



Essential Virtual SAN (VSAN)

Administrator's Guide to VMware® Virtual SAN

Cormac Hogan
Duncan Epping

SECOND EDITION



FREE SAMPLE CHAPTER

SHARE WITH OTHERS



Essential Virtual SAN (VSAN)

Administrator's Guide to
VMware® Virtual SAN

Second Edition

VMware Press is the official publisher of VMware books and training materials, which provide guidance on the critical topics facing today's technology professionals and students. Enterprises, as well as small- and medium-sized organizations, adopt virtualization as a more agile way of scaling IT to meet business needs. VMware Press provides proven, technically accurate information that will help them meet their goals for customizing, building, and maintaining their virtual environment.

With books, certification and study guides, video training, and learning tools produced by world-class architects and IT experts, VMware Press helps IT professionals master a diverse range of topics on virtualization and cloud computing and is the official source of reference materials for preparing for the VMware Certified Professional examination.

VMware Press is also pleased to have localization partners that can publish its products into more than 42 languages, including, but not limited to, Chinese (Simplified), Chinese (Traditional), French, German, Greek, Hindi, Japanese, Korean, Polish, Russian, and Spanish.

For more information about VMware Press, please visit
<http://www.vmwarepress.com>.

This page intentionally left blank

Essential Virtual SAN (VSAN)

**Administrator's Guide to
VMware® Virtual SAN**

Second Edition

Cormac Hogan

Duncan Epping

vmware PRESS

Boston • Indianapolis • San Francisco
New York • Toronto • Montreal • London • Munich • Paris • Madrid
Capetown • Sydney • Tokyo • Singapore • Mexico City

Essential Virtual SAN (VSAN) Second Edition

Copyright © 2016 VMware, Inc

Published by Pearson Education, Inc.

Publishing as VMware Press

All rights reserved. Printed in the United States of America. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise.

ISBN-10: 0-13-451166-2

ISBN-13: 978-0-13-451166-5

Library of Congress Control Number: 2016903470

Printed in the United States of America

First Printing: May 2016

All terms mentioned in this book that are known to be trademarks or service marks have been appropriately capitalized. The publisher cannot attest to the accuracy of this information. Use of a term in this book should not be regarded as affecting the validity of any trademark or service mark.

VMware terms are trademarks or registered trademarks of VMware in the United States, other countries, or both.

Warning and Disclaimer

Every effort has been made to make this book as complete and as accurate as possible, but no warranty or fitness is implied. The information provided is on an as is basis. The authors, VMware Press, VMware, and the publisher shall have neither liability nor responsibility to any person or entity with respect to any loss or damages arising from the information contained in this book or from the use of any digital content or programs accompanying it.

The opinions expressed in this book belong to the author and are not necessarily those of VMware.

Special Sales

For information about buying this title in bulk quantities, or for special sales opportunities (which may include electronic versions; custom cover designs; and content particular to your business, training goals, marketing focus, or branding interests), please contact our corporate sales department at corpsales@pearsoned.com or (800) 382-3419.

For government sales inquiries, please contact governmentsales@pearsoned.com.

For questions about sales outside the United States, please contact intlcs@pearson.com

VMWARE PRESS PROGRAM MANAGER

Karl Childs

EXECUTIVE EDITOR

Mary Beth Ray

TECHNICAL EDITORS

Christian Dickmann

John Nicholson

SENIOR DEVELOPMENT EDITOR

Ellie Bru

MANAGING EDITOR

Sandra Schroeder

PROJECT EDITOR

Mandie Frank

COPY EDITOR

Cenveo® Publisher Services

PROOFREADER

Cenveo Publisher Services

INDEXER

Cenveo Publisher Services

EDITORIAL ASSISTANT

Vanessa Evans

DESIGNER

Chuