting through the Ethernet to a storage management device on a network segment that is intended for storage devices. This is shown in Figure 8-2.

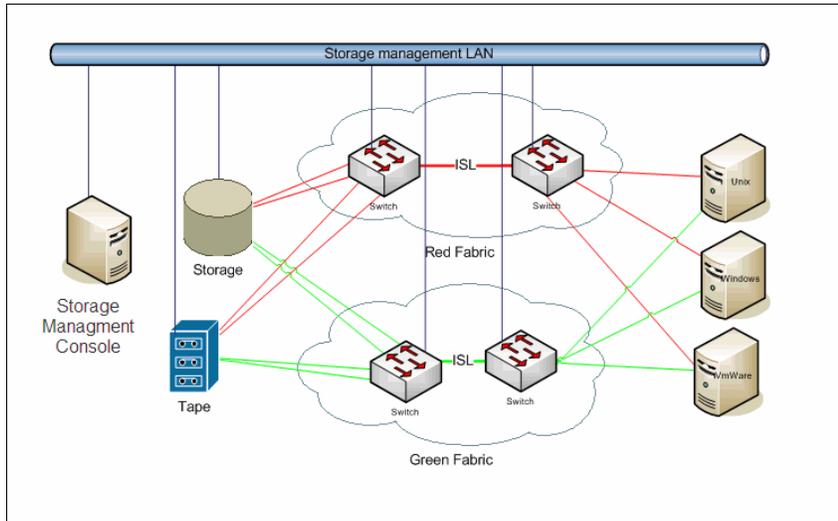


Figure 8-2 Storage Management network

The SAN storage level consists of the storage devices that integrate the SAN, such as disks, disk arrays, tapes, and tape libraries. As the configuration of a storage resource must be integrated with the configuration of the server's logical view of them, the SAN storage level management can also span both storage resources and servers.

8.1.3 Storage area network fault isolation and troubleshooting

In addition to providing tools for monitoring and configuring a SAN, one of the key benefits that a well-designed management mechanism can bring is the ability to efficiently detect, diagnose, and solve problems in a SAN.

There are many tools to collect the necessary data to perform problem determination and problem source identification (PD/PSI) in a SAN. Generally speaking, these tools offer the ability to:

- ▶ Monitor SAN health
- ▶ Report failures
- ▶ Monitor and identify storage devices
- ▶ Monitor the fabric for failures or imminent bottlenecks
- ▶ Interpret message and error logs
- ▶ Send SNMP traps or syslog messages

Although a well-designed management system can provide invaluable facilities, an easy-to-troubleshoot SAN still relies heavily on a good design, and on good documentation. In terms of PD/PSI, this means that configuration design information is understandable, available at any support level, and is always updated regarding the latest configuration. There must also be a database where all the information about connections, naming conventions, device serial numbers, WWN, zoning, system applications, and so on, is safely stored. Last,

but not least, there should be a responsible person that is in charge of maintaining this infrastructure, and monitoring the SAN health status.

8.2 Management interfaces and protocols

In this section, we present the main protocols and interfaces that have been developed to support management mechanisms.

8.2.1 Storage Networking Industry Association initiative

The Storage Networking Industry Association (SNIA) is using its Storage Management Initiative (SMI) to create and promote the adoption of a highly functional interoperable management interface for multivendor storage networking products. The SNIA strategic imperative is to have all storage managed by the SMI interface. The adoption of this interface allows the focus to switch to the development of value-add functionality. IBM is one of the industry vendors that is promoting the drive towards this vendor-neutral approach to SAN management.

In 1999, the SNIA and Distributed Management Task Force (DMTF) introduced open standards for managing storage devices. These standards use a common protocol that is called the *Common Information Model (CIM)* to enable interoperability. The web-based version of CIM (WBEM) uses XML to define CIM objects and process transactions within sessions. This standard proposes a CIM Object Manager (CIMOM) to manage CIM objects and interactions. CIM is used to define objects and their interactions. Management applications then use the CIM object model and XML over HTTP to provide for the management of storage devices. This enables central management by using open standards.

SNIA uses the xmlCIM protocol to describe storage management objects and their behavior. CIM allows management applications to communicate with devices by using object messaging encoded in xmlCIM.

The *Storage Management Interface Specification (SMI-S)* for SAN-based storage management provides basic device management, support for copy services, and virtualization. As defined by the standard, the CIM services are registered in a directory to make them available to device management applications and subsystems.

For more information about SMI-S, see this website:

<http://www.snia.org>

Open storage management with the Common Information Model

SAN management involves configuration, provisioning, logical volume assignment, zoning, and logical unit number (LUN) masking. Management also involves monitoring and optimizing performance, capacity, and availability. In addition, support for continuous availability and disaster recovery requires that device copy services are available as a viable failover and disaster recovery environment. Traditionally, each device provides a command-line interface (CLI) and a graphical user interface (GUI) to support these kinds of administrative tasks. Many devices also provide proprietary application programming interfaces (APIs) that allow other programs to access their internal capabilities.

For complex SAN environments, management applications are now available that make it easier to perform these kinds of administrative tasks over various devices.

The *Common Information Model (CIM)* interface and the *Storage Management Initiative Specification (SMI-S)* object model that is adopted by the Storage Networking Industry Association (SNIA), provide a standard model for accessing devices. This ability allows management applications and devices from various vendors to work with each other's products. This flexibility means that clients have more choice as to which devices work with their chosen management application, and which management applications they can use with their devices.

IBM embraces the concept of building open standards-based storage management solutions. IBM management applications are designed to work across multiple vendors' devices, and devices are being CIM-enabled to allow them to be controlled by other vendors' management applications.

Common Information Model Object Manager

The SMI-S standard designates that either a proxy or an embedded agent can be used to implement CIM. In each case, the CIM objects are supported by a *CIM Object Manager (CIMOM)*. External applications communicate with CIM via HTTP to exchange XML messages, which are used to configure and manage the device.

In a proxy configuration, the CIMOM runs outside of the device and can manage multiple devices. In this case, a *provider* component is installed into the CIMOM to enable the CIMOM to manage specific devices.

The providers adapt the CIMOM to work with different devices and subsystems. In this way, a single CIMOM installation can be used to access more than one device type, and more than one device of each type on a subsystem.

The CIMOM acts as a catcher for requests that are sent from storage management applications. The interactions between catcher and sender use the language and models that are defined by the SMI-S standard. This allows storage management applications, regardless of vendor, to query status and perform command and control using XML-based CIM interactions.

IBM developed its storage management solutions that are based on the CIMOM architecture, as shown in Figure 8-3.

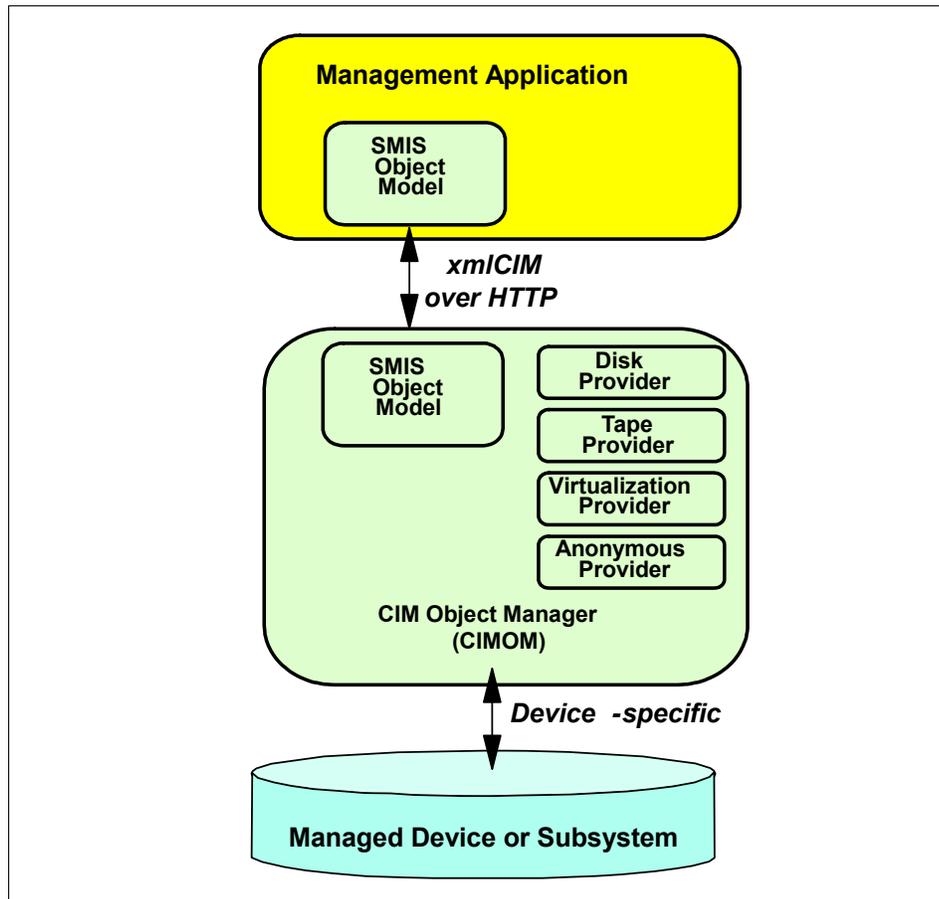


Figure 8-3 CIMOM component structure

8.2.2 Simple Network Management Protocol

Simple Network Management Protocol (SNMP), which is an IP-based protocol, has a set of commands for obtaining the status and setting the operational parameters of target devices. The SNMP management platform is called the *SNMP manager*, and the managed devices have the SNMP agent loaded. Management data is organized in a hierarchical data structure that is called the *Management Information Base (MIB)*. These MIBs are defined and sanctioned by various industry associations. The objective is for all vendors to create products in compliance with these MIBs so that inter-vendor interoperability at all levels can be achieved. If a vendor wants to include more device information that is not specified in a standard MIB, then that is usually done through MIB extensions.

This protocol is widely supported by LAN/WAN routers, gateways, hubs, and switches, and is the predominant protocol that is used for multivendor networks. Device status information (vendor, machine serial number, port type and status, traffic, errors, and so on) can be provided to an enterprise SNMP manager. A device can generate an alert by SNMP, in the event of an error condition. The device symbol, or icon, which is displayed on the SNMP manager console, can be made to turn red or yellow, or any warning color, and messages can be sent to the network operator.

Out-of-band developments

SNMP MIBs are being implemented for SAN fabric elements that allow out-of-band monitoring. The ANSI Fibre Channel Fabric Element MIB provides significant operational and configuration information about individual devices. The emerging Fibre Channel Management MIB provides more link table and switch zoning information that can be used to derive information about the physical and logical connections between individual devices.

8.2.3 Service Location Protocol

The *Service Location Protocol (SLP)* provides a flexible and scalable framework for providing hosts with access to information about the existence, location, and configuration of networked services. Traditionally, users had to find devices by knowing the name of a network host that is an alias for a network address. SLP eliminates the need for a user to know the name of a network host supporting a service. Rather, the user supplies the wanted type of service and a set of attributes that describe the service. Based on that description, the Service Location Protocol resolves the network address of the service for the user.

SLP provides a dynamic configuration mechanism for applications in local area networks. Applications are modeled as clients that need to find servers that are attached to any of the available networks within an enterprise. For cases where there are many different clients and services available, the protocol is adapted to use the nearby Directory Agents that offer a centralized repository for advertised services.

8.2.4 Vendor-specific mechanisms

These are some of the vendor-specific mechanisms that have been deployed by major SAN device providers.

Application programming interface

As you know, there are many SAN devices from many different vendors and everyone has their own management and configuration software. In addition, most of them can also be managed via a command-line interface (CLI) over a standard telnet connection, where an IP address is associated with the SAN device, or they can be managed by an RS-232 serial connection.

With different vendors and the many management and configuration software tools, we have a number of different products to evaluate, implement, and learn. In an ideal world, there would be one product to manage and configure all of the functions and features on the SAN platform.

Application programming interfaces (APIs) are one way to help this simplification become a reality. Some vendors make the API of their product available for other vendors to make it possible for common management in the SAN. This allows for the development of upper level management applications capable of interacting with multiple-vendor devices and offering the system administrator a single view of the SAN infrastructure.

Common Agent Services

Common Agent Services is a component that is designed to provide a way to deploy agent code across multiple user machines or application servers throughout an enterprise. The agents collect data from and perform operations on managed resources for Fabric Manager.

The Common Agent Services agent manager provides authentication and authorization and maintains a registry of configuration information about the agents and resource managers in

the SAN environment. The resource managers are the server components of products that manage agents that are deployed on the common agent. Management applications use the services of the agent manager to communicate securely with and to obtain information about the computer systems that are running the common agent software, referred to in this document, as the *agent*.

Common Agent Services also provide common agents to act as containers to host product agents and common services. The common agent provides remote deployment capability, shared machine resources, and secure connectivity.

Common Agent Services consists of these subcomponents:

- ▶ Agent manager

The agent manager is the server component of the Common Agent Services that provides functions that allow clients to get information about agents and resource managers. It enables secure connections between managed endpoints, maintains the database information about the endpoints and the software that is running on those endpoints, and processes queries against that database from resource managers. It also includes a registration service, which handles security certificates, registration, tracking of common agents and resource managers, and status collection and forwarding.

- ▶ Common agent

The common agent is a common container for all the subagents to run within. It enables multiple management applications to share resources when managing a system.

- ▶ Resource manager

Each product that uses Common Agent Services has its own resource manager and subagents. For example, Tivoli Provisioning Manager has a resource manager and subagents for software distribution and software inventory scanning.

Figure 8-4 shows the Common Agent topology.

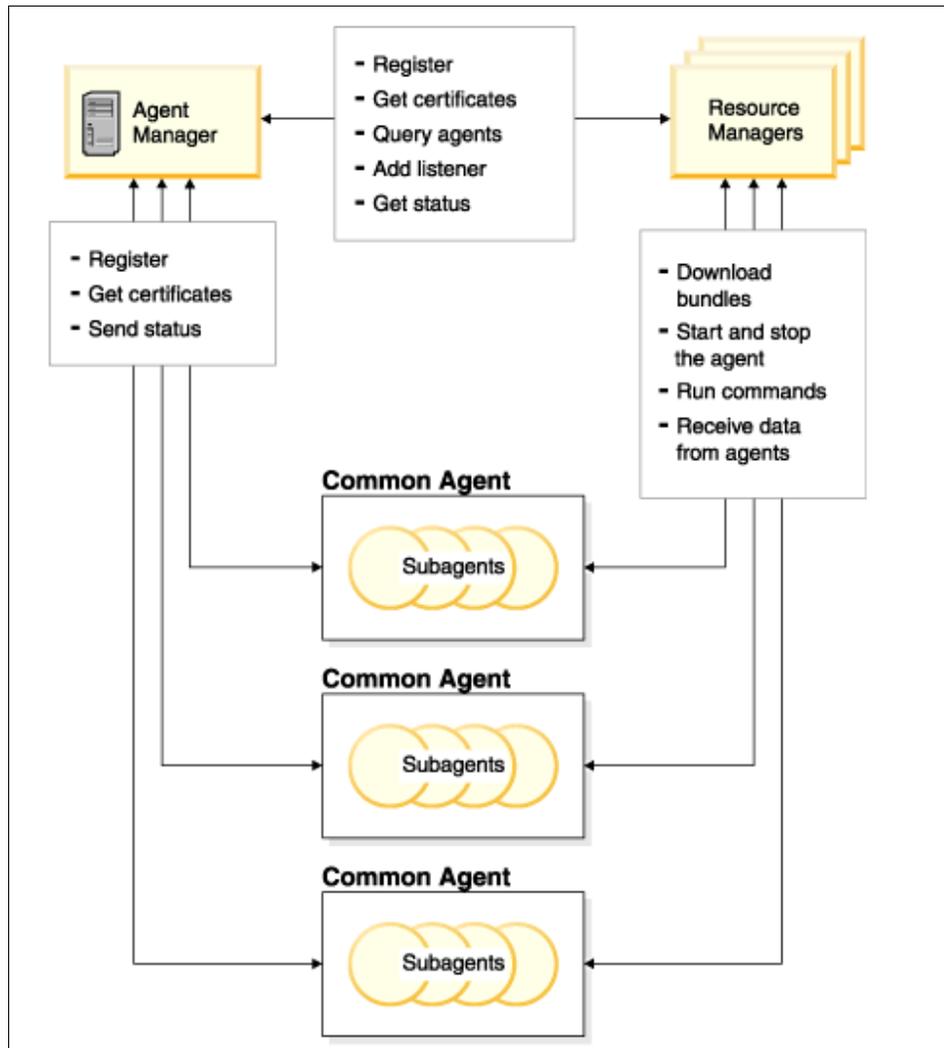


Figure 8-4 Common Agent Services

8.3 Management features

SAN management requirements are typified by having a common purpose, but are implemented in different fashions by the different vendors. Some prefer to use web browser interfaces, some prefer to use embedded agents, and some prefer to use the CLI, and some use a combination of all. There is no right or wrong way. Usually the selection of SAN components is based on a combination of what the hardware and software provide, not on the ease of use of the management solution. The high-level features of any SAN management solution are likely to include most of the following benefits:

- ▶ Cost effectiveness
- ▶ Open approach
- ▶ Device management
- ▶ Fabric management
- ▶ Pro-active monitoring
- ▶ Fault isolation and troubleshooting
- ▶ Centralized management
- ▶ Remote management
- ▶ Adherence to standards
- ▶ Resource management
- ▶ Secure access
- ▶ Standards compliant

8.3.1 Operations

When we describe management, it automatically includes the operational aspects of the environment. The SAN administrators are responsible for all configuration of the SAN switches.

It is common that the initial design and creation of a SAN environment includes only a handful of servers and few storage systems. However, the environment grows and new technology needs to be added. At this stage, it tends to get more complicated. That is why it is necessary to ensure that there is comprehensive documentation that documents all aspects of the environment, and it needs to be reviewed regularly to ensure that it is current.

Some of the standards and guidelines that need to be documented, include:

- ▶ Zoning standards:
 - How to create zones using preferred practices
 - Naming standards that are used in the SAN configuration
 - Aliases used
- ▶ Volume / LUN allocation standards:
 - Volume characteristics and their uses
 - Allocation rules
- ▶ Incident and problem guidelines: How to react in case of an incident.
- ▶ Roles and responsibilities: Roles and responsibilities within the team.
- ▶ SAN and storage installation preferred practices: Agreed process to install and configure the equipment.
- ▶ SAN and storage software and firmware upgrade roadmaps:
 - High-level overview of how to ensure that the environment is kept current
 - Change schedules

- ▶ Monitoring and performance guidelines: What is monitored and how are exceptions handled.

8.4 IBM Tivoli Storage Productivity Center

The *IBM Tivoli Storage Productivity Center* is an integrated hardware and software solution that provides a single point of entry for managing storage devices and other components of your data storage infrastructure. We describe a basic product explanation.

This product comes in several different varieties and we briefly describe them in this chapter.

- ▶ IBM Tivoli Storage Productivity Center Basic Edition
- ▶ IBM Tivoli Storage Productivity Center for Data
- ▶ IBM Tivoli Storage Productivity Center for Disk
- ▶ Tivoli Storage Productivity Center for Disk Select
- ▶ Tivoli Storage Productivity Center for Replication
- ▶ IBM Tivoli Storage Productivity Center Standard Edition
- ▶ IBM System Storage Productivity Center (SSPC)

For detailed information about each product, see this website:

<http://www.ibm.com/systems/storage/software/center/index.html>

The Tivoli Storage Productivity Center is an open storage infrastructure management solution that is designed to help:

- ▶ Reduce the effort of managing complex, heterogeneous storage infrastructures
- ▶ Improve storage capacity utilization
- ▶ Improve administrative efficiency

The Tivoli Storage Productivity Center provides reporting capabilities, identifying data usage and its location, and provisioning. It also provides a central point of control to move the data based on business needs to more appropriate online or auxiliary storage. The productivity center also centralizes the management of storage infrastructure capacity, performance, and availability.

Figure 8-5 shows the Tivoli Storage Productivity Center architecture overview.

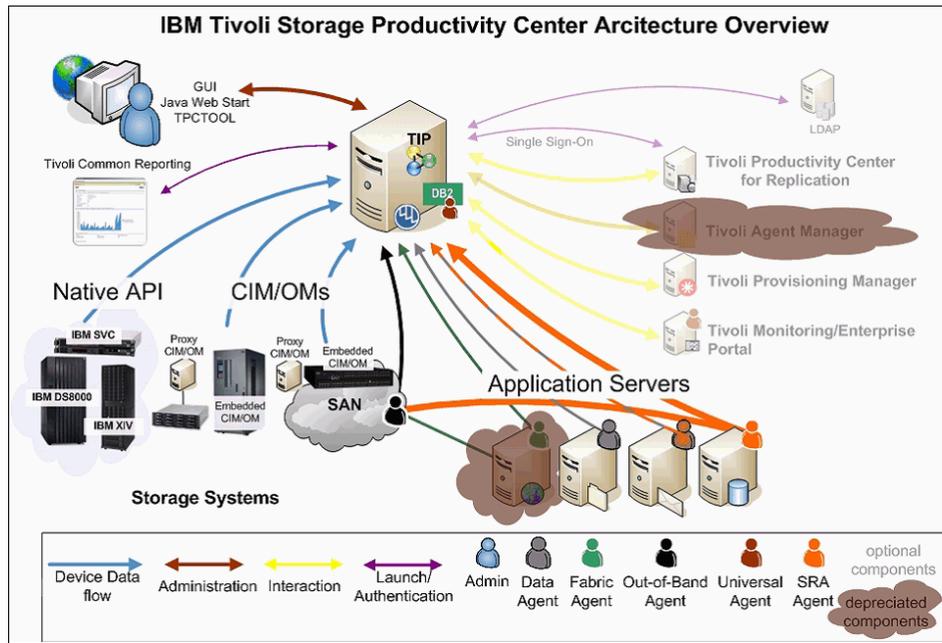


Figure 8-5 Tivoli Storage Productivity Center architecture

8.4.1 IBM Tivoli Storage Productivity Center for Data

IBM Tivoli Storage Productivity Center for Data provides over 400 enterprise-wide reports, monitoring and alerts, policy-based action and file system capacity automation in a heterogeneous environment. Tivoli Storage Productivity Center for Data is designed to help improve capacity utilization of file systems and databases and add intelligence to data protection and retention practices.

8.4.2 IBM Tivoli Storage Productivity Center for Disk

IBM Tivoli Storage Productivity Center for Disk is designed to provide storage device configuration and management from a single console. It includes performance capabilities to help monitor and manage performance, and measure service levels by storing received performance statistics into database tables for later use. Policy-based automation enables event action that is based on business policies. It sets performance thresholds for the devices that are based on selected performance metrics, generating alerts when those thresholds are exceeded. Tivoli Storage Productivity Center for Disk helps simplify the complexity of managing multiple SAN-attached storage devices.

8.4.3 IBM Tivoli Storage Productivity Center for Disk Select

IBM Tivoli Storage Productivity Center for Disk Select is designed to help reduce the complexity of managing storage devices by allowing administrators to configure, manage, and monitor performance of their entire storage infrastructure from a single console. Tivoli Storage Productivity Center for Disk Select provides the same features and functions as Tivoli Storage Productivity Center for Disk, but is limited to managing IBM System Storage DS3000, DS4000®, DS5000, and Storwize V7000 and IBM XIV® devices. It provides performance management, monitoring and reporting for these devices.

8.4.4 IBM Tivoli Storage Productivity Center Basic Edition

IBM Tivoli Storage Productivity Center Basic Edition is designed to provide device management services for IBM System Storage DS3000, DS4000, DS5000, and DS8000 products. This edition also provides services for IBM Storwize V7000, IBM SAN Volume Controller, IBM XIV, and heterogeneous storage environments. Productivity Center Basic Edition is a management option available with IBM Storage hardware acquisitions. This tool provides storage administrators a simple way to conduct device management for multiple storage arrays and SAN fabric components from a single integrated console.

This console is also the base of operations for the IBM Tivoli Storage Productivity Center Suite:

- ▶ Contains a storage topology viewer for a “big picture” perspective
- ▶ Offers asset and capacity reporting to improve storage utilization
- ▶ Assists with problem determination
- ▶ Can reduce storage complexity and improve interoperability
- ▶ Automates device discovery
- ▶ Extends existing device utilities
- ▶ Aids with server consolidation
- ▶ Storage Topology Viewer
- ▶ Ability to monitor, alert, report, and provision storage
- ▶ Status dashboard
- ▶ IBM System Storage DS8000 GUI integration with TPC Basic Edition

8.4.5 IBM Tivoli Storage Productivity Center Standard Edition

IBM Tivoli Storage Productivity Center Standard Edition is one of the industry’s most comprehensive storage resource management solutions that combines the consolidated benefits of the following three components as one bundle at a reduced price:

- ▶ Tivoli Storage Productivity Center for Data
- ▶ Tivoli Storage Productivity Center for Disk
- ▶ Tivoli Storage Productivity Center Basic Edition

In addition to the benefits and features of the Data, Disk, and Basic Editions, Tivoli Productivity Center Standard Edition offers more management, control, and performance reporting for the Fibre Channel SAN infrastructure.

Figure 8-6 shows the difference between the basic and standard edition of IBM Tivoli Storage Productivity Center.

Function	DS Storage Manager	SVC Admin Console	TPC Basic Edition	TPC Standard Edition
Storage Infrastructure Configuration/Status Reporting				
Device Discovery/Configuration	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Manage multiple DS8000s / SVCs from 1 User Interface	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Topology Viewer and Storage Health Management			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Provisioning, including Fabric zoning and Disk LUN assignment			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Configuration Management – Highlight configuration changes overtime periods, Best Practice recommendations, Storage configuration planning and recommendations, Security planner				<input checked="" type="checkbox"/>
Storage Reporting				
Basic Asset & Capacity Reporting			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Storage reporting on the relationships of computers, file systems and DS8000 LUNs/volumes			<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Capacity Analysis/Predictive Growth				<input checked="" type="checkbox"/>
Customized and Detailed Capacity Reporting – including Chargeback and Database Reporting				<input checked="" type="checkbox"/>
Performance Management				
Performance Reporting/Thresholds				<input checked="" type="checkbox"/>
Volume Performance Advisor – Recommend DS8000 configuration based on performance workloads				<input checked="" type="checkbox"/>
Fabric performance reporting and monitor				<input checked="" type="checkbox"/>

Figure 8-6 Basic and standard edition matrix

8.4.6 IBM Tivoli Storage Productivity Center for Replication

IBM Tivoli Storage Productivity Center for Replication improves time to value, improves continuous availability of open systems servers, and reduces the downtime of critical applications. It can manage the advanced copy services that are provided by the IBM System Storage DS8000, DS6000™ series, the IBM Enterprise Storage Server® (ESS) Model 800 and the IBM System Storage SAN Volume Controller

8.4.7 What is IBM System Storage Productivity Center?

The *IBM System Storage Productivity Center (SSPC)* is a consolidated storage-management solution that is designed as a new management console. It provides a single point of management by integrating the functionality of the Tivoli Productivity Center with the storage devices element managers in an easy-to-use user interface for management.

Developed as a solution to reduce complexity and improve overall interoperability between SAN components, SSPC also offers a storage topology viewer that acts as a dynamic graphical map of the overall SAN environment. This viewer includes information about the SAN device relationships, general system health, and detailed configuration data that can help storage administrators to better streamline operations.

SSPC is an appliance. It is important to understand that the operating system where SSPC runs requires regular maintenance no matter if it was originally sold by IBM or deployed on your own.

SSPC represents a simplification and an improvement in the way clients can manage their systems. Architecturally, SSPC is basically a Windows technology-based component of IBM System x hardware, delivered pre-loaded with the Tivoli Storage Productivity Center. It

also comes preinstalled with the Tivoli Storage Productivity Center Standard Edition software, which requires a separate license purchase.

8.4.8 What can be done from the System Storage Productivity Center?

The complete SSPC offers the following capabilities:

- ▶ Preinstalled and tested console: IBM designed and tested SSPC to support interoperability between server, software, and supported storage devices.
- ▶ IBM System Storage DS8000 GUI integration: With TPC V4.2.1, the DS Storage Manager GUI for the DS8000 is integrated with TPC for remote web access.
- ▶ IBM System Storage SAN Volume Controller (console and CIM agent V5.5): These management components of SAN Volume Controller are preinstalled on the SSPC along with Tivoli Storage Productivity Center Basic Edition, which together are designed to reduce the number of management servers.
- ▶ Automated device discovery: DS8000 and SAN Volume Controller storage devices can be automatically discovered and configured into Tivoli Storage Productivity Center for Replication environments. These devices are displayed in Tivoli Storage Productivity Center through a storage topology.
- ▶ Asset and capacity reporting: Tivoli Storage Productivity Center collects asset and capacity information from storage devices on the SAN, which can be kept for historical reporting, forecasting, and used for other tasks, such as analysis and provisioning.
- ▶ Advanced Topology Viewer: Provides a linked graphical and detailed view of the overall SAN, including device relationships and visual notifications.

8.5 Vendor management applications

Each vendor in the IBM SAN portfolio brings their own bespoke applications to manage and monitor the SAN. In the topics that follow, we give a high-level overview of each of them.

8.5.1 b-type

The *b-type* family switch management framework is designed to support the widest range of solutions, from the small workgroup SANs up to large enterprise SANs. The software that Brocade (the IBM valued original equipment manufacturer (OEM) partner) provides is called *Data Center Fabric Manager (DCFM)* and *Brocade Network Advisor (BNA)*. This software was added to the IBM portfolio as *IBM Data Center Fabric Manager* and *IBM Network Advisor*.

The following tools can be used with b-type SANs to centralize control and enable automation of repetitive administrative tasks:

- ▶ Web Tools:
A built-in web-based application that provides administration and management functions on a per switch basis.
- ▶ Data Center Fabric Manager (DCFM):
A client/server-based external application that centralizes management of IBM/Brocade multiprotocol fabrics within and across data centers, including support for FCoE and CEE.

► Fabric Watch:

A Fabric OS built-in tool that allows the monitoring of key switch elements: power supplies, fans, temperature, error counters, and so on.

► SNMP:

A feature that enables storage administrators to manage storage network performance, find, and solve storage network problems, and plan for storage network growth.

The following management interfaces allow you to monitor fabric topology, port status, physical status, and other information to aid in system debugging and performance analysis:

- Command-line interface (CLI) through a Telnet connection
- Advanced Web Tools
- SCSI Enclosure Services (SES)
- SNMP applications
- Management server

You can use all these management methods either in-band (Fibre Channel) or out-of-band (Ethernet), except for SES, which can be used for in-band only.

For more information about tools that can be used with b-type SANs, see this website:

<http://www.brocade.com/products/all/management-software/product-details/dcfm-enterprise/index.page>

8.5.2 Cisco

Fabric Manager and *Device Manager* are the centralized tools that are used to manage the Cisco SAN fabric and the devices that are connected to it. *Fabric Manager* can be used to manage fabric-wide settings such as zoning, but it can manage settings at an individual switch level as well.

Cisco product name changes: As of NX-OS 5.2, Cisco Fabric Manager and FMS will be known as Cisco Data Center Network Manager for SAN. Also, the LAN-focused Data Center Network Manager becomes Data Center Network Manager for LAN. *Data Center Network Manager (DCNM)* now refers to the converged product.

Cisco *DCNM* is advanced management software that provides comprehensive lifecycle management for the data center LAN and SAN.

Cisco DCNM Release 5.2 combines *Cisco Fabric Manager*, which previously managed SANs, and Cisco DCNM, which previously managed only LANs, into a unified product that can manage a converged data center fabric. As a part of the product merger in Cisco DCNM Release 5.2, the name Cisco DCNM for SAN replaces the name Cisco Fabric Manager. The name Cisco Fabric Manager still applies to Cisco Fabric Manager Release 5.0(x) and all earlier versions.

Cisco DCNM Release 5.2 supports the Cisco Nexus product family, Cisco MDS 9000 product family, Catalyst 6500 Series, and the Cisco UCS product family.

Fabric Manager provides high-level summary information about all the switches in a fabric, automatically launching the Web Tools interface when more detailed information is required. In addition, Fabric Manager provides improved performance monitoring over Web Tools alone.

Some of the capabilities of Fabric Manager are:

- ▶ Configures and manages the fabric on multiple efficient levels
- ▶ Intelligently groups multiple SAN objects and SAN management functions to provide ease and time efficiency in administering tasks
- ▶ Identifies, isolates, and manages SAN events across multiple switches and fabrics
- ▶ Provides drill-down capability to individual SAN components through tightly coupled Web Tools and Fabric Watch integration
- ▶ Discovers all SAN components and views: the real-time state of all fabrics
- ▶ Provides multi-fabric administration of secure Fabric OS SANs through a single encrypted console
- ▶ Monitors ISLs
- ▶ Manages switch licenses
- ▶ Performs fabric stamping

For more information about Fabric Manager, see this website:

<http://www.cisco.com/en/US/products/ps9369/index.html>

8.6 SAN multipathing software

In a well-designed SAN, your device is accessed by the host application over more than one path to potentially obtain better performance. This solution is also done to facilitate recovery in the case of controller, adapter, SFP, cable, or switch failure.

Multipathing software provides the SAN with an improved level of fault-tolerance and performance because it provides more than one physical path between the server and storage.

Traditionally, multipathing software would be supplied by each vendor to support its storage arrays. Currently, there is also the added option to use the multipathing software that often is embedded in the operating system itself. This approach has led to a server-centric approach to multipathing and is independent of the storage array itself. An approach such as this is often easier to implement from a testing and migration viewpoint.

Difference between Storage and a SAN: It is important to understand the key difference between a SAN and storage, though sometimes they are referred to as one.

Storage is where you keep your data.

SAN is the network that the data travels through between your server and storage.

Figure 8-7 shows an example of a dual fabric environment, where hosts have multipathing software and can access the storage if a path fails, or if a fabric fails.

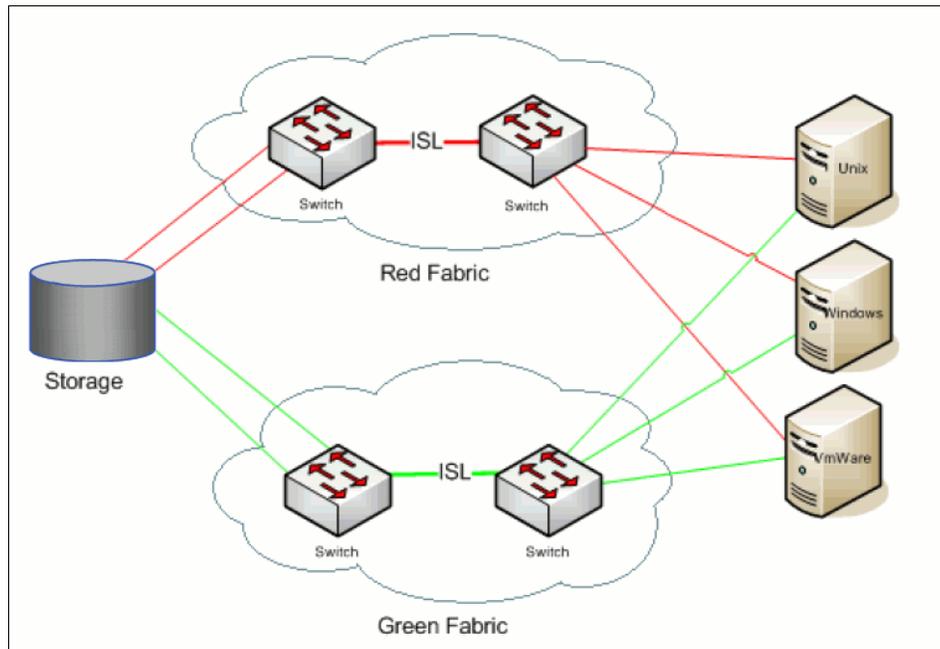


Figure 8-7 SAN overview

In the case of IBM, the *IBM Subsystem Device Driver (SDD)* is used and has the following benefits:

- ▶ Enhances data availability
- ▶ Dynamic input/output (I/O) load-balancing across multiple paths
- ▶ Automatic path failover protection
- ▶ Concurrent download of licensed machine code

When you think about how many paths should be configured to each volume, never exceed the supported level that is given by the storage device. When you implement zoning to a storage device, you must decide on how many paths to have. Detailed information about multipath drivers for IBM storage is at the following website:

http://www.ibm.com/support/docview.wss?rs=540&context=ST52G7&q=ssg1*&uid=ssg1S7000303&loc=en_US&cs

This is an example of how many paths you will have under different scenarios. In Figure 8-7, there are servers that are connected to the SAN with two HBAs, and they access their volumes through two storage ports on the storage device. This access is controlled by zoning, and by this action, gives them four working paths for their volumes: two from the Red Fabric and two from the Green Fabric for each server.

Figure 8-8 indicates that there is a single path failure, as indicated by the STOP sign.

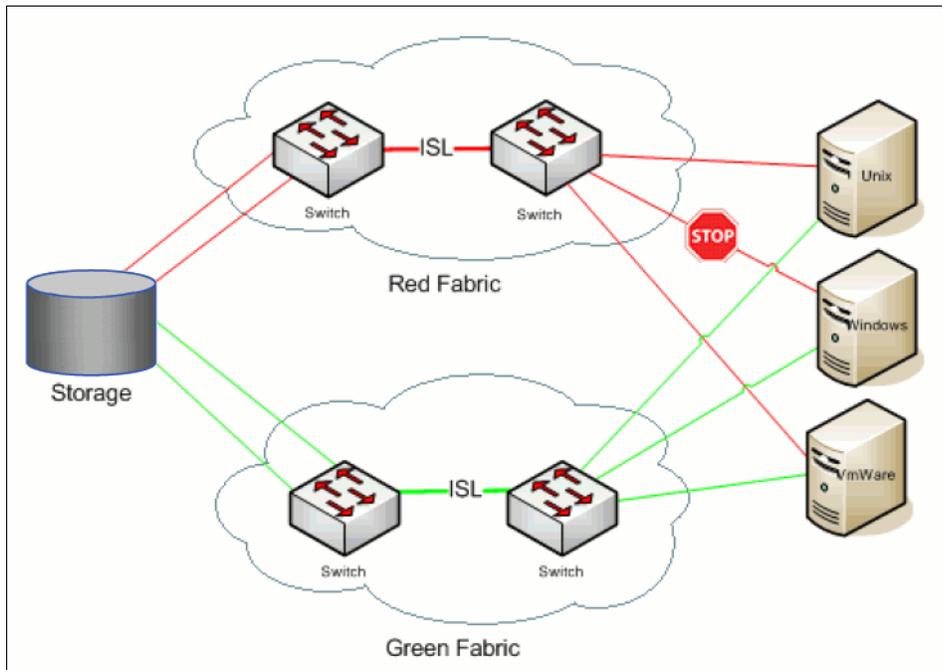


Figure 8-8 HBA failure in a single server

In Figure 8-8, the Windows server lost connectivity to the SAN and it does not have access any longer to the Red Fabric that is leaving. However, we do have working paths through the Green Fabric. All other servers that are running, are running without any issues.

Figure 8-9 shows us that a switch in the Red Fabric failed.

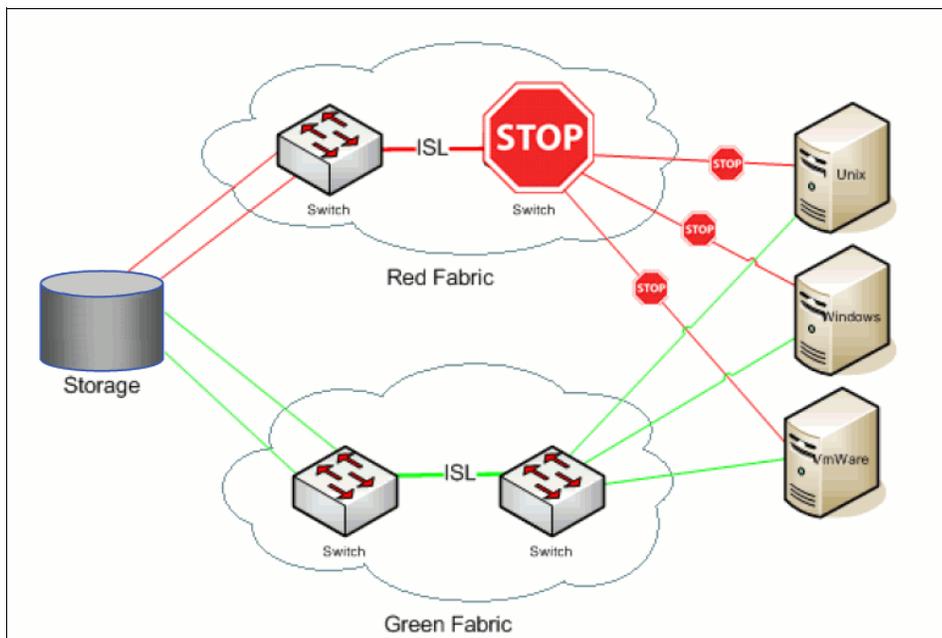


Figure 8-9 A switch is not functional and affects all attached devices

In Figure 8-9 on page 186, we have no access to that switch from our servers. The servers still have working paths through the Green Fabric.

Figure 8-10 shows that a link from the storage device to a switch, failed in the Red Fabric.

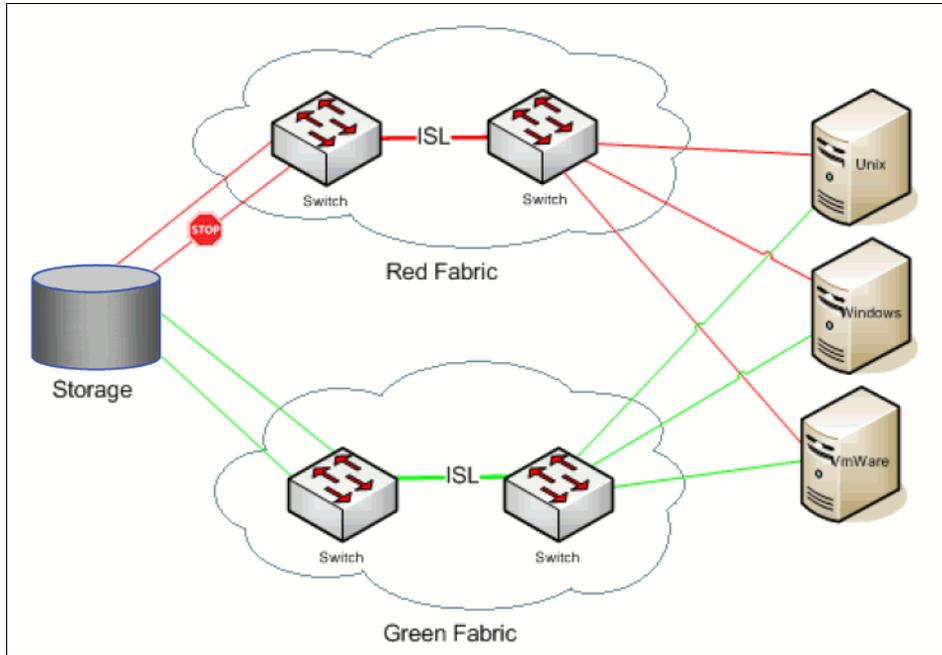


Figure 8-10 Storage device with a single failed connection to a switch

In Figure 8-10, the storage device lost one of four connections. One connection to the Red Fabric is not functional. Therefore, all servers that are using the same storage port now see three working paths out of the four possible. All servers that are zoned to the failed storage port are affected.

Figure 8-11 shows the storage device losing access to the Red Fabric.

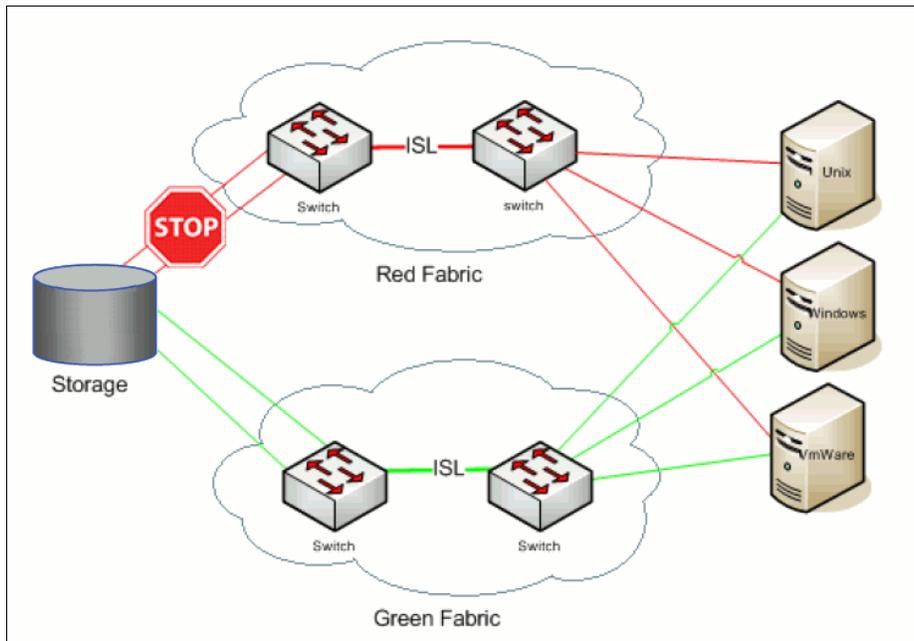


Figure 8-11 Storage device that lost two out of the four connections to the SAN

In Figure 8-11, our storage device lost access to the Red Fabric. All devices in the Red Fabric are running normally, it is only these two specific storage ports that failed. This scenario leaves our servers with two working paths through the Green Fabric. This configuration affects all servers that are zoned to these storage ports.

Figure 8-12 shows the storage device offline.

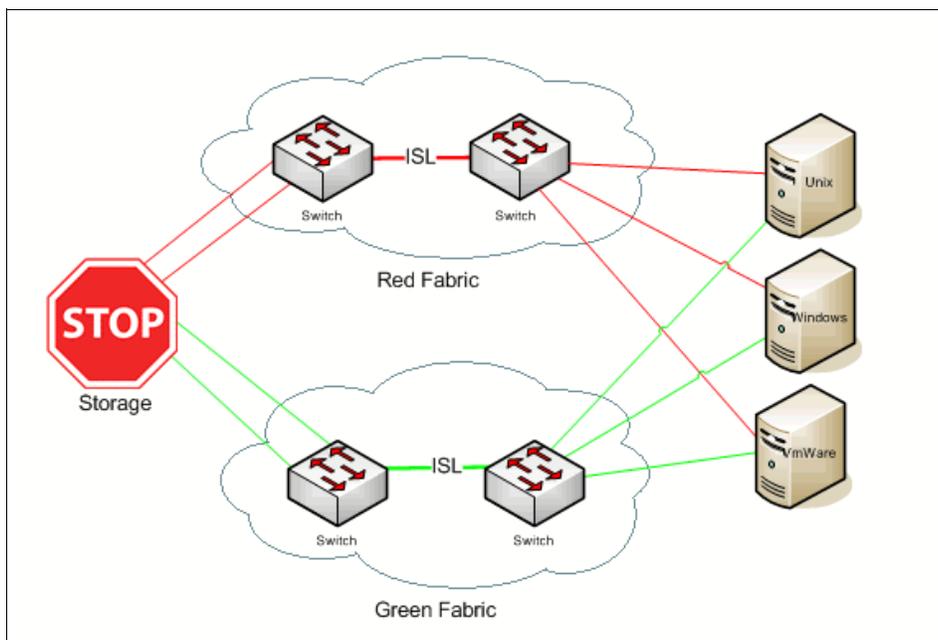


Figure 8-12 The storage device is offline; no working connections

In Figure 8-12 on page 188, we lost our storage device. No paths to any volumes on this device are available. No data is accessible. All servers that are zoned to this storage device are severely affected and have no access to this storage device.

If a correctly installed and supported version of a multipath driver was installed on the servers, we would have “survived” all scenarios except the last one, in which there would have been minimum impact.



Security

In this chapter, we provide an overview of the need for security, the different techniques that are available, and some of the key points to be aware of.

9.1 Security in the storage area network (SAN)

Security is always a major concern for networked systems administrators and users. Even for specialized networked infrastructures, such as SANs, special care must be taken so that information does not get corrupted, either accidentally or deliberately, or fall into the wrong hands. And we also must ensure that at a fabric level the correct security is in place; for example, to ensure that a user does not inadvertently change the configuration incorrectly.

Now that SANs have “broken” the traditional direct-attached storage paradigm of servers being cabled directly to servers, the inherent security that this provides is lost. The SAN and its resources might be shared by many users and many departments. The SAN might be shared by different operating systems that have differing ideas as to who owns what storage. To protect the privacy and safeguard the storage, SAN vendors came up with a segmentation feature to overcome this consideration. This feature is called *zoning*.

The fabric itself enforces the separation of data so that only those users that are intended to have access are able to communicate with the data which is intended for them.

Zoning, however, does not provide security in that sense; it implements only the means of segregation (isolation). The real security issue is the vulnerability when the data itself must travel outside of the data center, and over long distances. This type of travel often involves transmission over networks that are owned by different carriers.

We must look at security from two different angles: For data-in-flight, as explained in 9.4.2, “Data-in-flight” on page 198, and for data-at-rest, explained in 9.4.3, “Data-at-rest” on page 199.

More often than not, data is not encrypted when it is sent from the source to a target. Therefore, any information is readable with the correct tools, even though it is slightly more complicated than simply eavesdropping on a telephone line. Because all the data is sent at a block level with the Fibre Channel protocol (meaning that all data that is sent is squeezed into the Fibre Channel frame before sending), “sniffing” a frame or two might give you 2112 bytes of data. As an example of the difficulty, this amount would be similar to 1/333,000 of a normal CD or 13 milliseconds of a CD spanning 74 minutes. Obviously this comparison does not give you much information without putting it in the right context.

There is more concern if the whole Fibre Channel port or disk volumes and arrays are mirrored, or tapes that contain information end up in the wrong hands. However, to tamper with information from a SAN is not something that just happens, it is something that takes a concerted effort.

The storage architect and administrators must understand that in a SAN environment, often with a combination of diverse operating systems and vendor storage devices, some combination of technologies is required. This mixture ensures that the SAN is secure from unauthorized systems and users, whether accidental or deliberate.

In the discussions that follow, we briefly explore some of the technologies and their associated methodologies that can be used to ensure data integrity, and to protect and manage the fabric. Each technology has advantages and disadvantages. And each must be considered based on a well thought-out SAN security strategy, which is developed during the SAN design phase.

9.2 Security principles

It is a well-known fact that “a chain is only as strong as its weakest link”, and when describing computer security, the same concept applies. There is no point in locking all of the doors and then leaving a window open. A secure, networked infrastructure must protect information at many levels or layers, and have no single point of failure.

The levels of defense must be complementary, and work with each other. If you have a SAN, or any other network, that crumbles after a single penetration, then this level of defense is not a recipe for success.

There are a number of unique entities that must be given consideration in any environment. We describe some of the most important ones in the topics that follow.

9.2.1 Access control

Access control can be performed both with *authentication* and *authorization* techniques:

- Authentication** Means that the secure system must challenge the user (usually with a password) so that this user is identified.
- Authorization** After identifying a user, the system is able to know what this user is allowed to access and what they are not.

As in any IT environment, including SAN, access to information and to the configuration or management tools, must be restricted. Access must be granted to only those individuals that need to have access and are authorized to make changes. Any configuration or management software is typically protected with several levels of security. Levels usually start with a user ID and password that must be assigned appropriately to personnel based on their skill level and responsibility.

9.2.2 Auditing and accounting

An audit trail must be maintained for auditing and troubleshooting purposes, especially when you create a *root cause analysis (RCA)* after an incident occurs. Inspect and archive logs regularly.

9.2.3 Data security

Whether we describe data-at-rest or data-in-flight, data security consists of both data *confidentiality* and *integrity*:

- Data confidentiality** The system must guarantee that the information cannot be accessed by unauthorized people, that it remains confidential, and is only available for authorized personnel. As shown in the next section, confidentiality is usually accomplished by using data *encryption*.
- Data integrity** The system must guarantee that the data is stored or processed within its boundaries and that it is not altered or tampered with in any way.

The data security and integrity requirement aims to guarantee that data from one application or system does not become overlaid, corrupted, or otherwise destroyed. This requirement applies whether data is intentionally destroyed or destroyed by accident, either by other applications or systems. This requirement might involve some form of authorization, and the ability to fence off the data from one system from another system.

This data security necessity must be balanced with the requirement for the expansion of SANs to enterprise-wide environments, with a particular emphasis on multi-platform connectivity. True cross-platform data sharing solutions, as opposed to data partitioning solutions, are also a requirement. Security and access control also must be improved to guarantee data integrity.

In the topics that follow, we overview some of the common approaches to securing data that are encountered in the SAN environment. This list is not meant to be an in-depth description. It is merely an attempt to acquaint you with the technology and terminology that is likely to be encountered when a discussion on SAN security occurs.

9.2.4 Securing a fabric

In this section, some of the current methods for securing a SAN fabric are presented.

Fibre Channel Authentication Protocol

The Switch Link Authentication Protocol (SLAP/FC-SW-3) establishes a region of trust between switches. For an end-to-end solution to be effective, this region of trust must extend throughout the SAN, which requires the participation of fabric-connected devices, such as host bus adapters (HBAs). The joint initiative between Brocade and Emulex establishes Fibre Channel Authentication Protocol (FCAP) as the next-generation implementation of SLAP. Clients gain the assurance that a region of trust extends over the entire domain.

FCAP was incorporated into its fabric switch architecture and proposed the specification as a standard to ANSI T11 (as part of FC-SP). FCAP is a Public Key Infrastructure (PKI)-based cryptographic authentication mechanism for establishing a common region of trust among the various entities (such as switches and HBAs) in a SAN. A central, trusted third party serves as a guarantor to establish this trust. With FCAP, certificate exchange takes place among the switches and edge devices in the fabric to create a region of trust that consists of switches and HBAs.

The fabric authorization database is a list of the WWNs and associated information like domain IDs of the switches that are authorized to join the fabric.

The fabric authentication database is a list of the set of parameters that allows the authentication of a switch within a fabric. An entry of the authentication database holds at least the switch worldwide name (WWN), authentication mechanism Identifier, and a list of appropriate authentication parameters.

Zoning

Initially, SANs did not have any zoning. It was an any-to-any communication, and there was no real access control mechanism to protect storage that was used by one host from being accessed by another host. When SANs grew, this drawback became a security risk as SANs became more complex and were running more vital parts of the business. To mitigate the risk of unwanted cross communication, zoning was invented to isolate communication to devices within the same zone.

Persistent binding

Server-level access control is called *persistent binding*. Persistent binding uses configuration information that is stored on the server, and is implemented through the HBA driver of the server. This process binds a server device name to a specific Fibre Channel storage volume or logical unit number (LUN), through a specific HBA and storage port WWN. Or, put in more technical terms, it is a host-centric way to direct an operating system to assign certain Small Computer System Interface (SCSI) target IDs and LUNs.

Logical unit number masking

One approach to securing storage devices from hosts that want to take over already assigned resources, is logical unit number (LUN) masking. Every storage device offers its resources to the hosts with LUNs. For example, each partition in the storage server has its own LUN. If the host (server) wants to access the storage, it must request access to the LUN in the storage device. The purpose of LUN masking is to control access to the LUNs. The storage device itself accepts or rejects access requests from different hosts. The user defines which hosts can access which LUN with the storage device control program. Whenever the host accesses a particular LUN, the storage device checks its access list for that LUN. The device allows or disallows the host to gain access to the LUN.

Port binding

To provide a higher level of security, you can also use *port binding* to bind a particular device (as represented by a WWN) to a specific port that does not allow any other device to plug into the port. This device then assumes the role of the device that was there. The reason for this is that the “rogue” device that was inserted has a different WWN in which the port was bound.

Role-based access control

A *role-based access control feature (RBAC)* is available in most SAN devices today. By using RBAC, you can control user access and user authority in a simple way. RBAC allows you to provide users with access or permission to run tasks that are only within their skill set or job role.

Normally there are three definitions for RBAC:

- ▶ Role assignment
- ▶ Role authorization
- ▶ Permission authorization

Usually each role can contain multiple users and each user can be part of multiple roles. For example, if role1 users are only allowed access to configuration commands, and role2 users are only allowed access to debug commands, then if John belongs to both role1 and role2, he can access configuration and debug commands.

These predefined roles in a SAN environment are important to ensure that correct login and access is defined for each user.

9.2.5 Zoning, masking, and binding

Although zoning, masking, or binding are not classed as security products or mechanisms, combining all of their functionality together can make the SAN more secure than it would be without them.

9.3 Data security

These data security standards propose to secure Fibre Channel (FC) traffic between all FC ports and the domain controller.

The following methods are used for data security standards:

- ▶ FCPAP refers to Secure Remote Password Protocol (SRP), RFC 2945.
- ▶ DH-CHAP refers to Challenge Handshake Authentication Protocol (CHAP), RFC 1994.
- ▶ FCSec refers to IP Security (IPSec), RFC 2406.

The focus of the FCSec is to provide authentication of these entities:

- Node-to-node
- Node-to-switch
- Switch-to-switch

An additional function that might be possible to implement is *frame level encryption*.

The ability to perform switch-to-switch authentication in FC-SP enables a new concept in Fibre Channel: the secure *fabric*. Only switches that are authorized and properly authenticated are allowed to join the fabric.

Authentication in the secure fabric is twofold. The fabric wants to verify the identity of each new switch before it joins the fabric, and the switch that is wanting to join the fabric wants to verify that it is connected to the right fabric. Each switch needs a list of the worldwide names (WWNs) of the switches that are authorized to join the fabric. The switch also needs a set of parameters that are used to verify the identity of the other switches that belong to the fabric.

Manual configuration of such information within all the switches of the fabric is possible, but not advisable in larger fabrics. And there is the need of a mechanism to manage and distribute information about authorization and authentication across the fabric.

9.4 Storage area network encryption

What is data encryption, and symmetric and asymmetric encryption; or in-flight data, or data-at-rest? In the topics that follow, this terminology is explained and helps you to understand the fundamentals in encryption and key management.

9.4.1 Basic encryption definition

One of the first questions to answer is: “Do I need encryption”? In this section, we describe basic encryption, cryptographic terms, and ideas on how you can protect your data.

Encryption is one of the simple ways to secure your data. If the data is stolen, lost, or acquired in any way, it cannot be read without the correct encryption key.

Encryption has been used to exchange information in a secure and confidential way for many centuries. Encryption transforms data that is unprotected (plain or clear text) into encrypted data, or *ciphertext*, by using a key. It is difficult to “break” ciphertext to change it back to clear text without the associated encryption key.

There are two main types of encryption: *symmetric encryption* and *asymmetric encryption* (also called *public-key encryption*):

Symmetric When the same secret password, or key, is used to encrypt a message and decrypt the corresponding cipher text

Asymmetric When one key is used to encrypt a message and another to decrypt the corresponding cipher text

A *symmetric cryptosystem* follows a fairly straightforward philosophy: Two parties can securely communicate if both use the same *cryptographic algorithm* and possess the same secret key to encrypt and decrypt messages. This algorithm is the simplest and most efficient way of implementing secure communication, if the participating parties are able to securely exchange secret keys (or passwords).

Figure 9-1 illustrates symmetric encryption.

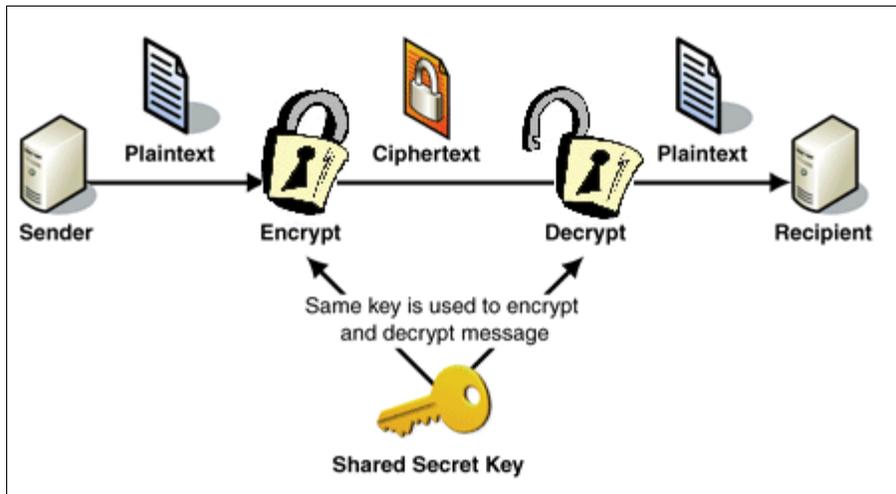


Figure 9-1 Symmetric cryptography

An *asymmetric* (or *public-key*) cryptosystem is a cryptographic system that uses a pair of unique keys, usually referred to as *public keys* and *private keys*. Each individual is assigned a pair of these keys to encrypt and decrypt information. A message encrypted by one of these keys can be decrypted only by the other key and vice versa:

- ▶ One of these keys is called a *public key* because it is made available to others for use when they encrypt information that is sent to an individual. For example, people can use a person's public key to encrypt information they want to send to that person. Similarly, people can use the user's public key to decrypt information that is sent by that person.
- ▶ The other key is called *private key* because it is accessible only to its owner. The individual can use the private key to decrypt any messages encrypted with the public key. Similarly, the individual can use the private key to encrypt messages so that the messages can be decrypted only with the corresponding public key.

This means that exchanging keys is not a security concern. An analogy to public-key encryption is that of a locked mailbox with a mail slot. The mail slot is exposed and accessible to the public; its location (the street address) is in essence the public key. Anyone knowing the street address can go to the door and drop a written message through the slot; however, only the person who possesses the key can open the mailbox and read the message.

Figure 9-2 illustrates the asymmetric cryptography process.

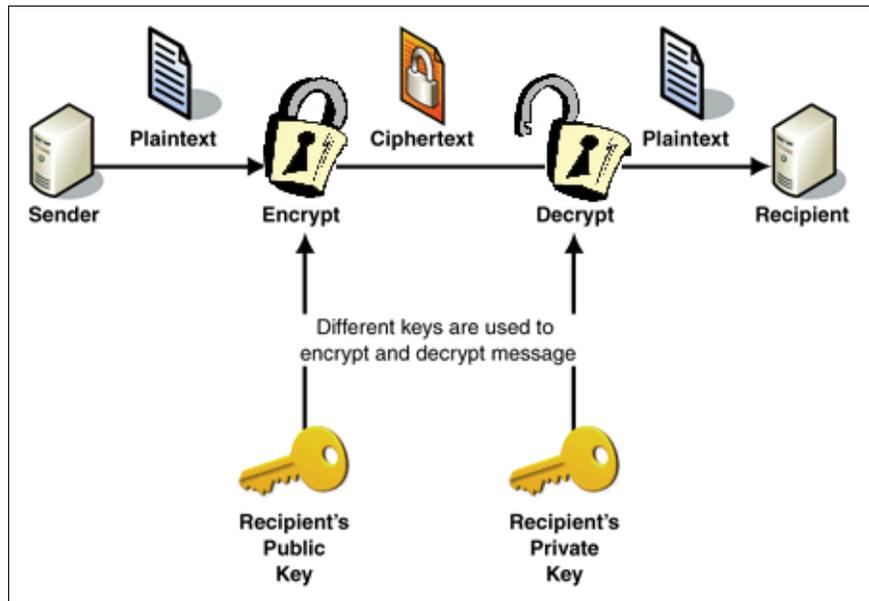


Figure 9-2 Asymmetric cryptography

The main disadvantage of public-key encryption when compared to symmetric encryption is that it demands much higher computing power to be performed as efficiently. For this reason, most of the current security systems use public-key mechanisms as a way to securely exchange symmetric encryption keys between parties that then use symmetric encryption for their communication. In this case, the exchanged symmetric secret key (or password) is called a *session key*.

However, there is still an issue regarding public-key cryptosystems: When you initially receive a public key from someone for the first time, how do you know that this individual is really who they claim to be? If “spoofing” the identity of someone is so easy, how do you knowingly exchange public keys? The answer is to use a *digital certificate*. A digital certificate is a digital document that is issued by a trusted institution that vouches for the identity and key ownership of an individual. The certificate guarantees authenticity and integrity.

In the next sections, some of the most common encryption algorithms and tools are presented along with terminology explanations.

9.4.2 Data-in-flight

Also known as *data-in-motion*, this term generically refers to protecting information any time that the data leaves its primary location. For example, when data is transmitted from the source across any type of network to a target. To secure this transmission, we use technologies such as Secure Sockets Layer (SSL), Virtual Private Network (VPN), and IP Security (IPSec) to assure data confidentiality. Then, we use other technologies such as digital certificates, message authentication codes, and keyed hashes to ensure data integrity. *Data-in-flight* is also information (data) that leaves the data center through, for example, an open network or leased dark fiber.

All of these areas can be addressed with encryption-based technologies.

9.4.3 Data-at-rest

Protecting data as it resides on the storage media, disk, or tape, is typically referred to as *data-at-rest*.

If encryption is used as part of the strategy for the protection of data-at-rest, this protection also indirectly addresses the issue of displayed tape media. This issue is addressed because, even if tapes fall into the wrong hands, the data that is stored on them is unreadable without the correct key. These security measures assume that you have enacted the appropriate key management techniques.

To gain the needed security level, you build layers of security on your SAN. You first increase the level of difficulty for an unauthorized user to even gain access to the data. You then compound that with the fact that private data is not stored in human-readable form.

9.4.4 Digital certificates

If you are using one of these encryption methods, you must also be certain that the person or machine you are sending to is the correct one. When you initially receive a public key from someone for the first time, how do you know that this individual is really the person that they claim to be? If “spoofing” the identity of someone is so easy, how do you knowingly exchange public keys? The answer is to use a digital certificate. A digital certificate is a digital document that is issued by a trusted institution that vouches for the identity and key ownership of an individual: it guarantees authenticity and integrity.

There are trusted institutions all over the world that generate trusted certificates. We use this type of mechanism also for the first time using a certificate that is generated by our switch. For more details, see 9.4.6, “Key management considerations and security standards” on page 200.

9.4.5 Encryption algorithm

After you decide that encryption is a must, you must also be aware that there are several encryption schemes to choose from. The most popular encryption schemes in use today include the following algorithms:

- ▶ 3DES
- ▶ DES
- ▶ AES
- ▶ RSA
- ▶ ECC
- ▶ Diffie-Hellman
- ▶ DSA
- ▶ SHA

For more information about IBM System Storage Data Encryption, see *IBM System Storage Data Encryption*, SG24-7797. For an example of how IBM implements encryption on the IBM System Storage SAN Volume Controller, see *Implementing the Storwize V7000 and the IBM System Storage SAN32B-E4 Encryption Switch*, SG24-7977.

If we look at the security aspect on its own, we have been focusing on establishing a perimeter of defense around system assets. Although securing access to our environments continues to be an important part of security, the typical business cannot afford to lock down its entire enterprise.

Open networks are now commonly used to connect clients, partners, employees, suppliers, and their data. While this offers significant advantages, it raises concerns about how a business protects its information assets and complies with industry and legislative requirements for data privacy and accountability. By using data encryption as a part of the solution, much of these concerns can be mitigated, as explained in next section.

9.4.6 Key management considerations and security standards

An encryption algorithm requires a key to transform the data. All cryptographic algorithms, at least the reputable ones, are in the public domain. Therefore, it is the key that controls access to the data. We cannot emphasize enough that you must safeguard the key to protect the data. A good tool for that purpose is *IBM Tivoli Key Lifecycle Management*, which we briefly describe in the next section.

IBM Tivoli Key Lifecycle Management

Because of the nature, security, and accessibility of encryption, data that is encrypted is dependent on the security of, and accessibility to, the decryption key. The disclosure of a decryption key to an unauthorized agent (individual person or system component) creates a security exposure in such a way that the unauthorized agent would also have access to the ciphertext that is generated with the associated encryption key.

Furthermore, if all copies of the decryption key are lost (whether intentionally or accidentally), no feasible way exists to decrypt the associated ciphertext, and the data that is contained in the ciphertext is said to have been cryptographically erased. If the only copies of certain data are cryptographically erased, then access to that data is permanently lost for all practical purposes.

This problem is why the security and accessibility characteristics of encrypted data can create considerations for you that do not exist with storage devices that do not contain encrypted data.

The primary reason for using encryption is that data is kept secure from disclosure and is kept from others that do not have sufficient authority. At the same time, data must be accessible to any agent that has both the authority and the requirement to gain access.

Two security considerations are important in this context:

- ▶ *Key security*

To preserve the security of encryption keys, the implementation must ensure that no one individual (system or person) has access to all the information that is required to determine the encryption key.

- ▶ *Key availability*

To preserve the access to encryption keys, redundancy can be provided by having multiple independent key servers that have redundant communication paths to encrypting devices. This ensures that the backup of each key server's data is maintained. Failure of any one key server or any one network, does not prevent devices from obtaining access to the data keys that are needed to provide access to the data.

The sensitivity of possessing and maintaining encryption keys, and the complexity of managing the number of encryption keys in a typical environment, results in a client requirement for a *key server*. A key server is integrated with encrypting products to resolve most of the security and usability issues that are associated with key management for encrypted devices. However, you must still be sufficiently aware of how these products interact to provide appropriate management of the computer environment.

Master key: Even with a key server, generally at least one encryption key, normally called the *master key (MK)*, must be maintained manually. For example, this is the key that manages access to all other encryption keys. This master key is a key that encrypts the data that is used by the key server to exchange keys.

Fundamentally, Tivoli Key Lifecycle Manager works by allowing administrators to connect with storage devices and then create and manage keystores. These stores are secure repositories of keys and certificate information that are used to encrypt and decrypt data, or to use existing keystores already in place. Over the course of the key lifecycle, all management functions, including creation, importation, distribution, backup, and archiving, are easily accomplished. These functions can be done by using the lifecycle manager's graphic interface, which can be accessed by using any standard browser on the network. Tivoli Key Lifecycle Manager thus serves as a central point of control, unifying key management even when different classes of storage devices are involved. For more information about Tivoli Key Lifecycle Manager, see this website:

<http://www-01.ibm.com/software/tivoli/beat/10212008.html>

There are two security standards that are important to ensuring the integrity of encryption products: FIPS 140 and Common Criteria. The official title for the standard Federal Information Processing Standard 140 (FIPS-140) is Security Requirements for Cryptographic Modules. FIPS 140-2 stands for the second revision of the standard and was released in 2001. Common Criteria has seven Evaluation Assurance Levels (EALs), which were defined in 1999. Together, these standards support a small industry for certifying security products and ensuring the integrity of encryption systems.

9.4.7 b-type encryption methods

B-type encryption devices from IBM are used to encrypt data at rest on a storage media and, starting with FOS 7.0, with 16 Gbps E_Ports in-flight encryption. When we describe storage media, this media can be either disk or tape.

In-flight encryption

The in-flight encryption and compression feature of Fabric OS allows frames to be encrypted or compressed at the egress point of an inter-switch link (ISL) between two b-type switches, and then to be decrypted or extracted at the ingress point of the ISL. This feature uses port-based encryption and compression. It is supported on 16 Gbps E_Ports only.

The purpose of encryption is to provide security for frames while they are in flight between two switches. The purpose of compression is for better bandwidth use on the ISLs, especially over long distance. An average compression ratio of 2:1 is provided. Frames are never left in an encrypted or compressed state when delivered to an end device, and both ends of the ISL must terminate at 16 Gbps ports.

Figure 9-3 shows the b-type in-flight encryption architecture.

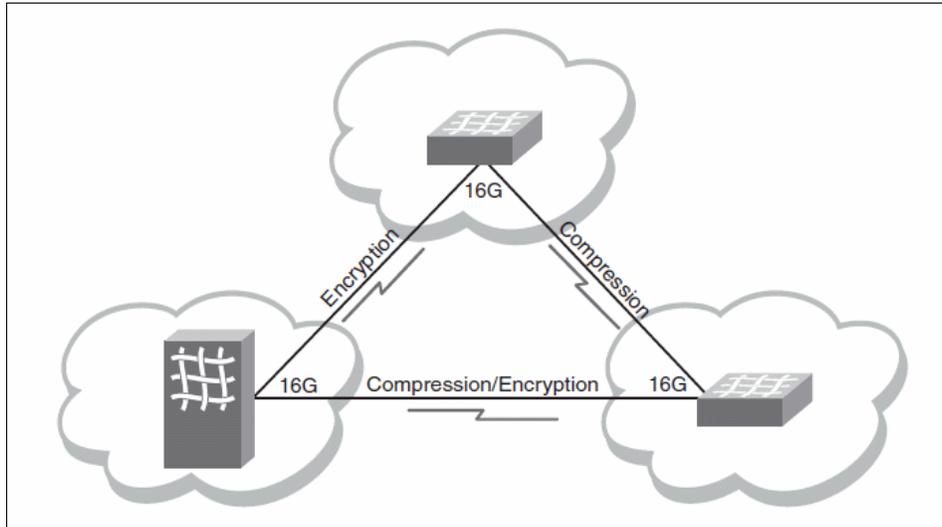


Figure 9-3 In-flight architecture

Encryption at rest

B-type fabric-based encryption solutions work transparently with heterogeneous servers, tape libraries, and storage subsystems. Although host-based encryption works only for a specified operating system and storage-based encryption works only for a specific vendor, b-type products are deployed in the core of the fabric to encrypt Fibre Channel-based traffic. Users deploy b-type encryption solutions via either the FS8-18 Encryption Blade or the 2U, rack-mounted IBM SAN32B-E4 Encryption Switch.

The Device Encryption Key (DEK) is important. Because it is needed to encrypt and decrypt the data, it must be random and 256 bits in length. B-type encryption devices use a True Random Number Generator (TRNG) to generate each DEK. For encrypting data destined for a disk drive, one DEK is associated with one logical unit number (LUN). The Institute of Electrical and Electronic Engineers 1619 (IEEE 1619) standard on encryption algorithms for disk drives is known as *AES256-XTS*. The encrypted data from the *AES256-XTS* algorithm is the same length as the unencrypted data. Therefore, the b-type encryption device can encrypt the data, block by block, without expanding the size of the data. The key management is done by using external software such as Tivoli Key Lifecycle Manager.

Figure 9-4 shows a simple b-type encryption setup.

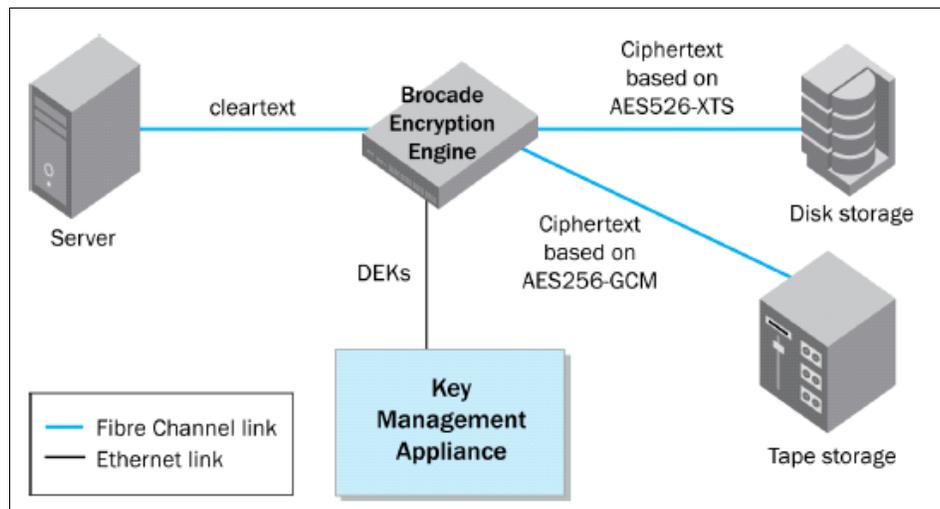


Figure 9-4 B-type encryption and key management

9.4.8 Cisco encryption methods

Cisco has two methods of encrypting SAN information: In-flight encryption and storage media encryption. For more information about both of these methods, see these websites:

- ▶ http://www.cisco.com/en/US/prod/collateral/ps4159/ps6409/ps5990/white_paper_c11-545124.html
- ▶ http://www.cisco.com/en/US/prod/collateral/ps4159/ps6409/ps6028/ps8502/product_data_sheet0900aecd8068ed59.html

In-flight encryption

Cisco TrustSec Fibre Channel Link Encryption is an extension of the FC-SP standard and uses the existing FC-SP architecture. Fibre Channel data that is traveling between E_Ports on 8 Gbps modules is encrypted. Cisco uses the 128-bit Advanced Encryption Standard (AES) encryption algorithm and enables either AES-Galois/Counter Mode (AES-GCM) or AES-Galois Message Authentication Code (AES-GMAC). AES-GCM encrypts and authenticates frames, and AES-GMAC authenticates only the frames that are being passed between the two peers. Encryption is performed at line rate by encapsulating frames at egress with encryption by using the GCM authentication mode with 128-bit AES encryption. At ingress, frames are decrypted and authenticated with integrity checks.

There are two primary use cases for Cisco TrustSec Fibre Channel Link Encryption. In the first use case, clients are communicating outside the data center over native Fibre Channel (for example, dark fiber, Coarse Wavelength-Division Multiplexing (CWDM) or Dense Wavelength-Division Multiplexing (DWDM)). In the second use case, encryption is performed within the data center for security-focused clients such as defense and intelligence services.

Figure 9-5 shows Cisco TrustSec Fibre Channel Link Encryption.

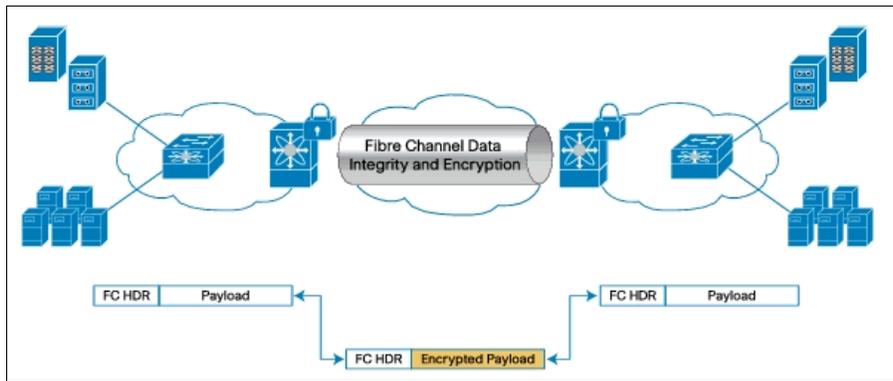


Figure 9-5 Cisco TrustSec encryption

Encryption at rest

Cisco uses Storage Media Encryption (SME) which protects data at rest on heterogeneous tape drives, virtual tape libraries (VTLs), and disk arrays in a SAN environment by using highly secure IEEE Advanced Encryption Standard (AES) algorithms.

Encryption is performed as a transparent Fibre Channel fabric service, which greatly simplifies the deployment and management of sensitive data on SAN-attached storage devices. Storage in any virtual SAN (VSAN) can make full use of Cisco SME. Secure lifecycle key management is included, with essential features such as key archival, shredding, automatic key replication across data centers, high-availability deployments, and export and import for single- and multiple-site environments. Provisioning and key management for Cisco SME are both integrated into Cisco Fabric Manager and Data Center Network Manager (DCNM); no additional software is required for key management.

Figure 9-6 shows the SME architecture.

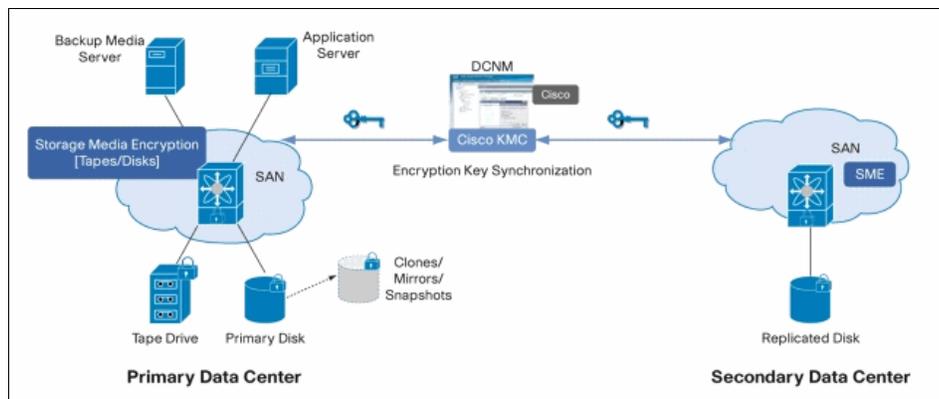


Figure 9-6 SME architecture

9.5 Encryption standards and algorithms

The following list defines some of the most popular encryption algorithms in use today:

- AES** *Advanced Encryption Standard (AES)* is a symmetric 128-bit block data encryption technique that was developed by Belgian cryptographers: Joan Daemen and Vincent Rijmen. The US government adopted the algorithm as its encryption technique in October 2000, replacing the DES encryption that it used. AES works at multiple network layers simultaneously. The National Institute of Standards and Technology (NIST) of the US Department of Commerce selected the algorithm, called *Rijndael* (pronounced Rhine Dahl or Rain Doll), out of a group of five algorithms under consideration. AES is the first publicly accessible and open cipher that is approved by the National Security Agency (NSA) for top secret information.
- RSA** The *RSA* algorithm involves three steps: key generation, encryption, and decryption. This algorithm was created in 1977 by Ron Rivest, Adi Shamir, and Len Adleman at MIT; the letters RSA are the initials of their surnames. It was the first algorithm that is known to be suitable for digital signing and data encryption, and one of the first great advances in public key cryptography. RSA is still widely used in electronic commerce protocols, and is believed to be secure given sufficiently long keys and the use of up-to-date implementations.
- ECC** *Elliptic curve cryptography* is an approach to public-key cryptography that is based on the mathematics of elliptic curves over finite fields. The use of elliptic curves in cryptography was suggested independently by Neal Koblitz and Victor S. Miller in 1985. Elliptic curves are also used in several integer factorization algorithms that have applications in cryptography. An example of such an algorithm is Lenstra elliptic curve factorization, but this use of elliptic curves is *not* usually referred to as elliptic curve cryptography.
- Diffie-Hellman** The *Diffie-Hellman (D-H)* key exchange is a cryptographic protocol which allows two parties that have no prior knowledge of each other to jointly establish a shared secret key over an insecure communications channel. This key can then be used to encrypt subsequent communications by using a symmetric key cipher.
- DSA** The *Digital Signature Algorithm (DSA)* is a United States Federal Government standard for digital signatures. It was proposed by the National Institute of Standards and Technology (NIST) in August 1991 for use in their Digital Signature Standard (DSS), specified in FIPS 186, adopted in 1993. A minor revision was issued in 1996 as FIPS 186-1, and the standard was expanded further in 2000 as FIPS 186-2, and again in 2009 as FIPS 186-3. DSA is covered by US Patent 5,231,668, filed 26 July 1991, and attributed to David W. Kravitz, a former NSA employee.
- SHA** The *Secure Hash Algorithm (SHA)* family is a set of related cryptographic hash functions. The most commonly used function in the family, *SHA-1*, is employed in a large variety of popular security applications and protocols, including TLS, SSL, PGP, SSH, S/MIME, and IPSec. The algorithm was also used on the Nintendo Wii gaming console for signature verification when booting occurs.

9.6 Security common practices

As we mentioned before, you might have the most sophisticated security system installed in your house; but, it is not worth anything if you leave the window open. At a high level, consider practicing the following security preferred practices at a minimum:

- ▶ Change default configurations and passwords often.
- ▶ Check and double check configuration changes to ensure that only the data that is supposed to be accessed, can be accessed.
- ▶ Management of devices usually takes a *telnet* form, with encrypted management protocols being used.
- ▶ Remote access often relies on unsecured networks. Ensure that the network is secure and that some form of protection is in place to guarantee only those people with the correct authority are allowed to connect.
- ▶ Ensure that the operating systems that are connected are as secure as they can be. If the operating systems are connected to an internal and external LAN, ensure that this connection cannot be used. Access might be gained by using loose configurations.
- ▶ Assign the correct roles to administrators.
- ▶ Ensure that the devices are in physically secure locations.
- ▶ Ensure that the passwords are changed if the administrator leaves. Also, ensure that passwords are changed regularly.

Finally, the SAN security strategy in its entirety must be periodically addressed as the SAN infrastructure develops, and as new technologies emerge and are introduced into the environment.

These safeguards do not guarantee that your information is 100% secure, but they will go some way in ensuring that all but the most ardent “thieves” are kept out.



Solutions

The added value of a storage area network (SAN) lies in the use of its technology to provide tangible and desirable benefits to the business. These benefits are provided by the use of fast, secure, reliable, and highly available networking solutions. Benefits range from increased availability and flexibility to more functionality that can reduce application downtime.

In this chapter, we provide a description of general SAN applications, and the types of components that are required to implement them. There is far more complexity than is presented here. For instance, this text does not cover how to choose one switch over another, or how many inter-switch links (ISLs) are necessary for a specific SAN design. These strategic decisions must be always considered by experienced IT architects, and that is beyond the intended scope of this book. We introduce the basic principles and key considerations that must be taken into account to choose an optimal solution for SAN deployments.

10.1 Introduction

During the last few years and with the continued development of the communication and computing technologies and products, SANs have evolved and are getting much more complex. Now, we are not referring to just a simple fiber-optic connection between SAN devices. Examples of these devices include: SAN switches, routers, tape drives, disk device subsystems, and target host systems that use standard Fibre Channel host bus adapters (HBAs). Technology has moved way beyond those solutions and continues to do so.

Today, businesses are looking for solutions that enable them to increase the data transfer rate within the most complex data centers. Businesses also want solutions that provide high availability of managed applications and systems, implement data security, and provide storage efficiency. At the same time, businesses want to reduce the associated costs and power consumption.

Organizations must find a smooth, effective, and cost-efficient way to upgrade their current or traditional SAN infrastructure. The upgraded infrastructure provides a less complex, and more powerful and flexible data center of the next generation.

There are many categories in which SAN solutions can be classified. We chose to classify ours as such: infrastructure simplification, business continuity, and information lifecycle management (ILM). In the topics that follow, we describe the use of basic SAN design patterns to build solutions for different requirements. These requirements range from simple data movement techniques that are frequently employed as a way to improve business continuity, up to sophisticated storage pooling techniques that are used to simplify complex infrastructures.

Before SAN solutions and requirements are described, we present some basic principles to be considered when you plan a SAN implementation or upgrade.

10.2 Basic solution principles

A number of important decisions must be made by the system architect, either when a new SAN is being designed, or when an existing SAN is being expanded. Such decisions usually refer to the choice of the connectivity technology, the preferred practices for adding capacity to a SAN, or the more suitable technology for achieving data integration. This section describes some of these aspects.

10.2.1 Connectivity

Connecting servers to storage devices through a SAN fabric is often the first step that is taken in a phased SAN implementation. Fibre Channel attachments have the following benefits:

- ▶ Improved performance by running Small Computer System Interface (SCSI) over Fibre Channel
- ▶ Extended connection distances (sometimes called *remote storage*)
- ▶ Enhanced addressability

Many implementations of Fibre Channel technology are simple configurations that remove some of the restrictions of the existing storage environments, and allow you to build one common physical infrastructure. The SAN uses common cabling to the storage and other peripheral devices. The handling of separate sets of cables, such as OEMI, ESCON, SCSI

single-ended, SCSI differential, SCSI LVD, and others, have caused the IT organization management much trauma as it attempted to treat each of these differently. One of the biggest issues is the special handling that is needed to circumvent the various distance limitations.

Installations without SANs commonly use SCSI cables to attach to their storage. SCSI has many restrictions such as limited speed, only a few devices that can be attached, and severe distance limitations. Running SCSI over Fibre Channel helps to alleviate these restrictions. SCSI over Fibre Channel helps improve performance and enables more flexible addressability and much greater attachment distances when compared to a normal SCSI attachment.

A key requirement of this type of increased connectivity is providing consistent management interfaces for configuration, monitoring, and management of these SAN components. This type of connectivity allows companies to begin to reap the benefits of Fibre Channel technology, while also protecting their current storage investments.

The flexibility and simplification of the SAN infrastructure can be dramatically improved by using Fibre Channel over Ethernet (FCoE), which evolved over the last few years. This enablement can easily replace dedicated switching solutions for LAN and SAN with a single device that is able to transfer both types of data: IP packets and storage data. We call these deployments, *Converged Networks*. In the following topics, we briefly present the basic migration steps to convergency.

10.2.2 Adding capacity

The addition of storage capacity to one or more servers might be facilitated while the device is connected to a SAN. Depending on the SAN configuration and the server operating system, it might be possible to add or remove devices without stopping and restarting the server.

If new storage devices are attached to a section of a SAN with loop topology (mainly tape drives), the *loop initialization primitive (LIP)* might affect the operation of other devices that are on the loop. This setback might be overcome by slowing down the operating system activity to all of the devices on that particular loop, before you attach the new device. This setback is far less of a problem with the latest generation of loop-capable switches. If storage devices are attached to a SAN by a switch, then the use of the switch and management software makes it possible to make the devices available to any system that is connected to the SAN.

10.2.3 Data movement and copy

Data movement solutions require that data be moved between similar or dissimilar storage devices. Today, data movement or replication is performed by the server or multiple servers. The server reads data from the source device, perhaps transmitting the data across a LAN or WAN to another server. Then, the data is written to the destination device. This task ties up server processor cycles and causes the data to travel twice over the SAN. The data travels one time from the source device to a server, and then a second time from a server to a destination device.

The objective of SAN data movement solutions is to avoid copying data through the server, and across a LAN or WAN. This practice frees up server processor cycles and LAN or WAN bandwidth. Today, this data replication can be accomplished in a SAN by using intelligent tools and utilities and between data centers that use, for example, FCoE protocol on a WAN.

The following sections list some of the available copy services functions.

Data migration

One of the critical tasks for a SAN administrator is to move data between two independent SAN infrastructures. The administrator might move data from an old storage system that is being discontinued, to the new enterprise and highly performing disk system. There are basically two scenarios: When SANs are independent and cannot be interconnected together, even if they reside in the same data center; and when the disk systems can be cross-connected through SAN switches.

Data replication over storage area networks

In this scenario, we are able to interconnect both storage devices (both SANs) together and migrate data directly from an old to the new storage box. This step is completed without any interruption of service or performance affect on the application or host server. This type of migration is what we consider, a *block-level data copy*. In this type of migration, storage systems do not analyze the data on disks, they just split it into blocks and copy the data that has changed or been modified. Many storage vendors, including IBM, offer replication services for their disk storage systems as an optional feature of service delivery, usually as a part of a backup and recovery solution. Copy services can be even further extended to long distances through a WAN to fulfill disaster recovery requirements or just to make application services highly available across geographies.

Figure 10-1 demonstrates how this data (logical unit number (LUN)) migration works, in principle.

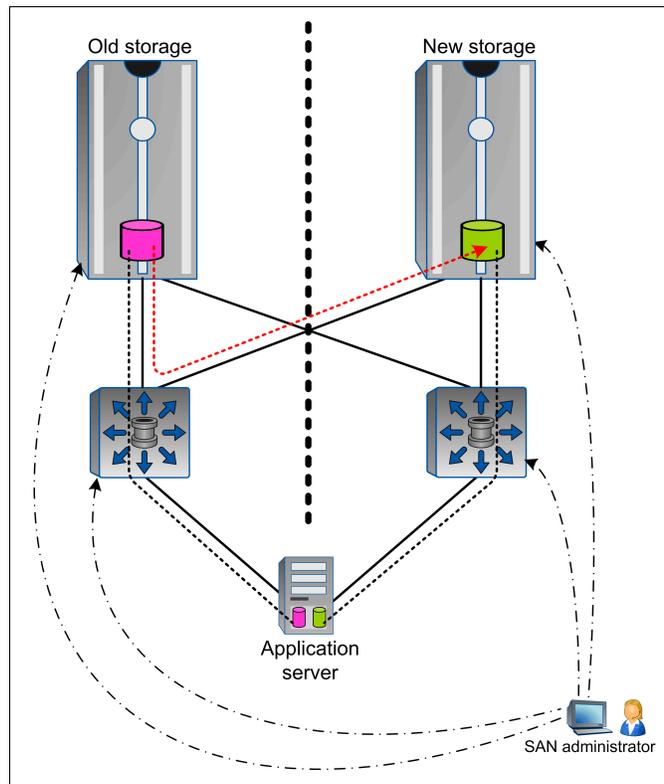


Figure 10-1 SAN-based replication of LUNs

In Figure 10-1, the storage administrator is challenged to migrate data to the newly deployed, highly performing disk storage system without interruption to the most critical SAP

applications of the client. Luckily, we can manage both source and target storage systems. These systems are configured to communicate together through SAN switches. Disk copy services are able to replicate specific LUNs from the old to the new storage devices and, most importantly, without any performance affect to the SAP application.

In addition, this procedure is often used to prepare a standby application server that is connected to the replicated disk LUNs. Or, this procedure is used just to replace the old server hardware where the SAP application is running, all with the minimum outage necessary to switch the application over to the prepared server.

Host-based data migration

Host-based migration of storage data is the option that is used when the storage administrator is not able to establish a connection between the source and target disk storage system. This type of migration usually happens in data centers with two independent SANs. In most cases, each of these SANs is managed by a different team of administrators or even by different vendors.

The principle of the migration is shown in Figure 10-2.

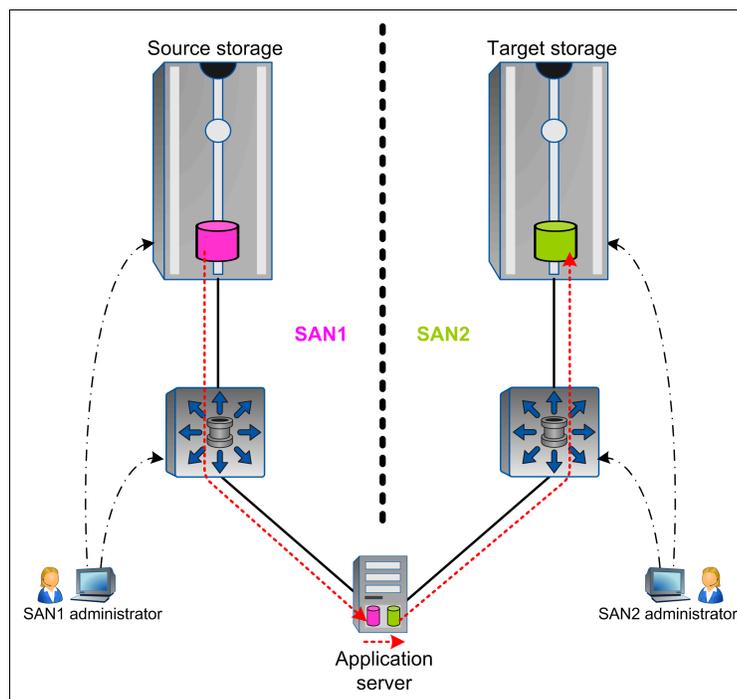


Figure 10-2 Host-based migration of data

The application server is connected to both SANs by use of independent host bus adapters (HBAs). Application owners and SAN2 administrators analyze the current disk structure that is assigned from the source storage system. The same disk capacity is to be assigned by the SAN2 administrator to the application server. The application or system owner then migrates the data from the source disk to the target disk. This migration is done manually by using the operating system functions (the application is offline), or disk mirroring must be enabled. When the data is synchronized between the source and target disks, the mirror can be broken, source disks can be unassigned, and the source storage system can be disconnected from the application server. The disadvantage of this solution is a significant I/O operation on the source and target LUNs that can potentially affect the performance of critical applications.

Remote data copy and migration

Remote copy or *data migration* is a business requirement that is used to protect data from disasters, or to migrate data from one location to avoid application downtime for planned outages such as hardware or software maintenance. Another challenge of remote copy services is to provide a highly available or fault-tolerant infrastructure for business critical systems and applications across data centers, typically over long distances, sometimes even continents.

Remote copy solutions are either *synchronous* or *asynchronous*, and they require different levels of automation to guarantee data consistency across disks and disk subsystems. Remote copy solutions are implemented only for disks at a physical or logical volume data block level. There are complete solutions from various vendors to support data migration projects to optimally schedule and use client network resources and to eliminate the affect on critical production environments. Products such as these help clients efficiently and effectively migrate the whole SAN data volumes from small remote data centers to the central one across a WAN without interruption to the service.

In the future, with more advanced storage management techniques such as outboard hierarchical storage management and file pooling, remote copy solutions would need to be implemented at the file level. This implies more data to be copied, and requires more advanced technologies to guarantee data consistency across files, disks, and tape in multi-server heterogeneous environments. Data center networking infrastructure is required to support various data transfer protocols to support these requirements. Examples of these interfaces include: FCoE, Converged Enhanced Ethernet (CEE), or simple iSCSI.

Real-time snapshot copy

Another outboard copy service that is enabled by Fibre Channel technology is *real-time snapshot* (also known as *T0* or *time=zero*) copy. This service is the process of taking an online snapshot, or freezing the data (databases, files, or volumes) at a certain time. This process allows the applications to update the original data while the frozen copy is duplicated. With the flexibility and extensibility that Fibre Channel brings, these snapshot copies can be made to either local or remote storage devices. The requirement for this type of function is driven by the need for 24x7 availability of key database systems. This solution is optimal in homogeneous infrastructures that consist of the devices from a single vendor.

10.2.4 Upgrading to faster speeds

One of the other considerations of any SAN environment is how newer, faster technology is to be introduced. Both 8 Gbps Fibre Channel and 10 GbE products already have a significant footprint in the market and participate in data center networking. We are now seeing vendors move forward with even faster technologies and products such as 16 Gbps Fibre Channel ports and HBAs. For most applications, this faster technology does not mean that they can immediately benefit. Applications that have random or “bursty” I/O might not necessarily gain any advantage. Only those applications and systems that stream large amounts of data are likely to see the most immediate benefits. One place that makes sense for 16 Gbps to be used is the inter-switch link (ISL). This scenario has two advantages: The increased speed between switches is the obvious one; the other advantage is that it might be possible to have fewer ISLs with the increased bandwidth. Having fewer ISLs means that it might be possible to reassess ISLs and use them to attach hosts or storage.

Another consideration that must be taken into account is the cost factor. IT architects and investors must evaluate their current SAN solutions in data centers and make strategic decisions to determine if it is beneficial to continue with the upgrade to a dedicated Fibre Channel solution that is running 16 Gbps devices. Or, the architects and investors must

determine if it is the right time to consider an upgrade to converged networks and use, for example, FCoE. There are many products that are available on the market that support such transformations and transitions and protect the investments of the clients for the future.

10.3 Infrastructure simplification

High on the list of critical business requirements is the need for IT infrastructures to better support business integration and transformation efforts. At the heart of these data center efforts is often the simplification and streamlining of core storage provisioning services and storage networking.

Viewed in the broadest sense, infrastructure simplification represents an optimized view and evolutionary approach (or the next logical step beyond basic server consolidation and virtualization) for companies on the verge of becoming true on-demand businesses. That is, businesses that are highly competitive in the market.

Is your IT infrastructure a complex set of disparate, server-specific, siloed applications that are operating across an endless area of servers (that is: transaction processing servers, database servers, tiered application servers, data gateways, human resource servers, accounting servers, manufacturing servers, engineering servers, email servers, web servers, and so on)? If so, then you must be able to answer questions such as: “Where can we deploy the next application?” “Where can we physically put the next server?” “How can we extend our storage resources?” “How can we connect more virtual or physical servers?” Or, just “Is there a simpler way to manage all of these servers?” We try to answer all of these questions in the following topics.

10.3.1 Where does the complexity come from?

A SAN, in theory, is a simple thing. It is a path from a server to a common storage resource. Therefore, where did all the complexity come from?

Limited budgets and short-sighted strategic thinking push IT organizations into looking for short-term solutions to pain points. When a new application or project becomes available, the easy, inexpensive option is to add another low-cost server. Because this “server sprawl”, or proliferation of UNIX and Windows Intel servers is an attractive short-term solution, the infrastructure costs to support these inexpensive servers often exceeds the purchase price of the server.

Now, storage systems are also added to the sprawl. Every server has two or four HBAs and a share of the consolidated storage. As more servers are added, we run out of SAN ports, so we add another switch, and then another, and finally another. Now we have “SAN sprawl” with a complex interlinked fabric that is difficult to maintain or change.

To make things more difficult, the servers are probably purchased from multiple vendors, with decisions made on cost, suitability to a specific application, or merely the personal preference of someone. The servers of different vendors are tested on specific SAN configurations. Every server producer has its own interoperability matrix or list of SAN configurations that the vendor tested, and that particular vendor supports. It might be difficult for a SAN administrator to find the appropriate devices and configurations that work together smoothly.

10.3.2 Storage pooling

Before SANs, the concept of the physical pooling of devices in a common area of the computing center was often not possible, and when it was possible, it required expensive and unique extension technology. By introducing a network between the servers and the storage resources, this problem is minimized. Hardware interconnections become common across all servers and devices. For example, common trunk cables can be used for all servers, storage, and switches.

This section briefly describes the two main types of storage device pooling: *disk pooling* and *tape pooling*.

Disk pooling

Disk pooling allows multiple servers to use a common pool of SAN-attached disk storage devices. Disk storage resources are pooled within a disk subsystem or across multiple IBM and non-IBM disk subsystems. And, capacity is assigned to independent file systems supported by the operating systems on the servers. The servers are potentially a heterogeneous mix of UNIX, Microsoft Windows, and even IBM z/OS.

Storage can be dynamically added to the disk pool and assigned to any SAN-attached server when and where it is needed. This function provides efficient access to shared disk resources without a level of indirection that is associated with a separate file server. This scenario is possible because storage is effectively *directly attached* to all the servers, and efficiencies of scalability result from consolidation of storage capacity.

When storage is added, *zoning* can be used to restrict access to the added capacity. Because many devices (or LUNs) can be attached to a single port, access can be further restricted by using LUN-masking. This masking means being able to specify who can access a specific device or LUN.

Attaching and detaching storage devices can be done under the control of a common administrative interface. Storage capacity can be added without stopping the server, and can be immediately made available to the applications.

Figure 10-3 shows an example of disk storage pooling across two servers.

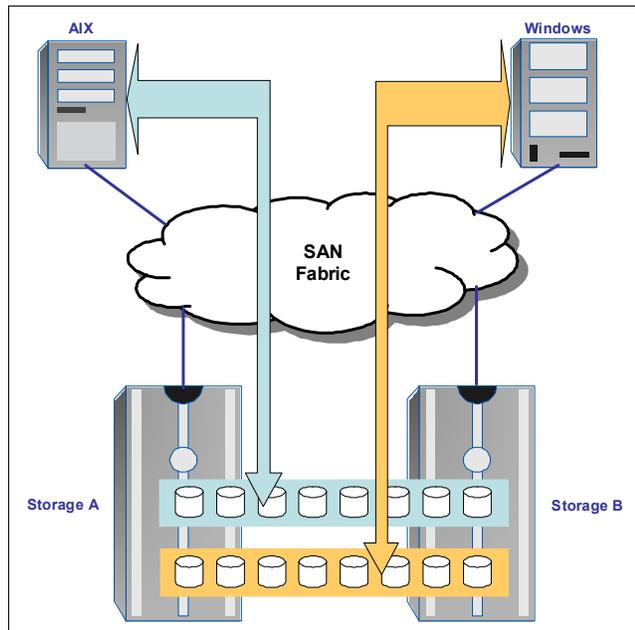


Figure 10-3 Disk pooling concept

In Figure 10-3, one server is assigned a pool of disks that are formatted to the requirements of the file system, and the second server is assigned another pool of disks, possibly in another format. The third pool might be the space that is not yet allocated or can be a pre-formatted disk for future use. Again, all the changes in the disk structure can be done dynamically, without any interruption to the service.

Tape pooling

Tape pooling addresses the problem faced today in an open systems environment in which multiple servers are unable to share tape resources across multiple hosts. Older methods of sharing a device between hosts consist of either manually switching the tape device from one host to the other, or writing applications that communicate with connected servers through distributed programming.

Tape pooling allows applications on one or more servers to share tape drives, libraries, and cartridges in a SAN environment in an automated, secure manner. With a SAN infrastructure, each host can directly address the tape device as though it is connected to all of the hosts.

Tape drives, libraries, and cartridges are owned by either a central manager (tape library manager) or a peer-to-peer management implementation. These devices are dynamically allocated and reallocated to systems (tape library clients) as required, based on demand. Tape pooling allows for resource sharing, automation, improved tape management, and added security for tape media.

Software is required to manage the assignment and locking of the tape devices to serialize tape access. Tape pooling is an efficient and cost effective way of sharing expensive tape resources, such as automated tape libraries. At any particular instant in time, a tape drive can be owned by one system only.

This concept of tape resource sharing and pooling is proven in medium to enterprise backup and archive solutions that use, for example, IBM Tivoli Storage Manager with SAN-attached IBM tape libraries.

Logical volume partitioning

At first sight an individual might ask: “How will logical volume partitioning make my infrastructure simpler? It looks as if we are creating more and more pieces to manage in my storage”. Conceptually, this thought is correct, but the benefit of *logical volume partitioning* is to address the need for maximum volume capacity and to effectively use it within target systems. It is essentially a way of dividing the capacity of a single storage server into multiple pieces. The storage subsystems are connected to multiple servers, and storage capacity is partitioned among the various subsystems.

Logical disk volumes are defined within the storage subsystem and assigned to servers. The logical disk is addressable from the server. A logical disk might be a subset or superset of disks that are only addressable by the subsystem itself. A logical disk volume can also be defined as subsets of several physical disks (striping). The capacity of a disk volume is set when defined. For example, two logical disks, with different capacities (for example, 50 GB and 150 GB) might be created from a single 300 GB hardware addressable disk. Each of these two disks is assigned to a different server, leaving 100 GB of unassigned capacity. A single 2000 GB logical disk might also be created from multiple real disks that exist in different storage subsystems. The underlying storage controller must have the necessary logic to manage the volume grouping and guarantee access securely to the data.

The function of a storage controller can be further used by some of the storage virtualization engines, such as the IBM SAN Volume Controller. This engine, when compared to environments that do not use this controller, offers even better and more scalability and virtualization of storage resources. The SAN Volume Controller provides these benefits with less management effort and clearer visibility to the target host systems.

Figure 10-4 shows multiple servers that are accessing logical volumes that were created by using the different alternatives that are previously mentioned. (The logical volume, *Unallocated volume*, is not assigned to any server).

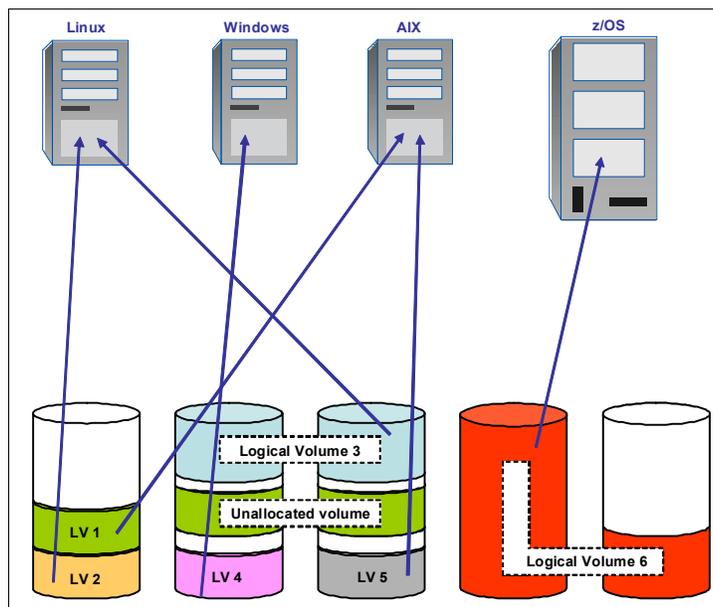


Figure 10-4 Conceptual model of logical volume partitioning

10.3.3 Consolidation

We can improve scalability, security, and manageability by enabling devices in separate SAN fabrics to communicate without merging fabrics into a single, large SAN fabric. This capability enables clients to initially deploy separate SAN solutions at the departmental and data center levels and then to consolidate them into large enterprise SAN solutions. This consolidation occurs as their experience and requirements grow and change. This type of solution is also known as *Data Center Bridging*.

Clients deploy multiple SAN islands for different applications with different fabric switch solutions. The growing availability of iSCSI server capabilities creates the opportunity for low-cost iSCSI server integration and storage consolidation. Additionally, depending on the choice of router, they provide *Fibre Channel over IP (FCIP)* or iFCP capability.

The available multiprotocol SAN routers provide an iSCSI Gateway Service to integrate low-cost Ethernet-connected servers to existing SAN infrastructures. It also provides Fibre Channel, FC-FC Routing Service to interconnect multiple SAN islands without requiring the fabrics to merge into a single large SAN.

Figure 10-5 shows an example of using a multiprotocol router and converged core switch to extend SAN capabilities across the long distances or just over the metropolitan areas.

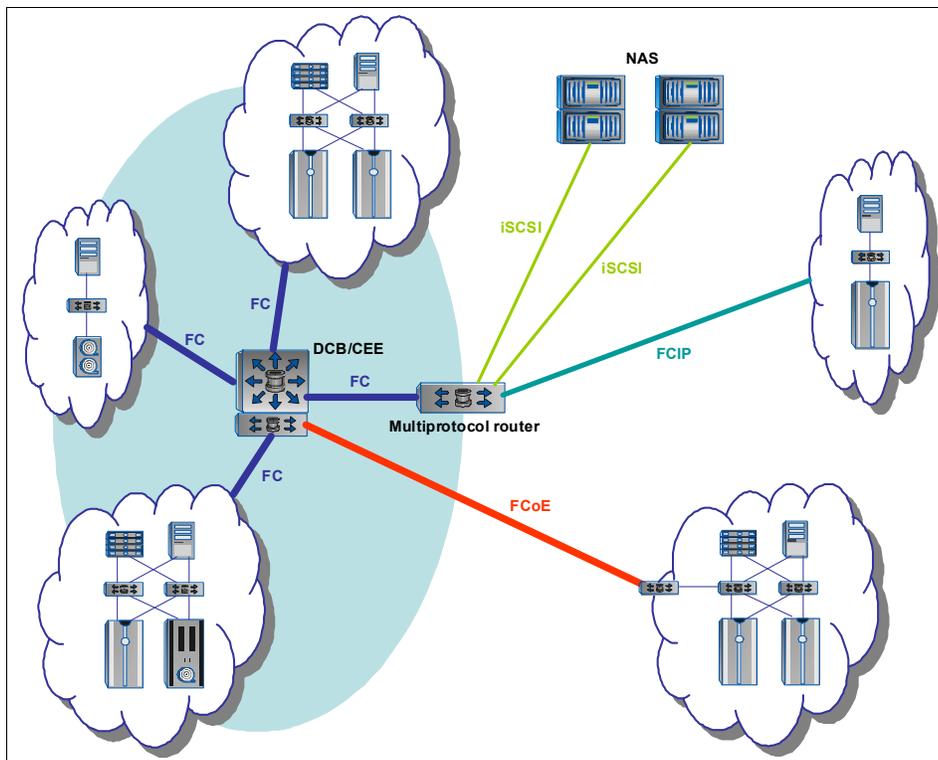


Figure 10-5 The concept of SAN consolidation

A multiprotocol capable router solution brings a number of benefits to the marketplace. In our example, there are discrete SAN islands, and number of different protocols involved. To merge these SAN fabrics, it would involve a number of disruptive and potentially expensive actions:

- ▶ Downtime
- ▶ Purchase of more switches and ports
- ▶ Purchase of HBAs

- ▶ Migration costs
- ▶ Configuration costs
- ▶ Purchase of more licenses
- ▶ Ongoing maintenance

However, by installing a multiprotocol router or core FCoE-enabled switch or director, there are many advantages:

- ▶ Least disruptive method
- ▶ No need to purchase extra HBAs
- ▶ Minimum number of ports to connect to the router
- ▶ No expensive downtime
- ▶ No expensive migration costs
- ▶ No ongoing maintenance costs other than router
- ▶ Support of other protocols
- ▶ Increases return on investment (ROI) by consolidating resources
- ▶ Can be used to isolate the SAN environment to be more secure

There are more benefits that the router and core switch can provide. In this example, an FC-FC routing service that negates the need for a costly SAN fabric merge exercise, the advantages are apparent, and real. The router can also be used to provide the following benefits:

- ▶ Device connectivity across multiple SANs for infrastructure simplification
- ▶ Tape-backup consolidation for information lifecycle management (ILM)
- ▶ Long-distance SAN extension for business continuity
- ▶ Low-cost server connectivity to SAN resources

10.3.4 Migration to a converged network

Medium and enterprise data centers usually run multiple separate networks. These networks include an Ethernet network for client to server and server to server communications, and a Fibre Channel SAN for the same type of connections. To support various types of networks, data centers use separate redundant interface modules for each network: Ethernet network interface cards (NICs) and Fibre Channel interfaces (HBAs) in their servers, and redundant pairs of switches at each layer in the network architecture. Use of parallel infrastructures increase capital costs, makes data center management more difficult, and diminishes business flexibility.

The principle of consolidation of both independent networks to share a single, integrated networking infrastructure relies on utilization of FCoE and helps address these challenges efficiently and effectively. In the following topics, we briefly describe how to upgrade your current infrastructure to a converged network in three principal steps. The prerequisite of the converged network is lossless 10 Gbps (GoE) Ethernet, inline with the Data Center Bridging standards (DCB).

Figure 10-6 on page 219 presents the concept of the migration to convergency.

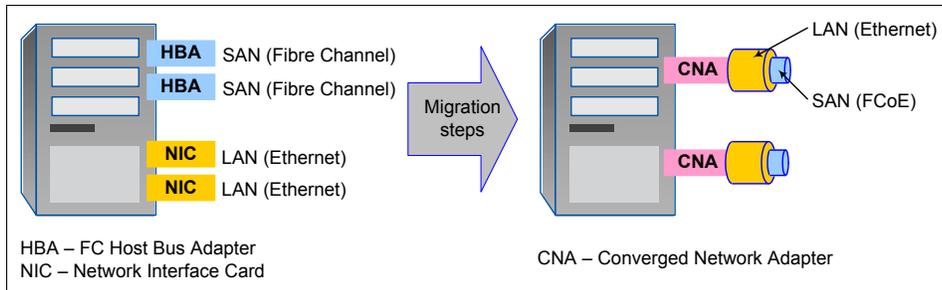


Figure 10-6 Conceptual model of migration to a converged network

The following list provides a summary of the key benefits of upgrading to a converged network:

- ▶ Reduced capital expenditures by 15 - 30%, depending on current infrastructure
- ▶ Reduced operational and management costs by 30 - 50%
- ▶ Improved network and storage efficiencies
- ▶ Increased assets and storage utilization
- ▶ Improved flexibility and business agility
- ▶ Reduced power consumption by 15 - 20%

The following three steps detail the process of migrating to a converged network:

1. Access layer convergence

Assume that we have separate adapters for a 1 Gbps or 10 Gbps Ethernet communication (2 - 8 adapters for each server) and FC HBAs for storage networking (2 - 6 adapters, usually dual-port). In this step, we are going to replace these combinations by using Converged Network Adapters (CNAs). See Figure 10-7 on page 220.

Fabric deployment: For illustration purposes, only a single fabric data center solution is presented in all of the figures. In real data centers, dual-fabric deployment is essential.

Additionally, we must install a switch that supports both protocols, IP and FCoE, usually as a top-of-rack (TOR) device. Therefore, the TOR device that supports DCB standards and multiple protocols can continue to work with the existing environment. The device can also separate the network traffic from the data storage traffic and direct each of them to the correct part of the data center networking infrastructure. All the segmentation of the traffic is done at the access layer; this is the first step of the overall process.

Figure 10-7 shows the first step in the migration process: Access layer convergence.

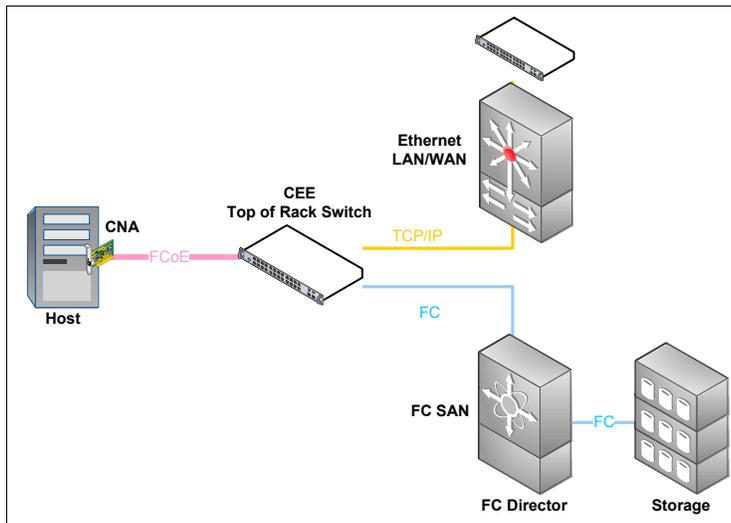


Figure 10-7 Access layer convergence

2. Fabric convergence

The second step is to implement more core types of switches that support data center bridging and converged network protocols. Therefore, rather than implementing a converged network on TOR switches or blades, we move this function to the core directors. There is a second stage of the development of DCB standards that introduces Multi-Hop bridging because there are different solutions from each of the vendors of the SAN networking products. See Figure 10-8 for details.

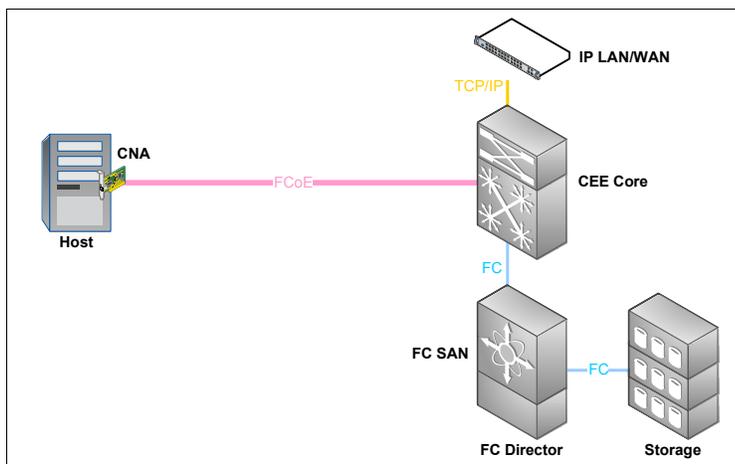


Figure 10-8 Fabric convergence

3. Storage convergence

For the final step of the migration, we implement native FCoE-enabled storage devices. Now, there are various vendors with midrange to enterprise disk storage systems that already offer FCoE. This step enables clients to migrate the current FC-attached storage data to the FCoE-enabled storage system and disconnect the original FC core and edge switches. This step dramatically reduces the requirements for operation and management of the infrastructure, reduces the power consumption, and simplifies the complexity of the

network (rack space and cabling). Figure 10-9 shows the final status of the converged network infrastructure.

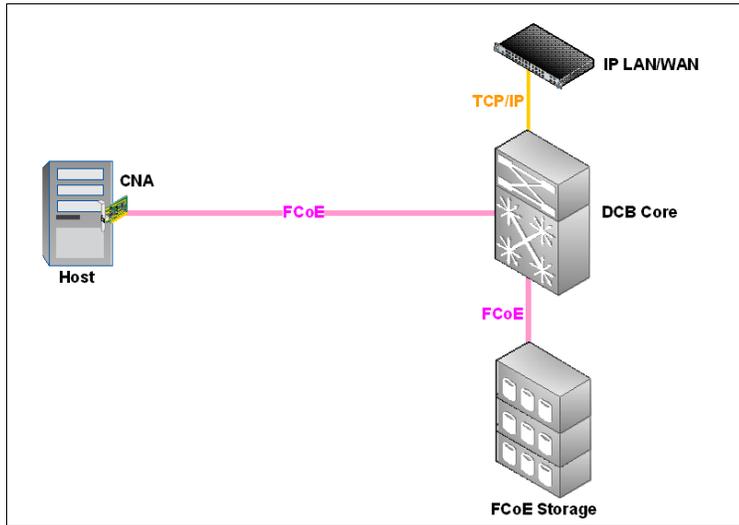


Figure 10-9 Storage convergence

FCoE in the converged network offers several benefits and advantages over existing approaches to I/O consolidation:

- Compatibility with existing Fibre Channel SANs by preserving well-known Fibre Channel concepts. Examples of these concepts include: Virtual SANs (VSANs), worldwide names (WWNs), FC IDs (FCIDs), multipathing, and zoning to servers and storage arrays.
- A high level of performance, comparable to the performance of current Ethernet and Fibre Channel networks. These networks are achieved by using a hardware-based Ethernet network infrastructure that is not limited by the overhead of higher-layer TCP/IP protocols.
- The exceptional scalability of Ethernet at the highest available speeds (1, 10, and 40 GoE, and eventually 100 GoE).
- Simplified operations and management (no change to the management infrastructure currently deployed in SANs).

10.4 Business continuity and disaster recovery

On-demand businesses rely on their IT systems to conduct business. Everything must be working all the time. Failure truly is not an option these days. A sound and comprehensive business continuity strategy that encompasses high availability, near-continuous operations, fault-tolerant systems, and disaster recovery is essential.

Today, data protection of multiple network or SAN-attached servers is performed according to one of two backup and recovery paradigms: Local backup and recovery, or network backup and recovery.

The *local backup and recovery* solution has the advantage of speed because the data does not travel over the network. However, with a local backup and recovery approach, there are costs for overhead (because local devices must be acquired for each server, and are thus

difficult to use efficiently). There are also costs for management overhead (because of the need to support multiple tape drives, libraries, and mount operations).

The *network backup and recovery* approach which uses shared tape libraries and tape drives, the best over SAN, is cost-effective. This approach is efficient because it allows centralization of storage devices that use one or more network-attached devices. This centralization shortens the ROI because the installed devices are used efficiently. One tape library can be shared across many servers. Management of a network backup and recovery environment is often simpler than the local backup and recovery environment. This is true because it eliminates the potential need to perform manual tape mount operations on multiple servers.

SANs combine the best of both approaches. This benefit is accomplished by central management of backup and recovery. It is also accomplished by assigning one or more tape devices to each server and by using Fibre Channel protocols to transfer data directly from the disk device to the tape device, or vice versa, over the SAN.

Another popular topic in this category is instant business continuity if there is a device failure. Clients that run the most critical applications cannot afford to wait until their data is restored to a fixed or standby server or devices from backups. Server or application clusters allow clients to continue their business with minimal outage or even without any disruption. We are referring to highly available or fault-tolerant systems.

In the following sections, we describe these approaches in more detail.

10.4.1 Clustering and high availability

SAN architecture naturally allows multiple systems (target hosts) to access the same disk resources in the medium or enterprise disk storage systems, even concurrently. This feature enables specific applications (AIX IBM HACMP™, IBM PowerHA®, Windows Cluster Services, and so on) that run on the hosts to introduce highly available or fault-tolerant application systems. These systems assure that in a case of a single host failure, the application is automatically (without any manual administrator intervention) moved over to the backup cluster host (high availability: short outage). Or, the failed host is isolated from application processing and the workload is balanced between other working cluster nodes (fault-tolerant: no disruption to service).

The conceptual scheme of a highly available cluster solution is shown in Figure 10-10.

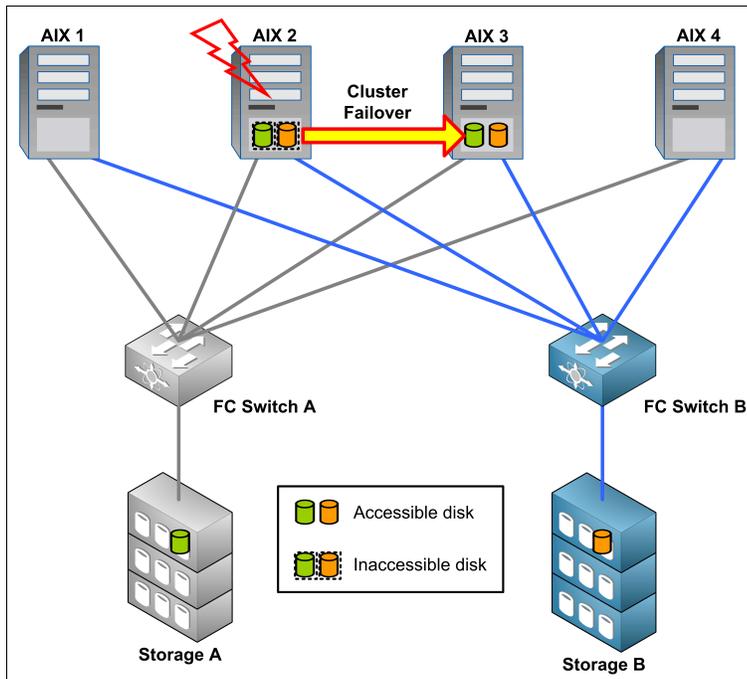


Figure 10-10 Highly available cluster system

In Figure 10-10, an application is running on system AIX2 and it is managing mirrored disks from both of the storage systems (green and orange). SAN zoning allows both cluster nodes (AIX2 and AIX3) to operate the same set of disks and the cluster has one primary cluster node active (an active-passive cluster). At the time of the AIX2 failure, cluster services that are running on AIX3 recognize the failure and automatically move all the application resources to AIX3. The disk sets are activated there and the application is started in the correct sequence.

Figure 10-11 explains the configuration of a fault-tolerant clustered environment.

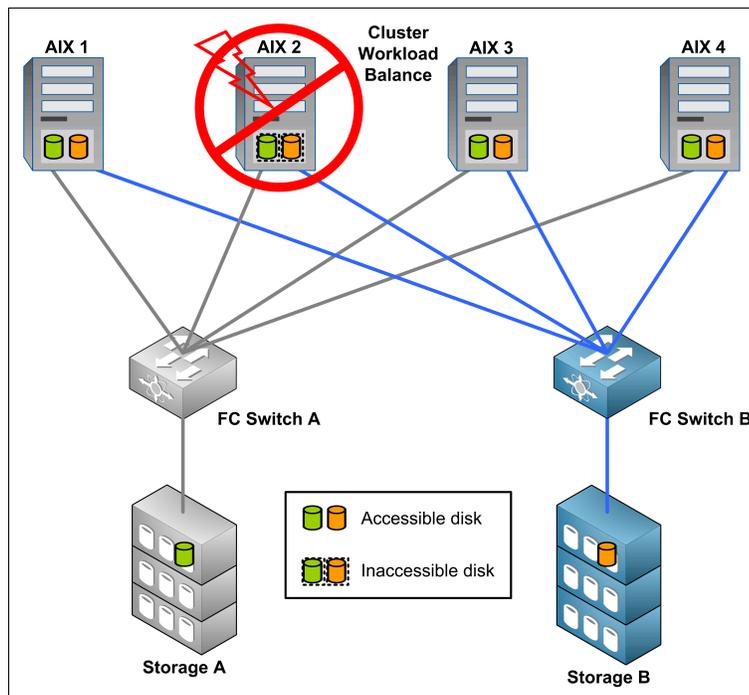


Figure 10-11 Fault-tolerant cluster system

SAN zoned disks are available to all four AIX host systems, the clusters are active and the master application works concurrently on them with workload balancing. If there is an AIX2 system failure, the cluster application automatically deactivates the assigned disks and redistributes the workload between the remaining active cluster nodes. There is no interruption to the business. A configuration such as this is costly and usually is only employed for business critical applications, such as banking systems and air traffic control systems.

10.4.2 LAN-free data movement

The network backup and recovery paradigm implies that data flows from the backup and recovery client (usually a file or database server) to the centralized backup and recovery server. Or, the data flows between the backup and recovery servers, over the Ethernet network. The same applies for the archive or storage management applications. Often the network connection is the bottleneck for data throughput especially in the case of large database systems. This is because of the network connection bandwidth limitations. The SAN offers the advantage to offload the backup data out of the LAN.

Tape drive and tape library sharing

A basic requirement for LAN-free backup and recovery is the ability to share tape drives and tape libraries between the central backup tape repository and backup and recovery clients with large database files. Systems with a higher number of small files still use the network for data transportation because they cannot benefit from a LAN-free environment.

In the tape drive and tape library sharing approach, the backup and recovery server or client that requests a backup copy to be copied to or from tape, reads or writes the data directly to the tape device by using SCSI commands. This approach bypasses the network transport's latency and network protocol path length; therefore, it can offer improved backup and

recovery speeds in cases where the network is the constraining factor. The data is read from the source device and written directly to the destination device. The central backup and recovery server controls only the tape mount operations and stores the references (metadata) into its embedded database system.

Figure 10-12 shows an example of tape library sharing and LAN-free backups.

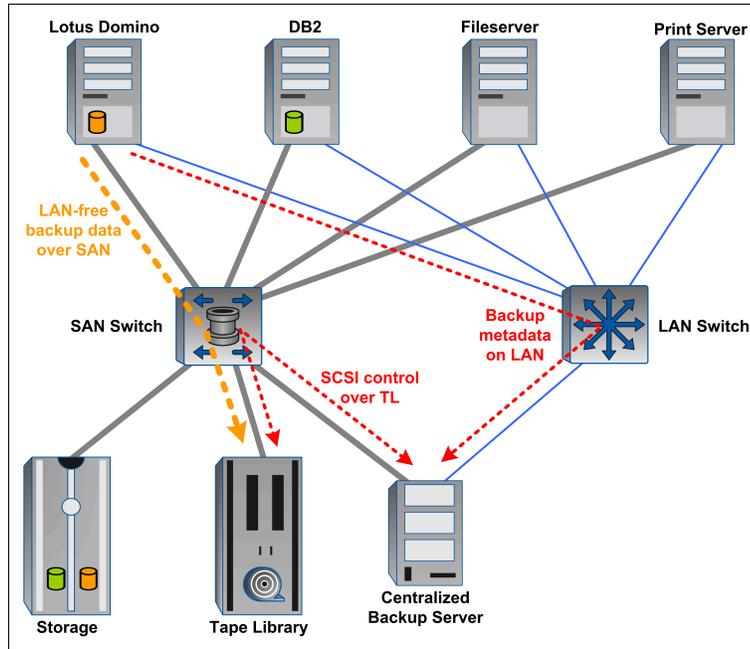


Figure 10-12 LAN-free and LAN-based backups

IBM Lotus® Domino and IBM DB2 database systems benefit from a greater performance of backups over Fibre Channel, directly to tapes. However, small servers with a higher number of files still continue to back up to a LAN or WAN.

IBM offers enterprises a centralized backup and recovery solution that supports various platforms and database systems. IBM Tivoli Storage Manager (ITSM) and its component ITSM for SANs, enables clients to perform online backups and archives of large application systems directly to tape over a SAN, and without significant affect on performance.

10.4.3 Disaster backup and recovery

A SAN can facilitate disaster backup solutions because of the greater flexibility that is allowed in connecting storage devices to servers. Backup solutions are also simplified because of the greater distances that are supported when compared to SCSI restrictions. It is now possible to perform extended distance backups for disaster recovery within a campus or city, as shown in Figure 10-13 on page 226.

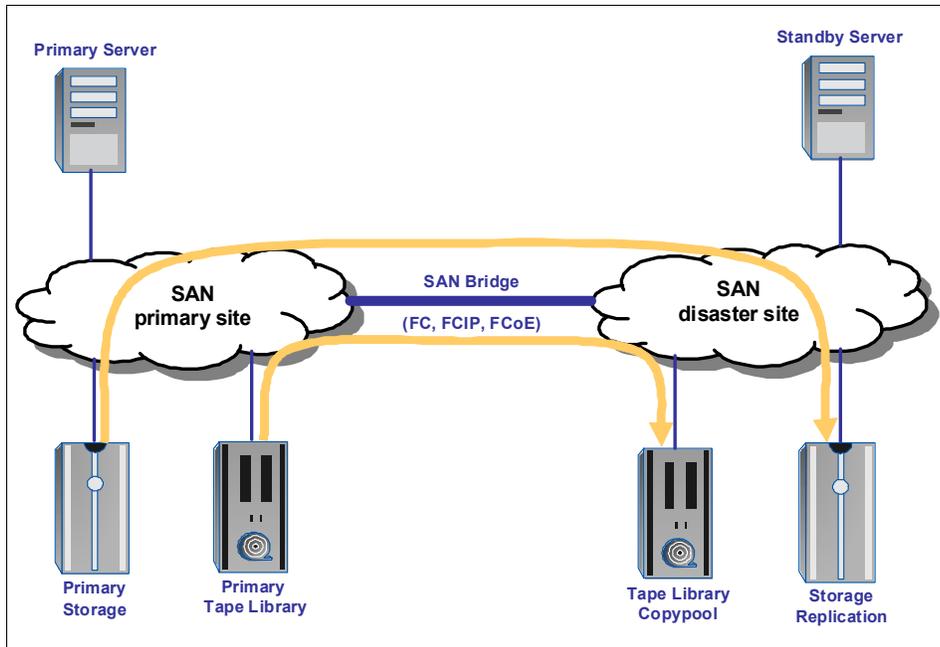


Figure 10-13 Disaster backup at a remote site using SAN bridging

When longer distances are required, SANs must be connected using gateways and WANs. One of the solutions is FCoE.

Depending on business requirements, disaster protection deployments might use copy services that are implemented in disk subsystems and tape libraries (that might be achieved by utilization of SAN services), SAN copy services, and most likely a combination of both.

10.5 Information lifecycle management

Information lifecycle management (ILM) is a process for managing information through its lifecycle, from conception until disposal, in a manner that optimizes storage and access at the lowest cost.

ILM is not just hardware or software, it includes processes and policies to manage the information. It is designed upon the recognition that different types of information can have different values at different points in their lifecycle. Predicting storage needs and controlling costs can be especially challenging as the business grows.

The overall objectives of managing information with ILM are to help reduce the total cost of ownership (TCO) and help implement data retention and compliance policies. To effectively implement ILM, owners of the data need to determine how information is created, how it ages, how it is modified, and when it can safely be deleted. ILM segments the data according to value, which can help create an economical balance and sustainable strategy to align storage costs with businesses objectives and information value.

10.5.1 Information lifecycle management

To manage the data lifecycle and to make your business ready for on-demand services, there are four main elements that can address your business in an ILM-structured environment:

- ▶ Tiered storage management
- ▶ Long-term data retention and archiving
- ▶ Data lifecycle management
- ▶ Policy-based archive management

In the next sections, we describe each of these elements in more detail.

10.5.2 Tiered storage management

Most organizations today seek a storage solution that can help them manage data more efficiently. They want to reduce the costs of storing large and growing amounts of data and files and to maintain business continuity. Through tiered storage, you can reduce overall disk-storage costs by providing the following benefits:

- ▶ Reducing overall disk-storage costs by allocating the most recent and most critical business data to higher performance disk storage. Costs can also be reduced by moving older and less critical business data to lower-cost disk storage.
- ▶ Speeding business processes by providing high-performance access to the most recent and most frequently accessed data.
- ▶ Reducing administrative tasks and human errors. Older data can be moved to lower-cost disk storage automatically and transparently.

Typical storage environment

Storage environments typically have multiple tiers of *data value*, such as application data that is needed daily, and archive data that is accessed infrequently. However, typical storage configurations offer only a single tier of storage, as shown in Figure 10-14, which limits the ability to optimize cost and performance.

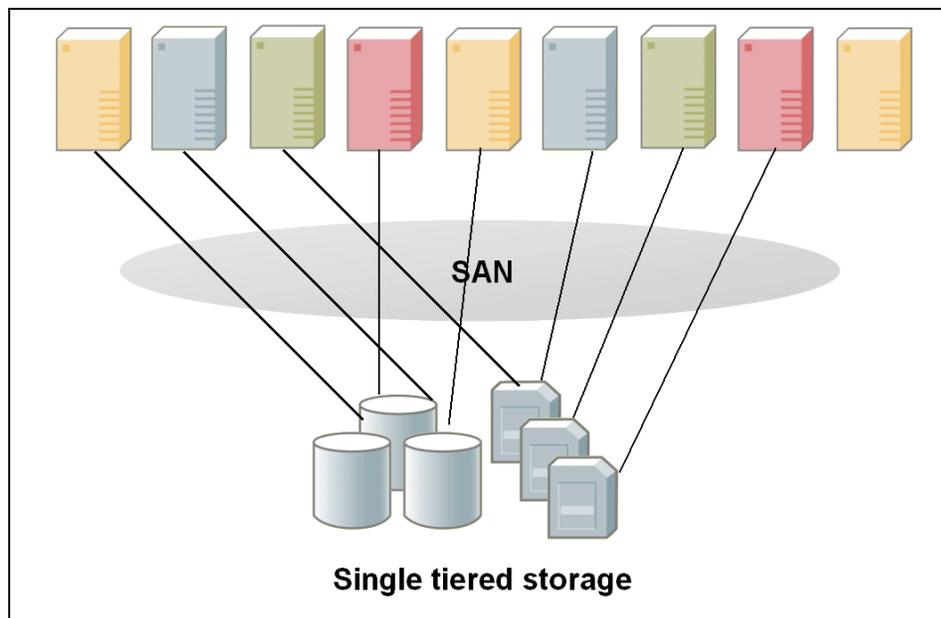


Figure 10-14 Traditional non-tiered storage environment

Multi-tiered storage environment

A tiered storage environment that uses the SAN infrastructure affords the flexibility to align storage cost with the changing value of information. The tiers are related to data value. The most critical data is allocated to higher performance disk storage, while less critical business data is allocated to lower-cost disk storage.

Each storage tier provides different performance metrics and disaster recovery capabilities. Creating classes and storage device groups is an important step to configure a tiered storage ILM environment.

Figure 10-15 shows a multi-tiered storage environment.

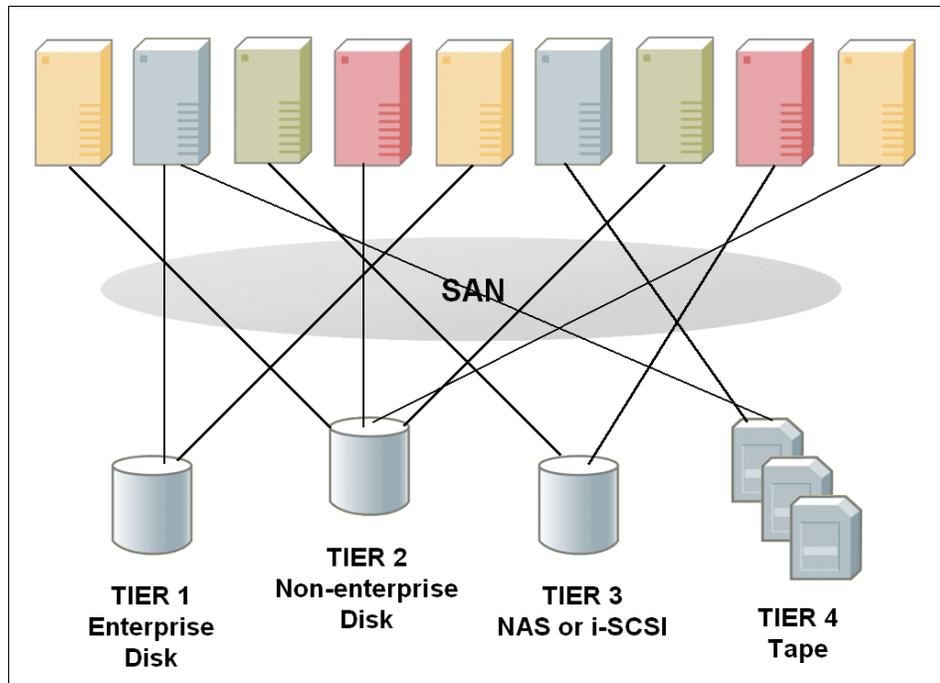


Figure 10-15 ILM tiered storage environment

An IBM ILM solution in a tiered storage environment is designed with the following factors in mind:

- ▶ Reduces the TCO of managing information. It can help optimize data costs and management, freeing expensive disk storage for the most valuable information.
- ▶ Segments data according to value. Segmenting can help create an economical balance and sustainable strategy to align storage costs with business objectives and information value.
- ▶ Helps to make decisions about moving, retaining, and deleting data, because ILM solutions are closely tied to applications.
- ▶ Manage information and determine how it is managed based on content, rather than migrating data that is based on technical specifications. This approach can help result in more responsive management, and offers you the ability to retain or delete information in accordance with business rules.
- ▶ Provide the framework for a comprehensive Enterprise Content Management strategy.

10.5.3 Long-term data retention

There is a rapidly growing class of data that is best described by how it is managed rather than the arrangement of its bits. The most important attribute of this data is its retention period, hence it is called *retention managed data*, and it is typically kept in an archive or a repository. In the past, it has been variously known as archive data, fixed content data, reference data, unstructured data, and other terms that imply its read-only nature. It is often measured in terabytes and is kept for long periods of time, sometimes forever.

In addition to the sheer growth of data, laws and regulations that govern the storage and secure the retention of business and client information, are increasingly becoming part of the business landscape. These obstacles make data retention a major challenge to any institution. An example of this challenge is the Sarbanes-Oxley Act, enacted in the US in 2002.

Businesses must comply with these laws and regulations. Regulated information can include email, instant messages, business transactions, accounting records, contracts, or insurance claims processing, all of which can have different retention periods. These periods can be for two years, seven years, or can be retained forever. Data is an asset when it must be kept; however, data kept past its mandated retention period might also become a liability. Furthermore, the retention period can change because of factors such as litigation. All of these factors mandate tight coordination and the need for ILM.

Not only are there numerous state and governmental regulations that must be met for data storage, there are also industry-specific and company-specific ones. And these regulations are constantly being updated and amended. Organizations must develop a strategy to ensure that the correct information is kept for the right period of time, and is readily accessible when it must be retrieved at the request of regulators or auditors.

It is easy to envision the exponential growth in data storage that results from these regulations and the accompanying requirement for a means of managing this data. Overall, the management and control of retention-managed data is a significant challenge for the IT industry when you take into account factors such as cost, latency, bandwidth, integration, security, and privacy.

10.5.4 Data lifecycle management

At its core, the process of ILM moves data up and down a path of tiered storage resources. These resources include high-performance, high-capacity disk arrays; lower-cost disk arrays such as Serial Advanced Technology Attachment (SATA); tape libraries; and permanent archival media, where appropriate. Yet ILM involves more than just data movement; it encompasses scheduled deletion and regulatory compliance as well. Because decisions about moving, retaining, and deleting data are closely tied to application use of data, ILM solutions are usually closely tied to applications.

ILM has the potential to provide the framework with a comprehensive information management strategy, and to help ensure that information is stored on the most cost-effective media. This framework helps enable administrators to use tiered and virtual storage, and to process automation. By migrating unused data off more costly, high-performance disks, ILM is designed to help by performing the following functions:

- ▶ Reduce the cost to manage and retain data.
- ▶ Improve application performance.
- ▶ Reduce backup windows and ease system upgrades.
- ▶ Streamline data management.
- ▶ Allow the enterprise to respond to demand in real time.

- ▶ Support a sustainable storage management strategy.
- ▶ Scale as the business grows.

ILM is designed to recognize that different types of information can have different value at different points in their lifecycle.

As shown in Figure 10-16, data can be allocated to a specific storage level aligned to its cost, with policies that define when and where data is moved.

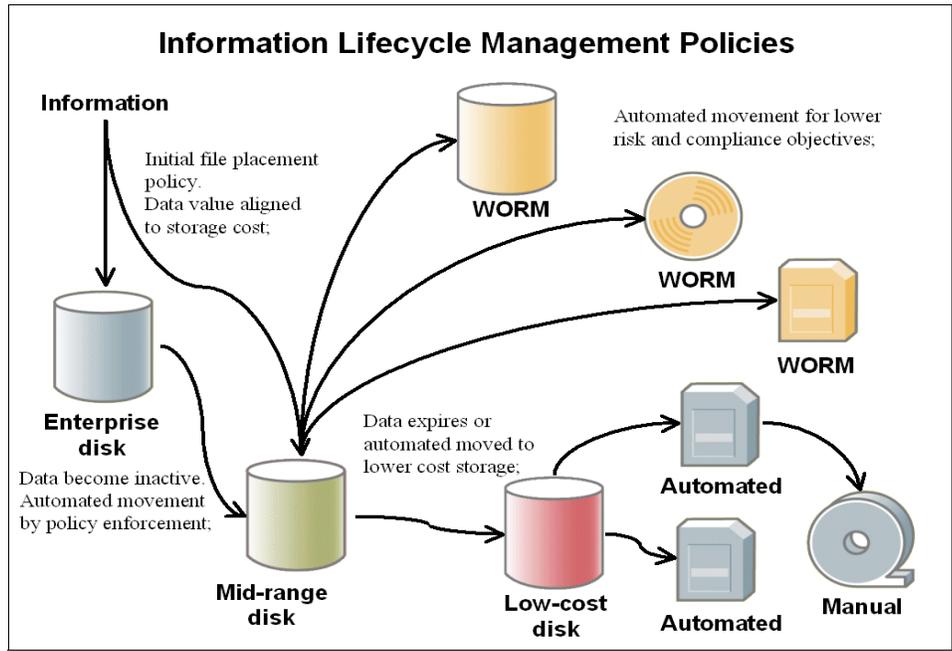


Figure 10-16 ILM policies

10.5.5 Policy-based archive management

Businesses of all sizes upgrade to electronic-business solutions and a new way of doing business. These organizations already have mountains of data and content that have been captured, stored, and distributed across the enterprise. This wealth of information provides a unique opportunity. By incorporating these assets into electronic-business solutions, and at the same time, delivering newly generated information media to their employees and clients, a business can reduce costs and information redundancy. Organizations can also use the potential profit-making aspects of their information assets.

Growth of information in corporate databases such as enterprise resource planning (ERP) and email systems makes organizations consider moving unused data off the high-cost disks. To effectively manage information in corporate databases, businesses must take the following steps:

- ▶ Identify database data that is no longer being regularly accessed and move it to an archive, where it remains available.
- ▶ Define and manage what to archive, when to archive, and how to archive from the mail or database system to the back-end archive management system.

Database archive solutions can help improve performance for online databases, reduce backup times, and improve application upgrade times.

Email archiving solutions are designed to reduce the size of corporate email systems by moving email attachments and messages to an archive from which they can easily be recovered if needed. This action helps reduce the need for user management of email, improves the performance of email systems, and supports the retention and deletion of email.

The way to archive is to migrate and store all information assets into an electronic-business enabled content manager. ERP databases and email solutions generate large volumes of information and data objects that can be stored in content management archives. An archive solution allows you to free system resources, while maintaining access to the stored objects for later reference. Allowing the archive to manage and migrate data objects gives a solution the ability to have ready access to newly created information that carries a higher value. And at the same time, you are still able to retrieve data that is archived on less expensive media.

More information is available in *ILM Library: Information Lifecycle Management Best Practices Guide*, SG24-7251.



Storage area networks and green data centers

System storage networking products and their deployment in large enterprise data centers significantly participate in the rapid growth of floorspace, power, and cooling resources.

In this chapter, we briefly introduce the concepts of a green data center strategy and how a storage area network (SAN) and IBM System Networking align with the green goal. In addition, we also describe the IBM smarter data center that facilitates the evolution of energy efficient operations.

11.1 Data center constraints

Many data centers are running out of power and space. They cannot add more servers because they reached either their power or space limits, or perhaps they reached the limit of their cooling capacity.

In addition, environmental concerns are becoming priorities because they can also impede the ability of a company to grow. Clients all over the world prefer to purchase products and services from companies that have a sustainable approach to the environment. Clients also want products and services that are able to meet any targets that might be imposed on them, whether from inside their company or from outside, in the form of government legislation.

Because environmental sustainability is a business imperative currently, many data center clients are looking at ways to save energy and cut costs so that their company can continue to grow. However, it is also a time to consider transformation in spending, not just cutting costs. Only smarter investments in technology and perhaps a different way of thinking is needed to achieve green efficiency in data centers.

Data centers must provide flexibility to respond quickly to future unknowns in business requirements, technology, and computing models. They need to adapt to be more cost-effective, for both capital expenditure (CAPEX) and operational expenditure (OPEX). Additionally, they require active monitoring and management capabilities to provide the operational insights to meet the required availability, resiliency, capacity planning, and energy efficiency.

The IT architects need to consider four key factors that drive the efficiency of data center operations:

- ▶ **Energy cost.** The cost of a kilowatt of electricity rose only slightly in recent years, but the cost of operating servers has increased significantly. The context around this paradox is that the energy consumption of the servers is increasing exponentially faster than the utility cost. Rising demand has accelerated the adoption of virtualization technologies and increased virtual image densities. Hence, this increased demand drives total server energy consumption higher, while the amortized cost of operation per workload is decreasing.
- ▶ **Power capacity.** Some companies cannot deploy more servers because more electrical power is not available. Many suppliers, especially those in crowded urban areas, are telling clients that power feeds are at capacity limits and that they have no more power to sell. New server, storage, and networking products give better performance at lower prices, but can also be power hungry. The effort to overcome a power supply threshold is a huge investment.
- ▶ **Cooling capacity.** Many data centers are now 10 - 20 years old, and the cooling facilities are not adapted to the present needs. Traditional cooling methods allowed for 2-3 kW of cooling per rack. Today's requirements are 20-30 kW per rack. Heat density is many times past the design point of the data center.
- ▶ **Space limitation.** Each time a new project or application comes online, new images, servers, or storage subsystems are added. Therefore, the space utilization is growing exponentially because of business requirements. When images, servers, and storage cannot be added, except by building another data center, growth becomes expensive.

11.1.1 Energy flow in data center

To understand how to reduce energy consumption, you need to understand where and how the energy is used. You can study energy use in a data center by taking three different views:

- ▶ How energy is distributed between IT equipment (servers, storage, network devices) and supporting facilities (power, cooling, and lighting)
- ▶ How energy is distributed between the different components of the IT equipment (processor, memory, disk, and so on)
- ▶ How the energy allocated to IT resources is used to produce business results (Are idle resources that are powered on, spending energy without any effect?)

Figure 11-1 shows how energy is used by several components of a typical non-optimized data center. Each component is divided into two portions: IT equipment (servers, storage, network) and the infrastructure around it that supports the IT equipment (chillers, humidifiers, air conditioners, power distribution units, uninterruptible power supply (UPS), lights, and so on).

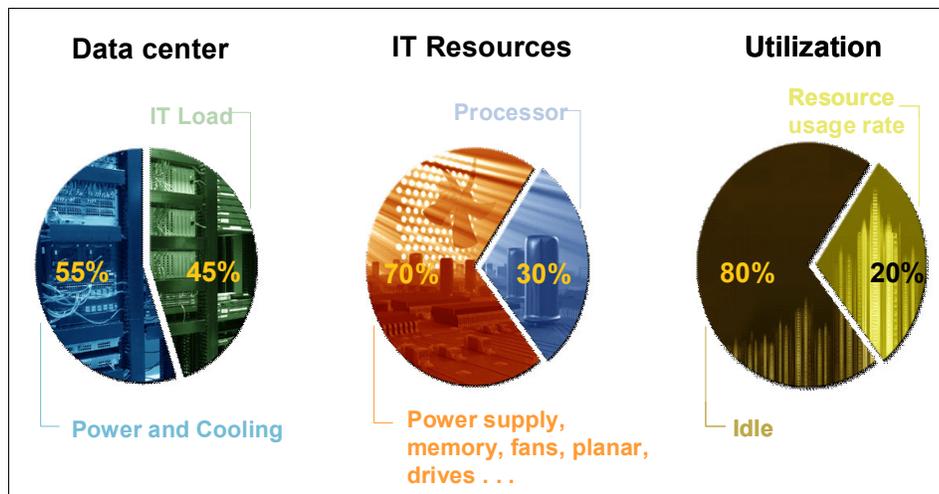


Figure 11-1 Energy usage in a typical data center¹

In typical data centers, the IT equipment does not use 55% of the overall energy that is brought into the data center. Therefore, this portion of the energy is not producing calculations, data storage, and so on. The concept of a green data center is able to eliminate this waste and reduce such inefficiency.

Energy conversion: Basic laws of thermodynamics state that energy cannot be created or destroyed, it changes only in form. The efficiency of this conversion is less than 100% (in a real world, much less than 100%).

Solution designers and IT architects must also consider the energy consumption of the components at the IT equipment level. For example, in a typical server, the processor uses only 30% of the energy and the remainder of the system uses 70%. Therefore, efficient hardware design is crucial. Features like virtualization of physical servers can help to change this ratio to a more reasonable value.

Finally, companies need to consider the use of IT resources in the data center. A typical server utilization rate is around 20%. Underutilized systems can be a significant issue because much energy is expended on non-business purposes, thus wasting a major

¹ Data source: Creating Energy-Efficient Data Centers, US Department of Energy

investment. Again, server virtualization, consolidation, and addressed provisioning of IT resources help to use the entire capacity of your IT equipment.

Data centers must become immensely more efficient to meet their needs while they keep costs in check as the demand for and price of resources continue to rise. But the realization of this efficiency requires a deep and pervasive transformation in how data centers are designed, managed, operated, populated, and billed. These aspects mandate a unified and coordinated effort across organizational and functional boundaries toward a common set of goals.

In the following text, we introduce the concept of green data centers and how IBM supports the migration to the concept of data centers of the next generation that are effective, cost-efficient, and environment friendly.

11.2 Data center optimization

To enable your data center to become more effective, less power-consuming, and more cost-efficient in terms of infrastructure management and operation, IT architects must consider two components of the migration strategy:

- ▶ **Optimization of the site and facilities.** Includes data center cooling, heating, ventilation, air conditioning (HVAC), UPS, and power distribution to the site and within the data center. A standby power supply or alternative power sources must also be considered.
- ▶ **Optimization of the IT equipment.** Pertains to IT equipment in the data center that generates business value to the clients, such as servers (physical and virtual), disk and tape storage devices, and networking products.

Applying innovative technologies within the data center can yield more computing power per kilowatt. The IT equipment continues to become more energy efficient. Technology evolution and innovation outpace the life expectancy of data center equipment. Therefore, many companies are finding that replacing older IT equipment with newer models can significantly reduce overall power and cooling requirements and free up valuable floor space. For example, IBM studies demonstrate that blade servers reduce power and cooling requirements 25 - 40% over 1U technologies. Replacing equipment before it is fully depreciated might seem unwise. However, the advantages that new models can offer (lower energy consumption and two to three times more computing power than older models), combined with potential space, power, and cooling recoveries, are usually enough to offset any lost asset value.

11.2.1 Strategic considerations

The strategy of moving towards having a green data center and the overall cost effective IT infrastructure consists of four major suggested steps:

- ▶ Centralization
 - Consolidate many small remote centers into fewer
 - Reduce infrastructure complexity
 - Improve facility management
 - Reduce staffing requirements
 - Improve management cost
- ▶ Physical consolidation
 - Consolidate many servers into fewer on physical resource boundaries
 - Reduce system management complexity
 - Reduce physical footprint of servers in the data center

- ▶ Virtualization
 - Remove physical resource boundaries
 - Increase hardware utilization
 - Allocate less than physical boundary
 - Reduce software license costs
- ▶ Application integration
 - Migrate many applications to fewer, more powerful server images
 - Simplify IT environment
 - Reduce operational resources
 - Improve application-specific tuning and monitoring

11.3 Green storage

Computer systems are not the only candidates for energy savings. As the amount of managed data grows over the years exponentially, one of the top candidates is also storage systems within data centers. Each component of the storage system has power and cooling requirements.

Published studies show that the proportion of energy that is used by storage disk systems and storage networking products varies 15 - 25% of the overall energy consumption of the typical data center. This number significantly increases because there are continuously growing requirements for storage space.

However, no matter what efficiency improvements are made, active (spinning) disk drives still use energy if they are powered on. Consequently, the most energy-intensive strategy for data storage is to keep all of the data of the organization on active disks. While this process provides the best access performance, it is not the most environmentally-friendly approach, nor is it normally an absolute requirement.

Green storage technologies occupy less raw storage capacity to store the same amount of native valuable client data. Therefore, the energy consumption per gigabyte of raw capacity falls accordingly.

The storage strategy for green data centers includes the following elements:

- ▶ Information lifecycle management (ILM)
- ▶ Consolidation and virtualization
- ▶ On-demand storage provisioning
- ▶ Hierarchical storage and storage tiering
- ▶ Compression and deduplication

In the following sections, we briefly describe each of them.

11.3.1 Information lifecycle management

Information lifecycle management (ILM) is a process for managing information through its lifecycle, from conception until disposal, in a manner that optimizes storage and access at the lowest cost.

ILM is not just hardware or software, it includes processes and policies to manage the information. It is designed upon the recognition that different types of information can have different values at different points in their lifecycle. Predicting storage needs and controlling costs can be especially challenging as the business grows. Although the total value of stored

information increases overall, historically, not all data is created equal, and the value of that data to business operations fluctuates over time.

This trend is shown in Figure 11-2, and is commonly referred to as the *data lifecycle*. The existence of the data lifecycle means that all data cannot be treated the same.

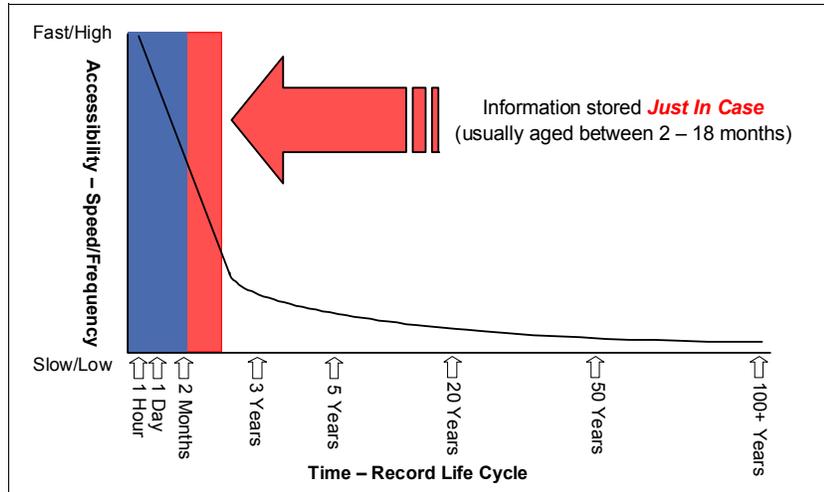


Figure 11-2 Data lifecycle

However, infrequently accessed or inactive data can become suddenly valuable again as events occur, or as new business initiatives or projects are taken on. Historically, the requirement to retain information results in a *buy more storage* mentality. However, this approach serves to only increase overall operational costs and complexity, and increases the demand for hard-to-find qualified personnel.

Typically, only around 20% of the information is active and frequently accessed by users. The remaining 80% is either inactive or even obsolete. See the Figure 11-3 for details.

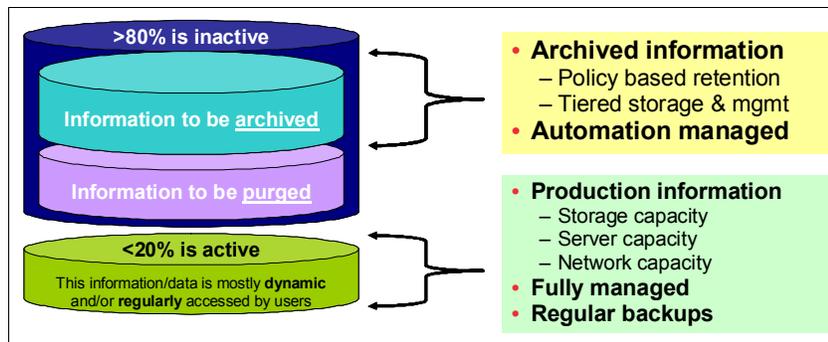


Figure 11-3 Usage of data

The automated identification of the storage resources in an infrastructure and an analysis of how effectively those resources are being used, is the crucial part of the ILM. File-system and file-level evaluation uncovers categories of files that, if deleted or archived, can potentially represent significant reductions in the amount of data that must be stored, backed up, and managed. The key position in the ILM process is the automated control through policies that are customizable with actions that can include centralized alerting, distributed responsibility, and fully automated response. This process also includes deletion of data.

See more details in *ILM Library: Information Lifecycle Management Best Practices Guide*, SG24-7251.

11.3.2 Storage consolidation and virtualization

As the need for data storage continues to spiral upward, traditional physical approaches to storage management become increasingly problematic. Physically expanding the storage environment can be costly, time-consuming, and disruptive. These drawbacks are compounded when expansion must be done again and again in response to ever-growing storage demands. Yet, manually improving storage utilization to control growth can be challenging. Physical infrastructures can also be inflexible at a time when businesses must be able to make even more rapid changes to stay competitive.

The alternative is a centralized, consolidated storage pool of disk devices that are easy to manage and are transparent to be provisioned to the target host systems. Going further, the consolidated or centralized storage can be virtualized, where storage virtualization software presents a view of storage resources to servers that is different from the actual physical hardware in use.

Figure 11-4 shows storage consolidation.

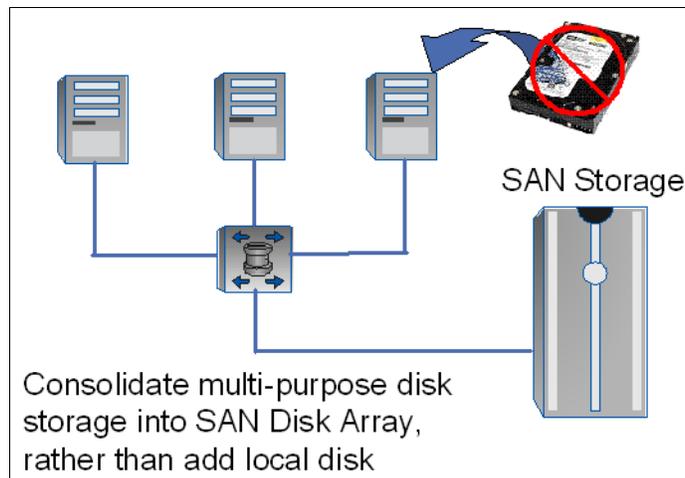


Figure 11-4 Storage consolidation

This logical view can hide undesirable characteristics of storage while it presents storage in a more convenient manner for applications. For example, storage virtualization might present storage capacity as a consolidated whole, hiding the actual physical boxes that contain the storage.

In this way, storage becomes a logical pool of resources that exists virtually, regardless of where the actual physical storage resources are in the larger information infrastructure. These software-defined virtual resources are easier and less disruptive to change and manage than hardware-based physical storage devices, since they do not involve moving equipment or making physical connections. As a result, they can respond more flexibly and dynamically to changing business needs. Similarly, the flexibility that is afforded by virtual resources makes it easier to match storage to business requirements.

Virtualization offers significant business and IT advantages over traditional approaches to storage. Storage virtualization can help organizations in the following ways:

- ▶ Reduce data center complexity and improve IT productivity by managing multiple physical resources as fewer virtual resources.
- ▶ Flexibly meet rapidly changing demands by dynamically adjusting storage resources across the information infrastructure.
- ▶ Reduce capital and facility costs by creating virtual resources instead of adding more physical devices.
- ▶ Improve utilization of storage resources by sharing available capacity and deploying storage on demand only as it is needed.
- ▶ Deploy tiers of different storage types to help optimize storage capability while controlling cost and power and cooling requirements.

The value of a virtualized infrastructure is in the increased flexibility that is created by having pools of system resources on which to draw and in the improved access to information that is afforded by a shared infrastructure. There is also value in the lower total cost of ownership that comes with decreased management costs, increased asset utilization, and the ability to link infrastructure performance to specific business goals.

For more information about how IBM Storage Virtualization solutions can help your organization meet its storage challenges, see the *IBM Information Infrastructure Solutions Handbook*, SG24-7814, or see this website:

<http://www.ibm.com/systems/storage/virtualization/>

11.3.3 On-demand storage provisioning

The provisioning of SAN-attached storage capacity to a server can be a time consuming and cumbersome process. The task requires skilled storage administrators. And the complexity of the task can restrict the ability of an IT department to respond quickly to requests to provision new capacity. However, there is a solution to this issue through automation. An on-demand storage provisioning solution monitors the current disk usage of specified target host systems and applications and allocates more disk capacity for the period of business need.

End-to-end storage provisioning is the term that is applied to the whole set of steps that are required to provision usable storage capacity to a server. Provisioning covers the configuration of all the elements in the chain. This process includes the steps from carving out a new volume on a storage subsystem, through creating a file system at the host and making it available to the users or applications.

Typically, this process involves a storage administrator that uses several different tools, each focused on the specific task at hand, or the tasks are spread across several people. This spread of tasks and tools creates many inefficiencies in the provisioning process, which affect the ability of IT departments to respond quickly to changing business demands. The resulting complexity and high degree of coordination that is required can also lead to errors and can possibly affect the systems and application availability.

Automation of the end-to-end storage provisioning process by using workflow automation can significantly simplify this task of provisioning storage capacity. Each individual step is automated and the rules for preferred practices around zoning, device configuration, and path selection can be applied automatically. The benefits are increased responsiveness to business requirements, lower administration costs, and higher application availability.

11.3.4 Hierarchical storage and tiering

Companies continue to deploy storage systems that deliver many different classes of storage that range from high performance and high cost to high capacity and low cost. Through the deployment of SANs, many of these storage assets are now physically connected to servers that run many different types of applications and create many kinds of information. Finally, with the arrival of network-resident storage services for distributed management of volume, files, and data replication, IT managers have more control. The managers can intelligently provision, reallocate, and protect storage assets to meet the needs of many different applications across the network, instead of device by device.

In a tiered storage environment, data is classified and assigned dynamically to different tiers. For example, we can use expensive fast-performing storage components to store often-accessed and mission-critical files, in contrast to using cheaper storage for less used non-critical data. In conclusion, it improves efficiency and saves costs. We can identify the following typical storage tiers, which are categorized based on performance and cost for each gigabyte:

- ▶ High-performing SAN-attached disk systems (SSD, SAS)
- ▶ Medium-performing SAN-attached disks (SAS, SATA)
- ▶ Network-attached storage systems (NAS)
- ▶ Tape storage and other media with sequential access

Each level of storage tier can be assigned manually by a storage administrator or data can be moved automatically between tiers, which is based on migration policies. The conceptual model of storage tiering is shown in Figure 11-5.

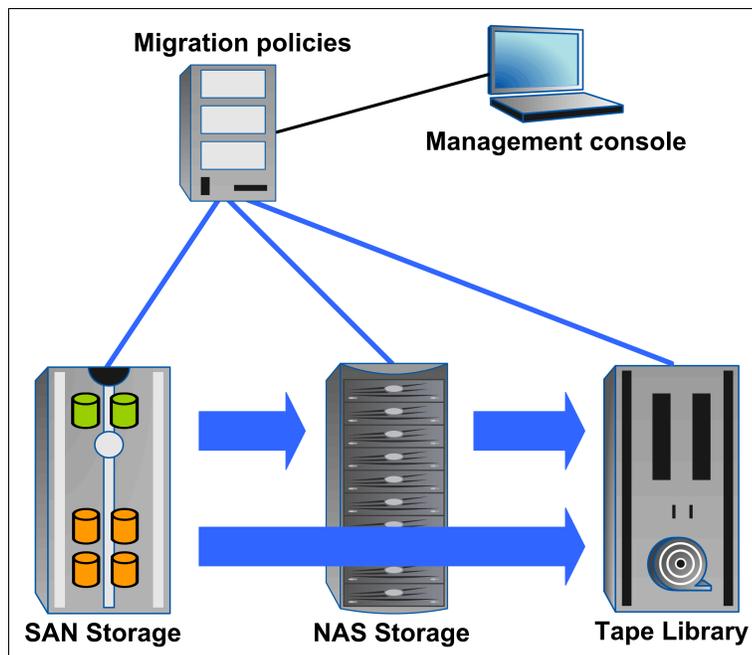


Figure 11-5 Principle of tiered storage

IBM offers various tools and utilities for storage tiering and hierarchical management for different scenarios. Tools range from those such as IBM Easy Tier, which is used in enterprise disk storage systems, all the way through to the IBM Tivoli Storage Manager (ITSM) for Hierarchical Storage Management for Windows, and ITSM for Space Management for the AIX/Linux platform. Specific positions in tiered management have the

IBM Global Parallel File System (GPFS™) that allows data migration between different level of storage sources as well.

11.3.5 Data compression and deduplication

Business data growth rates will continue to increase rapidly in the coming years. Likewise, retention and retrieval requirements for new and existing data will expand, driving still more data to disk storage. As the amount of disk-based data continues to grow, there is an ever-increasing focus on improving data storage efficiencies across the information infrastructure.

Data reduction is a tactic which can decrease the disk storage and network bandwidth that is required, lower total cost of ownership (TCO) for storage infrastructures, optimize use of existing storage assets, and improve data recovery infrastructure efficiency. Compression and deduplication and other forms of data reduction are features that can exist within multiple components of the information infrastructure.

Compression immediately reduces the amount of required physical storage across all storage tiers. This solution, which supports online compression of existing data, allows storage administrators to gain back free disk space in the existing storage system. This solution must be done without having to change any administrative processes or enforce users to clean up or archive data. The benefits to the business are immediate because the capital expense of upgrading the storage system is delayed. As data is stored in compressed format at the primary storage system, all other storage tiers and the transports in between them observe the same benefits. Replicas, backup images, and replication links all require less expenditure after the implementation of compression at the source.

After compression is applied to the stored data, the required power and cooling for each unit of storage is reduced. This reduction is possible because more logical data is stored on the same amount of physical storage. In addition, within a particular storage system, more data can be stored; therefore, the overall rack unit requirements are lowered. Figure 11-6 on page 243 shows the typical compression rates that can be achieved with specific IBM products.

The exact compression ratio depends on the nature of the data. IBM sees compression ratios as high as 90 percent in certain Oracle database configurations and about 50% with PDF files. As always, compression ratios vary by data type and how the data is used.

In contrast to compression, the *data deduplication* mechanism identifies identical chunks of data within a storage container. This process keeps only one copy of each chunk, while all the other logically identical chunks are pointed to this chunk. There are various implementations of this method. One option is in-line deduplication and the other one is post-processing. Each chunk of data must be identified in a way that is easily comparable. Chunks are processed by using a parity calculation or cryptographic hash function. This processing gives the chunks shorter identifiers that are known as *hash values*, *digital signatures*, or *fingerprints*. These fingerprints can be stored in an index or catalog where they can be compared quickly with other fingerprints to find matching chunks.

Figure 11-6 on page 243 shows the typical compression rates that can be achieved with specific IBM products.

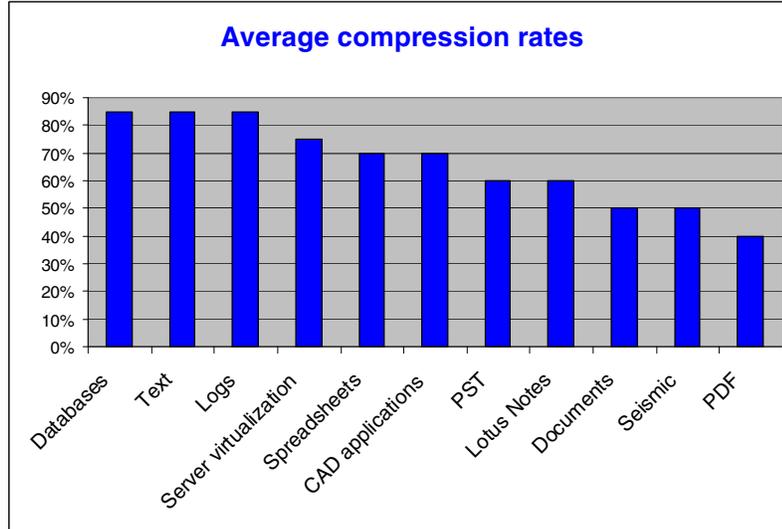


Figure 11-6 The average compression rates

Data deduplication processing can occur on the client, on an infrastructure server, or on the storage system. Each option has factors to consider:

- ▶ **Client-based deduplication:** This process reduces the amount of data that is being transferred over the network to the storage system. But, there are often requirements for extra CPU and disk I/O processing on the client side.
- ▶ **Server-based deduplication:** This process allows you to deduplicate the data of multiple clients at a scheduled time. But, this process requires extra CPU and disk I/O processing on the infrastructure server (for example, IBM Tivoli Storage Manager server).
- ▶ **Storage-based deduplication:** This process occurs at the disk storage device level, where the data is stored. This type of deduplication is generally transparent to the clients and servers. It uses CPU and disk I/O on the storage system.

For more information about data compression, deduplication, and concrete solutions from IBM, see *Introduction to IBM Real-time Compression Appliances*, SG24-7953 and *Implementing IBM Storage Data Deduplication Solutions*, SG24-7888.



The IBM product portfolio

This chapter guides you through the IBM System Storage storage area network (SAN) components that IBM offers through its marketing channels, either as original equipment manufacturer (OEM) products or as an authorized reseller.

In addition to the typical Fibre Channel-attached products, we also provide brief descriptions of storage and virtualization devices, even though they are usually not classified as SAN devices.

12.1 Classification of IBM storage area network products

To stay competitive in the global marketplace, the right people must have the right access to the right information at the right time to be effective, creative, and highly innovative. IBM offers a comprehensive portfolio of SAN switches, storage, software, services, and solutions to reliably bring information to people in a cost effective way. IBM provides flexible, scalable, and open standards-based business-class and global enterprise-class storage networking solutions for the on-demand world.

IBM helps you to align your storage investment with the value of the information by using a wide range of tiered storage options, policy-based automation, and intelligent information management solutions. The IBM System Storage portfolio offers the broadest range of storage solutions in the industry (including disk, tape, SAN, software, financial, and services offerings). This portfolio enables companies to create long-term solutions that can be tailored to their business needs. IBM System Storage tiered disk, tape, and SAN switch solutions provide a wide variety of choices to align and move data to cost-optimized storage. This process is based on policies, matching the storage solution with the service level requirements, and the value of the data in the growing environments.

IBM enables their clients to confidently protect strategic information assets and to efficiently comply with regulatory and security requirements with the unrivaled breadth of storage solutions from IBM. IBM SAN directors and routers provide metro and global connectivity between sites over Internet Protocol, IP networks. IBM SAN extension solutions include data compression and encryption services to help protect your data in flight between secure data centers.

IBM solutions are optimized for the unique needs of midsize organizations, large enterprises, cloud computing providers, and others. Clients can get just what they need, saving time and money. A key benefit of selecting IBM for your next information infrastructure project is access to a broad portfolio of outstanding products and services. IBM offers highly rated, patented technology that delivers unique value.

In this chapter, we do not provide an in-depth analysis into all of the technical details of each product. The intention of this chapter and book is to introduce the principles and basic components of the SAN environments in a reasonable extent, in a way that is easy to understand and follow.

We also do not cover the SAN networking products for IBM BladeCenter® technologies. That technology is covered in the IBM Redbooks publication, *IBM BladeCenter Products and Technology*, SG24-7523.

We can identify the following categories of SAN products that IBM offers:

- ▶ SAN Fibre Channel networking
 - Entry level SAN switches
 - Midrange SAN switches
 - Enterprise SAN directors
 - Multiprotocol routers
- ▶ Storage device subsystems
 - Entry level disk systems
 - Midrange disk systems
 - Enterprise disk systems

- ▶ Tape storage systems
 - Fibre Channel tape drives
 - Autoloaders and entry level tape libraries
 - Midrange tape libraries
 - Enterprise tape libraries
- ▶ Storage virtualization
 - Disk storage virtualization
 - Tape virtualization
 - SAN products for cloud computing
- ▶ IP-based data center networking for SAN environments
 - Hardware products for converged networks
 - Software solutions for virtual fabrics

For more technical details and information about other IBM storage products, see the *IBM System Storage Solutions Handbook*, SG24-5250. And for a comprehensive description of each product and its market position, see the IBM storage website:

<http://www.ibm.com/systems/storage/>

12.2 SAN Fibre Channel networking

This section provides brief information about IBM products for Fibre Channel-based (optical) data center networking solutions, starting from entry level SAN switches to midrange switches, and up to enterprise SAN directors and multiprotocol routers.

For more information about the latest IBM SAN products, see this website:

<http://www-03.ibm.com/systems/networking/switches/san/index.html>

12.2.1 Entry SAN switches

Entry level SAN switches represent easy-to-use preconfigured data center networking solutions for small and medium business (SMB) environments. IBM offers two products that fall into this category:

- ▶ IBM System Storage SAN24B-4 Express
- ▶ Cisco MDS 9124 Express for IBM System Storage

Here we summarize both of them.

IBM System Storage SAN24B-4 Express

This system provides high-performance, scalable, and simple-to-use fabric switching with 8, 16, or 24 ports that operate at 8, 4, 2, or 1 gigabits per second (Gbps) (depending on which optical transceiver is used). This system is for servers that run Microsoft Windows, IBM AIX, UNIX, and Linux operating systems, server clustering, infrastructure simplification, and business continuity solutions. The SAN24B-4 Express includes a *EZSwitchSetup wizard*, which is an embedded setup tool that is designed to guide novice users through switch setup, often in less than 5 minutes. Figure 12-1 shows the front view of the SAN24B-4 switch.



Figure 12-1 Front view of the IBM System Storage SAN24B-4 Express switch

A single SAN24B-4 Express switch can serve as the cornerstone of a SAN for individuals that want to obtain the benefits of storage consolidation and are just beginning to implement Fibre Channel storage systems. Such an entry level configuration can consist of one or two Fibre Channel links to a disk storage array or to a Linear Tape Open (LTO) tape drive. An entry level 8-port storage consolidation solution can support up to seven servers with a single path to either disk or tape. The Ports-on-Demand feature is designed to enable a base switch to grow to 16 and 24 ports to support more servers and more storage devices without taking the switch offline. A high-availability solution can be created with redundant switches. This capability is ideal for server clustering environments.

Such a configuration can support 6 - 22 servers, each with dual Fibre Channel adapters cross-connected to redundant SAN24B-4 Express switches. These switches are cross-connected to a dual controller storage system. The focus of the SAN24B-4 Express is as the foundation of small to medium sized SANs. However, it can be configured to participate as a full member in an extended fabric configuration with other members of the IBM System Storage and former TotalStorage SAN b-type and m-type families. This capability helps

provide investment protection as SAN requirements evolve and grow over time. The IBM System Storage SAN24B-4 Express switch provides the following features and characteristics:

- ▶ Efficient 1U design with 8, 16, or 24 ports configuration on demand
- ▶ Auto-sensing 8, 4, 2, or 1 Gbps ports that enable high performance and improved utilization while providing easy installation and management
- ▶ Hot swappable small form factor pluggables (SFPs)
- ▶ Inter-switch link (ISL) trunking for up to eight ports provides a total bandwidth of 128 Gbps
- ▶ Provides Fibre Channel interfaces E_port, F_port, FL_port, and M_port
- ▶ Advanced zoning (hardware-enforced) helps to protect against non-secure, unauthorized, and unauthenticated network and management access and worldwide name (WWN) spoofing
- ▶ Hot firmware activation enables fast firmware upgrades that eliminate disruption to the existing fabric
- ▶ Compatibility with an earlier version with IBM b-type and m-type
- ▶ Optional as-needed licensed features:
 - Adaptive Networking
 - Advance Performance Monitor
 - Extended Fabric
 - Fabric Watch
 - Trunking activation
 - Server Application Optimization (SAO)

Cisco MDS 9124 Express for IBM System Storage

This system provides high-performance, scalable, and simple-to-use fabric switching with 8, 16, or 24 ports that are operating at 1, 2 and 4 Gbps. This system is for servers that run Microsoft Windows, UNIX, Linux, Novell NetWare, and IBM OS/400® operating systems, server clustering, infrastructure simplification, and business continuity solutions. The switch includes a replaceable power supply, virtual SAN, Cisco Fabric Manager, and redundant power supply feature that is designed to simplify setup and ongoing maintenance for Cisco MDS 9000 users. The front view of the multilayer switch is shown in Figure 12-2.



Figure 12-2 Front view of the Cisco MDS 9124 Multilayer fabric switch

A single Cisco MDS 9124 switch can serve as an initial building block for a SAN for businesses that want to obtain the benefits of storage consolidation and are just beginning to implement Fibre Channel storage systems. An entry level configuration, for example, might consist of one or two Fibre Channel links to a disk storage array, or to an LTO tape drive. An entry level, 8-port storage consolidation solution can support up to seven servers with a single path to either disk or tape. The On-Demand Port Activation feature is designed to enable a base switch to grow from 8 to 24 ports, in eight port increments, to support more servers and more storage devices without taking the switch offline.

Higher availability solutions can be created by using multiple Cisco MDS 9124 switches. Such implementations would be well-suited to server clustering environments. Such a configuration

can support 6 - 22 servers, each with dual Fibre Channel adapters cross-connected to redundant 9124 switches, which are cross-connected to a dual controller storage system. The following list provides the main features and available options:

- ▶ Efficient 1U design with 8, 16, or 24 ports configuration on demand
- ▶ Auto-sensing 4, 2, or 1 Gbps ports, available with shortwave, 4 km longwave, or 10 km longwave SFPs, all hot swappable
- ▶ Provides Fibre Channel interfaces E_port, F_port, and FL_port
- ▶ Simple SAN configuration that uses Cisco Fabric Manager
- ▶ Replaceable power supply and redundant cooling
- ▶ Optional features include:
 - Hot-swap redundant power supply
 - Cisco MDS 9000 Enterprise package activation
 - Cisco MDS 9000 Fabric Manager Server package activation

Discontinued entry SAN switches

The following products and devices are withdrawn from marketing by IBM. However, these products are still commonly seen in many data centers of small or medium business solutions that rely on SAN networking:

- ▶ IBM System Storage SAN10Q-2
- ▶ IBM TotalStorage SAN16B-2
- ▶ IBM TotalStorage SAN16M-2

For more information about currently available entry IBM SAN switches, see this website:

<http://www.ibm.com/systems/storage/san/entry/index.html>

12.2.2 Midrange SAN switches

The IBM Midrange SAN switches provide scalable and affordable small and medium business (SMB) and enterprise solutions for storage networking:

- ▶ Cost conscious SMB clients with limited technical skills
- ▶ Integrated, scalable, high-availability IBM virtualization family solutions
- ▶ Heterogeneous Windows, Linux, IBM iSeries®, UNIX, and Mainframe servers
- ▶ IBM xSeries®, iSeries, IBM pSeries®, and zSeries Server sales channels
- ▶ Support the IBM System Storage Virtualization family of products, System Storage, and former TotalStorage devices and disk subsystems, IBM Tivoli Storage Manager, SAN Manager, SRM, and Multiple Device Manager
- ▶ Integrated solutions at affordable prices with worldwide IBM support and IBM TotalStorage Solution Center (TSSC) services

The category of midrange SAN switches includes the following products:

- ▶ IBM System Storage SAN40B-4
- ▶ IBM System Storage SAN80B-4
- ▶ IBM System Storage SAN48B-5
- ▶ Cisco MDS 9148 for IBM System Storage
- ▶ IBM System Storage SAN32B-E4

IBM System Storage SAN40B-4

A compact, high-performance, easy-to-install Fibre Channel (FC) SAN switch which enables multiple servers to connect to external disk and tape systems. The SAN40B-4 supports most common operating systems and connects to most common servers and external storage devices. It supports server virtualization with virtual data paths across a single FC link. This system can connect to other SAN switches, routers, and directors to form a multi-switch fabric for increased connectivity and scalability of clients' SAN infrastructure. The front view of the SAN40B-4 switch (model 2498-B40) is shown in Figure 12-3.



Figure 12-3 IBM System Storage SAN40B-4 switch

The IBM System Storage SAN40B-4 SAN fabric switch provides 24, 32, or 40 active ports and is designed for high performance with 8 Gbps link speeds and compatibility with an earlier version to support links that run at 4, 2, or 1 Gbps link speeds. High availability features make it suitable for use as a core switch in midrange environments or as an edge-switch in enterprise environments where a wide range of SAN infrastructure simplification and business continuity configurations are possible. IBM Power Systems, IBM System x, and IBM System z and many non-IBM disk and tape devices are supported in many common operating system environments. Optional features provide specialized distance extension, dynamic routing between separate or heterogeneous fabrics, link trunking, IBM Fibre Channel connection (FICON), Server Application Optimization (SAO), performance monitoring, and advanced security capabilities.

The IBM System Storage SAN40B-4 switch provides the following features and characteristics:

- ▶ Efficient 1U design with 24, 32, or 40 ports configuration on demand
- ▶ Auto-sensing 8, 4, 2, and 1 Gbps ports that enable high performance and improved utilization while providing easy installation and management
- ▶ Hot swappable small form factor pluggables (SFPs) and redundant power supply and cooling
- ▶ Inter-switch link (ISL) trunking for up to eight ports provides a total bandwidth of 128 Gbps
- ▶ Provides Fibre Channel interfaces E_port, F_port, FL_port, M_port, and optional EX_port
- ▶ N_port ID virtualization enables host images behind identical host bus adapters (HBAs) to connect to an F_port that uses a unique N_port ID
- ▶ Advanced zoning (hardware-enforced) helps protect against non-secure, unauthorized, and unauthenticated network and management access and worldwide name (WWN) spoofing
- ▶ Hot firmware activation enables fast firmware upgrades that eliminate disruption to the existing fabric
- ▶ Compatibility with an earlier version with IBM b-type and m-type

IBM System Storage SAN80B-4

This system is a compact, high-performance, easy-to-install Fibre Channel (FC) SAN switch that enables multiple servers to connect to external disk and tape systems. It extends the basic model SAN40B-4 by an additional 40 8 Gbps ports. The SAN80B-4 (model 2498-B80) supports most common operating systems and connects to most common servers and external storage devices. Supports server virtualization with virtual data paths across a single FC link. Can connect to other SAN switches, routers, and directors to form a multi-switch fabric for increased connectivity. See Figure 12-4 for the front view of SAN80B-4.



Figure 12-4 Front view of the IBM System Storage SAN80B-4

The IBM System Storage SAN80B-4 SAN fabric switch provides 48, 64, or 80 active ports. This switch is designed for high performance with 8 Gbps link speeds and compatibility with an earlier version to support links that run at 4, 2, and 1 Gbps link speeds. High availability features make it suitable for use as a core switch in midrange environments or as an edge switch in enterprise environments. IBM Power Systems, System x, System z, and many non-IBM disk and tape devices are supported in many common operating system (OS) environments. Optional features provide specialized distance extension, dynamic routing between separate or heterogeneous fabrics, link trunking, IBM FICON, Server Application Optimization (SAO), performance monitoring and advanced security capabilities.

IBM System Storage SAN80B-4 demonstrates the same functions and features as the basic model of midrange switches SAN40B-4, and the same operational capabilities, just with the following dissimilarities:

- ▶ Efficient 2U design with 48, 64, and 80 ports configuration on demand
- ▶ Significantly reduced power consumption by single chassis and fewer field-replaceable units (FRUs)
- ▶ 16-port activation licensed feature

IBM System Storage SAN48B-5

The SAN48B-5 Switch is designed to meet the demands of hyper-scale, private cloud storage environments by delivering 16 Gbps Fibre Channel technology and capabilities that support highly virtualized environments.

The SAN48B-5 (2498-F48) delivers SAN technology within a flexible, simple, and easy-to-use solution. In addition to providing scalability, SAN48B-5 can address demanding Reliability, Availability, and Serviceability (RAS) requirements to help minimize downtime to support mission-critical cloud environments.

The front view of the IBM System Storage SAN48B-5 16 Gbps Fibre Channel switch is shown in Figure 12-5.

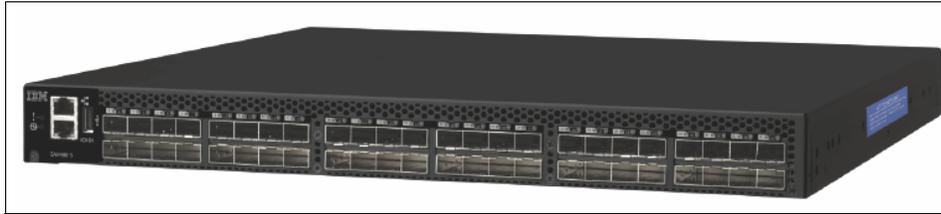


Figure 12-5 IBM System Storage SAN48B-5 fabric switch

The following list includes the product highlights of the System Storage SAN48B-5 fabric switch:

- ▶ Performance of 16 Gbps with up to 48 ports in an energy-efficient, 1U enclosure
- ▶ Speed of 2, 4, 8, 10, or 16 Gbps on all ports which produce an aggregate 768 Gbps full-duplex throughput
- ▶ High-performance of 128 Gbps and resilient frame-based trunking
- ▶ 10 Gbps Fibre Channel integration on the same port for dense wavelength division multiplexing (DWDM) metro connectivity on the same switch
- ▶ In-flight Data Compression and Encryption for efficient link utilization and security
- ▶ Redundant, hot-swap components and non-disruptive software upgrades
- ▶ Diagnostic Port (D-Port) feature for physical media diagnostic, troubleshooting, and verification services
- ▶ Multi-tenancy in cloud environments through Virtual Fabrics, Integrated Routing, quality of service (QoS), and fabric-based zoning features

Cisco MDS 9148 for IBM System Storage

The Cisco MDS 9148 for IBM System Storage Multilayer Fabric Switch (2417-C48) is designed to provide an affordable, highly capable, and scalable storage networking solution for small, midrange, and large enterprise clients. The system can be used as part of SAN solutions from simple single switch configurations to larger multi-switch configurations in support of fabric connectivity and advanced business continuity capabilities.

As seen in Figure 12-6, the switch is designed to offer outstanding value in a compact one-rack-unit (1U) form factor. With the ability to expand from 16 to 48 ports in eight-port increments, the Cisco MDS 9148 can be used as the foundation for small, stand-alone SANs, as a top-of-rack switch, or as an edge switch in larger core-edge SAN infrastructures.



Figure 12-6 Cisco MDS 9148 for IBM System Storage

The Cisco MDS 9148 Multilayer Fabric Switch is designed to support quick configuration with zero-touch plug-and-play features and task wizards that allow it to be deployed quickly and easily in networks of any size. Powered by Cisco MDS 9000 NX-OS Software, it includes advanced storage networking features and functions. This switch is compatible with Cisco

MDS 9000 Series Multilayer Directors and Switches, providing transparent, end-to-end service delivery in core-edge deployments.

Fabric connectivity capabilities can be the basis for infrastructure simplification solutions for IBM System i, System p, and System x servers and storage consolidation and high-availability server clustering with IBM System Storage disk arrays. Business continuity capabilities can help businesses protect valuable data with IBM System Storage tape subsystems and IBM Tivoli Storage Manager data protection software.

The Cisco MDS 9148 for IBM System Storage switch incorporates the following features:

- ▶ Flexibility and scalability of up to 48 auto-sensing 1/2/4/8 Gbps ports in 1U standard rack-mountable or stand-alone unit
- ▶ Intelligent storage networking services that are powered by Cisco MDS 9000 NX-OS management software, which enables Virtual SANs, Inter-VSAN routing (IVR), Port Channels, QoS, and security
- ▶ Support for virtual environments with full N_port ID virtualization
- ▶ Simplified storage management and diagnostic procedures that are enabled by Cisco Fabric Manager, FC ping, FC traceroute, Switched Port Analyzer (SPAN), and Cisco Fabric Analyzer
- ▶ Support of 4/8 Gbps shortwave, 4 Gbps 4 km longwave, or 4/8 Gbps 10 km longwave SFPs
- ▶ Redundant power supplies and cooling

IBM System Storage SAN32B-E4 Encryption switch

Data is one of the most highly valued resources in a competitive business environment. Protecting that data, controlling access to it, and verifying its authenticity while maintaining its availability are priorities in our security conscious world. Increasing regulatory requirements are also helping to drive the need for the adequate security of data. Encryption is a powerful and widely used technology that helps protect data from loss and inadvertent or deliberate compromise.

The IBM System Storage SAN32B-E4 Encryption Switch (2488-E32) is a high performance stand-alone device that is designed for protecting data-at-rest in mission critical environments. In addition to helping IT organizations achieve compliance with regulatory mandates and meeting industry standards for data confidentiality, the SAN32B-E4 Encryption Switch also protects them against potential litigation and liability that follows a reported breach. The front view of the Encryption switch is shown in Figure 12-7.



Figure 12-7 IBM System Storage SAN32B-E4 Encryption switch

For data center fabric security, IBM provides advanced encryption services for SANs with the IBM System Storage SAN32B-E4 Encryption Switch. The switch is a high speed, highly reliable hardware device that delivers fabric-based encryption services to protect data assets

either selectively or on a comprehensive basis. The 8 Gbps SAN32B-E4 Fibre Channel Encryption Switch scales nondisruptively, providing from 48 up to 96 Gbps of encryption processing power to meet the needs of the most demanding environments with flexible, on-demand performance. It also provides compression services at speeds up to 48 Gbps for tape storage systems. Moreover, it is tightly integrated with one of the industry-leading, enterprise class key management systems, the IBM Tivoli Key Lifecycle Manager. This solution can scale to support key lifecycle services across distributed environments.

At a glance, the IBM System Storage SAN32B-E4 Encryption Switch provides the following features:

- ▶ 32 autosensing 8/4/2/1 Gbps active ports in base configuration
- ▶ Choice of 8 Gbps shortwave, longwave, or extended distance SFPs
- ▶ Fabric-based data-at-rest encryption with 48 Gbps disk and tape encryption processing. Transparent, online encryption of *cleartext* logical unit numbers (LUNs) and rekeying of encrypted LUNs with no disruption
- ▶ Provides Fibre Channel interfaces E_port, F_port, FL_port, M_port, and optional EX_port
- ▶ Redundant power supplies and cooling
- ▶ Tight integration with Tivoli Key Lifecycle Manager with support of multi-vendor disk storage subsystems

Discontinued midrange SAN switches

IBM withdrew the following products from marketing, nevertheless these switches might still be available by using the wide network of IBM Business Partners. Also, you can likely find them in many small or medium business SAN solutions or client data centers.

- ▶ IBM TotalStorage SAN32B-2
- ▶ IBM TotalStorage SAN32B-2 Express
- ▶ IBM TotalStorage SAN32M-2
- ▶ IBM TotalStorage SAN32M-2 Express
- ▶ IBM System Storage SAN64B-2
- ▶ Cisco MDS 9120 and 9140 Multilayer fabric switches
- ▶ Cisco MDS 9216i and 9216A Multilayer fabric switches
- ▶ Cisco MDS 9020 Fabric switch

For more information about IBM midrange SAN switches, see this website:

<http://www.ibm.com/systems/storage/san/midrange/index.html>

12.2.3 Enterprise SAN directors

The IBM Enterprise SAN directors provide the data center networking infrastructure with the following benefits:

- ▶ The highest availability and scalability, and intelligent software to simplify the management of complex, integrated enterprise SANs
- ▶ Heterogeneous Windows, Linux, IBM System i, UNIX, and Mainframe servers
IBM DSxxxx, FASTT, ESS, LTO, and ETS storage
- ▶ Supports the IBM System Storage Virtualization family of products and storage systems, IBM Tivoli Storage Manager, SAN Manager Storage Resource Manager, and Multiple Device Manager
- ▶ Offers customized solutions with competitive prices, worldwide IBM support, and IBM Global Services and IBM financial services

IBM offers the following enterprise SAN directors through its marketing channels:

- ▶ IBM System Storage SAN384B-2
- ▶ IBM System Storage SAN768B-2
- ▶ Cisco MDS 9506 for IBM System Storage
- ▶ Cisco MDS 9509 for IBM System Storage
- ▶ Cisco MDS 9513 for IBM System Storage

IBM System Storage SAN384B-2 and SAN768B-2

The IBM System Storage SAN768B-2 and IBM System Storage SAN384B-2 fabric backbones are highly robust network switching platforms that are designed for evolving enterprise data centers. Each system combines breakthrough performance, scalability, and energy efficiency with long-term investment protection.

Supporting open systems and IBM System z environments, these platforms address data growth and server virtualization challenges to:

- ▶ Enable server, SAN, and data center consolidation
- ▶ Minimize disruption and risk
- ▶ Reduce infrastructure and administrative costs

Built for large enterprise networks, the SAN768B-2 has eight vertical blade slots to provide up to 384 16 Gbps or 512 8 Gbps FC ports. The SAN384B-2 is ideal for midsize core or edge deployments, providing four horizontal blade slots and up to 192 16 Gbps or 256 8 Gbps FC ports. The flexible blade architecture also supports FCoE, fabric-based encryption, SAN extension advanced functionality for high performance servers, I/O consolidation, data protection, and disaster recovery solutions.

The SAN768B-2 and SAN384B-2 are extremely efficient at reducing power consumption, cooling, and the carbon footprint in data centers. Although these switches provide exceptional performance and scale, these networking backbones use less than 1 watt per Gbps. As members of the IBM System Storage family of b-type SAN products, the SAN768B-2 and the SAN384B-2 are designed to participate in fabrics that contain other b-type and m-type devices that are manufactured by Brocade. This versatile hardware can serve as the backbone in a complex fabric and provide connections to other b-type and m-type directors, switches, and routers.

The SAN768B-2 and SAN384B-2 backbones use Brocade Fabric OS (FOS), which provides several characteristic features. Some of these features include Bottleneck Detection, Top Talkers (part of Advanced Performance Monitoring), and Adaptive Networking, which is a suite of tools that include Ingress Rate Limiting, Traffic Isolation, and QoS. Managed through the IBM System Storage Data Center Fabric Manager (DCFM) or the command-line interface (CLI), these advanced capabilities help optimize fabric behavior and application performance.

Both products are shown in Figure 12-8.

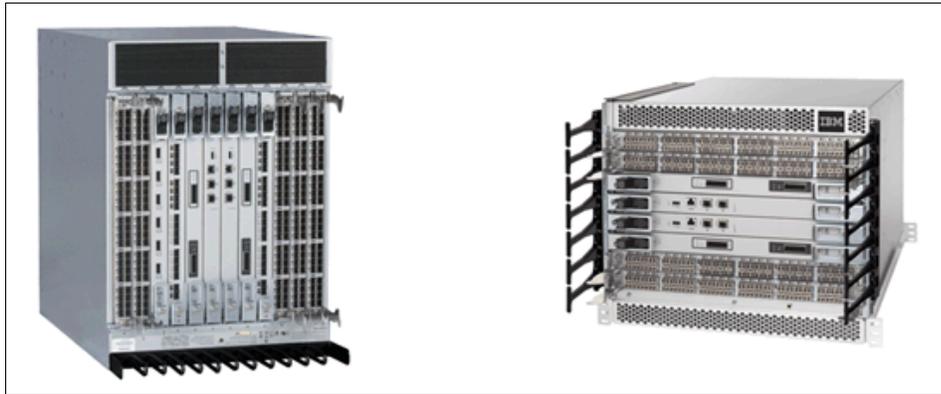


Figure 12-8 IBM System Storage SAN768B-2 (left) and SAN384B-2 (right)

The following list provides the blades that are available for both models:

- ▶ 16 Gbps 32-port or 48-port FC blades
- ▶ 8 Gbps 64-port FC blades
- ▶ 8 Gbps 16-port FC encryption blade
- ▶ 10 Gbps six-port FC blade
- ▶ FCoE 10 GbE 24-port blade that supports twenty-four 10 Gbps CEE/FCoE ports
- ▶ 8 Gbps extension blade that supports twelve 8 Gbps FC ports and ten 1 Gbps GbE ports, or two optional 10 Gbps GbE ports

The following list provides highlights of the features and functions of both products:

- ▶ Redundant control processor modules, power supplies, and cooling
- ▶ Auto-sensing 16/8/4/2 Gbps E_Port, F_Port, FL_Port, EX_port, and 10 Gbps E_Port Fibre Channel interfaces and 10 Gbps converged ethernet ports
- ▶ 16/8/4/2 FICON interfaces and 1 GbE or 10 GbE FCIP ethernet interfaces
- ▶ Management protocols that include HTTP, SNMP v1/v3 (FE MIB, FC Management MIB), Telnet
- ▶ Management and operational software includes Auditing, Syslog, Web Tools, Fabric Watch, IBM System Storage Data Center Fabric Manager (DCFM), and CLI
- ▶ Full compatibility with an earlier version with IBM System Storage and IBM TotalStorage b-type and m-type SAN directors, switches, and routers; other directors, switches, and routers that are manufactured by Brocade
- ▶ Fabric-based data-at-rest encryption for disk array LUNs, heterogeneous tape drives, and virtual tape libraries to enforce data confidentiality and privacy requirements (with selected blades); support for Tivoli Key Lifecycle Manager
- ▶ Advanced features such as ISL trunking, Server Application Optimization (SAO), Virtual Fabrics, Adaptive Networking Services, Advanced Performance Monitoring, support for extended fabrics up to 10 km

Cisco MDS 9506 and MDS 9509 for IBM System Storage

The Cisco MDS 9506 and MDS 9509 Multilayer Directors for IBM System Storage, support 1, 2, 4, 8 and 10 Gbps Fibre Channel switch connectivity and intelligent network services. These features help improve the security, performance, and manageability that is required to consolidate geographically dispersed storage devices into a large enterprise SAN. Administrators can use these directors to help address the needs for high performance and

reliability in SAN environments that range from small workgroups to large, integrated global enterprise SANs.

Both, the Cisco MDS 9506 and MDS 9509 (see Figure 12-9) for IBM System Storage use two Supervisor-2 Modules that are designed for high availability and performance. The Supervisor-2 Module combines an intelligent control module and a high performance crossbar switch fabric in a single unit. It uses Fabric Shortest Path First (FSPF) multipath routing, which provides intelligence to load balance across a maximum of 16 equal-cost paths and to dynamically reroute traffic if a switch fails.

Each Supervisor-2 Module provides the necessary crossbar bandwidth to deliver full system performance in the MDS 9506 director with up to four Fibre Channel switching modules (8 for MDS 9509). It is designed to provide that loss or removal of a single crossbar module has no affect on system performance. Fibre Channel switching modules are designed to optimize performance, flexibility, and density. The Cisco MDS 9506 Multilayer Director requires a minimum of one and allows a maximum of four switching modules, while MDS 9509 uses 1 - 16 modules. These modules are available in either 12-, 24-, and 48-port 4 Gbps configurations, allowing the Cisco MDS 9506 to support 12 - 192 Fibre Channel ports per chassis (336 in MDS 9509). Optionally, a four-port 10 Gbps Fibre Channel module is available for high performance ISL connections over metro optical networks.



Figure 12-9 Cisco 9506 (left) and MDS 9509 (right) Multilayer Directors

Advanced traffic management capabilities are integrated into the switching modules to help simplify deployment and to optimize performance across a large fabric. The PortChannel capability allows users to aggregate up to 16 physical 2 Gbps Inter-Switch Links into a single logical bundle, providing optimized bandwidth utilization across all links. The bundle might span any port from any 16-port switching module within the chassis, providing up to 32 Gbps throughput.

The following list provides the available switching modules:

- ▶ Twelve, 24, or 48 ports 4 Gbps Fibre Channel module
- ▶ Twenty-four or 48 ports 8 Gbps Fibre Channel module
- ▶ Four ports 8 Gbps with 44 ports 4 Gbps Fibre Channel module
- ▶ Four ports 10 Gbps Fibre Channel module

The following list provides the highlights of the two storage devices:

- ▶ Provides Fibre Channel throughput of up to 8 Gbps per port and up to 64 Gbps with each PortChannel ISL connection

- ▶ Offers scalability for 12 - 192 (336) Fibre Channel ports
- ▶ Offers 10 Gbps ISL ports for Inter-Data Center links over metro optical networks
- ▶ Offers Gigabit Ethernet IP, GbE ports for iSCSI, or FCIP connectivity over global networks
- ▶ Includes virtual SAN (VSAN) capability for SAN consolidation into virtual SAN islands on a single physical fabric
- ▶ Includes high-availability design with non-disruptive firmware upgrades
- ▶ Enterprise, SAN Extension over IP, Mainframe, Storage Services Enabler, and Fabric Manager Server Packages provide added intelligence and value

Cisco MDS 9513 for IBM System Storage

The Cisco MDS 9513 for IBM System Storage provides 12 - 528 Fibre Channel ports, with 4 and 8 Gbps support and a high-availability design. It offers 4 - 44 10 Gbps ports for ISL connectivity across metro optical networks. It includes VSAN capability for SAN consolidation into virtual SAN “islands” on a single physical fabric. The Cisco MDS 9513 provides network security features for large enterprise SAN deployment. The director also offers intelligent networking services to help simplify mainframe FICON and Fibre Channel SAN management and reduce total cost of ownership (TCO). See the front view of the chassis in Figure 12-10.



Figure 12-10 Front view of the Cisco MDS 9513

The Cisco MDS 9513 Multilayer Director uses two Supervisor-2 Modules, which are designed to support high availability. The Supervisor-2 Module is designed to provide industry-leading scalability, intelligent SAN services, non-disruptive software upgrades, graceful process restart and failover, and redundant operation. Dual crossbar-switching fabric modules provide a total internal switching bandwidth of 2.4 Tbps for inter-connection of up to 11 Fibre Channel switching modules.

Fibre Channel switching modules improve performance, flexibility, and density. The Cisco MDS 9513 for IBM System Storage requires a minimum of one Fibre Channel switching module and allows a maximum of 11. These modules are available in 12-, 24-, or 48-port

4 and 8 Gbps configurations, which enable the Cisco MDS 9513 to support 12 - 528 Fibre Channel ports per chassis. Optionally, a four-port 10 Gbps Fibre Channel module is available for high-performance ISL connections over metro optical networks.

The following list includes the highlights of the Cisco MDS 9513:

- ▶ Supports Fibre Channel throughput of up to 8 Gbps per port and up to 64 Gbps with each PortChannel ISL connection
- ▶ Offers Gigabit Ethernet (GbE) IP ports for iSCSI or FCIP connectivity over global networks
- ▶ Offers scalability for 12 - 528 1, 2, 4, and 8 Gbps Fibre Channel ports
- ▶ High-availability design with support for non-disruptive firmware upgrades Includes VSAN capability for SAN consolidation into virtual SAN islands on a single physical fabric
- ▶ Enterprise, SAN Extension over IP, Mainframe, Storage Services Enabler, and Fabric Manager Server Packages provide added intelligence and value

Discontinued enterprise SAN directors

IBM withdrew the following products from marketing. However, you can still find them in medium business or large strategic data centers and client environments:

- ▶ IBM TotalStorage SAN Director M14
- ▶ IBM TotalStorage SAN140M
- ▶ IBM TotalStorage SANC40M
- ▶ IBM TotalStorage SAN256M
- ▶ IBM TotalStorage SAN256B
- ▶ IBM System Storage SAN384B
- ▶ IBM System Storage SAN768B

For more information about IBM SAN directors, see this website:

<http://www.ibm.com/systems/storage/san/enterprise/>

12.2.4 Multiprotocol routers

IBM currently has two multiprotocol routers in its portfolio:

- ▶ IBM System Storage SAN06M-R
- ▶ Cisco MDS 9222i for IBM System Storage

IBM System Storage SAN06M-R multiprotocol router

This is an entry-level multiprotocol extension router that is designed to connect two SANs over a wide distance by using the Internet as the interconnection fabric (upgrades to enterprise level functions are available). Intended to support business continuity solutions between supported servers at one site and support servers or IBM System Storage disk or tape devices at a distant location.

SAN06M-R delivers high performance with 8 Gbps FC ports and hardware assisted traffic processing for line-rate performance. It uses existing Internet, IP-based infrastructures for metro and global SAN extension for business continuity solutions. Up to eight virtual FCIP tunnels are available to help maximize scalability and utilization of metropolitan area network (MAN) or wide area network (WAN) resources.

The router provides hardware-based compression, large window sizes, and selective acknowledgement of IP packets that are designed to optimize performance of SAN extension over IP networks. The front view of the device is shown in Figure 12-11.



Figure 12-11 Front view of the SAN06M-R

The following list provides a summary of the router and its highlights:

- ▶ 1U 19" packaging that is designed for rack-mount or table-top
- ▶ Designed for high-performance with up to 8 Gbps FC autosensing ports and up to 1 Gbps Ethernet (GbE) ports
- ▶ Supports either 8, 4, and 2; or 4, 2, and 1 Gbps FC link speeds
- ▶ Shortwave and longwave SFPs can be intermixed in the same router
- ▶ Hardware-based compression with extensive buffering
- ▶ Optional FICON with CUP and FICON Accelerator enable support for enterprise class environments
- ▶ SAN isolation from Internet, WAN or MAN failures.

Cisco MDS 9222i for IBM System Storage

The Cisco MDS 9222i for IBM System Storage is designed to address the needs of medium sized businesses and large enterprises with a wide range of SAN capabilities. It can be used as a cost effective high performance SAN extension over IP router switch for midrange SMB clients in support of IT simplification and business continuity solutions. It can also be used as a remote site router switch for device aggregation and SAN extension over IP to data center directors for large enterprise clients.

Business continuity solutions include data protection with IBM System Storage tape libraries and devices and IBM Tivoli Storage Manager data protection software; and disaster protection with IBM System Storage disk metro and global mirroring disaster recovery solutions. The front view of the MDS 9222i is shown in Figure 12-12.



Figure 12-12 Front view of the Cisco MDS 9222i

The following list provides product highlights of the Cisco MDS 9222i:

- ▶ Base switch includes eighteen 4 Gbps Fibre Channel ports with two shortwave SFPs; and four GbE IP ports with SAN Extension over IP
- ▶ Optional 8 Gbps FC switch ports support optical SFP+ transceivers and 10 Gbps ISL ports support optical X2 transceivers for edge switch attachment to Cisco MDS 9500 directors
- ▶ SAN extension with high performance FCIP acceleration and hardware-based compression capabilities and security with hardware-based encryption
- ▶ Multiservice design for high performance business continuity solutions with Windows, UNIX, Linux, NetWare, IBM OS/400, and IBM z/OS servers
- ▶ Includes VSAN capability for SAN consolidation into virtual SAN islands on a single physical fabric

Discontinued IBM multiprotocol routers

The following products and devices are no longer available through the IBM marketing channels:

- ▶ IBM TotalStorage04M-R
- ▶ IBM TotalStorage SAN16B-R
- ▶ IBM TotalStorage SAN16M-R
- ▶ IBM System Storage SAN18B-R

For more information about each of these products, see this website:

<http://www.ibm.com/systems/storage/san/routers/index.html>

12.3 IBM System Storage Disk Systems

The IBM System Storage Disk Systems products and offerings provide compelling storage solutions with superior value for all levels of business. In the following section, we briefly describe SAN-attached storage disk device subsystems as a key component of every data center to keep data available in an effective and cost-efficient way.

We do not describe network-attached storage (NAS) because that is not the primary objective of this publication. SAN disk systems for storage virtualization and cloud computing are described in 12.5, “Storage virtualization and cloud computing” on page 278.

We do not provide detailed information about expansion units (EXP) of each of the disk subsystems. For more information about enclosures, we refer you to IBM resources.

12.3.1 Entry level disk systems

Designed to deliver advanced functionality at a breakthrough price, these systems provide an exceptional solution for workgroup storage applications such as email, file, print, and web servers. Other features include collaborative databases and remote boot for diskless servers.

We describe the following product: IBM System Storage DS3500 Express.

Enclosures and expansion units (not described in this section) include the following systems:

- ▶ IBM System Storage EXP2500 Express
- ▶ IBM System Storage EXP3000 Express
- ▶ IBM System Storage EXP3512 Expansion Enclosure
- ▶ IBM System Storage EXP3524 Expansion Enclosure

IBM System Storage DS3500 Express

The IBM System Storage DS3500 Express Storage™ Systems are the newest addition to the IBM System Storage DS3000 series family of entry disk storage systems. The DS3500 delivers affordable, entry-level configurations for small and medium businesses in compact 2U, 19-inch rack mount enclosures. This solution provides the flexibility to scale in capacity, performance, host interfaces, and advanced functions as your business grows or requirements change.

The DS3500 combines next-generation controller technology with the latest, high-performance host interface technologies to deliver new levels of performance to the DS3000 series. With the DS3500, you choose the initial system configuration that matches your performance requirements and budget:

- ▶ Single controller for an entry solution with low initial investment
- ▶ Dual controllers for higher performance
- ▶ Dual controller with the Turbo Performance option for the best results

The DS3500 is available in two models (see Figure 12-13):

- ▶ DS3512 Express with 12 3.5-inch 6 Gbps SAS attached disk
- ▶ DS3534 Express with 24 2.5-inch 6 Gbps SAS attached disk



Figure 12-13 IBM System Storage DS3512 (top) and DS3524 (bottom)

The following list includes product highlights of the DS3512 and DS3524 devices:

- ▶ Compact 2U 19-inch rack mountable chassis
- ▶ Single or dual controllers, 6 Gbps SAS attached 3.5-inch (DS3512) or 2.5-inch (DS3524) disks
- ▶ Support for 192 disk drives through the attachment of EXP3500 expansion units (EXP3512 with 12 3.5-inch or EXP3524 with 24 2.5-inch disk bays)
- ▶ Eight Gbps FC, 6 Gbps SAS, 1 Gbps, and 10 Gbps iSCSI host connectivity available
- ▶ Support for up to 128 storage partitions with RAID levels 0, 1, 3, 5, 6, 10
- ▶ Dual redundant hot-swappable power supplies and cooling fans

Discontinued entry disk systems

IBM withdrew the following product portfolio from marketing:

- ▶ IBM TotalStorage DS3100
- ▶ IBM TotalStorage DS3200
- ▶ IBM TotalStorage DS3300
- ▶ IBM System Storage DS3400
- ▶ IBM TotalStorage EXP24 Expansion Unit
- ▶ IBM TotalStorage EXP420 Expansion Unit

For more information about each product, particularly for storage expansion units that are not covered in the previous sections, see this website:

<http://www.ibm.com/systems/storage/disk/entry/index.html>

12.3.2 Midrange disk systems

Throughout this section, we provide a brief overview of disk storage device systems that IBM offers for small and medium business solutions. Again, we do not describe storage enclosures and expansion units in detail.

The following disk systems are described in this section:

- ▶ IBM System Storage DS5020 Express
- ▶ IBM System Storage DS5000
- ▶ IBM System Storage DS3950 Express

The following list provides the enclosures that are available:

- ▶ IBM System Storage EXP5060 Expansion unit for DS5000 family
- ▶ IBM System Storage EXP520 Expansion unit for DS5020 Express
- ▶ IBM System Storage EXP395 Expansion unit for DS3950

IBM System Storage DS5020 Express

Optimized data management requires storage solutions with high data availability, strong storage management capabilities, and powerful performance features. The IBM System Storage DS5020 Express is designed to provide lower total cost of ownership (TCO), high performance, robust functionality, and unparalleled ease of use.

Additionally, auto-negotiating 8 Gbps Fibre Channel interfaces allow the DS5020 Express to integrate seamlessly into an existing 2 Gbps or 4 Gbps infrastructure. This integration offers investment protection going forward to when the SAN inevitably becomes 8 Gbps.

Aside from 8 Gbps FC connections, the DS5020 Express offers an optional: 1 Gbps iSCSI interface for less demanding applications and lower-cost implementation, up to 67.2 TB of Fibre Channel or FC-SAS physical storage capacity, 224 TB of SATA physical storage capacity, and 33.6 TB of solid-state drive (SSD) physical storage capacity. Other features of this device include power system management, data management, and data protection features.

See Figure 12-14 on page 265 for the front view of the DS5020 Express.

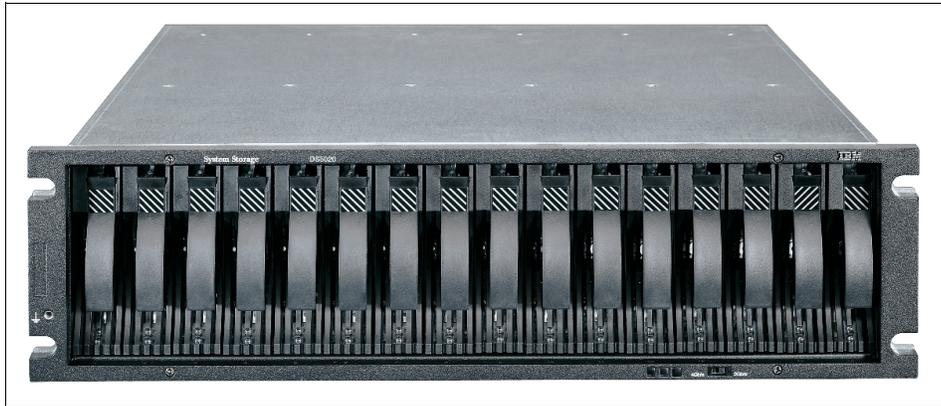


Figure 12-14 Top front view of the DS5020 Express model

The following list provides the product highlights of the DS5020 Express:

- ▶ High-performance 8 Gbps FC connections and 1 Gbps iSCSI
- ▶ Up to 224 TB of physical storage capacity with 112 (using six EXP520 expansion units) 2-TB SATA disk drives in up to 128 storage partitions
- ▶ Transparent management of DS3000 and DS5000 products that use DS Storage Manager
- ▶ Support for intermixing Fibre Channel/FC-SAS/SED/SATA/SSD drives enables cost-effective tiered storage
- ▶ Redundant, hot-swappable power supplies and cooling

IBM System Storage DS5000

The DS5000 series storage systems (DS5100 and DS5300) are designed to meet the demanding open-systems requirements of today and tomorrow, while establishing a new standard for lifecycle longevity with field-replaceable host interface cards. Seventh generation architecture delivers relentless performance, real reliability, multidimensional scalability, and unprecedented investment protection.

The DS5000 storage systems are equally adept at supporting transactional applications, such as databases and online transaction processing (OLTP); throughput-intensive applications, such as high-performance computing (HPC) and rich media; and concurrent workloads for consolidation and virtualization. With relentless performance and superior reliability and availability, DS5000 series storage systems can support the most demanding service level agreements (SLAs) for the most common operating systems, including Microsoft Windows, UNIX, and Linux. And when requirements change, you can add or replace host interfaces, grow capacity, add cache, and reconfigure the system as needed.

Figure 12-15 shows the top front view of IBM System Storage DS5100 as a member of the DS5000 family of IBM storage products.



Figure 12-15 Top front view of the DS5100

The following list provides the product highlights of the IBM System Storage DS5100:

- ▶ Efficient, compact 4U packaging that is designed for a 19-inch rack
- ▶ Easy-to-use, easy-to-configure management interface that is able to manage the DS3000, DS4000, and DS5000 series storage systems
- ▶ Scalable up to 448 drives by using the EXP5000 enclosure and up to 960 TB of high-density storage with the EXP5060 enclosure
- ▶ Support for intermixing drive types (FC, FC-SAS, SED, SATA, and SSD) and host interfaces (8 Gbps FC and 1/10 Gbps iSCSI)
- ▶ Two performance levels (base: DS5100 and high: DS5300) with the ability to field-upgrade performance levels
- ▶ Designed to support high availability with hot-swappable components and non-disruptive firmware upgrades

IBM System Storage DS3950 Express

As part of the DS series, the DS3950 Express offers high performance 8 Gbps capable Fibre Channel connections, an optional 1 Gbps iSCSI interface for less demanding applications, and lower-cost implementation. This improved performance provides up to 67.2 TB of Fibre Channel or FC-SAS physical storage capacity, up to 224 TB of SATA physical storage capacity, and powerful system management, data management, and data protection features. The DS3950 Express is designed to expand from workgroup to enterprise-wide capability with up to six Fibre Channel expansion units with the EXP395 Expansion Unit.

The design of the DS3950 avoids over-configuration to provide an affordable entry-point device. The device offers seamless “pay-as-you-grow” scalability as requirements change, which categorizes the product to the midrange category. Its efficient storage utilization lowers the raw capacity requirement and support that is needed for intermixing high performance and high capacity drives.

The IBM System Storage DS3950 Express model is shown in Figure 12-16.

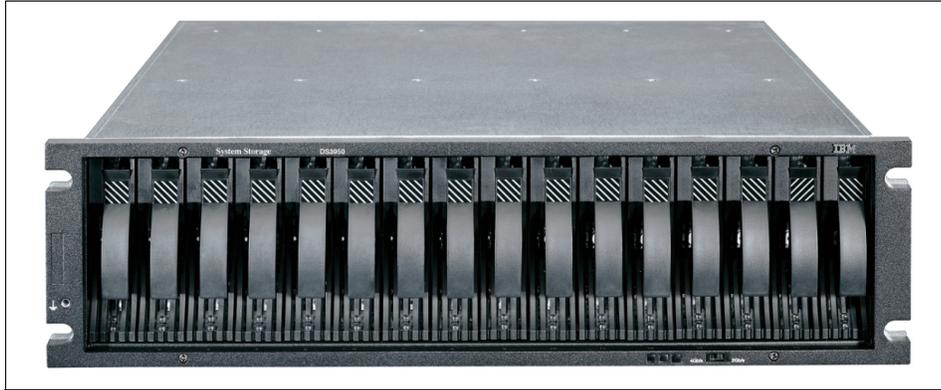


Figure 12-16 Top front view of the DS3950 Express

The following list provides highlights of the DS3950 Express:

- ▶ Support for intermixing drive types (FC, FC-SAS, and SATA) and host interfaces (8 Gbps FC and 1 Gbps iSCSI)
- ▶ Support for up to 112 disk drive modules with the attachment of six EXP395 expansion units that gives up to 224 TB of physical storage capacity, in up to 128 storage partitions
- ▶ Fully integrated replication features such as IBM FlashCopy®, Enhanced Remote Mirror, and VolumeCopy
- ▶ Hot-swappable redundant power supplies and cooling fans in a standard 19-inch rack-mountable unit

Discontinued midrange disk products

IBM withdrew the following products from marketing. However, clients can still find them in medium business or large strategic data centers and client environments:

- ▶ IBM TotalStorage DS4200
- ▶ IBM TotalStorage DS4700
- ▶ IBM TotalStorage DS4800
- ▶ IBM TotalStorage EXP4000 Expansion unit

For more information about each product, and particularly for storage expansion units that are not covered in the previous text, see this website:

<http://www.ibm.com/systems/storage/disk/midrange/index.html>

12.3.3 Enterprise disk systems

With their high capacity, scalability, broad server support, and virtualization features, the enterprise class storage systems are well suited for simplifying the storage environment. These systems consolidate data from multiple storage systems on a single system.

In this section, we describe these high-end, enterprise class IBM disk device subsystems:

- ▶ IBM System Storage DS8000 family
- ▶ IBM System Storage DSC3700

IBM System Storage DS8000

The IBM System Storage DS8000 offers high-performance, high-capacity, secure storage systems that are designed to deliver resiliency and total value for the most demanding, heterogeneous storage environments.

The DS8000 series is built on powerful and market-proven IBM POWER® microprocessors in dual two-way or dual four-way shared symmetric multiprocessor (SMP) complexes. The latest DS8800 model (Figure 12-17) introduces a new generation of hardware and is the fastest disk system in the IBM storage portfolio. The DS8800 provides dual IBM POWER6® controllers, 8 Gbps host and device adapters, and 6 Gbps SAS disk drives. This model delivers a new level of performance in a design that condenses more data in a smaller footprint. Compared to the POWER5 processor in the DS8700 models, the POWER6 processor offers up to over a 40% performance improvement. These upgrades, along with the new 2.5-inch SAS drives, deliver a dramatic boost in data throughput in a more compact footprint. This design is a powerful combination that is aimed at taming the most demanding enterprise applications.



Figure 12-17 Three frames of DS8800 represent the extraordinary scalability

Despite their distinct hardware components, both the DS8800 and DS8700 models now ship with a common microcode built on over 10 years of market proven reliability. This common microcode not only reinforces the IBM commitment to superior reliability, it also enables IBM to deliver new functions on both the DS8800 and DS8700 hardware platforms concurrently. Certainly, DS8700 clients appreciate the investment protection for their deployments. And, all DS8000 clients continue to enjoy the interoperability of most remote mirror and copy functions across older DS8300, DS8100, and IBM TotalStorage Enterprise Storage Server (ESS) models.

The DS8000 series also supports various major server platforms, including IBM z/OS, IBM z/VM®, Linux on IBM System z, IBM i, IBM OS/400, IBM i5/OS™, and IBM AIX operating systems. Other platforms that are supported include: Linux, HP-UX, Sun Solaris, Novell NetWare, VMware, and Microsoft Windows environments, among many others. With

such broad platform support, the DS8000 series can easily accommodate a wide array of applications and their distinct service levels.

IBM System Storage DS8700 and DS8800 disk device subsystems represent simply the world class performance to help optimize responsiveness in an on-demand world, with a minimum of *five-nines* (99.999%) availability.

The following list provides the most important highlights and features of the DS8000 series:

- ▶ Dual symmetric multiprocessing (SMP) processor complexes
- ▶ Four Gbps FC-attached disk drives (DS8700) and 6 Gbps SAS (DS8800)
- ▶ Support for 2 - 32 host adapters and up to 128 FC/FICON 8 Gbps host ports
- ▶ Minimum of eight drives and a maximum of 1056 drives scales up to 2048 TB of physical storage capacity
- ▶ Up to 384 GB cache memory with innovative caching algorithms
- ▶ Storage Pool Striping to automatically avoid disk hotspots
- ▶ I/O Priority Manager aligns quality of service levels to separate application workloads in the system
- ▶ IBM FlashCopy, Metro Mirror, Global Mirror, Metro/Global Mirror, and Global Copy provide flexible replication services
- ▶ The IBM System Storage Easy Tier feature automatically helps optimize solid-state storage (SSD) deployments in multitier systems
- ▶ Full Disk Encryption drive options for advanced protection of data at-rest

IBM System Storage DCS3700

The IBM System Storage DCS3700 is designed to meet the storage needs of highly scalable, data streaming applications in high performance computing environments.

Developed with density in mind, the DCS3700 storage system delivers up to 120 TB of physical capacity in a slim, 4U form factor. With the attachment of up to two DCS3700 Expansion Units, the physical capacity can be scaled up to 360 TB, all within a 12U space.

The DCS3700 (Figure 12-18 on page 270) features the latest technologies, including 6 Gbps SAS and 8 Gbps FC host interfaces, along with 6 Gbps SAS drives. The DCS3700 is equally adept at delivering throughput to bandwidth-intensive applications and I/O operations to transactional applications, such as databases and Microsoft Exchange. This capability is why the product is classified as an enterprise disk storage system, even if the dimensions are not typical for this category.

The DCS3700 is ideally suited for high performance streaming applications, such as rich media, financial markets, telecommunications, weather modeling, and others that need rigorous bandwidth requirements.

Figure 12-18 shows an image of the IBM System Storage DCS3700.



Figure 12-18 IBM System Storage DCS3700

The following list includes the key benefits and highlights of the IBM System Storage DCS3700:

- ▶ Six Gbps SAS high-density storage system that delivers scalable capacity at an affordable price point
- ▶ Multilevel data protection with IBM FlashCopy, Volume Copy, and Remote Mirroring across FC
- ▶ Mixed host interfaces support direct-attached storage (DAS) and SAN tiering to reduce overall operation and acquisition costs
- ▶ Investment protection and cost-effective backup and recovery with remote mirror across FC host ports and compatibility with DS3500, DS5000, and DS4000
- ▶ Up to 180 drives per system with the attachment of two DCS3700 Expansion Units (60 drives per enclosure) with 20 drives minimum drive quantity per enclosure

Discontinued enterprise disk systems

IBM no longer offers the following products through its marketing channels. However, these systems are still seen in large and strategic hosting data centers across geographies:

- ▶ IBM TotalStorage Enterprise Storage Server (ESS)
- ▶ IBM System Storage DS8100
- ▶ IBM System Storage DS8300

For more information about IBM high-end, enterprise disk systems, see this website:

<http://www.ibm.com/systems/storage/disk/enterprise/index.html>

12.4 IBM Tape Storage Systems

Tape systems were traditionally associated with the mainframe computer market. This association is because they represented an essential element in mainframe systems architecture since the early 1950s, as a cost-effective way to store large amounts of data. In contrast, the midrange and client/server computer market has made limited use of tape technology until recently.

However, over the past few years, growth in the demand for data storage and reliable backup and archiving solutions has greatly increased the need to provide manageable and cost-effective tape library products. The value of using tape for backup purposes has only gradually become obvious and important in these environments.

In this section, we briefly guide you through the IBM Fibre Channel attached tape drives, autoloaders, and automated tape libraries. We do not describe the current tape technology (mainly, Linear Tape Open (LTO)) because this topic is widely described in other publications.

For more information the current tape technology, see the *IBM System Storage Tape Library Guide for Open Systems*, SG24-5946. The guide is available at this website:

<http://www.ibm.com/systems/storage/product/tape.html>

12.4.1 Fibre Channel tape drives

Client data is vital to business operations. IBM offers entry-level and enterprise tape products that are designed to provide backup and protection of client data at an appealing cost to a company's budget. In this section, we give an overview of the following tape drives:

- ▶ IBM System Storage TS2230, TS2240, TS2250 half-height LTO tape drives
- ▶ IBM System Storage TS2340 and TS2350 full-height LTO tape drives
- ▶ IBM System Storage TS1120, TS1130, TS1140 for 3590/3592 cartridges

We intend to only describe Fibre Channel-connected tape drives (that is, 3/6 Gbps SAS attached and LVD/HVD SCSI attached) because this book is focused on SAN-related networking technology.

IBM Ultrium LTO Half-Height tape drives

IBM Ultrium Linear Tape Open (LTO) Half-Height external tape drives include the following models:

- ▶ IBM System Storage TS2230 Tape Drive Express for LTO-3 cartridges
- ▶ IBM System Storage TS2240 Tape Drive Express for LTO-4 cartridges
- ▶ IBM System Storage TS2250 Tape Drive Express for LTO-5 cartridges

TS2230 and TS2240 tape drives are suited for handling backup, save, and restore, and archival data storage needs with higher capacity and higher data transfer rate than previous generation. Additionally, TS2250 tape drives continue to support tape encryption as the previous generation of LTO4. They all support Write-Once-Read-Many (WORM) tape cartridges for compliance purposes.

Figure 12-19 shows all three models of IBM Ultrium LTO Half-Height external tape drives.



Figure 12-19 IBM TS2230 (upper left), TS2240 (lower left), and TS2250 (right)

Half-High tape drives: IBM External Half-High tape drives use a 6 Gbps SAS connection only. However, there is an available FC-attached model, which is dedicated for internal installation into entry tape libraries and enclosures. These versions are not considered as external tape drives. Therefore, they are not covered in detail in this section.

The following list provides highlights of the IBM External Half-High tape drives:

- ▶ External stand-alone or rack-mountable, easy-to-install units
- ▶ Standard 6 Gbps SAS interface for external models; 8 Gbps FC (4 Gbps LTO3/LTO4) interface for internal half-height LTO5 tape drives
- ▶ Support for Write Once Read Many (WORM) cartridges and hardware compression; generation LTO5 offers tape partitioning on selected platforms
- ▶ Native capacity of 1.5 TB (3.0 TB compressed) for LTO5 tape cartridges (800 GB/1.6 TB for LTO4, 400 GB/800 GB for LTO3)
- ▶ Maximum data rate 140 MBps (LTO5), 120 MBps (LTO4), 80 MBps (LTO3)
- ▶ Transparent installation and configuration across various supported platforms that use dedicated IBM tape device drivers

IBM Ultrium LTO Full-Height tape drives

IBM Ultrium LTO Full-Height external tape drives incorporate industry-proved IBM Linear Tape Open (LTO) full sized tape drives of the series 3580. They are delivered in compact stand-alone or rack mountable enclosures easy-to-install. IBM offers two models of external tape drives in this category:

- ▶ IBM System Storage TS2340 Tape Drive Express for LTO-4 cartridges
- ▶ IBM System Storage TS2350 Tape Drive Express for LTO-5 cartridges

Similar to the half-height external tape drives, the full-height tape drives are offered by IBM as internal modules for installation in midrange or enterprise tape libraries that use either a 3/6 Gbps SAS or 8/4 Gbps FC interface. Both of the devices are presented in Figure 12-20.



Figure 12-20 Full-Height IBM external tape drives TS2340 (left) and TS2350 (right)

The following list provides a brief summary of the features and benefits of both the products:

- ▶ External stand-alone or rack-mountable, easy-to-install devices
- ▶ TS2340: 3 Gbps or 320 MBps SCSI interface. TS2350: 6 Gbps SAS port. Internal tape drives type 3580: 4 Gbps FC interface for LTO4, 8 Gbps FC-attached LTO5 tape drives
- ▶ Native capacity of 1.5 TB (3.0 TB compressed) for LTO5 tape cartridges (800 GB/1.6 TB for LTO4)
- ▶ Maximum data transfer rate 140 MBps (LTO5), up to 120 MBps (LTO4)

- ▶ Available WORM cartridges and hardware compression, generation LTO5 offers tape partitioning on selected platforms

IBM System Storage 3592 tape drives

IBM System Storage 3592 series of tape drives offer a design that is focused on high capacity and performance, and high reliability for storing your mission critical data. This series includes the following available IBM tape drives, as shown in Figure 12-21:

- ▶ IBM System Storage TS1120
- ▶ IBM System Storage TS1130
- ▶ IBM System Storage TS1140

Whether configured as a stand-alone drive or part of an automated tape library, these models offer fast access to data and high capacity in a single drive. These features help to reduce the complexity of your tape infrastructure. They have the same form factor as its predecessors and machine type 3592. See 12.5, “Storage virtualization and cloud computing” on page 278 for details.



Figure 12-21 IBM TS1120 (left), TS1130 (center), and TS1140 (right)

To help optimize drive utilization and reduce infrastructure requirements, these three models of tape drives can be shared among supported open system hosts on a SAN or on IBM FICON mainframe hosts when they are attached to an IBM System Storage Tape Controller for System z.

The following list provides product highlights of the three models of IBM tape drives:

- ▶ High performance with data transfer rate up to 650 MBps with compression
- ▶ Flexible media, including short and long length cartridges, rewritable and WORM formats and media partitioning on given platforms
- ▶ Four TB capacity by using JC/JY media, 1.6 TB by using JB/JX media, and 500 GB by using JK media (TS1140)
- ▶ One TB capacity by using JB/JX media, 640 GB by using JA/JW media, and 128 GB by using JJ/JR media (TS1130)
- ▶ Capacity of 700 GB by using JB/JX media, 300/500 GB by using JA/JW media, and 60/100 GB capacity by using JJ/JR media (TS1120)
- ▶ IBM Power Systems, System i, System p, System z, and System x support

12.4.2 Autoloaders and entry tape libraries

Single tape drive autoloaders and entry-class tape libraries are suited to handle backup, save, and restore and for archival data storage needs in small to medium size environments. They benefit from LTO technology because they typically incorporate IBM half-height tape drives that are connected over Fibre Channel or a 3/6 Gbps SAS interface.

In the following section, we describe only the model that is able to accommodate FC-attached full-height tape drives: the IBM System Storage TS3200 Tape Library Express.

We do not describe other available IBM entry tape automation products because they do not offer an FC interface and are directly attached to the host systems:

- ▶ IBM System Storage TS2900 Tape Autoloader
- ▶ IBM System Storage TS3100 Tape Library Express

Technical details of the TS2900 and TS3100 devices (including TS3200) can be found at the following website:

<http://www.ibm.com/systems/storage/tape/entry/index.html>

IBM System Storage TS3200 Tape Library Express

The TS3200 Express Model and its storage management applications are designed to address capacity, performance, data protection, reliability, affordability, and application requirements. It is designed as a functionally rich, high capacity, entry level tape storage solution that incorporates the LTO Ultrium tape technology. The IBM TS3200 Express model is an excellent solution for large capacity or high performance tape backup with or without random access. The TS3200 is also an excellent choice for tape automation for IBM Power Systems, IBM System x, and other open systems.

The IBM TS3200 Express (Figure 12-22) is designed to support the newest generation of LTO with up to two IBM Ultrium 5 full-height tape drives or up to four IBM Ultrium 5 half-height tape drives, as well as LTO generations 3 and 4 tape drives using a 4U form factor.



Figure 12-22 Top front view of the IBM TS3200 Express

IBM System Storage TS3200 Tape Library Express provides these features:

- ▶ Supports the LTO-5, LTO-4, or LTO-3 tape drives, Low Voltage Differential (LVD) 320 MBps SCSI, 8 Gbps FC, and 6 Gbps SAS attachments
- ▶ Up to two full-height or four half-height IBM tape drives
- ▶ Sequential or random access mode with a standard barcode reader
- ▶ 4U rack mountable or stand-alone form factor with up to 48 cartridge slots provides native capacity up to 72 TB (144 TB with LTO-5 compression)
- ▶ Remote library management through a transparent, intuitive web interface

12.4.3 Midrange tape libraries

Whether a small/medium size business is expanding operations or experiencing rapid data growth, IBM midrange tape libraries are designed to help meet the client needs of data backup, archive, and management. These tape products are designed to offer reliability, performance, and flexibility for today and the future. This section introduces the IBM tape automation product: IBM System Storage TS3310 Tape Library.

We provide the typical features of this midrange tape library. For further technical details, including all of the available configuration scenarios, see the *IBM System Storage Tape Library Guide for Open Systems*, SG24-5946; see the most current updates at this website:

<http://www.ibm.com/systems/storage/tape/ts3310/index.html>

The following midrange product is withdrawn from IBM offerings: IBM System Storage TS3400 Tape Library.

IBM System Storage TS3310

The IBM System Storage TS3310 Tape Library is a modular, scalable tape library that is designed to address the tape storage needs of rapidly growing companies who find themselves space and resource constrained with tape backup and other tape applications.

Designed around a 5U high modular base library unit, the TS3310 can scale vertically with expansion for LTO tape cartridges, drives, and redundant power supplies. Each expansion module (9U) contains 92 physical LTO cartridge storage cells and space for up to four LTO-5, LTO-4, and LTO-3 tape drives. Additionally, the module has space for up to two power supply modules - with one redundant. See Figure 12-23 for the configuration of L5B base module and one expansion unit E9U.



Figure 12-23 Base unit and enclosure of the TS3310

The TS3310 supports either the standard or WORM LTO data cartridge and continued support for encryption of data with LTO-4 and LTO-5 tape drives. IBM Tivoli Key Lifecycle Manager is required for encryption key management with Ultrium 5 drives.

Product highlights of the IBM System Storage TS3310 include the following features:

- ▶ Modular, scalable tape library that is designed to grow up to the capacity of 409 LTO cartridges and 18 LTO-5, LTO-4, and LTO-3 tape drives
- ▶ Form factor from 5U (base module) up to 41U (base module plus four expansion units) in standard 19-inch rack mountable or stand-alone configuration
- ▶ Intuitive web-based remote management
- ▶ Hot-swap tape drives and power supplies
- ▶ Available datapath and control path failover for redundant host connectivity
- ▶ Support for a wide range of systems which include IBM System p, System x, System i, AS/400®, IBM RS/6000®, Intel, Hewlett-Packard (HP), and Sun

12.4.4 Enterprise tape libraries

The enterprise tape libraries provide storage solutions for large, unattended storage operations from today's midrange up to the enterprise (z/OS and Open Systems) environments. We describe the following model of the enterprise tape automation product: IBM System Storage TS3500 Tape Library.

The IBM TS3500 is available as a universal, effective tape solution for either Open Systems or System z environments; therefore, IBM no longer offers the following product: IBM TotalStorage 3494 Tape Library

Comprehensive information about the TS3500, including deployment in a z/OS environment, is available: *IBM TS3500 Tape Library with System z Attachment A Practical Guide to Enterprise Tape Drives and TS3500 Tape Automation*, SG24-6789. See this website:

<http://www.ibm.com/systems/storage/tape/ts3500/index.html>

IBM System Storage TS3500

The TS3500 Tape Library offers outstanding retrieval performance with typical cartridge move times of less than 3 seconds. This performance is possible by combining reliable, automated tape handling and storage with reliable, high-performance IBM LTO Ultrium and 3592 tape drives.

The TS3500 Tape Library can be partitioned into multiple logical libraries. This feature makes it an excellent choice for consolidating tape workloads from multiple heterogeneous Open Systems servers. It also enables support for System z attachment in the same library.

In addition, the TS3500 Tape Library provides outstanding reliability and redundancy. These benefits are possible through the provision of redundant power supplies in each frame, an optional second cartridge accessor, service bays for nondisruptive maintenance, control and data path failover, and dual grippers within each cartridge accessor. Both library and drive firmware can now be upgraded nondisruptively, that is, without interrupting the normal operations of the library.

The IBM TS3500 Tape Library (Machine Type 3584) is a modular, highly scalable tape library that consists of frames that house tape drives and cartridge storage slots. You can install a single-frame base library (Figure 12-24 on page 277) and expand it to 16 frames, tailoring the library to match your system capacity requirements.

Up to 12 IBM Ultrium tape drives can be installed in a single frame. All generations of LTO tape drives are supported by TS3500. As the enterprise tape libraries are being implemented

in the most complex data centers and client backup solutions that use SAN, most of tape drive attachments are made based on FC interface (up to 8 Gbps with Ultrium LTO-5).

Figure 12-24 shows the base frame of the TS3500.



Figure 12-24 Base frame of the TS3500

The following list provides the benefits of the IBM TS3500 in enterprise SAN environments:

- ▶ One base frame, and up to 15 expansion frames per library; up to 15 libraries that are interconnected per complex
- ▶ Up to 12 FC attached drives per frame (192 per library, 2700 per complex)
- ▶ IBM 3592 JA/JJ/JB/JC and JW/JR/JX/JY (WORM) cartridges or IBM LTO Ultrium 5, 4, 3, 2, 1 cartridges
- ▶ Up to 60 PB compressed with IBM Ultrium 5 cartridges per library, up to 900 PB compressed per complex,
- ▶ Up to 180 PB compressed with 3592 extended capacity cartridges per library, up to 2.7 EB compressed per complex
- ▶ Remote management by using a web browser
- ▶ Advanced Library Management System (ALMS), SNMP functionality, multipath architecture, and persistent WWN

12.5 Storage virtualization and cloud computing

In many IT departments, increased user demand has led to haphazard storage growth, resulting in sprawling, heterogeneous storage environments. These environments make it difficult to achieve optimal utilization and to provision storage capacity for new users and applications. Storage virtualization can put an end to these problems. It enables companies to logically aggregate disk storage so capacity can be efficiently allocated across applications and users.

Virtualization solutions help take the cost and complexity out of IT infrastructures. In this section, we describe two important topics: disk and tape virtualization, and which IBM products participate in this process. Finally, we briefly explain what benefits bring IBM storage products to the cloud computing and which of them are the key building blocks of the *IBM Smart Business Storage Cloud*.

12.5.1 Disk storage virtualization

IBM disk storage virtualization products provide simplified and centralized management of clients' medium to enterprise storage infrastructures that consist of different disk systems, even from different vendors. In this section, we describe the following IBM disk systems that enable storage virtualization:

- ▶ IBM System Storage SAN Volume Controller
- ▶ IBM XIV Storage System
- ▶ IBM Storwize V7000 Unified

For more details about each of these products, visit:

<http://www.ibm.com/systems/storage/virtualization>

IBM System Storage SAN Volume Controller

SANs enable companies to share homogeneous storage resources across the enterprise. But for many clients, information resources are spread over various locations and storage environments, often with products from different vendors, who supply everything from mainframes to notebooks. To achieve higher utilization of resources, clients now need to share their storage resources from all their environments, regardless of the vendor. IBM System Storage SAN Volume Controller contributes towards this goal of a solution that can help strengthen existing SANs by increasing storage capacity, efficiency, uptime, administrator productivity, and functionality.

IBM System Storage SAN Volume Controller Software is delivered preinstalled on SAN Volume Controller Storage Engines so it is quickly ready for implementation when the engines are attached to your SAN. SAN Volume Controller Storage Engines are based on proven IBM System x server technology and are always deployed in redundant pairs (see Figure 12-25 on page 279), which are designed to deliver high availability.



Figure 12-25 IBM SAN Volume Controller in a clustered pair

SAN Volume Controller is designed to take control of existing storage, retaining all of your existing information. This ability helps the speed and provides simplified implementation, while helping to minimize the need for more storage. When the SAN Volume Controller is implemented, you can change the configuration quickly and easily as needed.

SAN Volume Controller is designed to support nondisruptive data migration between storage systems. In addition, SAN Volume Controller helps make storage potentially available to all attached servers, greatly increasing the flexibility for using for example VMware vMotion. Without SAN Volume Controller, use of vMotion might be limited by the storage that is being dedicated to specific servers.

Because SAN Volume Controller is displayed to servers as a single type of storage, virtual server provisioning is also simplified because only a single driver type is needed in server images. This feature also simplifies administration of those server images. Similarly, SAN Volume Controller eases the replacement of storage or the movement of data from one storage type to another because these changes do not require changes to server images. Without SAN Volume Controller, changes of the storage type might require disruptive changes to the server images.

The SAN Volume Controller is based on the IBM System x3550 with Intel Xeon 5600 2.5 GHz processor and 24 GB of memory cache. The device incorporates four 8 Gbps FC ports, two 1 Gbps (optionally, an additional two 10 Gbps) iSCSI ports. The controller contains redundant power supplies in a standard 19-inch rack-mountable enclosure.

Briefly summarized, the IBM System Storage SAN Volume Controller is designed with the following functions in mind:

- ▶ Combines storage capacity from multiple vendors for centralized management. And, increases storage utilization by providing more flexible access to storage assets.
- ▶ Improves administrator productivity by enabling management of pooled storage from a single interface.
- ▶ Reduces downtime by insulating host applications from changes to the physical storage infrastructure.
- ▶ Enables a tiered storage environment to match the cost of storage to the value of data.
- ▶ Supports data migration among storage systems, without interruption to applications.
- ▶ Supports consolidated disaster recovery site servicing for more than one production location.

IBM XIV Storage System

IBM XIV Storage System is a proven, high-end disk storage device, which is designed for growth with an unmatched ease of use. IBM XIV eliminates the complexity of managing enterprise storage. It never compromises performance for reliability, providing consistent high performance without manual tuning for even the most demanding application workloads, while keeping TCO incredibly low. Its grid architecture delivers virtual storage that optimizes performance in virtualized environments and integrates seamlessly with cloud technologies to provide the agility that clients need to handle growth.

The IBM XIV Storage System is offered in two models, both based on the same proven XIV architecture and all-inclusive pricing approach:

- ▶ **IBM XIV** features powerful storage for most application needs, with excellent price-to-performance advantages. This model offers client-acclaimed value in handling a mix of diverse workloads at low TCO. It incorporates 1 or 2 TB SATA disk drives up to 161 TB usable capacity (15 modules).
- ▶ **IBM XIV Gen3** (Figure 12-26) features state-of-the-art hardware that can do more, and do more faster. Rely on it for your ultra-demanding performance objectives, including business intelligence, archiving, large email setups, data warehousing, and OLTP workloads, as well as ever-changing virtualized and cloud environments. The 2 TB SAS-attached disk drives are used.



Figure 12-26 Third generation of IBM XIV Storage System

Several architectural features contribute to the unique performance profile of the XIV system:

- ▶ **Massive parallelism in a fully distributed architecture:** XIV is based on a distributed architecture of interconnected modules that include multi-core processors, large cache, and high-density disk drives.
- ▶ **Distributed data:** The system stores data by breaking it down into 1-megabyte chunks that are called *partitions*, each mirrored for redundancy to another module.
- ▶ **Distributed bandwidth within modules:** Aggressive prefetching is enabled by the large cache-to-disk bandwidth that is available within each module.
- ▶ **Load balancing:** The system distributes the application load across all system modules uniformly.

- ▶ **High performance during disk rebuild:** XIV rebuilds a failed disk at unprecedented speed because of a distributed rebuild mechanism that engages all disks in the system in the rebuild process.

IBM XIV integrates easily with virtualization, email, database, analytics, and data protection solutions from Microsoft, IBM, SAP, Oracle, SAS, VMware, Hyper-V, and Symantec. The XIV Gen3 model gives applications a tremendous performance boost, helping clients meet increasing demands with fewer servers and networks. The XIV series plays a key role in IBM end-to-end dynamic infrastructure solutions, integrating seamlessly with IBM ProtecTIER, Scale Out Network Attached Storage, SAN Volume Controller, Storwize V7000, and Tivoli products.

IBM Storwize V7000 Unified

IBM Storwize V7000 Unified is a virtualized storage system. The system complements virtualized server environments that provide unmatched performance, availability, advanced functions, and a highly scalable capacity that has not been seen before in midrange disk systems. Storwize V7000 Unified is a powerful midrange disk system that has been easy to use and to enable rapid deployment without more resources.

Figure 12-27 shows the rack installation of both types of Storwize V7000 Unified bays: 12 3.5-inch disks and 24 2.5-inch disk drives.



Figure 12-27 IBM Storwize V7000 and its rack-mountable disk bays

Storwize V7000 Unified consolidates block and file workloads into a single storage system for simplicity of management and reduced cost. It offers greater efficiency and flexibility through built-in solid-state drive (SSD) optimization and thin provisioning technologies. Storwize V7000 Unified advanced functions also enable non-disruptive migration of data from existing storage, simplifying implementation and minimizing disruption to users. In addition, it enables you to virtualize and reuse existing disk systems, supporting a greater potential return on investment (ROI). The three key functions help to provide a single point of control to support improved storage efficiency:

- ▶ **Consolidation** of storage resources by efficient scaling, improves productivity and reduces cost.
- ▶ **Virtualization** of storage infrastructure can optimize expenditures, resources, and capabilities. New support for VMware vStorage APIs enables Storwize V7000 to take on

some storage-related tasks that were previously performed by VMware. This benefit helps to improve efficiency and frees up server resources for other more mission-critical tasks.

- ▶ **Tiering** optimizes storage by enabling data to be in a way that can improve system performance, reduce costs, and simplify information management. Using IBM System Storage Easy Tier technology, Storwize V7000 can use SSDs confidently, effectively, and economically by automatically and dynamically moving only the appropriate data to the SSDs in the system, which is based on performance monitoring.

The following list provides the IBM Storwize V7000 product highlights:

- ▶ Control enclosure supports up to 240 TB of physical capacity by using 12 standard 2U expansion units with SAS attached disk drives.
- ▶ Two control enclosures can be clustered for high availability and optimal performance.
- ▶ Host attachments that use eight 8 Gbps, four 1 Gbps, or optionally 10 Gbps iSCSI ports.
- ▶ File module provides attachment to 1 Gbps and 10 Gbps NAS environments.
- ▶ Incorporated SSDs support business applications that need to grow dynamically with effective space utilization (thin provisioning). Easy Tier automatically migrates frequently used data to SSDs.
- ▶ Support of *IBM Advanced Copy Services* such as IBM FlashCopy, Metro Mirror, and Global Mirror. IBM Tivoli FlashCopy Manager shortens backup and recovery times.

12.5.2 Tape storage virtualization

To maintain continuous business operations, address regulatory requirements, and archive business records, clients need an infrastructure that enables them to manage their data from online application storage to offline, permanent archive media. Tape is a key part of both the backup and archive lifecycle. Tape still provides the lowest TCO alternative for securely storing long-term archives for record keeping and disaster recovery. As data centers and data stores grow, tape operations can become more complex. This growth can lead to increased backup and restore times and higher management involvement and costs. Tape virtualization solutions can help to get over these constraints.

IBM offers the following tape storage virtualization products:

- ▶ IBM System Storage TS7610 ProtecTIER Deduplication Appliance Express
- ▶ IBM System Storage TS7650 ProtecTIER Deduplication Appliance
- ▶ IBM System Storage TS7650G ProtecTIER Deduplication Gateway
- ▶ IBM System Storage TS7680 ProtecTIER Deduplication Gateway for System z
- ▶ IBM Virtualization Engine TS7700 Family

IBM withdrew these tape virtualization products from marketing:

- ▶ IBM Virtualization Engine TS7520
- ▶ IBM Virtualization Engine TS7530
- ▶ IBM Virtual Tape Server 3494 - B10
- ▶ IBM Virtual Tape Server 3494 - B20

For more information about tape storage virtualization, see this IBM Redbooks publication: *Implementing IBM Storage Data Deduplication Solutions*, SG24-7888, and *TS7680 Deduplication ProtecTIER Gateway for System z*, SG24-7796. See this website for a product summary:

<http://www.ibm.com/systems/storage/tape/virtualization/index.html>

IBM ProtecTIER deduplication appliances

IBM ProtecTIER deduplication solutions feature revolutionary and patented HyperFactor data deduplication technology. This technology provides enterprise class performance, scalability, and proven enterprise level data integrity to meet the disk-based data protection needs of the enterprise data center, down to midmarket environments. These benefits are realized while also gaining significant infrastructure cost reductions.

The IBM System Storage TS7610 ProtecTIER Deduplication Appliance Express provides fast, reliable, and easy-to-deploy backup and recovery for midsize IT environments. This solution has a preconfigured repository and can be configured with either a Virtual Tape Library (VTL) or Symantec OpenStorage (OST) interface. Available in two configuration options (4 and 5.4TBs), the TS7610 provides capacity, price, and performance features; as well as, reliability, availability, and serviceability (RAS) features that are required by midsize clients.

The IBM System Storage TS7650 ProtecTIER Deduplication Appliance is designed to improve backup and recovery operations. The system is available in four preconfigured solutions that range from 7 TBs - 36 TB in a cluster. This integrated solution makes it easy to harness the power of deduplication without making radical changes to the existing environment. Again, TS7650 has a preconfigured repository and can be deployed with either VTL or OST.

Figure 12-28 presents both models of the IBM ProtecTIER Deduplication Appliances.



Figure 12-28 IBM ProtecTIER Deduplication Appliances: TS7650 and TS7610 Express

The following list provides the key benefits of the TS7650 and the TS7610 Express:

- ▶ Improves backup and recovery performance with high-speed disk-based data protection that enables more efficient, reliable storage of valuable data, while optimizing storage infrastructure and reducing TCO.
- ▶ Simple to install, manage, and maintain data protection for midsize IT environments.
- ▶ Provides VTL or OST interface support.
- ▶ Up to 100/500 MBps or more of inline data deduplication performance.
- ▶ Up to 25 times or more storage capacity reduction.

- ▶ Emulation of up to 4/12 virtual libraries, 64/256 virtual drives, and 8192/128000 virtual cartridges.
- ▶ Capacity to store up to 135/900 TBs or more of backup data on a single 5.4/36 TB appliance.

IBM ProtecTIER Deduplication Gateways

The IBM ProtecTIER Deduplication Appliances offer complete preconfigured deduplication solutions with installed storage capacity for medium and enterprise data centers. However, *IBM ProtecTIER Deduplication Gateways* require more SAN-attached storage systems to hold and protect business critical and valuable backup data. IBM introduces two deduplication products, one for Open Systems: IBM TS7650G ProtecTIER Deduplication Gateway, and one for the mainframe environment: IBM TS7680 ProtecTIER Deduplication Gateway for System z.

Both models of deduplication gateways offer high-performance inline data deduplication, high-availability clustering with Global Deduplication, and flexibility to support up to 1 petabyte (PB) of physical storage capacity on both IBM and non-IBM supported storage systems. Figure 12-29 shows the front panel of both, identically looking ProtecTIER Deduplication Gateways.



Figure 12-29 IBM TS7650G ProtecTIER Deduplication Gateway

The following list provides the benefits and key highlights of ProtecTIER gateways:

- ▶ High-speed backups with up to 2000 MBps (7.2 TB/hr) or more sustained inline deduplication backup performance and up to 2800 MBps (10 TB/hr) or more sustained recovery performance.
- ▶ Easily scalable to provide up to 25-PB backup storage capacity, up to 512 virtual tape drives, and 1 million logical volumes per two node clusters.
- ▶ Non-hash-based approach avoids the possibility of data loss because of a hash collision.
- ▶ ProtecTIER Native Replication uses deduplication technology in the disk repositories at both the primary and secondary sites to lower bandwidth requirements.
- ▶ Inline deduplication enables replication to occur concurrently with backup operations to increase responsiveness and ability to restore data quickly when needed.
- ▶ Intuitive graphical user interface-based (GUI-based) management, and monitoring tools

IBM Virtualization Engine TS7700

The *IBM Virtualization Engine TS7700* is a family of mainframe virtual tape solutions that are designed to optimize tape processing. With one solution, the implementation of a fully integrated tiered storage hierarchy of disk and tape takes advantage of the benefits of both technologies. These benefits help to enhance performance and provide the capacity that is needed for today's tape processing requirements. Deploying this innovative subsystem can help reduce batch processing time, TCO, and management involvement.

A TS7700 (Figure 12-30) can help improve the efficiency of mainframe tape operations by efficiently using disk storage, tape capacity, and tape speed, and by providing many tape addresses. These benefits help make the TS7700 a suitable repository for local and remote backups and archival data. Two models are available for purchase.

The *TS7720 Virtualization Engine* provides high capacity for workloads that are cache friendly because of their rapid recall requirements. The TS7720 features 2-TB SATA disk drives with RAID 6 to allow clients to scale their solution to meet the needs of growing workloads without affecting application availability.

The *TS7740 Virtualization Engine* supports attachment to and uses the performance and capacity of the IBM System Storage TS1130 and TS1120 Tape Drives. Or, the TS7740 uses the IBM TotalStorage 3592 Model J1A Tape Drive that is installed in an IBM System Storage TS3500 Tape Library or an IBM TotalStorage 3494 Tape Library. Support for these tape drives can help to reduce the number of cartridges and the size of the library. This reduction occurs by allowing storage of up to 3 TB on a single 3592 JB cartridge, assuming 3:1 compression.

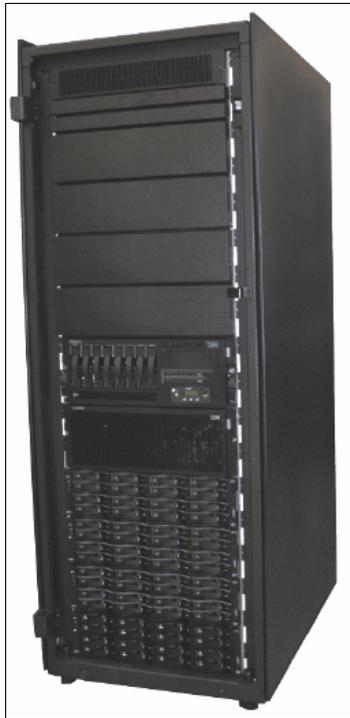


Figure 12-30 IBM TS7700

12.5.3 Storage systems for cloud computing

IBM offers three types of cloud solutions, for storage and other services:

- ▶ Smart Business on IBM SmartCloud™ are standardized services that are provided by IBM on a pay-per-use basis.
- ▶ Smart Business Cloud services are private cloud services, behind the firewall of a client, which are built or run by IBM.
- ▶ Smart Business Systems are purpose built, integrated service delivery platforms.

IBM Smart Business Storage Cloud

IBM Information Infrastructure includes next generation virtualized storage and storage management products that can support the demands of cloud computing. Cloud computing applications are typically deployed in a virtualized environment with a strong security model. The following IBM products are key building blocks of IBM Smart Business Storage Cloud:

- ▶ **IBM System Storage SAN Volume Controller:** Virtualize IBM and non-IBM storage to enable resource pooling, thin provisioning, and simplified management. See “IBM System Storage SAN Volume Controller” on page 278 for details.
- ▶ **IBM XIV Storage System:** Automates and virtualizes data management, and dramatically simplifies systems management to help tame dynamic workloads of clients. The XIV is introduced in “IBM XIV Storage System” on page 280.
- ▶ **IBM Scale Out Network Attached Storage:** Offers an extreme scale-out capability for large storage infrastructures that require high availability. Scale Out Network Attached Storage also delivers computing services that make the supporting technology almost invisible. It enables applications and services to be uncoupled from the underlying infrastructure, enabling businesses to adjust to change quickly. As a result, Scale Out Network Attached Storage can easily integrate with the strategies of your organization to develop a more dynamic enterprise. The comprehensive information about Scale Out Network Attached Storage is included in *IBM Scale Out Network Attached Storage: Architecture, Planning, and Implementation Basics*, SG24-7875
- ▶ **IBM Tivoli Storage Productivity Center:** Manages virtualized storage and generates usage reports. Automates storage performance and capacity management. Detailed guidance through this robust family of IBM Tivoli software products is available in the *IBM Tivoli Storage Productivity Center V4.2 Release Guide*, SG24-7894. Its benefits for SAN infrastructure are available in *SAN Storage Performance Management Using TotalStorage Productivity Center*, SG24-7364.
- ▶ **Media encryption:** This encryption is essential for cloud applications because it provides a strong security model with minimal overhead. This encoding also provides a data shredding capability that costs almost nothing to start (just delete the centralized encryption key). Drive-level media encryption and centralized key management is available for midrange and enterprise disk and tape. IBM Tivoli Key Lifecycle Manager is the core IBM application for simple, strong, and centralized encryption key management.

12.6 IP-based networking for SAN environments

Ethernet networking within complex data centers, especially 10 Gbps Ethernet (10 GbE), provides a significant number of benefits to the storage networks:

- ▶ Simplifies storage management
- ▶ Increases utilization of assets in the storage networks
- ▶ Improves flexibility of storage environments to adopt new solutions
- ▶ Reduces cost of infrastructure by using consolidation of assets
- ▶ Improves efficiency of storage and ethernet network

Ethernet protocols such as iSCSI or FCoE enable clients to start thinking about migration from typical SAN Fibre Channel networking to consolidated storage and Ethernet networking on single network. IBM offers various products to support and enable converged networking for flexibility and efficiency of asset utilization in complex data centers or just to connect remote offices with all the benefits of SAN (Data Center Bridging).

In the following text, we briefly describe two types of IBM products:

- ▶ **Hardware offerings** include IBM devices that provide a consolidated networking solution that uses Converged Network Adapters, which transport both SAN and ethernet LAN data.
- ▶ **Software solutions** use IBM Virtual Fabric emulation on given Emulex network adapters and specific IBM System Networking products.

We do not describe the range of IBM Ethernet products in this book. For more information about the Ethernet approach of IBM, see these websites:

IBM System Networking RackSwitch™

<http://www.ibm.com/systems/networking/switches/rack.html>

IBM System Networking BladeCenter

<http://www-03.ibm.com/systems/networking/switches/bladecenter.html>

IBM Distributed Virtual Switch

<http://www-03.ibm.com/systems/networking/switches/virtual/index.html>

12.7 Hardware solutions for network convergence

In this section, we introduce products that offer converged SAN and LAN networking by utilization of *Converged Network Adapters*:

- ▶ Cisco Nexus 5000 for IBM System Storage
- ▶ IBM System Networking RackSwitch G8124
- ▶ IBM System Networking RackSwitch G8264

The following device is withdrawn from marketing: IBM Converged Switch B32.

Cisco Nexus 5000 for IBM System Storage

Cisco Nexus 5000 switches for IBM System Storage are designed for data center environments with technology that supports consistent low latency Ethernet solutions, with front to back cooling, and with network ports in the rear. This configuration brings switching into close proximity with the servers, making cable runs short and simple. The switch family is highly serviceable, with optional redundant, hot-pluggable power supplies and fan modules. It uses data center class Cisco NX-OS to support high reliability and ease of management.

The switch family, by using cut-through architecture, supports line-rate 10 Gigabit Ethernet on all ports while maintaining consistent low latency. This support is independent of the packet size and services that are enabled. The product family supports IEEE Data Center Bridging and Converged Enhanced Ethernet (CEE) capabilities that can help increase the reliability, efficiency, and scalability of Ethernet networks. These features allow the switch to support multiple traffic classes over an Ethernet fabric, thus enabling consolidation of LAN, SAN, and cluster environments. Its ability to connect FCoE to native Fibre Channel, protects existing storage system investments while dramatically simplifying in-rack cabling.

In addition to supporting standard 10 Gigabit Ethernet network interface cards (NICs) on servers, the Cisco Nexus 5000 switches (Figure 12-31 on page 288) integrate with multifunction adapters called *converged network adapters (CNAs)* that combine the functions of Ethernet NICs and Fibre Channel *host bus adapters (HBAs)*. This integration makes the transition to a single, unified network fabric that is consistent with existing practices, management software, and operating system drivers. The switch family is compatible with integrated transceivers and Twinax cabling solutions to help deliver cost-effective connectivity

for 10 Gigabit Ethernet to the servers at the rack level. This compatibility reduces or eliminates the need for expensive optical transceivers.



Figure 12-31 Cisco Nexus 5000 switches for IBM System Storage

The following list provides product highlights of the *Cisco Nexus 5000* switches for IBM System Storage:

- ▶ Designed as a 1U (Cisco Nexus 5010 28 port switch, Cisco Nexus 5548P, and 5548UP, with up to 48 ports) and as a 2U (Cisco Nexus 5020 56 port switch and Cisco Nexus 5596UP, with up to 96 ports) 19-inch rack mountable or stand-alone enclosure.
- ▶ Expansion modules include eight 1, 2, 4, and 8 Gbps FC ports; four 10 GbE and four 1, 2, 4, and 8 Gbps FC ports; six 10 GbE ports.
- ▶ Ten GbE ports are capable of transporting both storage and LAN traffic, which eliminates the need for separate server SAN and LAN adapters and cables.
- ▶ Consistent management is provided through consistency of both Cisco NX-OS Software and Cisco MDS 9000 SAN-OS Software management models and tools.
- ▶ IEEE Data Center Bridging features for lossless transmission, priority flow control, and enhanced transmission selection.
- ▶ Enterprise-class availability features such as hot-swappable, field replaceable, redundant power supplies, redundant fan modules, and port expansion modules.

For more information about switches, see this website:

<http://www.ibm.com/systems/networking/hardware/ethernet/c-type/nexus/>

12.7.1 IBM Virtual Fabric solution

IBM Virtual Fabric solution for IBM System z uses IBM System Networking *convergence ready* products and specific *Emulex Virtual Fabric adapters*. These products are not classified as the typical hardware convergence solutions that we described in Section 12.7, “Hardware solutions for network convergence” on page 287.

This new and innovative solution that is based on Emulex adapters and IBM System Networking Rack Switch products is different from other vNIC solutions. The difference is the fact that it carves up dedicated pipes between the adapter and the switch. This solution is built on industry standards that provide maximum performance in both directions, while allowing for pipes to be allocated at any speed 1 - 10 GbE.

Using a single 10 GbE dual-port adapter and creating virtual pipes to lessen the number of upstream switch ports, helps to drive out costs and complexity in the IT infrastructure. This scenario is possible by requiring up to 75% fewer adapters, cables, and upstream switch ports. It is also important to note that this function can be used across multiple application

environments, not just virtualization. When compared to using multiple 1 GbE ports, clients were able to see the following advantages:

- ▶ Potential of over 40% acquisition cost savings
- ▶ Up to 75% reduction in power consumption
- ▶ Simpler management with less cabling and fewer components to manage
- ▶ Easy integration into existing client setups (virtual or non-virtual)

In the following sections, we briefly describe IBM System Networking rack switches that participate in the IBM Virtual Fabric solution. For more information about this offering, see this website:

<http://www.ibm.com/systems/x/options/networking/virtualfabric/>

IBM System Networking RackSwitch G8124

The *IBM System Networking RackSwitch G8124* is a 10 GbE switch that is designed for the data centers, which provides a virtual, cooler, and easier network solution. The G8124 offers twenty-four 10 GbE ports in a high density, 1U footprint. Designed with top performance in mind, the RackSwitch G8124 provides line-rate, high-bandwidth switching, filtering, and traffic queuing. These features without delaying data and with large data-center grade buffers to keep traffic moving.

The G8124 is virtual, which provides rack-level virtualization of networking interfaces. IBM VMready® software enables movement of virtual machines, which provides matching movement of VLAN assignments, access control lists (ACLs), and other networking and security settings. VMready works with all leading VM providers, such as VMware, Citrix, Xen, and Microsoft. The G8124 also supports *Virtual Fabric*, which allows for the carving up of a physical network interface card (NIC) into 2 - 8 virtual NICs (vNICs). This action creates a virtual pipe between the adapter and the switch for improved performance, availability, and security, while it reduces cost and complexity.

Low latency that is offered by the G8124 makes it ideal for latency sensitive applications, such as high-performance computing clusters and financial applications. The G8124 also supports the newest protocols that include Converged Enhanced Ethernet (CEE) and Data Center Bridging for support of FCoE and can be used for NAS or iSCSI.

IBM System Networking RackSwitch G8124 is shown in Figure 12-32.



Figure 12-32 Front view of IBM System Networking RackSwitch G8124

The following list provides product highlights and benefits for the System Networking RackSwitch G8124:

- ▶ Optimal for high-performance computing and applications that require high bandwidth and low latency.
- ▶ All ports are non-blocking 10 Gigabit Ethernet with deterministic latency of 680 nanoseconds.
- ▶ VMready helps reduce configuration complexity and improves security levels in virtualized environments.
- ▶ Virtual Fabric capability allows for the carving up of a physical NIC into multiple virtual NICs (with Emulex adapters).

- ▶ Variable speed fans automatically adjust as needed, helping to reduce energy consumption.
- ▶ Seamless, standards-based integration into existing Cisco switches, and other networks help reduce downtime and the learning curve.

IBM System Networking RackSwitch G8264

IBM System Networking RackSwitch G8264 is a 10 and 40 GbE switch that is designed for the data center, which provides speed, intelligence, and interoperability on proven platforms.

The RackSwitch G8264 offers up to 64x10 GbE and up to four 40 GbE ports: 1.28 Tbps in a 1U footprint. Designed with top performance in mind, the RackSwitch G8264 provides line-rate, high-bandwidth switching, filtering, and traffic queuing without delaying data. Large data center grade buffers keep traffic moving. Redundant power and fans, along with numerous high availability features, enable the RackSwitch G8264 to be available for business sensitive traffic.

The low latency that is offered by the G8264 (Figure 12-33) makes it ideal for latency sensitive applications such as high performance computing clusters and financial applications. The G8264 supports the newest protocols that include: Data Center Bridging/Converged Enhanced Ethernet (DCB/CEE) for support of FCoE.

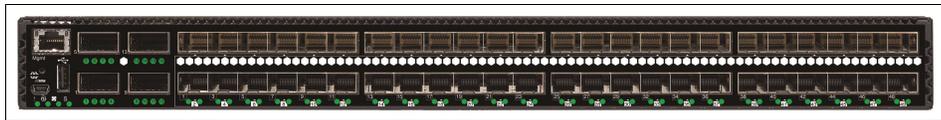


Figure 12-33 Front view of IBM System Networking RackSwitch G8264

The following list provides the product and business benefits of the System Networking RackSwitch G8264:

- ▶ Optimized for high-performance computing (HPC) and other applications that require high bandwidth and low latency.
- ▶ Software is based on Internet standards for optimal interoperability with Cisco or other vendor networks.
- ▶ VMready and Virtual Fabric for virtualized networks (with dedicated Emulex adapters for IBM System x).
- ▶ Forty-eight 10 GbE SFP+ ports, four 40 GbE QSFP+ ports, up to sixty-four 10 GbE SFP+ ports with optional breakout cables.
- ▶ Hot-swappable redundant power supplies and fans.

12.8 IBM Flex System networking

IBM Flex System™ offers intelligent, integrated, and flexible network architecture that can fit with your existing or future environment. These high performance Ethernet offerings that are coupled with on-demand scalability, offer an easy way to scale as IT requirements grow.

IBM Flex System Fabric is:

- ▶ Integrated: Helps manage discrete aspects of the data center as an integrated system through the built-in management appliance.
- ▶ Optimized: High performance scalable offerings with available 1 Gb, 10 Gb, and 40 Gb uplinks allow easy integration with an existing network. Provides simple and cost effective scalability for future growth.

- ▶ Automated: Automated provisioning and setup of both the physical and virtual network.

To meet today's complex and ever-changing business demands, you need a solid foundation of server, storage, networking, and management resources that is simple to deploy, yet can quickly, and automatically, adapt to changing conditions. You also need access to, and the ability to take advantage of, broad expertise and proven preferred practices in systems management, applications, hardware maintenance, and more. The *IBM PureFlex™ System* combines advanced IBM hardware and software along with patterns of expertise and integrates them into optimized solutions that are easy to deploy.

The network resources in an IBM PureFlex System are tightly integrated into the system to support virtualization and simple, integrated management. You can move from managing a physical network to managing a logical network in a virtualized environment, supporting business services instead of network components. With integrated management tools based on open standards, these resources are easy to provision and deploy so that you can reduce the cost of managing your virtual fabric. You have fewer elements to manage, but still get port and bandwidth flexibility with highly scalable switches. With scalable components, you can buy a base product and purchase and enable more ports without adding new hardware.

12.8.1 IBM Flex System Fabric EN4093 10Gb Scalable Switch

The *IBM Flex System Fabric EN4093 10Gb Scalable Switch* is a 10 Gb 64-port upgradeable midrange to high-end switch module, offering Layer 2/3 switching that is designed to install within the I/O module bays of the Enterprise Chassis. The switch has the following features:

- ▶ Up to 42 internal 10 Gb ports
- ▶ Up to 14 external 10 Gb uplink ports (SFP+ connectors)
- ▶ Up to two external 40 Gb uplink ports (QSFP+ connectors)

The switch is considered suited for clients with the following needs:

- ▶ Build a 10 Gb infrastructure
- ▶ Implement a virtualized environment
- ▶ Require investment protection for 40 Gb uplinks
- ▶ Want to reduce TCO, improve performance, while maintaining high levels of availability and security
- ▶ Want to avoid oversubscription (Traffic from multiple internal ports attempting to pass through a lower quantity of external ports, leading to congestion and performance affect)

The EN4093 10Gb Scalable Switch is shown in Figure 12-34.



Figure 12-34 IBM Flex System Fabric EN4093 10Gb Scalable Switch

As listed in Table 12-1, the switch is initially licensed with 14 10 Gb internal ports that are enabled and ten 10 Gb external uplink ports enabled. Further ports can be enabled, including the two 40 Gb external uplink ports with the Upgrade 1 and Upgrade 2 license options. Upgrade 1 must be applied before Upgrade 2 can be applied.

Table 12-1 IBM Flex System Fabric EN4093 10Gb Scalable Switch part numbers and port upgrades

Part number	Feature code ^a	Product description	Total ports enabled		
			Internal	10 Gb uplink	40 Gb uplink
49Y4270	A0TB / 3593	IBM Flex System Fabric EN4093 10Gb Scalable Switch <ul style="list-style-type: none"> ▶ 10x external 10 Gb uplinks ▶ 14x internal 10 Gb ports 	14	10	0
49Y4798	A1EL / 3596	IBM Flex System Fabric EN4093 10Gb Scalable Switch (Upgrade 1) <ul style="list-style-type: none"> ▶ Adds 2x external 40 Gb uplinks ▶ Adds 14x internal 10 Gb ports 	28	10	2
88Y6037	A1EM / 3597	IBM Flex System Fabric EN4093 10Gb Scalable Switch (Upgrade 2) (requires Upgrade 1): <ul style="list-style-type: none"> ▶ Adds 4x external 10 Gb uplinks ▶ Add 14x internal 10 Gb ports 	42	14	2

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

The key components on the front of the switch are shown in Figure 12-35.

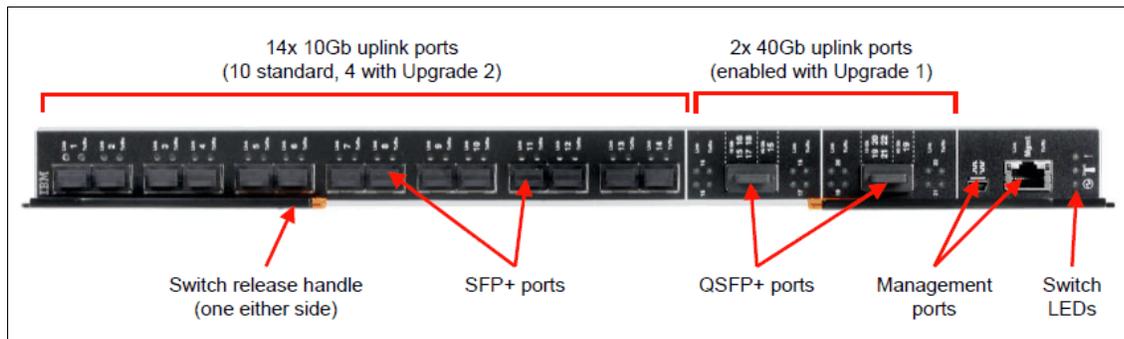


Figure 12-35 IBM Flex System Fabric EN4093 10Gb Scalable Switch

Each upgrade license enables more internal ports. To take full advantage of those ports, each compute node needs the appropriate I/O adapter installed:

- ▶ The base switch requires a two-port Ethernet adapter (one port of the adapter goes to each of two switches).
- ▶ Upgrade 1 requires a four-port Ethernet adapter (two ports of the adapter to each switch).
- ▶ Upgrade 2 requires a six-port Ethernet adapter (three ports to each switch).

Upgrade 2: Adding Upgrade 2 enables an additional 14 internal ports, delivering the ability to have 42 internal ports: three ports that are connected to each of the 14 compute nodes in the chassis. To take full advantage of all 42 internal ports, a 6-port adapter is required, but this type of adapter is currently not available.

Upgrade 2 still provides a benefit even with a 4-port adapter because this upgrade enables four extra external 10 Gb uplinks as well.

The rear of the switch has 14 SFP+ module ports and two QSFP+ module ports. The QSFP+ ports can be used to provide either two 40 Gb uplinks or eight 10 Gb ports, using one of the supported QSFP+ to 4x 10 Gb SFP+ cables that are listed in Table 12-2. This cable splits a single 40 Gb QSFP+ port into 4 SFP+ 10 Gb ports.

For management of the switch, there is a mini USB port and also an Ethernet management port provided.

The supported SFP+ and QSFP+ modules and cables for the switch are listed in Table 12-2.

Table 12-2 Supported SFP+ modules and cables

Part number	Feature code ^a	Description
Serial console cables		
90Y9338	A2RR / None	IBM Flex System Management Serial Access Cable Kit
SFP transceivers - 1 GbE		
81Y1618	3268 / EB29	IBM SFP RJ-45 Transceiver (does not support 10/100 Mbps)
81Y1622	3269 / EB2A	IBM SFP SX Transceiver
90Y9424	A1PN / None	IBM SFP LX Transceiver
SFP+ transceivers - 10 GbE		
46C3447	5053 / None	IBM SFP+ SR Transceiver
90Y9412	A1PM / None	IBM SFP+ LR Transceiver
44W4408	4942 / 3382	10GBase-SR SFP+ (MMFiber) transceiver
SFP+ Direct-attach copper (DAC) cables - 10 GbE		
90Y9427	A1PH / ECB4	1 m IBM Passive DAC SFP+
90Y9430	A1PJ / ECB5	3 m IBM Passive DAC SFP+
90Y9433	A1PK / None	5 m IBM Passive DAC SFP+
QSFP+ transceiver and cables - 40 GbE		
49Y7884	A1DR / EB27	IBM QSFP+ 40GBASE-SR Transceiver (Requires either cable 90Y3519 or cable 90Y3521)
90Y3519	A1MM / None	10 m IBM MTP Fiber Optical Cable (requires transceiver 49Y7884)
90Y3521	A1MN / None	30 m IBM MTP Fiber Optical Cable (requires transceiver 49Y7884)
QSFP+ breakout cables - 40 GbE to 4x10 GbE		
49Y7886	A1DL / EB24	1 m 40 Gb QSFP+ to 4 x 10Gb SFP+ Cable
49Y7887	A1DM / EB25	3 m 40 Gb QSFP+ to 4 x 10Gb SFP+ Cable

Part number	Feature code ^a	Description
49Y7888	A1DN / EB26	5 m 40 Gb QSFP+ to 4 x 10Gb SFP+ Cable
QSFP+ Direct-attach copper (DAC) cables - 40 GbE		
49Y7890	A1DP / None	1 m QSFP+ to QSFP+ DAC
49Y7891	A1DQ / None	3 m QSFP+ to QSFP+ DAC

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

The EN4093 10Gb Scalable Switch has the following features and specifications:

- ▶ Internal ports
 - Forty-two internal full-duplex 10 Gigabit ports (14 ports are enabled by default. Optional FoD licenses are required to activate the remaining 28 ports).
 - Two internal full-duplex 1 GbE ports connected to the chassis management module.
- ▶ External ports
 - Fourteen ports for 1 Gb or 10 Gb Ethernet SFP+ transceivers (support for 1000BASE-SX, 1000BASE-LX, 1000BASE-T, 10GBASE-SR, or 10GBASE-LR) or SFP+ copper direct-attach cables (DAC). Ten ports are enabled by default. An optional FoD license is required to activate the remaining four ports. SFP+ modules and DAC cables are not included and must be purchased separately.
 - Two ports for 40 Gb Ethernet QSFP+ transceivers or QSFP+ DACs (ports are disabled by default. An optional FoD license is required to activate them). QSFP+ modules and DAC cables are not included and must be purchased separately.
 - One RS-232 serial port (mini-USB connector) that provides an additional means to configure the switch module.
- ▶ Scalability and performance
 - 40 Gb Ethernet ports for extreme uplink bandwidth and performance
 - Fixed-speed external 10 Gb Ethernet ports to use 10 Gb core infrastructure
 - Autosensing 10/1000/1000 external Gigabit Ethernet ports for bandwidth optimization
 - Non-blocking architecture with wire-speed forwarding of traffic and aggregated throughput of 1.28 Tbps
 - Media access control (MAC) address learning: automatic update, support of up to 128,000 MAC addresses
 - Up to 128 IP interfaces per switch
 - Static and LACP (IEEE 802.3ad) link aggregation, up to 220 Gb of total uplink bandwidth per switch, up to 64 trunk groups, up to 16 ports per group
 - Support for jumbo frames (up to 9,216 bytes)
 - Broadcast/multicast storm control
 - IGMP snooping to limit flooding of IP multicast traffic
 - IGMP filtering to control multicast traffic for hosts that participate in multicast groups
 - Configurable traffic distribution schemes over trunk links that are based on source and destination IP or MAC addresses or both
 - Fast port forwarding and fast uplink convergence for rapid STP convergence

- ▶ Availability and redundancy
 - Virtual Router Redundancy Protocol (VRRP) for Layer 3 router redundancy
 - IEEE 802.1D STP for providing L2 redundancy
 - IEEE 802.1s Multiple STP (MSTP) for topology optimization, up to 32 STP instances are supported by single switch
 - IEEE 802.1w Rapid STP (RSTP) provides rapid STP convergence for critical delay-sensitive traffic like voice or video
 - Per-VLAN Rapid STP (PVRST) enhancements
 - Layer 2 Trunk Failover to support active/standby configurations of network adapter teaming on compute nodes
 - Hot Links provides basic link redundancy with fast recovery for network topologies that require Spanning Tree to be turned off
- ▶ VLAN support
 - Up to 1024 VLANs supported per switch, with VLAN numbers that range 1 - 4095 (4095 is used for the management module connection only)
 - 802.1Q VLAN tagging support on all ports
 - Private VLANs
- ▶ Security
 - VLAN-based, MAC-based, and IP-based ACLs
 - 802.1x port-based authentication
 - Multiple user IDs and passwords
 - User access control
 - Radius, TACACS+ and LDAP authentication and authorization
- ▶ Quality of service (QoS)
 - Support for IEEE 802.1p, IP ToS/DSCP, and ACL-based (MAC/IP source and destination addresses, VLANs) traffic classification and processing
 - Traffic shaping and remarking based on defined policies
 - Eight Weighted Round Robin (WRR) priority queues per port for processing qualified traffic
- ▶ IP v4 Layer 3 functions
 - Host management
 - IP forwarding
 - IP filtering with ACLs, up to 896 ACLs supported
 - VRRP for router redundancy
 - Support for up to 128 static routes
 - Routing protocol support (RIP v1, RIP v2, OSPF v2, BGP-4), up to 2048 entries in a routing table
 - Support for DHCP Relay
 - Support for IGMP snooping and IGMP relay
 - Support for Protocol Independent Multicast (PIM) in Sparse Mode (PIM-SM) and Dense Mode (PIM-DM).
- ▶ IP v6 Layer 3 functions
 - IPv6 host management (except default switch management IP address)

- IPv6 forwarding
- Up to 128 static routes
- Support for OSPF v3 routing protocol
- IPv6 filtering with ACLs
- ▶ Virtualization
 - Virtual Fabric with vNIC (virtual NICs)
 - 802.1Qbg Edge Virtual Bridging (EVB)
 - VMready
- ▶ Converged Enhanced Ethernet
 - Priority-Based Flow Control (PFC) (IEEE 802.1Qbb) extends 802.3x standard flow control to allow the switch to pause traffic that is based on the 802.1p priority value in the VLAN tag of each packet.
 - Enhanced Transmission Selection (ETS) (IEEE 802.1Qaz) provides a method for allocating link bandwidth that is based on the 802.1p priority value in the VLAN tag of each packet.
 - Data Center Bridging Capability Exchange Protocol (DCBX) (IEEE 802.1AB) allows neighboring network devices to exchange information about their capabilities.
- ▶ Manageability
 - Simple Network Management Protocol (SNMP V1, V2, and V3)
 - HTTP browser GUI
 - Telnet interface for CLI
 - SSH
 - Serial interface for CLI
 - Scriptable CLI
 - Firmware image update (TFTP and FTP)
 - Network Time Protocol (NTP) for switch clock synchronization
- ▶ Monitoring
 - Switch LEDs for external port status and switch module status indication
 - Remote Monitoring (RMON) agent to collect statistics and proactively monitor switch performance
 - Port mirroring for analyzing network traffic passing through the switch
 - Change tracking and remote logging with syslog feature
 - Support for sFLOW agent for monitoring traffic in data networks (separate sFLOW analyzer that is required elsewhere)
 - POST diagnostics

For more information, see the IBM Redbooks Product Guide: *IBM Flex System Fabric EN4093 10Gb Scalable Switch*, available at this website:

<http://www.redbooks.ibm.com/abstracts/tips0864.html?Open>

12.8.2 IBM Flex System EN4091 10Gb Ethernet Pass-thru

The *EN4091 10Gb Ethernet Pass-thru* module offers a one-for-one connection between a single node bay and an I/O module uplink. It has no management interface and can support both 1 Gb and 10 Gb dual-port adapters that are installed in the compute nodes. If quad-port adapters are installed in the compute nodes, only the first two ports will have access to the ports of the pass-thru module.

The necessary 1 GbE or 10 GbE module (SFP, SFP+ or DAC) must also be installed in the external ports of the pass-thru. This installation is required to support the wanted speed (1 Gb or 10 Gb) and medium (fiber optic or copper) for the adapter ports on the compute nodes.

The IBM Flex System EN4091 10Gb Ethernet Pass-thru is shown in Figure 12-36.



Figure 12-36 IBM Flex System EN4091 10Gb Ethernet Pass-thru

The ordering part number and feature codes are listed in Table 12-3.

Table 12-3 EN4091 10Gb Ethernet Pass-thru part number and feature codes

Part number	Feature code ^a	Product name
88Y6043	A1QV / 3700	IBM Flex System EN4091 10Gb Ethernet Pass-thru

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

The EN4091 10Gb Ethernet Pass-thru has the following specifications:

- ▶ **Internal ports**
Fourteen internal full-duplex Ethernet ports that can operate at 1 Gb or 10 Gb speeds.
- ▶ **External ports**
Fourteen ports for 1 Gb or 10 Gb Ethernet SFP+ transceivers (support for 1000BASE-SX, 1000BASE-LX, 1000BASE-T, 10GBASE-SR, or 10GBASE-LR) or SFP+ copper direct-attach cables (DAC). SFP+ modules and DAC cables are not included and must be purchased separately.
- ▶ Unmanaged device that has no internal Ethernet management port; however, it is able to provide its vital product data (VPD) to the secure management network in the Chassis Management Module.
- ▶ Supports 10 Gb Ethernet signaling for CEE, FCoE, and other Ethernet based transport protocols.
- ▶ Allows direct connection from the 10 Gb Ethernet adapters that are installed in compute nodes in a chassis to an externally located top-of-rack (TOR) switch or other external device.

Four-port adapters: The EN4091 10Gb Ethernet Pass-thru has only 14 internal ports. As a result, only two ports on each compute node are enabled, one for each of two pass-thru modules that is installed in the chassis. If four-port adapters are installed in the compute nodes, ports 3 and 4 on those adapters are not enabled.

There are standard 3 I/O module status LEDs. Each port has link and activity LEDs.

Table 12-4 lists the supported transceivers and direct-attach copper (DAC) cables.

Table 12-4 IBM Flex System EN4091 10Gb Ethernet Pass-thru part numbers and feature codes

Part number	Feature codes ^a	Description
SFP+ transceivers - 10 GbE		
44W4408	4942 / 3282	10 GbE 850 nm Fiber SFP+ Transceiver (SR)
46C3447	5053 / None	IBM SFP+ SR Transceiver
90Y9412	A1PM / None	IBM SFP+ LR Transceiver
SFP transceivers - 1 GbE		
81Y1622	3269 / EB2A	IBM SFP SX Transceiver
81Y1618	3268 / EB29	IBM SFP RJ45 Transceiver
90Y9424	A1PN / None	IBM SFP LX Transceiver
Direct-attach copper (DAC) cables		
81Y8295	A18M / EN01	1 m 10GE Twinax Act Copper SFP+ DAC (active)
81Y8296	A18N / EN02	3 m 10GE Twinax Act Copper SFP+ DAC (active)
81Y8297	A18P / EN03	5 m 10GE Twinax Act Copper SFP+ DAC (active)
95Y0323	A25A / None	1 m IBM Active DAC SFP+ Cable
95Y0326	A25B / None	3 m IBM Active DAC SFP+ Cable
95Y0329	A25C / None	5 m IBM Active DAC SFP+ Cable

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

For more information, see the IBM Redbooks Product Guide: *IBM Flex System EN4091 10Gb Ethernet Pass-thru*, available from this website:

<http://www.redbooks.ibm.com/abstracts/tips0865.html?Open>

12.8.3 IBM Flex System EN2092 1Gb Ethernet Scalable Switch

The *EN2092 1Gb Ethernet Switch* provides support for L2/L3 switching and routing. The switch provides the following features:

- ▶ Up to 28 internal 1 Gb ports
- ▶ Up to 20 external 1 Gb ports (RJ45 connectors)
- ▶ Up to four external 10 Gb uplink ports (SFP+ connectors)

The switch is shown in Figure 12-37.



Figure 12-37 IBM Flex System EN2092 1Gb Ethernet Scalable Switch

As listed in Table 12-5, the switch comes standard with 14 internal and 10 external Gigabit Ethernet ports enabled. Further ports can be enabled, including the four external 10 Gb uplink ports. Upgrade 1 and the 10 Gb Uplinks upgrade can be applied in either order.

Table 12-5 IBM Flex System EN2092 1Gb Ethernet Scalable Switch part numbers and port upgrades

Part number	Feature code ^a	Product description
49Y4294	A0TF / 3598	IBM Flex System EN2092 1Gb Ethernet Scalable Switch <ul style="list-style-type: none"> ▶ 14 internal 1 Gb ports ▶ 10 external 1 Gb ports
90Y3562	A1QW / 3594	IBM Flex System EN2092 1Gb Ethernet Scalable Switch (Upgrade 1) <ul style="list-style-type: none"> ▶ Adds 14 internal 1 Gb ports ▶ Adds 10 external 1 Gb ports
49Y4298	A1EN / 3599	IBM Flex System EN2092 1Gb Ethernet Scalable Switch (10 Gb Uplinks) <ul style="list-style-type: none"> ▶ Adds four external 10 Gb uplinks

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

The key components on the front of the switch are shown in Figure 12-38.

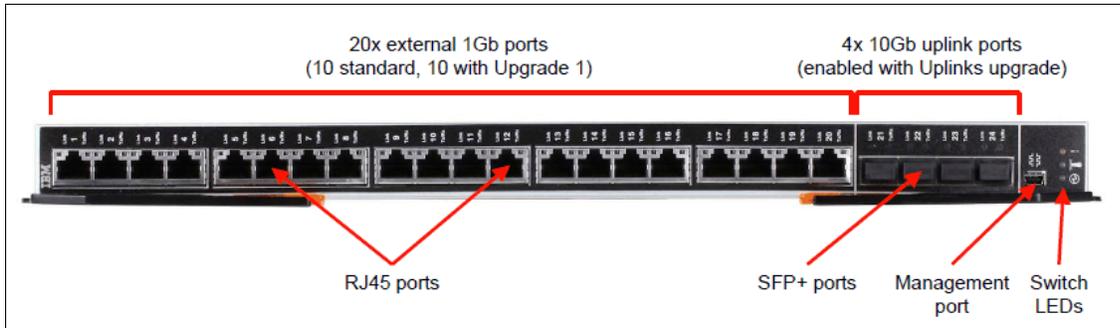


Figure 12-38 IBM Flex System EN2092 1Gb Ethernet Scalable Switch

The standard switch has 14 internal ports and the Upgrade 1 license enables 14 more internal ports. To take full advantage of those ports, each compute node needs the appropriate I/O adapter installed:

- ▶ The base switch requires a two-port Ethernet adapter that is installed in each compute node (one port of the adapter goes to each of two switches)
- ▶ Upgrade 1 requires a four-port Ethernet adapter that is installed in each compute node (two ports of the adapter to each switch)

The standard has 10 external ports enabled. Additional external ports are enabled with license upgrades and can be installed in either order:

- Upgrade 1 enables 10 more ports for a total of 20 ports
- Uplinks Upgrade enables the four 10 Gb SFP+ ports

This switch is considered ideal for clients who have the following needs:

- ▶ Still use 1 Gb as their networking infrastructure
- ▶ Deploy virtualization and require multiple 1 Gb ports
- ▶ Want investment protection for 10 Gb uplinks
- ▶ Looking to reduce TCO, improve performance, while maintaining high levels of availability and security
- ▶ Looking to avoid oversubscription (multiple internal ports attempting to pass through a lower quantity of external ports, leading to congestion or performance affect).

The switch has three switch status LEDs and one mini-USB serial port connector for console management.

Uplink Ports 1 - 20 are RJ45 and the 4 x 10Gb uplink ports are SFP+. The switch supports either SFP+ modules or DAC cables. The supported SFP+ modules and DAC cables for the switch are listed in Table 12-6.

Table 12-6 IBM Flex System EN2092 1Gb Ethernet Scalable Switch SFP+ and DAC cables

Part number	Feature code ^a	Description
SFP transceivers		
81Y1622	3269 / EB2A	IBM SFP SX Transceiver
81Y1618	3268 / EB29	IBM SFP RJ45 Transceiver
90Y9424	A1PN / None	IBM SFP LX Transceiver
SFP+ transceivers		
44W4408	4942 / 3282	10 GbE 850 nm Fiber SFP+ Transceiver (SR)
46C3447	5053 / None	IBM SFP+ SR Transceiver
90Y9412	A1PM / None	IBM SFP+ LR Transceiver
DAC cables		
90Y9427	A1PH / None	1 m IBM Passive DAC SFP+
90Y9430	A1PJ / ECB5	3 m IBM Passive DAC SFP+
90Y9433	A1PK / None	5 m IBM Passive DAC SFP+

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

The EN2092 1 Gb Ethernet Scalable Switch has the following features and specifications:

- ▶ Internal ports
 - Twenty-eight internal full-duplex Gigabit ports (14 ports are enabled by default. An optional FoD license is required to activate another 14 ports.)
 - Two internal full-duplex 1 GbE ports that are connected to the chassis management module
- ▶ External ports
 - Four ports for 1 Gb or 10 Gb Ethernet SFP+ transceivers (support for 1000BASE-SX, 1000BASE-LX, 1000BASE-T, 10GBASE-SR, or 10GBASE-LR) or SFP+ copper direct-attach cables (DAC). These ports are disabled by default. An optional FoD license is required to activate them. SFP+ modules are not included and must be purchased separately.
 - Twenty external 10/100/1000 1000BASE-T Gigabit Ethernet ports with RJ-45 connectors (10 ports are enabled by default. An optional FoD license is required to activate another 10 ports).
 - One RS-232 serial port (mini-USB connector) that provides an additional means to configure the switch module.
- ▶ Scalability and performance
 - Fixed-speed external 10 Gb Ethernet ports for maximum uplink bandwidth
 - Autosensing 10/1000/1000 external Gigabit Ethernet ports for bandwidth optimization
 - Non-blocking architecture with wire-speed forwarding of traffic
 - Media access control (MAC) address learning: automatic update, support of up to 32,000 MAC addresses
 - Up to 128 IP interfaces per switch
 - Static and LACP (IEEE 802.3ad) link aggregation, up to 60 Gb of total uplink bandwidth per switch, up to 64 trunk groups, up to 16 ports per group
 - Support for jumbo frames (up to 9,216 bytes)
 - Broadcast/multicast storm control
 - IGMP snooping for limit flooding of IP multicast traffic
 - IGMP filtering to control multicast traffic for hosts participating in multicast groups
 - Configurable traffic distribution schemes over trunk links that are based on source/destination IP or MAC addresses or both
 - Fast port forwarding and fast uplink convergence for rapid STP convergence
- ▶ Availability and redundancy
 - Virtual Router Redundancy Protocol (VRRP) for Layer 3 router redundancy
 - IEEE 802.1D STP for providing L2 redundancy
 - IEEE 802.1s Multiple STP (MSTP) for topology optimization, up to 32 STP instances that are supported by a single switch
 - IEEE 802.1w Rapid STP (RSTP) (provides rapid STP convergence for critical delay-sensitive traffic like voice or video)
 - Per-VLAN Rapid STP (PVRST) enhancements
 - Layer 2 Trunk Failover to support active/standby configurations of network adapter teaming on compute nodes

- Hot Links provides basic link redundancy with fast recovery for network topologies that require Spanning Tree to be turned off
- ▶ VLAN support
 - Up to 1024 VLANs supported per switch, with VLAN numbers that range 1 - 4095 (4095 is used for the management module connection only)
 - 802.1Q VLAN tagging support on all ports
 - Private VLANs
- ▶ Security
 - VLAN-based, MAC-based, and IP-based ACLs
 - 802.1x port-based authentication
 - Multiple user IDs and passwords
 - User access control
 - Radius, TACACS+ and LDAP authentication and authorization
- ▶ Quality of service (QoS)
 - Support for IEEE 802.1p, IP ToS/DSCP, and ACL-based (MAC/IP source and destination addresses, VLANs) traffic classification and processing
 - Traffic shaping and remarking based on defined policies
 - Eight Weighted Round Robin (WRR) priority queues per port for processing qualified traffic
- ▶ IP v4 Layer 3 functions
 - Host management
 - IP forwarding
 - IP filtering with ACLs, up to 896 ACLs supported
 - VRRP for router redundancy
 - Support for up to 128 static routes
 - Routing protocol support (RIP v1, RIP v2, OSPF v2, BGP-4), up to 2048 entries in a routing table
 - Support for DHCP Relay
 - Support for IGMP snooping and IGMP relay
 - Support for Protocol Independent Multicast (PIM) in Sparse Mode (PIM-SM) and Dense Mode (PIM-DM)
- ▶ IP v6 Layer 3 functions
 - IPv6 host management (except default switch management IP address)
 - IPv6 forwarding
 - Up to 128 static routes
 - Support for OSPF v3 routing protocol
 - IPv6 filtering with ACLs
- ▶ Virtualization
 - VMready
- ▶ Manageability
 - Simple Network Management Protocol (SNMP V1, V2, and V3)
 - HTTP browser GUI
 - Telnet interface for CLI

- SSH
- Serial interface for CLI
- Scriptable CLI
- Firmware image update (TFTP and FTP)
- Network Time Protocol (NTP) for switch clock synchronization
- ▶ Monitoring
 - Switch LEDs for external port status and switch module status indication
 - Remote Monitoring (RMON) agent to collect statistics and proactively monitor switch performance
 - Port mirroring for analyzing network traffic passing through the switch
 - Change tracking and remote logging with the syslog feature
 - Support for the sFLOW agent for monitoring traffic in data networks (separate sFLOW analyzer that is required elsewhere)
 - POST diagnostics

For more information, see the IBM Redbooks Product Guide: *IBM Flex System EN2092 1Gb Ethernet Scalable Switch*, available at this website:

<http://www.redbooks.ibm.com/abstracts/tips0861.html?Open>

12.8.4 IBM Flex System FC5022 16Gb SAN Scalable Switch

The *IBM Flex System FC5022 16Gb SAN Scalable Switch* is a high-density, 48-port 16 Gbps Fibre Channel switch that is used in the Enterprise Chassis. The switch provides 28 internal ports to compute nodes by way of the midplane, and 20 external SFP+ ports. These SAN switch modules deliver an embedded option for IBM Flex System users that are deploying storage area networks in their enterprise. They offer end-to-end 16 Gb and 8 Gb connectivity.

The N_Port Virtualization mode streamlines the infrastructure by reducing the number of domains to manage while it enables the ability to add or move servers without an affect to the SAN. Monitoring is simplified by an integrated management appliance, or clients that use end-to-end Brocade SAN can use the Brocade management tools.

Figure 12-39 shows the IBM Flex System FC5022 16Gb SAN Scalable Switch.

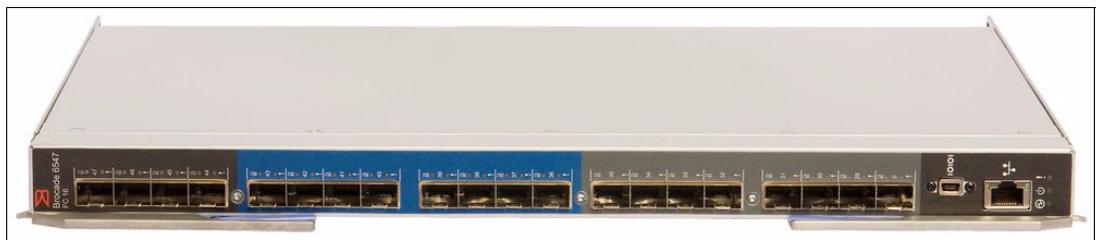


Figure 12-39 IBM Flex System FC5022 16Gb SAN Scalable Switch

Two versions are available as listed in Table 12-7 on page 304: a 12-port switch module and a 24-port switch with the Enterprise Switch Bundle (ESB) software. The port count can be applied to internal or external ports by using a feature that is called *Dynamic Ports on Demand (DPOD)*.

Table 12-7 IBM Flex System FC5022 16Gb SAN Scalable Switch part numbers

Part number	Feature codes ^a	Description	Ports enabled
88Y6374	A1EH / 3770	IBM Flex System FC5022 16Gb SAN Scalable Switch	12
90Y9356	A1EJ / 3771	IBM Flex System FC5022 24-port 16Gb ESB SAN Scalable Switch	24

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

Table 12-8 provides a feature comparison by model for FC5022 switches.

Table 12-8 Feature comparison by model

Feature	FC5022 16Gb ESB Switch (90Y9356)	FC5022 16Gb SAN Scalable Switch (88Y6374)
Number of active ports	24	12
Full fabric	Included	Included
Access Gateway	Included	Included
Advanced zoning	Included	Included
Enhanced Group Management	Included	Included
ISL Trunking	Included	Not available
Adaptive Networking	Included	Not available
Advanced Performance Monitoring	Included	Not available
Fabric Watch	Included	Not available
Extended Fabrics	Included	Not available
Server Application Optimization	Included	Not available

With Dynamic Ports on Demand (DPOD), ports are licensed as they come online. With the FC5022 16Gb SAN Scalable Switch, the first 12 ports that are reporting (on a first-come, first-served basis) on boot-up are assigned licenses. These 12 ports might be any combination of external or internal Fibre Channel (FC) ports. After all licenses are assigned, you can manually move those licenses from one port to another. Because this move is dynamic, no defined ports are reserved except ports 0 and 29. The FC5022 16Gb ESB Switch has the same behavior, the only difference is the number of ports.

The part number for the switch includes the following items:

- ▶ One IBM Flex System FC5022 16Gb SAN Scalable Switch or IBM Flex System FC5022 24-port 16Gb ESB SAN Scalable Switch
- ▶ Important Notices flyer
- ▶ Warranty flyer
- ▶ Documentation CD-ROM

The switch does not include a serial management cable. However, IBM Flex System Management Serial Access Cable, 90Y9338, is supported and contains two cables: a mini-USB-to-RJ45 serial cable, and a mini-USB-to-DB9 serial cable. Either of these cables can be used to connect to the switch locally for configuration tasks and firmware updates.

Transceivers

The switch comes without SFP+, they must be ordered separately to provide outside connectivity. Table 12-9 lists supported SFP+ options.

Table 12-9 Supported SFP+ transceivers

Part number	Feature code ^a	Description
88Y6416	5084 / 5370	Brocade 8 Gb SFP+ SW Optical Transceiver
88Y6393	A22R / 5371	Brocade 16 Gb SFP+ Optical Transceiver

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

Benefits

The switches offer the following key benefits:

- ▶ **Exceptional price and performance for growing SAN workloads**
 The FC5022 16Gb SAN Scalable Switch delivers exceptional price and performance for growing SAN workloads through a combination of market-leading 1,600 MB/sec throughput per port and an affordable high-density form factor. The 48 FC ports produce an aggregate 768 Gbps full-duplex throughput, plus any external eight ports can be trunked for 128 Gbps ISLs. The 16 Gbps port technology dramatically reduces the number of ports and associated optics and cabling that is required through 8/4 Gbps consolidation. Therefore, the cost savings and simplification benefits are substantial.
- ▶ **Accelerating fabric deployment and serviceability with diagnostic ports**
 Diagnostic Ports (D_Ports) are a new port type that is supported by the FC5022 16Gb SAN Scalable Switch. This device enables administrators to quickly identify and isolate 16 Gbps optics, as well as port and cable problems, reducing fabric deployment and diagnostic times. If the optical media is found to be the source of the problem, it can be transparently replaced because 16 Gbps optics are hot-pluggable.
- ▶ **A building block for virtualized, private cloud storage**
 The FC5022 16Gb SAN Scalable Switch supports multi-tenancy in cloud environments through VM-aware end-to-end visibility and monitoring, QoS, and fabric-based advanced zoning features. The FC5022 16Gb SAN Scalable Switch enables secure distance extension to virtual private or hybrid clouds with dark fiber support, and in-flight encryption and data compression. Internal fault-tolerant and enterprise-class RAS features help minimize downtime to support mission-critical cloud environments.
- ▶ **Simplified and optimized interconnect with Brocade Access Gateway**
 The FC5022 16Gb SAN Scalable Switch can be deployed as a full-fabric switch or as a Brocade Access Gateway, which simplifies fabric topologies and heterogeneous fabric connectivity. Access Gateway mode uses N_Port ID Virtualization (NPIV) switch standards to present physical and virtual servers directly to the core of SAN fabrics. This makes it transparent to the SAN fabric, greatly reducing management of the network edge.
- ▶ **Maximizing investments**
 To help optimize technology investments, IBM offers a single point of serviceability that is backed by industry-renowned education, support, and training. In addition, the

IBM 16/8 Gbps SAN Scalable Switch is in the IBM ServerProven® program, enabling compatibility among various IBM and partner products. IBM recognizes that clients deserve the most innovative, expert integrated systems solutions.

Features and specifications

The FC5022 16Gb SAN Scalable Switches have the following features and specifications:

- ▶ Internal ports
 - Twenty-eight internal full-duplex 16 Gb FC ports (up to 14 internal ports can be activated with Port-on-Demand feature; remaining ports are reserved for future use).
 - Internal ports operate as F_ports (fabric ports) in native mode or in access gateway mode.
 - Two internal full-duplex 1 GbE ports that are connected to the chassis management module.
- ▶ External ports
 - Twenty external ports for 16 Gb SFP+ or 8 Gb SFP+ transceivers supporting 4 Gb, 8 Gb, and 16 Gb port speeds (SFP+ modules are not included and must be purchased separately. See Table 12-2 on page 293). Ports are activated with Port-on-Demand feature.
 - External ports can operate as F_ports (fabric ports), FL_ports (fabric loop ports), or E_ports (expansion ports) in native mode or as N_Ports (Node Ports) in access gateway mode.
 - One external 1 GbE port (1000BASE-T) with RJ-45 connector for switch configuration and management.
 - One RS-232 serial port (mini-USB connector) that provides an additional means to configure the switch module.
- ▶ Access gateway mode (N_Port ID Virtualization - NPIV) support
- ▶ Power-on self-test diagnostics and status reporting
- ▶ ISL Trunking (licensable), which allows up to eight ports (at 16, 8, or 4 Gbps speeds) to combine to form a single, logical ISL with a speed of up to 128 Gbps (256 Gbps full duplex). This configuration allows for optimal bandwidth utilization, automatic path failover, and load balancing.
- ▶ Brocade Fabric OS (FOS), which delivers distributed intelligence throughout the network and enables a wide range of value-added applications, such as Brocade Advanced Web Tools and Brocade Advanced Fabric Services (on certain models)
- ▶ Supports up to 768 Gbps I/O bandwidth
- ▶ 420 million frames switches per second, 0.7 microseconds latency
- ▶ 8,192 buffers for up to 3,750 km extended distance at 4 Gbps FC (Extended Fabrics license required)
- ▶ In-flight 64 Gbps Fibre Channel compression and decompression support on up to two external ports (no license required)
- ▶ In-flight 32 Gbps encryption and decryption on up to two external ports (no license required)
- ▶ 48 Virtual Channels (VCs) per port
- ▶ Port mirroring to monitor ingress or egress traffic from any port within the switch
- ▶ Two I2C connections able to interface to redundant management modules
- ▶ Hot pluggable: up to four hot pluggable switches per chassis

- ▶ Single fuse circuit
- ▶ Four temperature sensors
- ▶ Managed with Brocade Web Tools
- ▶ Supports a minimum of 128 domains in Native mode and Interoperability mode
- ▶ Nondisruptive code load in Native mode and Access Gateway mode
- ▶ 255 N_port logins per physical port
- ▶ D_port support on external ports
- ▶ Class 2 and Class 3 frames
- ▶ SNMP v1 and v3 support
- ▶ SSH v2 support
- ▶ SSL support
- ▶ NTP client support (NTP V3)
- ▶ FTP support for firmware upgrades
- ▶ SNMP/MIB monitoring functionality that is contained within the Ethernet Control MIB-II (RFC1213-MIB)
- ▶ End-to-end optics and link validation
- ▶ Sends switch events and syslogs to the Chassis Management Module (CMM)
- ▶ Traps identify cold start, warm start, link up/link down and authentication failure events
- ▶ Support for IPv4 and IPv6 on the management ports

The FC5022 16Gb SAN Scalable Switches come standard with the following software features:

- ▶ Brocade Full Fabric mode: Enables high performance 16 Gb or 8 Gb fabric switching
- ▶ Brocade Access Gateway mode: uses NPIV to connect to any fabric without adding switch domains to reduce management complexity
- ▶ Dynamic Path Selection: enables exchange-based load balancing across multiple Inter-Switch Links for superior performance
- ▶ Brocade Advanced Zoning: segments a SAN into virtual private SANs to increase security and availability
- ▶ Brocade Enhanced Group Management: enables centralized and simplified management of Brocade fabrics through IBM Network Advisor

Enterprise Switch Bundle (ESB) software licenses

The *IBM Flex System FC5022 24-port 16Gb ESB SAN Scalable Switch* includes a complete set of licensed features. These features maximize performance, ensure availability, and simplify management for the most demanding applications and expanding virtualization environments.

This switch comes with 24 port licenses that can be applied to either internal or external links on this switch.

This switch also includes the following Enterprise Switch Bundle (ESB) software licenses:

- ▶ Brocade Extended Fabrics
 - Provides up to 1000 km of switches fabric connectivity over long distances.

- ▶ Brocade ISL Trunking
Can aggregate multiple physical links into one logical link for enhanced network performance and fault tolerance.
- ▶ Brocade Advanced Performance Monitoring
Enables performance monitoring of networked storage resources. This license includes the TopTalkers feature.
- ▶ Brocade Fabric Watch
Monitors mission-critical switch operations. Fabric Watch now includes new Port Fencing capabilities.
- ▶ Adaptive Networking
Adaptive Networking provides a rich set of capabilities to the data center or virtual server environments. It ensures high priority connections to obtain the bandwidth necessary for optimum performance, even in congested environments. It optimizes data traffic movement within the fabric by Ingress Rate Limiting, QoS, and Traffic Isolation Zones.
- ▶ Server Application Optimization (SAO)
This license optimizes overall application performance for physical servers and virtual machines. SAO, when deployed with Brocade Fibre Channel HBAs, extends Brocade Virtual Channel (VC) technology from fabric to the server infrastructure. This configuration delivers application-level, fine-grain QoS management to the HBAs and related server applications.

Supported Fibre Channel standards

The switches support the following Fibre Channel standards:

- ▶ FC-AL-2 INCITS 332: 1999
- ▶ FC-GS-5 ANSI INCITS 427 (includes the following standard): FC-GS-4 ANSI INCITS 387: 2004
- ▶ FC-IFR INCITS 1745-D, revision 1.03 (under development)
- ▶ FC-SW-4 INCITS 418:2006 (includes the following standards):
 - FC-SW-3 INCITS 384: 2004
 - FC-VI INCITS 357: 2002
 - FC-TAPE INCITS TR-24: 1999
- ▶ FC-DA INCITS TR-36: 2004 (includes the following standards):
 - FC-FLA INCITS TR-20: 1998
 - FC-PLDA INCITS TR-19: 1998
 - FC-MI-2 ANSI/INCITS TR-39-2005
 - FC-PI INCITS 352: 2002
 - FC-PI-2 INCITS 404: 2005
- ▶ FC-PI-4 INCITS 1647-D, revision 7.1 (under development)
- ▶ FC-PI-5 INCITS 479: 2011
- ▶ FC-FS-2 ANSI/INCITS 424:2006 (includes the following standards):
 - FC-FS INCITS 373: 2003
 - FC-LS INCITS 433: 2007
- ▶ FC-BB-3 INCITS 414: 2006 (includes the following standards):
 - FC-BB-2 INCITS 372: 2003

- FC-SB-3 INCITS 374: 2003 (replaces FC-SB ANSI X3.271: 1996; FC-SB-2 INCITS 374: 2001)
- RFC 2625 IP and ARP Over FC
- RFC 2837 Fabric Element MIB
- MIB-FA INCITS TR-32: 2003
- FCP-2 INCITS 350: 2003 (replaces FCP ANSI X3.269: 1996)
- ▶ SNIA Storage Management Initiative Specification (SMI-S) Version 1.2 (includes the following standards):
 - SNIA Storage Management Initiative Specification (SMI-S) Version 1.03 ISO standard IS24775-2006. Replaces (ANSI INCITS 388: 2004)
 - SNIA Storage Management Initiative Specification (SMI-S) Version 1.1.0
 - SNIA Storage Management Initiative Specification (SMI-S) Version 1.2.0

For more information, see the IBM Redbooks Product Guide: *IBM Flex System FC5022 16Gb SAN Scalable Switch*, available at this website:

<http://www.redbooks.ibm.com/abstracts/tips0870.html?Open>

12.8.5 IBM Flex System FC3171 8Gb SAN Switch

This 8 Gb SAN switch from QLogic is a full-fabric Fibre Channel switch module that can be converted to a pass-thru module when configured in transparent mode.



Figure 12-40 IBM Flex System FC3171 8Gb SAN Switch

The I/O module has 14 internal ports and six external ports. All ports are licensed on the switch because there are no port licensing requirements. Ordering information is listed in Table 12-10.

Table 12-10 FC3171 8Gb SAN Switch

Part number	Feature code ^a	Product name
69Y1930	A0TD / 3595	IBM Flex System FC3171 8Gb SAN Switch

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

There are no SFPs supplied as standard, the SFP modules and cables that are listed in Table 12-11 on page 310 are supported.

Table 12-11 FC3171 8Gb SAN Switch supported SFP modules and cables

Part number	Feature codes ^a	Description
44X1964	5075 / 3286	IBM 8 Gb SFP+ SW Optical Transceiver
39R6475	4804 / 3238	4 Gb SFP Transceiver Option

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

It is possible to reconfigure the FC3171 8Gb SAN Switch to become a pass-thru module. This step is possible by using the switch graphical user interface (GUI) or command-line interface (CLI). The module can then be converted back to a full function SAN switch at some future date. The switch requires a reset when turning the transparent mode on or off.

The switch can be configured by either the command line, or via QuickTools:

- ▶ *Command Line.* Access the switch by the console port through the Chassis Management Module or through the Ethernet port. This method requires a basic understanding of the CLI commands.
- ▶ *QuickTools.* Requires a current version of the JRE on your workstation before pointing a web browser to the IP address of the switch. The IP address of the switch must be configured. QuickTools does not require a license and code is included.

On this switch when in full fabric mode, access to all of the Fibre Channel security features is provided. Security includes additional services that are available, such as Secure Socket Layer (SSL) and Secure Shell (SSH). In addition, RADIUS servers might be used for device and user authentication. When SSL/SSH is enabled, the security features are available to be configured. This allows the SAN administrator to configure which devices are allowed to log in to the Full Fabric Switch module, by creating security sets with security groups. These groups are configured on a per switch basis. The security features are not available when in the pass-thru mode.

The following list provides the FC3171 8Gb SAN Switch specifications and standards:

- ▶ Fibre Channel standards:
 - C-PH version 4.3
 - FC-PH-2
 - FC-PH-3
 - FC-AL version 4.5
 - FC-AL-2 Rev 7.0
 - FC-FLA
 - FC-GS-3
 - FC-FG
 - FC-PLDA
 - FC-Tape
 - FC-VI
 - FC-SW-2
 - Fibre Channel Element MIB RFC 2837
 - Fibre Alliance MIB version 4.0
- ▶ Fibre Channel protocols:
 - Fibre Channel service classes: class 2 and class 3
 - Operation modes: Fibre Channel class 2 and class 3, connectionless

- ▶ External port type:
 - Full fabric mode: Generic loop port (GL_port)
 - Transparent mode: Transparent fabric port (TF_port)
- ▶ Internal port type:
 - Full fabric mode: Fabric port (F_port)
 - Transparent mode: Transparent host port/NPIV mode (TH_port)
 - Support for up to 44 host NPIV logins
- ▶ Port characteristics:
 - External ports are automatically detected and self- configuring
 - Port LEDs illuminate at startup
 - Number of Fibre Channel ports: 6 external ports and 14 internal ports
 - Scalability: Up to 239 switches maximum depending on your configuration
 - Buffer credits: 16 buffer credits per port
 - Maximum frame size: 2148 bytes (2112-byte payload)
 - Standards-based FC FC-SW2 Interoperability
 - Support for up to a 255-to-1 port-mapping ratio
 - Media type: Small form-factor pluggable plus (SFP+) module
- ▶ 2 Gb specifications
 - 2 Gb fabric port speed: 1.0625 or 2.125 Gbps (gigabits per second)
 - 2 Gb fabric latency: Less than 0.4 msec
 - 2 Gb fabric aggregate bandwidth: 80 Gbps at full duplex
- ▶ 4 Gb specifications
 - 4 Gb switch speed: 4.250 Gbps
 - 4 Gb switch fabric point-to-point: 4 Gbps at full duplex
 - 4 Gb switch fabric aggregate bandwidth: 160 Gbps at full duplex
- ▶ 8 Gb specifications
 - 8 Gb switch speed: 8.5 Gbps
 - 8 Gb switch fabric point-to-point: 8 Gbps at full duplex
 - 8 Gb switch fabric aggregate bandwidth: 320 Gbps at full duplex
- ▶ Nonblocking architecture to prevent latency
- ▶ System processor: IBM PowerPC®

For more information, see the IBM Redbooks Product Guide: *IBM Flex System FC3171 8Gb SAN Switch*, available from this website:

<http://www.redbooks.ibm.com/abstracts/tips0866.html?Open>

12.8.6 IBM Flex System FC3171 8Gb SAN Pass-thru

The *IBM Flex System FC3171 8Gb SAN Pass-thru I/O* module is an 8 Gbps Fibre Channel pass-thru SAN module that has 14 internal ports and six external ports. It is shipped with all of the ports enabled.

Figure 12-41 shows the switch.



Figure 12-41 IBM Flex System FC3171 8Gb SAN Pass-thru

Ordering information is listed in Table 12-12.

Table 12-12 FC3171 8Gb SAN Pass-thru part number

Part number	Feature code ^a	Description
69Y1934	A0TJ / 3591	IBM Flex System FC3171 8Gb SAN Pass-thru

a. The first feature code listed is for configurations ordered through IBM System x sales channels. The second feature code is for configurations ordered through the IBM Power Systems channel.

Future requirements: If there is a potential future requirement to enable full fabric capability, then this switch should not be purchased and instead the FC3171 8Gb SAN Switch should be considered.

There are no SFPs supplied with the switch and must be ordered separately. Supported transceivers and fiber optic cables are listed in Table 12-13.

Table 12-13 FC3171 8Gb SAN Pass-thru supported modules and cables

Part Number	Feature Code	Description
44X1964	5075 / 3286	IBM 8 Gb SFP+ SW Optical Transceiver
39R6475	4804 / 3238	4 Gb SFP Transceiver Option

The FC3171 8Gb SAN Pass-thru can be configured using either command line or QuickTools:

- ▶ **Command Line.** Access the module by the console port through the Chassis Management Module or through the Ethernet Port. This method requires a basic understanding of the CLI commands.
- ▶ **QuickTools.** Requires a current version of the JRE on your workstation before pointing a web browser to the modules IP address. The IP address of the module must be configured. QuickTools does not require a license and code is included.

The pass-thru module supports the following standards:

- ▶ **Fibre Channel standards:**
 - C-PH version 4.3
 - FC-PH-2
 - FC-PH-3
 - FC-AL version 4.5
 - FC-AL-2 Rev 7.0
 - FC-FLA
 - FC-GS-3
 - FC-FG

- FC-PLDA
- FC-Tape
- FC-VI
- FC-SW-2
- Fibre Channel Element MIB RFC 2837
- Fibre Alliance MIB version 4.0
- ▶ Fibre Channel protocols:
 - Fibre Channel service classes: class 2 and class 3
 - Operation modes: Fibre Channel class 2 and class 3, connectionless
- ▶ External port type: Transparent fabric port (TF_port)
- ▶ Internal port type: Transparent host port/NPIV mode (TH_port)
 - Support for up to 44 host NPIV logins
- ▶ Port characteristics:
 - External ports are automatically detected and self- configuring
 - Port LEDs illuminate at startup
 - Number of Fibre Channel ports: 6 external ports and 14 internal ports
 - Scalability: Up to 239 switches maximum depending on your configuration
 - Buffer credits: 16 buffer credits per port
 - Maximum frame size: 2148 bytes (2112-byte payload)
 - Standards-based FC FC-SW2 Interoperability
 - Support for up to a 255-to-1 port-mapping ratio
 - Media type: Small form-factor pluggable plus (SFP+) module
- ▶ Fabric point-to-point bandwidth: 2 Gbps or 8 Gbps at full duplex
- ▶ 2 Gb Specifications
 - 2 Gb fabric port speed: 1.0625 or 2.125 Gbps (gigabits per second)
 - 2 Gb fabric latency: Less than 0.4 msec
 - 2 Gb fabric aggregate bandwidth: 80 Gbps at full duplex
- ▶ 4 Gb Specifications
 - 4 Gb switch speed: 4.250 Gbps
 - 4 Gb switch fabric point-to-point: 4 Gbps at full duplex
 - 4 Gb switch fabric aggregate bandwidth: 160 Gbps at full duplex
- ▶ 8 Gb Specifications
 - 8 Gb switch speed: 8.5 Gbps
 - 8 Gb switch fabric point-to-point: 8 Gbps at full duplex
 - 8 Gb switch fabric aggregate bandwidth: 320 Gbps at full duplex
- ▶ System processor: PowerPC
- ▶ Maximum frame size: 2148 bytes (2112-byte payload)
- ▶ Nonblocking architecture to prevent latency

For more information, see the IBM Redbooks Product Guide: *IBM Flex System FC3171 8Gb SAN Pass-thru*, available from:

<http://www.redbooks.ibm.com/abstracts/tips0866.html?Open>



Certification

In this chapter, we provide an insight into some of the various professional certifications that exist and that are appropriate to the topics in this book.

13.1 Why certification?

A good question is: Why should individuals take the effort to certify when there is more than enough work to do anyway? The following benefits to an individual can be realized:

- ▶ Validates your skills and knowledge
- ▶ Gains peer recognition
- ▶ Potential to become more valuable to your company and in the marketplace
- ▶ Ratifies your skills as an industry professional

To an employer, the following benefits can be seen:

- ▶ Great way of benchmarking the skill level of the employee
- ▶ Gives confidence about the employee's ability to support storage networks
- ▶ Demonstrates standards-based, non-proprietary, and vendor-neutral storage concepts

13.2 IBM Professional Certification Program

Today's marketplace is both crowded and complex. Individuals and businesses that do not stay ahead of the curve, risk being left behind. To develop a solid, competitive advantage, and to remain ahead of that curve, technology specialists are turning to professional certification from IBM. The extensive IBM portfolio of integrated certifications includes servers, software, application, and solution skills. The certification process is designed to prepare you and your company to meet business initiatives with real solutions.

The IBM Professional Certification Program helps in laying the groundwork for your personal journey to become a world-class resource to your clients, colleagues, and company. The program does this by providing you with the appropriate skills and accreditation that are needed to succeed.

13.2.1 About the program

The IBM Professional Certification Program is both a journey and a destination. It is a business solution; a way for skilled IT professionals to demonstrate their expertise to the world. The certification validates your skills and demonstrates your proficiency in the latest IBM technology and solutions.

The certification requirements can be tough. It is a rigorous process that differentiates you from everyone else.

The following list provides the mission of the IBM Professional Certification Program:

- ▶ To provide a reliable, valid, and fair method of assessing skills and knowledge.
- ▶ To provide IBM with a method of building and validating the skills of individuals and organizations.
- ▶ To develop a loyal community of highly skilled certified professionals who recommend, sell, service, support, or use IBM products and solutions.

13.2.2 Certifications by product

IBM has a myriad of certification courses that cover software, hardware, and products.

For a complete list of the courses that are available, see this website:

<http://www-03.ibm.com/certify/certs/index.shtml>

13.2.3 Mastery tests

Mastery tests are used to verify the mastery of knowledge that is covered in a course or a defined set of learning materials. They are not certification tests, which are designed to validate skills that are needed in a specific job role. Rather, mastery tests help to assure an individual achieved a foundation of knowledge and understanding of a subject matter.

Mastery tests supplement certifications as a method used by IBM to evaluate knowledge of IBM sales and technical professionals. As with certifications, the successful completion of a mastery test might be required for participation in some IBM Business Partner activities.

For a complete list of the courses available, see this website:

http://www-03.ibm.com/certify/mastery_tests/index_bd.shtml

13.3 Storage Networking Industry Association certifications

The *Storage Networking Industry Association (SNIA)* provides vendor-neutral certifications. There are different certification options within SNIA. The certification program itself is called the *Storage Networking Certification Program (SNCP)*.

The SNCP provides a strong foundation of vendor-neutral, systems-level credentials that integrate with and complement individual vendor certifications.

The structure of the SNCP is enhanced to reflect the advancement and growth of storage networking technologies over the past few years. And the structure is refined to provide for expanded offerings in the future. Through evolving and enhancing the SNCP, the SNIA is establishing a uniform standard by which individual knowledge and skill sets can be judged.

Before the establishment of the SNIA SNCP, there was no single standard by which to measure a professional's knowledge of storage networking technologies. Through its certification program, the SNIA is working to establish open standards for storage networking certification that IT organizations can trust.

SNCP Foundations exam withdrawn: The SNCP Foundations exam (S10–101) is withdrawn.

13.3.1 SNIA Certified Storage Professional (SCSP)

Attaining this vendor-neutral credential demonstrates that a storage networking professional attained a foundation of knowledge and expertise in fundamental storage networking technologies and concepts. This credential is independent of any product-specific certifications.

13.3.2 SNIA Certified Storage Engineer (SCSE)

This vendor-neutral credential addresses storage managers and administrators that are responsible for effective management, administration, troubleshooting, and diagnosing a SAN. This certification includes the following specialties:

- ▶ Performance management
- ▶ Implementation of upgrades
- ▶ Installation of new SANs
- ▶ Backup and recovery
- ▶ Business continuance

13.3.3 SNIA Certified Storage Architect (SCSA)

This vendor-neutral credential is designed for storage architects and storage networking professionals who assess, plan, and design complex storage networking solutions.

This certification validates the abilities of an individual and allows them to use industry standards in their programs.

13.3.4 SNIA Certified Storage Networking Expert (SCSN-E)

This credential is a culmination of the preceding vendor-neutral technical expertise. Combined with one or more vendor certifications, this credential positions the individual to work in a multi-vendor landscape, for example, in a storage managed services, partner, or reseller environment.

13.3.5 SNIA Qualified Data Protection Associate

This vendor-neutral credential is designed for individuals that are seeking to gain expertise and validate their knowledge in the storage data protection area.

This credential signifies that the candidate can:

- ▶ Recognize and describe concepts of data protection, restoration, and recovery methods
- ▶ Assess data protection planning and strategies
- ▶ Use management tools and practices
- ▶ Evaluate data protection methods and practices
- ▶ Assess security/confidentiality
- ▶ Troubleshoot potential pain points of data protection and recovery

13.3.6 SNIA Qualified Storage Virtualization Associate

This vendor-neutral credential is designed for individuals that are seeking to gain expertise and validate their knowledge in storage virtualization.

This credential signifies that the candidate can:

- ▶ Define virtualization concepts
- ▶ Describe the benefits of virtualization
- ▶ Identify the potential pain points of virtualization
- ▶ Describe virtualization implementation strategies
- ▶ Explain administrative and management tasks that are required for virtualization

13.3.7 SNIA Qualified Storage Sales Professional

The SNIA Qualified Storage Sales Professional (SQSSP) credential is vendor neutral and is designed to validate storage concepts, terminology, and basic client need assessments.

In direct response to firms that request specialized training for their sales and marketing professionals, this vendor-neutral training and credential provides individuals with the technologies in today's complex data center.

This training and credential prepares sales professionals by providing them with a broad set of information about various different solutions in the storage ecosystem.

13.3.8 CompTIA Storage+ Powered by SNIA

CompTIA Storage+ Powered by SNIA is a vendor-neutral certification that validates the knowledge and skills that are required of IT storage professionals.

The CompTIA Storage+ Powered by SNIA certification exam covers the knowledge and skills that are required to configure basic networks to include archive, backup, and restoration technologies. Additionally, the successful candidate is able to understand the fundamentals of

business continuity, application workload, system integration, and storage and system administration. This professional must also be able to perform basic troubleshooting on connectivity issues and referencing documentation.

The exam is targeted toward IT storage professionals with at least 12 months of experience. Though it is not required, CompTIA A+, CompTIA Network+, or CompTIA Server+ certification is recommended. For more information, see this website:

<http://certification.comptia.org/getCertified/certifications/storage.aspx>

13.4 Brocade certifications

Brocade has a large track of certification exams. After the completion of three of the tracks that are outlined in the following lists, it is possible to achieve the top credential of *Brocade Distinguished Architect*.

Each of the individual tracks has the following exams that are associated with them:

Brocade Certified Professional FICON Track

- ▶ Brocade Accredited Data Center Specialist Exam 160-130
- ▶ Brocade Accredited FICON Specialist Exam 160-140
- ▶ Brocade Certified Fabric Administrator (BCFA) Exam 143-410
- ▶ Brocade Certified Fabric Professional (BCFP) Exam 143-070
- ▶ Brocade Certified Architect For FICON (BCAF) Exam 143-120

Brocade Certified Professional Data Center Track

- ▶ Brocade Accredited Server Connectivity Specialist Exam 160-020
- ▶ Brocade Accredited Data Center Specialist Exam 160-130
- ▶ Brocade Certified Fabric Administrator (BCFA) Exam 143-410
- ▶ Brocade Certified Fabric Professional (BCFP) Exam 143-070
- ▶ Brocade Certified SAN Manager (BCSM) Exam 143-360
- ▶ Brocade Certified Fabric Designer (BCFD) Exam 143-260

Brocade Certified Professional Internetworking Track

- ▶ Brocade Accredited Internetworking Specialist Exam 160-120
- ▶ Brocade Certified Network Engineer Exam 150-120
- ▶ Brocade Certified Layer 4-7 Engineer Exam 150-320
- ▶ Brocade Certified Network Professional Exam 150-220
- ▶ Brocade Certified Layer 4-7 Professional Exam 150-420
- ▶ Brocade Accredited WLAN Specialist Exam 160-170
- ▶ Brocade Certified Network Designer Exam 150-510

Brocade Certified Professional Converged networking Track

- ▶ Brocade Accredited FCoE Specialist Exam 160-160
- ▶ Brocade Certified Ethernet Fabric Engineer Exam 150-610
- ▶ Brocade Certified FCoE Professional (BCFCoEP) Exam 143-510

13.4.1 Brocade Accredited Server Connectivity Specialist

This certification is for professionals with understanding of basic Brocade server adapter concepts, and can demonstrate knowledge of installation, maintenance, and troubleshooting.

This certification requires successful completion of the following exam: Brocade Accredited Server Connectivity Specialist Exam 160-020.

13.4.2 Brocade Accredited Data Center Specialist

This certification is for professionals that have the understanding of basic Fibre Channel theory, terminology, hardware, and the various reasons and benefits of implementing a storage area network (SAN). This certification requires successful completion of the following exam: Brocade Accredited Data Center Specialist: Exam 160-130.

13.4.3 Brocade Accredited Fibre Channel connection (FICON) Specialist

This certification is for professionals who understand the basic mainframe terminology and the relationship between FICON and Open Systems. This professional can also recognize the functions of Field Management System (FMS) and IBM RMF™ and can identify Brocade hardware and software products that support FICON. This certification requires the successful completion of the following exam: Brocade Accredited FICON Specialist Exam 160-140.

13.4.4 Brocade Accredited FCoE Specialist

This certification is for professionals who have an understanding of FCoE, CEE, the FCoE Initialization Protocol, and the associated Brocade hardware. This certification requires the successful completion of the following exam: Brocade Accredited FCoE Specialist Exam 160-160.

13.4.5 Brocade Accredited Internetworking Specialist

This certification is for professionals who have an understanding of basic internetworking terminology, hardware, routing concepts, the OSI 7 Layer Model, and the TCP/IP protocol suite, which includes IP addressing. This certification requires the successful completion of the following exam: Brocade Accredited Internetworking Specialist Exam 160-120.

13.4.6 Brocade Accredited WLAN Specialist

This certification is for professionals with the ability to demonstrate knowledge of wireless concepts, and install, maintain, and troubleshoot a Brocade Mobility solution. This certification requires the successful completion of the following exam: Brocade Accredited WLAN Specialist Exam 160-170.

13.4.7 Brocade Certified Fabric Administrator (BCFA)

This certification is for beginners who have a basic Fibre Channel (FC) protocol understanding and can perform basic switch configurations and troubleshooting. This certification requires the successful completion of the following exam: Brocade Certified Fabric Administrator (BCFA): Exam 143-410.

13.4.8 Brocade Certified Fabric Professional (BCFP)

This certification is for professionals that have the understanding of advanced FC technologies like FC switching and routing in an extended environment. These professionals

are also able to configure, administer, and troubleshoot the FC router and extended fabrics. Candidates can also implement adaptive networking in a SAN. This certification requires the successful completion of the following exam: Brocade Certified Fabric Professional (BCFP): Exam 143-070.

Brocade beta exam: Brocade announced a new Brocade Certified Ethernet Fabric Professional beta exam. For more information about this exam, see this website:

<http://community.brocade.com/docs/DOC-2814>

13.4.9 Brocade Certified SAN Manager (BCSM)

This certification is for experts who are skilled with administering, configuring, and troubleshooting the Brocade products with the help of management tools. This credential also focuses on the implementation of security enhancements, monitoring, and alerting of Brocade SAN. This certification requires the successful completion of the following exam: Brocade Certified SAN Manager (BCSM) Exam 143-360.

13.4.10 Brocade Certified Fabric Designer (BCFD)

This certification is for skilled professionals with the ability to design Data Center Fabric and to provide the implementation plans. This credential requires design skills with various criteria such as reliability, availability, and scalability; and to also plan for integration of new devices into current infrastructure. This certification requires the successful completion of the following exam: Brocade Certified Fabric Designer (BCFD) Exam 143-260.

13.4.11 Brocade Certified Architect For FICON (BCAF)

This certification is for experts who have a good understanding of IBM System z I/O and are able to identify, design, implement, and support Brocade products for mainframe FICON requirements. This certification requires the successful completion of the following exam: Brocade Certified Architect For FICON (BCAF) Exam 143-120.

13.4.12 Brocade Certified FCoE Professional (BCFCoEP)

This certification is for professionals who are able to show skills on FCoE and CEE concepts. They are also able to design, implement, support, and troubleshoot Brocade's FCoE products. This certification requires the successful completion of the following exam: Brocade Certified FCoE Professional (BCFCoEP) Exam 143-510.

13.4.13 Brocade Certified Ethernet Fabric Engineer

This certification is for professionals who are able to demonstrate knowledge of Ethernet fabric concepts and Brocade Ethernet fabric products. These professionals must also be able to install, configure, manage, and troubleshoot Brocade Ethernet fabrics. This certification requires successful completion of the following exam: Brocade Certified Ethernet Fabric Engineer Exam 150-610.

13.4.14 Brocade Certified Network Engineer

This certification is for professionals who are able to install and maintain IP switching and routing (Layer 2/3) networks that are based on Brocade products. This certification requires successful completion of the following exam: Brocade Certified Network Engineer Exam 150-120.

13.4.15 Brocade Certified Layer 4-7 Engineer

This certification is for professionals with the ability to install, configure, maintain, and perform basic troubleshooting of Brocade Layer 4-7 application delivery products. This certification requires successful completion of the following exam: Brocade Certified Layer 4-7 Engineer Exam 150-320.

13.4.16 Brocade Certified Network Professional

This certification is for professionals with the ability to install, configure, maintain, and troubleshoot Brocade Ethernet switches and routers in complex environments. This certification requires successful completion of the following exam: Brocade Certified Network Professional Exam 150-220.

13.4.17 Brocade Certified Layer 4-7 Professional

This certification is for professionals with the ability to design, configure, administer, and troubleshoot complex implementations of Brocade Layer 4-7 application delivery solutions. This certification requires successful completion of the following exam: Brocade Certified Layer 4-7 Professional Exam 150-420.

13.4.18 Brocade Certified Network Designer

This certification is for professionals with the ability to design a campus or enterprise network using Brocade solutions. This certification requires successful completion of the following exam: Brocade Certified Network Designer 150-510.

For more information about certification on Brocade, see this website:

<http://www.brocade.com/education/certification-accreditation/index.page>

13.5 Cisco certification

Cisco has various certifications for different product categories. This section focuses on SAN and system networking. Cisco has five levels of general IT certification: Entry, Associate, Professional, Expert, and Architect.

Table 13-1 on page 324 lists various Cisco certifications paths and their corresponding exams for different levels from Entry to Expert.

Table 13-1 Cisco certification levels and paths

Certification paths	Entry level	Associate	Professional	Expert	Architect
Routing and Switching	CCENT	CCNA	CCNP	CCIE Routing and Switching	Cisco Certified Architect (CCAr)
Design	CCENT	CCNA and CCDA	CCDP	CCDE	
Network Security	CCENT	CCNA Security	CCSP and CCNP Security	CCIE Security	
Wireless	CCENT	CCNA Wireless	CCNP Wireless	CCIE Wireless	
Storage Networking	CCENT	CCNA	CCNP	CCIE Storage Networking	

13.5.1 Cisco Certified Entry Networking Technician (CCENT)

This is the entry level certification for professionals with the ability to install, operate, and troubleshoot a small enterprise branch network, including basic network security. This certification requires successful completion of the following exam: Interconnecting Cisco Networking Devices Part 1 640-822 ICND1.

13.5.2 Cisco Certified Network Associate (CCNA)

This certification is for professionals with the ability to install, configure, operate, and troubleshoot medium-size routed and switched networks. This certification includes the ability to implement and verify connections to remote sites in a WAN. This certification requires successful completion of the following exam: Cisco Certified Network Associate (CCNA 640-802).

13.5.3 Cisco Certified Network Associate Security (CCNA Security)

This certification is for professionals with the ability to: Secure Cisco networks, develop a security infrastructure, recognize threats and vulnerabilities to networks, and mitigate security threats. This certification requires successful completion of the following exam: Implementing Cisco IOS Network Security (IINS 640-553).

13.5.4 Cisco Certified Network Associate Wireless (CCNA Wireless)

This certification is for professionals with associate-level knowledge and skills to configure, implement, and support wireless LANs that use Cisco equipment. This certification requires successful completion of the following exam: Implementing Cisco Unified Wireless Networking Essentials (IUWNE 640-721).

13.5.5 Cisco Certified Design Associate (CCDA)

This certification is for professionals with the ability to design a Cisco converged network. This professional is also able to design routed and switched network infrastructures and services

that involve LAN, WAN, and broadband access. This certification requires successful completion of the following exam: Designing for Cisco Internetwork Solutions: Exam 640-864.

13.5.6 Cisco Certified Network Professional (CCNP)

This certification is for professionals with the ability to plan, implement, verify, and troubleshoot local and wide-area enterprise networks. These professionals must also be able to work collaboratively with specialists on advanced security, voice, wireless, and video solutions. This certification requires successful completion of these exams:

- ▶ The Implementing Cisco IP Routing (ROUTE 642-902)
- ▶ Implementing Cisco IP Switched Networks (SWITCH 642-813)
- ▶ Troubleshooting and Maintaining Cisco IP Networks (TSHOOT 642-832)

13.5.7 CCNP Security certification

This certification is for professionals that are responsible for security in routers, switches, networking devices, and appliances. This designation is also for professionals that are responsible for choosing, deploying, supporting, and troubleshooting: Firewalls, VPNs, and intrusion detection system/intrusion prevention system (IDS/IPS) solutions for their networking environments. This certification requires successful completion of the following exams:

- ▶ Securing Networks with Cisco Routers and Switches (SECURE) v1.0: Exam 642-637
- ▶ Deploying Cisco ASA Firewall Solutions (FIREWALL v1.0): Exam 642-617
- ▶ Deploying Cisco ASA VPN Solutions (VPN v1.0): Exam 642-647
- ▶ Implementing Cisco Intrusion Prevention System v7.0: Exam 642-627

13.5.8 CCNP Wireless certification

This certification is for professionals with good expertise in designing, implementing, and operating Cisco Wireless networks and mobility infrastructures. The professionals also must be able to assess and translate network business requirements into technical specifications which can be incorporated into successful installations. This certification requires successful completion of the following exams:

- ▶ Conducting Cisco Unified Wireless Site Survey (CUWSS): Exam 642-731
- ▶ Implementing Cisco Unified Wireless Voice Networks (IUWVN): Exam 642-741
- ▶ Implementing Cisco Unified Wireless Mobility Services (IUWMS): Exam 642-746
- ▶ Implementing Advanced Cisco Unified Wireless Security (IAUWS): Exam 642-736

13.5.9 Cisco Certified Design Professional (CCDP)

This certification is for professionals with advanced knowledge of network design concepts and principles. This professional can describe, design, and create advanced addressing and routing, security, network management, and a data center of multi-layered enterprise architectures that include virtual private networking and wireless domains. This certification requires successful completion of the following exams:

- ▶ Implementing Cisco IP Routing: Exam 642-902
- ▶ Implementing Cisco IP Switched Networks: Exam 642-813
- ▶ Designing Cisco Network Service Architectures: Exam 642-874

13.5.10 Cisco Certified Internetwork Expert (CCIE) - Routing and Switching

This certification is for professionals who are expert-level network engineers that plan, operate, and troubleshoot complex, converged network infrastructures. This certification requires successful completion of the following exams:

- ▶ CCIE Routing and Switching: Written Exam #350-001, v4.0
- ▶ CCIE Routing and Switching v4.0: Lab Exam

13.5.11 Cisco Certified Internetwork Expert (CCIE) - Security

This certification is for professionals who have the knowledge and skills to implement, maintain, and support extensive Cisco Network Security Solutions using the latest industry preferred practices and technologies. This certification requires successful completion of the following exams:

- ▶ CCIE Security: Written Exam 350-018
- ▶ CCIE Security: Lab Exam v3.0

13.5.12 Cisco Certified Internetwork Expert (CCIE) - Wireless

This certification is for professionals with a broad theoretical knowledge of wireless networking and a solid understanding of wireless local area networking (WLAN) technologies from Cisco. This certification requires successful completion of the following exams:

- ▶ CCIE Wireless: Exam 350-050
- ▶ CCIE Wireless: Lab Exam

13.5.13 Cisco Certified Design Expert (CCDE)

This certification is for professionals with the ability to design large enterprise networks and develop solutions which address planning, design, integration, optimization, operations, security, and support the infrastructure. This certification requires successful completion of the following exams:

- ▶ CCDE: Written Exam 352-001
- ▶ CCDE: Practical Exam

13.5.14 Cisco CCIE Storage Networking

This certification is for professionals with a good understanding of Fibre Channel protocols, Cisco products, management, and troubleshooting tools from entry level to high-end products and applications. This certification requires successful completion of the following exams:

- ▶ CCIE Storage Networking: Written Exam 350-04
- ▶ CCIE Storage Networking: Lab Exam

13.5.15 Cisco Certified Architect

CCDE certification is a prerequisite for this designation. Professionals who apply for this certification must appear for an interview with a Cisco board of members for the validation of the architect role, and during which a skills assessment is conducted.

13.5.16 Cisco specialization tracks

Apart from the general IT certification tracks, Cisco also offers certifications for various specializations. For example, for data center professionals, there are certifications with various specializations such as sales, design, and support. The following Cisco specialization tracks are offered:

Data Center Networking Infrastructure

- ▶ Cisco Data Center Networking Infrastructure Design Specialist
- ▶ Cisco Data Center Networking Infrastructure Sales Specialist
- ▶ Cisco Data Center Networking Infrastructure Support Specialist

Data Center Storage Networking

- ▶ Cisco Data Center Storage Networking Design Specialist
- ▶ Cisco Data Center Storage Networking Sales Specialist
- ▶ Cisco Data Center Storage Networking Support Specialist

For more information about Cisco certifications, see this website:

http://www.cisco.com/web/learning/1e3/learning_career_certifications_and_learning_paths_home.html

13.6 The Open Group certifications

The Open Group provides certification programs for people, products, and services that meet their standards. For enterprise architects and IT specialists, the certification programs provide a worldwide professional credential for knowledge, skills, and experience. For IT products, Open Group Certification Programs offer a worldwide guarantee of conformance.

13.6.1 The Open Group Certified IT Specialists (Open CITS)

This designation is a vendor-neutral, Open Group, global certification program for IT specialists. This certification applies to various domains of IT industry and not specific only to Storage domain. Depending upon the skill set and experience of the professional, there are three levels of certifications for IT specialists:

- ▶ Level 1 - Certified IT Specialist: Professionals who are able to perform as a contributing specialist with assistance or supervision, with a wide range of appropriate skills.
- ▶ Level 2 - Master Certified IT Specialist: Professionals who are able to perform independently as lead specialist, and take responsibility for the delivery of solutions.
- ▶ Level 3 - Distinguished IT Specialist: Professionals who deliver leadership, scope, depth, and breadth of impact.

13.6.2 The Open Group Certified Architect (Open CA)

This designation is for architect professionals in IT, business, and enterprise architecture. There are three levels of certifications in this designation:

- ▶ Level 1 - Certified: Professionals with the ability to perform with assistance or supervision and who have a wide range of appropriate skills as a contributing architect.
- ▶ Level 2 - Master: Professionals with the ability to perform independently and take responsibility for the delivery of systems and solutions as the lead architect.

- ▶ Level 3 - Distinguished: Professionals who have significant breadth and depth of impact on the business through the application of IT architecture.

13.6.3 The Open Group certification

For more information about the Open Group certification programs, see this website:

<http://www3.opengroup.org/certifications>

13.7 Juniper Networks Certification Program (JNCP)

Juniper Networks offers Junos-based and non-Junos based platform-specific certifications.

13.7.1 JNCP Junos-based certification tracks

JNCP offers three tracks of professional certifications that are based on the Junos platform.

The following Junos-based professional certification tracks are offered:

- ▶ Service Provider Routing and Switching
- ▶ Enterprise Routing and Switching
- ▶ Junos Security

Table 13-2 Junos tracks

Certification track	Associate JNCIA	Specialist JNCIS	Professional JNCIP	Expert JNCIE
Service Provider and Switching	JNCIA-Junos	JNCIS-SP	JNCIP-SP	JNCIE-SP
Enterprise Routing & Switching	JNCIA-Junos	JNCIS-ENT	JNCIP-ENT	JNCIE-ENT
Junos Security	JNCIA-Junos	JNCIS-SEC	JNCIP-SEC	JNCIE-SEC

13.7.2 Service Provider Routing and Switching track

This section describes the certifications in the Service Provider Routing and Switching track.

Juniper Networks Certified Internet Associate (JNCIA–Junos)

This credential is designed for experienced networking professionals with beginner to intermediate knowledge of networking. The written exam verifies the candidate's understanding of the Juniper Networks Junos OS, networking fundamentals, and basic routing and switching.

This certification requires successful completion of the following exam:
Juniper Networks Junos Associate (JNCIA-Junos) Exam code: JN0-101

Juniper Networks Certified Internet Specialist (JNCIS-SP)

This credential is designed for experienced networking professionals with beginner to intermediate knowledge of routing and switching implementations in Junos. The written exam

verifies the candidate's basic understanding of routing and switching technologies and related platform configuration and troubleshooting skills.

This certification requires successful completion of the following exams:

- ▶ Juniper Networks Certified Internet Associate-Junos (JNCIA-Junos)
(Acceptable substitutions: JNCIA-ER, JNCIA-EX, JNCIA-M, or JNCIS-M)
- ▶ Juniper Networks Certified Internet Specialist (JNCIS-SP) Exam code: JN0-360

Juniper Networks Certified Internet Professional (JNCIP–SP)

This credential is designed for experienced networking professionals with advanced knowledge of the Juniper Networks Junos OS. The written exam verifies the candidate's understanding of advanced routing technologies and related platform configuration and troubleshooting skills.

This certification requires successful completion of the following exams:

- ▶ Juniper Networks Certified Internet Specialist (JNCIS-SP) Exam code: JN0-360
(Acceptable substitution: JNCIS-M)
- ▶ Juniper Networks Certified Internet Professional (JNCIP-SP) Exam code: JN0-660

Juniper Networks Certified Internet Expert (JNCIE–SP)

At the pinnacle of the Service Provider Routing and Switching certification track is the one-day JNCIE-SP Lab Exam. This exam is designed to validate the ability of the networking professional to implement, troubleshoot, and maintain Juniper Networks service provider networks. The eight-hour format of this exam requires that candidates build a service provider network that consists of multiple MX series routers. Successful candidates perform system configuration on all devices and implement various: Protocols, policies, VPNs, HA capabilities, and Class of Services on 8 MX Series Ethernet Services Routers.

This certification requires successful completion of the following exams:

- ▶ Juniper Networks Certified Internet Professional (JNCIP-SP) Exam code: JN0-660
(Acceptable substitution: JNCIP-M)
- ▶ Juniper Networks Certified Internet Expert (JNCIE–SP) Lab Exam code: JPR-960

13.7.3 Enterprise Routing and Switching track

This section describes the certifications in the Enterprise Routing and Switching track.

Juniper Networks Certified Internet Specialist (JNCIS-ENT)

This credential is designed for experienced networking professionals with beginner to intermediate knowledge of routing and switching implementations in Junos. The written exam verifies the candidate's basic understanding of routing and switching technologies and related platform configuration and troubleshooting skills.

This certification requires successful completion of the following exams:

- ▶ JNCIA-Junos
(Acceptable substitutions: JNCIA-ER, JNCIA-EX, JNCIA-M, or JNCIS-M)
- ▶ Juniper Networks Certified Internet Specialist (JNCIS-ENT) Exam code: JN0-343

Juniper Networks Certified Internet Professional (JNCIP-ENT)

This credential is designed for experienced networking professionals with advanced knowledge of the Juniper Networks Junos OS. The written exam verifies the candidate's understanding of advanced enterprise routing and switching technologies, and related platform configuration and troubleshooting skills.

This certification requires successful completion of the following exams:

- ▶ Juniper Networks Certified Specialist Enterprise Routing and Switching (JNCIS-ENT) Exam code: JN0-343
- ▶ Juniper Networks Certified Internet Professional (JNCIP-ENT) Exam code: JN0-643

Juniper Networks Certified Internet Expert (JNCIE-ENT)

At the pinnacle of the Enterprise Routing and Switching certification track is the one-day JNCIE-ENT practical exam. This exam is designed to validate the ability of the networking professional to deploy, configure, manage, and troubleshoot Junos-based enterprise routing and switching platforms. Throughout this eight-hour practical exam, candidates build an enterprise network infrastructure that consists of multiple routers and switching devices. Successful candidates perform system configuration on all devices, and configure protocols and features like IPV6, OSPF V2, OSPF V3, BGP, MSDP, PIM, SSM, RSTP, LLDP, 802.1x, CoS, and routing policies.

This certification requires successful completion of the following exam:

Juniper Networks Certified Internet Expert (JNCIE-ENT) Exam code: JPR-943

13.7.4 Junos security track

This section describes the certifications in the Junos security track.

Juniper Networks Certified Internet Specialist (JNCIS-SEC)

This credential is designed for experienced networking professionals with intermediate knowledge of the Juniper Networks Junos software for SRX Series devices. The written exam verifies the candidate's understanding of security technologies and related platform configuration and troubleshooting skills.

This certification requires successful completion of the following exams:

- ▶ Juniper Networks Certified Associate Junos (JNCIA-Junos)
(Acceptable substitutions: JNCIA-ER, JNCIA-M, JNCIS-M, or JNCIA-EX)
- ▶ Juniper Networks Certified Internet Specialist (JNCIS-SEC) Exam code: JN0-332

Juniper Networks Certified Internet Professional (JNCIP-SEC)

This credential is designed for experienced networking professionals with advanced knowledge of the Juniper Networks Junos software for SRX Series devices. The written exam verifies the candidate's understanding of advanced security technologies and related platform configuration and troubleshooting skills.

This certification requires successful completion of the following exams:

- ▶ Juniper Networks Certified Specialist Security (JNCIS-SEC) Exam code: JN0-332
- ▶ Juniper Networks Certified Internet Professional (JNCIP-SEC) Exam code: JN0-632

Juniper Networks Certified Internet Expert (JNCIE-SEC)

At the pinnacle of the Junos Security certification track is the one-day JNCIE-SEC practical exam. This exam is designed to validate the ability of the networking professional to deploy, configure, manage, and troubleshoot JUNOS-based security platforms. Throughout this eight-hour practical exam, candidates build a secure enterprise network that consists of multiple firewall devices that are interconnected via IPSec VPNs. Successful candidates perform system configuration on all devices and configure secure management capabilities. Candidates also install complex policies and attack prevention features, high availability (HA) capabilities, and intrusion prevention system (IPS) features.

This certification requires successful completion of the following exams:

- ▶ Juniper Networks Certified Professional Security (JNCIP-SEC) Exam code: JN0-632
- ▶ Juniper Networks Certified Internet Expert (JNCIE-SEC) Exam code: JPR-932

13.8 Non-Junos certification tracks

Non-Junos platforms have certification tracks as shown in Table 13-3.

Table 13-3 Non-Junos track

Certification Track	Associate JNCIA	Specialist JNCIS	Professional JNCIP
E-Series	JNCIA-E	JNCIS-E	JNCIP-E
Firewall/ VPN	JNCIA-FWV	JNCIS-FWV	
SSL	JNCIA-SSL	JNCIS-SSL	
IDP	JNCIA-IDP		
Unified Access Control		JNCIS-AC	
WX Series	JNCIA-WX		

13.8.1 E-Series certification track

The Juniper Networks Certification Program (JNCP) E-series certification track is a multi-tiered program. The program allows participants to demonstrate competence with specific Juniper Networks technologies. Demonstration of the skills is done through a combination of written proficiency exams and hands-on configuration and troubleshooting exams. Successful candidates demonstrate thorough understanding of networking technology, in general; and understanding of the Juniper Networks E-series platforms and operating system, in particular.

Juniper Networks Certified Internet Associate (JNCIA- E)

This credential is designed for experienced networking professionals with beginner to intermediate knowledge of the Juniper Networks E-series platforms. The written exam verifies the candidate's basic understanding of Internet technology and related platform configuration and troubleshooting skills. JNCIA-E exam topics are based on the content of the Introduction to Juniper Networks Routers—E-series and E-series B-RAS Configuration Basics instructor-led training courses. This exam is not a prerequisite for the JNCIS-E exam.

This certification requires successful completion of the following exam:
Juniper Networks Certified Internet Associate (JNCIA- E) Exam code: JN0-120.

Juniper Networks Certified Internet Specialist (JNCIS-E)

The JNCIS-E exam is designed for networking professionals with advanced knowledge of, and experience with, the Juniper Networks E-series platforms. The JNCIS-E exam tests for a wider and deeper level of knowledge than does the JNCIA-E exam. Exam questions focus on the E-series platforms documentation set and on-the-job product experience. Questions also focus on the understanding of Internet technologies and design principles that are considered to be common knowledge at the Specialist level. Passing the JNCIS-E exam is a prerequisite for attempting the JNCIP-E practical exam.

This certification requires successful completion of the following exam:
Juniper Networks Certified Internet Specialist (JNCIS-E) Exam code: JN0-130.

Juniper Networks Certified Internet Professional (JNCIP-E)

The JNCIP-E exam is a one-day practical exam that is designed to validate the candidate's ability to successfully build an Internet service provider (ISP) consisting of multiple E-series virtual routers. This certification establishes the candidate's practical and theoretical knowledge of basic and advanced Internet technologies and the candidate's ability to effectively apply that knowledge in a hands-on environment. Candidates configure and troubleshoot routing scenarios using various protocols and technologies on E-series platforms.

This certification requires successful completion of the following exam:
Juniper Networks Certified Internet Professional (JNCIP-E) Exam code: CERT-JNCIP-E.

13.8.2 Firewall/VPN certification track

The Juniper Networks Certification Program (JNCP) Firewall/VPN certification track is a two-tiered program that allows participants to demonstrate competence with Juniper Networks Firewall with VPN products and the ScreenOS software.

Juniper Networks Certified Internet Associate (JNCIA-FWV)

This credential is designed for experienced networking professionals with beginner to intermediate knowledge of Juniper Firewall/VPN products and ScreenOS software. The written exam verifies the candidate's basic understanding of Internet and security technology and related device configuration. JNCIA-FWV exam topics are based on the content of the Configuring Juniper Networks Firewall/IPSec VPN Products instructor-led training course. This exam is not a prerequisite for the JNCIS-FWV certification.

This certification requires successful completion of the following exam:
Juniper Networks Certified Internet Associate (JNCIA-FWV) Exam code: JN0-522.

Juniper Networks Certified Internet Specialist (JNCIS-FWV)

The JNCIS-FWV is designed for networking professionals with advanced knowledge of, and experience with, Juniper Firewall/VPN products and ScreenOS software. The JNCIS-FWV exam tests for a wider and deeper level of knowledge than does the JNCIA-FWV exam. Sources of question content include all ScreenOS training courses, the Firewall/VPN and ScreenOS documentation set, and on-the-job product experience. Exam topics also include Internet technologies and design principles that are considered to be common knowledge at the Specialist level.

This certification requires successful completion of the following exam:
Juniper Networks Certified Internet Specialist (JNCIS-FWV) Exam code: JN0-532.

13.8.3 SSL certification track

The Juniper Networks Certification Program (JNCP) SSL certification track allows participants to demonstrate competence with Juniper Networks Secure Access products and their deployment.

Juniper Networks Certified Internet Associate (JNCIA-SSL)

Designed for experienced networking professionals with beginner-intermediate knowledge of the Juniper Networks Secure Access products and their deployment. JNCIA-SSL exam topics are based on the content of the Configuring Juniper Networks Secure Access instructor led training course.

This certification requires successful completion of the following exam:
Juniper Networks Certified Internet Associate (JNCIA-SSL) Exam code: JN0-562.

Juniper Networks Certified Internet Specialist (JNCIS-SSL)

Designed for experienced networking professionals with intermediate knowledge of the Juniper Networks Secure Access products and their deployment. JNCIS-SSL exam topics are based on the content of the Advanced Juniper Networks Secure Access instructor-led training course.

This certification requires successful completion of the following exam:
Juniper Networks Certified Internet Specialist (JNCIS-SSL) Exam code: JN0-570.

13.8.4 Intrusion Detection and Prevention (IDP) Track

The Juniper Networks Certification Program (JNCP) IDP certification track allows participants to demonstrate competence with Juniper Networks NetScreen IDP products and their deployment.

Juniper Networks Certified Internet Associate (JNCIA-IDP)

Designed for experienced networking professionals with beginner to intermediate knowledge of the Juniper Networks IDP products and their deployment. JNCIA-IDP exam topics are based on the content of the Implementing Intrusion Detection and Prevention (IIDP) instructor-led training course.

This certification requires successful completion of the following exam:
Juniper Networks Certified Internet Associate (JNCIA-IDP) Exam code: JN0-541.

13.8.5 Unified Access Control (UAC) Track

The Juniper Networks Certification Program (JNCP) Unified Access Control certification track allows participants to demonstrate competence with Juniper Networks Unified Access Control products and their deployment.

JNCIA-AC is now inactive: As of 13 April 2012, the JN0-141 AC, Associate (JNCIS-AC) exam became End of Life (EOL). The Juniper Networks Certified Internet Associate, AC (JNCIA-AC) certification and the Access Control Track is now inactive and unsupported. However, when earned, a JNCIA-AC credential is valid for two years.

The replacement exam is the JN0-314 Junos Pulse Access Control, Specialist (JNCIS-AC), which earns the Juniper Networks Certified Specialist Junos Pulse Access Control (JNCIS-AC) certification.

13.8.6 WX certification track

The Juniper Networks Certification Program (JNCP) WX certification track allows participants to demonstrate competence with Juniper Networks WAN Acceleration platforms and their deployment.

Juniper Networks Certified Internet Associate (JNCIA-WX)

Designed for experienced networking professionals with beginner to intermediate knowledge of the Juniper Networks WAN Acceleration (WX) and WAN Acceleration Cache (WXC) platforms and their deployment. JNCIA-WX exam topics are based on the content of the Operating Juniper Networks WX Application Acceleration Platforms (OJWX) instructor-led training course.

This certification requires successful completion of the following exam:
Juniper Networks Certified Internet Associate (JNCIA-WX) Exam code: JN0-311.

For more information about the Juniper certifications, see this website:

<http://www.juniper.net/us/en/training/certification/certification-tracks/>

Related publications

The publications listed in this section are considered particularly suitable for a more detailed discussion of the topics covered in this book.

IBM Redbooks

- ▶ *IBM TotalStorage: SAN Product, Design, and Optimization Guide*, SG24-6384
- ▶ *IBM TotalStorage SAN Volume Controller*, SG24-6423
- ▶ *Implementing an Open IBM SAN*, SG24-6116
- ▶ *Implementing the Cisco MDS 9000 in an Intermix FCP, FCIP, and FICON Environment*, SG24-6397
- ▶ *Introduction to SAN Distance Solutions*, SG24-6408
- ▶ *Introducing Hosts to the SAN Fabric*, SG24-6411
- ▶ *IP Storage Networking: IBM NAS and iSCSI Solutions*, SG24-6240
- ▶ *The IBM TotalStorage NAS Integration Guide*, SG24-6505
- ▶ *Implementing the IBM TotalStorage NAS 300G: High Speed Cross Platform Storage and Tivoli SANergy!*, SG24-6278
- ▶ *Using iSCSI Solutions' Planning and Implementation*, SG24-6291
- ▶ *IBM Storage Solutions for Server Consolidation*, SG24-5355
- ▶ *Implementing the Enterprise Storage Server in Your Environment*, SG24-5420
- ▶ *Implementing Linux with IBM Disk Storage*, SG24-6261
- ▶ *IBM Tape Solutions for Storage Area Networks and FICON*, SG24-5474
- ▶ *IBM Enterprise Storage Server*, SG24-5465
- ▶ *The IBM TotalStorage Solutions Handbook*, SG24-5250

The following publications from IBM Redbooks provide more information about IBM Flex System and IBM PureSystems. These are available from the following website:

<http://www.redbooks.ibm.com/portals/puresystems>

- ▶ *IBM PureFlex System and IBM Flex System Products & Technology*, SG24-7984
- ▶ *IBM Flex System p260 and p460 Panning and Implementation Guide*, SG24-7989
- ▶ *IBM Flex System Networking in an Enterprise Data Center*, REDP4834

Chassis and Compute Nodes:

- ▶ *IBM Flex System Enterprise Chassis*, TIPS0863
- ▶ *IBM Flex System p260 and p460 Compute Node*, TIPS0880
- ▶ *IBM Flex System x240 Compute Node*, TIPS0860
- ▶ *IBM Flex System Manager*, TIPS0862

Switches:

- ▶ *IBM Flex System EN2092 1Gb Ethernet Scalable Switch*, TIPS0861
- ▶ *IBM Flex System Fabric EN4093 10Gb Scalable Switch*, TIPS0864
- ▶ *IBM Flex System EN4091 10Gb Ethernet Pass-thru Module*, TIPS0865
- ▶ *IBM Flex System FC5022 16Gb SAN Scalable Switch and FC5022 24-port 16Gb ESB SAN Scalable Switch*, TIPS0870
- ▶ *IBM Flex System IB6131 InfiniBand Switch*, TIPS0871
- ▶ *IBM Flex System FC3171 8Gb SAN Switch and Pass-thru*, TIPS0866

Adapters:

- ▶ *IBM Flex System EN2024 4-port 1Gb Ethernet Adapter*, TIPS0845
- ▶ *IBM Flex System FC5022 2-port 16Gb FC Adapter*, TIPS0891
- ▶ *IBM Flex System CN4054 10Gb Virtual Fabric Adapter and EN4054 4-port 10Gb Ethernet Adapter*, TIPS0868
- ▶ *IBM Flex System FC3052 2-port 8Gb FC Adapter*, TIPS0869
- ▶ *ServeRAID M5115 SAS/SATA Controller for IBM Flex System*, TIPS0884
- ▶ *IBM Flex System IB6132 2-port FDR InfiniBand Adapter*, TIPS0872
- ▶ *IBM Flex System EN4132 2-port 10Gb Ethernet Adapter*, TIPS0873
- ▶ *IBM Flex System IB6132 2-port QDR InfiniBand Adapter*, TIPS0890
- ▶ *IBM Flex System FC3172 2-port 8Gb FC Adapter*, TIPS0867

You can search for, view, download, or order these documents and other Redbooks, Redpapers, Web Docs, drafts, and additional materials, at the following website:

ibm.com/redbooks

IBM Flex System education

The following courses are IBM educational offerings for IBM Flex System. Some course numbers and titles might have changed slightly after publication.

Course types: IBM courses prefixed with NGTxx are traditional, face-to-face classroom offerings. Courses prefixed with NGVxx are Instructor-Led Online (ILO) offerings. Courses prefixed with NGPxx are Self-paced Virtual Class (SPVC) offerings.

- ▶ *NGT10/NGV10/NGP10*, IBM Flex System - Introduction
- ▶ *NGT20/NGV20/NGP20*, IBM Flex System x240 Compute Node
- ▶ *NGT30/NGV30/NGP30*, IBM Flex System p260 and p460 Compute Nodes
- ▶ *NGT40/NGV40/NGP40*, IBM Flex System Manager Node
- ▶ *NGT50/NGV50/NGP50*, IBM Flex System Scalable Networking

For more information about these courses, and many other IBM System x educational offerings, see the global IBM Training website:

<http://www.ibm.com/training>

Referenced websites

These websites are also relevant as further information sources:

- ▶ IBM TotalStorage hardware, software, and solutions:
<http://www.storage.ibm.com>
- ▶ IBM System Storage: Storage area networks:
<http://www-03.ibm.com/servers/storage/san/>
- ▶ Brocade:
<http://www.brocade.com>
- ▶ Cisco:
<http://www.cisco.com>
- ▶ QLogic:
<http://www.qlogic.com>
- ▶ Emulex:
<http://www.emulex.com>
- ▶ Finisar:
<http://www.finisar.co>
- ▶ Tivoli:
<http://www.tivoli.com>
- ▶ IEEE:
<http://www.ieee.org>
- ▶ Storage Networking Industry Association:
<http://www.snia.org>
- ▶ Fibre Channel Industry Association:
<http://www.fibrechannel.com>
- ▶ SCSI Trade Association:
<http://www.scsita.org>
- ▶ Internet Engineering Task Force:
<http://www.ietf.org>
- ▶ American National Standards Institute:
<http://www.ansi.org>
- ▶ Technical Committee T10:
<http://www.t10.org>
- ▶ Technical Committee T11:
<http://www.t11.org>

Help from IBM

IBM Support and downloads

ibm.com/support

IBM Global Services

ibm.com/services



Redbooks

Introduction to Storage Area Networks and System Networking

(1.5" spine)
1.5" <-> 1.998"
789 <-> 1051 pages



Redbooks

Introduction to Storage Area Networks and System Networking

(1.0" spine)
0.875" <-> 1.498"
460 <-> 788 pages



Redbooks

Introduction to Storage Area Networks and System Networking

(0.5" spine)
0.475" <-> 0.873"
250 <-> 459 pages



Redbooks

Introduction to Storage Area Networks and System Networking

(0.2" spine)
0.17" <-> 0.473"
90 <-> 249 pages

(0.1" spine)
0.1" <-> 0.169"
53 <-> 89 pages



Redbooks

Introduction to Storage Area Networks and System Networking

(2.5" spine)
2.5" <-> nnn.n"
1315 <-> nnnn pages



Redbooks

Introduction to Storage Area Networks and System Networking

(2.0" spine)
2.0" <-> 2,498"
1052 <-> 1314 pages



Introduction to Storage Area Networks and System Networking



**Learn basic SAN and
System Networking
concepts**

**Introduce yourself
to the
business benefits**

**Discover the IBM
System Networking
portfolio**

The plethora of data that is created by the businesses of today is making storage a strategic investment priority for companies of all sizes. As storage takes precedence, three major initiatives emerge:

- ▶ Flatten and converge your network
IBM takes an open, standards-based approach to implement the latest advances in the flat, converged data center network designs of today. IBM System Networking solutions enable clients to deploy a high-speed, low-latency Unified Fabric Architecture.
- ▶ Optimize and automate virtualization
Advanced virtualization awareness reduces the cost and complexity of deploying physical and virtual data center infrastructure.
- ▶ Simplify management
IBM data center networks are easy to deploy, maintain, scale, and virtualize, delivering the foundation of consolidated operations for dynamic infrastructure management.

Welcome to the era of *Smarter Networking for Smarter Data Centers*.

The smarter data center with improved economics of IT can be achieved by connecting servers and storage with a high-speed and intelligent network fabric. A smarter data center that hosts IBM System Networking solutions can provide an environment that is smarter, faster, greener, open, and easy to manage.

This IBM Redbooks publication provides an introduction to the SAN and Ethernet networking, and how these networks help to achieve a smarter data center. This book is intended for people who are not very familiar with IT, or who are just starting out in the IT world.

**INTERNATIONAL
TECHNICAL
SUPPORT
ORGANIZATION**

**BUILDING TECHNICAL
INFORMATION BASED ON
PRACTICAL EXPERIENCE**

IBM Redbooks are developed by the IBM International Technical Support Organization. Experts from IBM, Customers and Partners from around the world create timely technical information based on realistic scenarios. Specific recommendations are provided to help you implement IT solutions more effectively in your environment.

**For more information:
ibm.com/redbooks**

SG24-5470-04

ISBN 0738437131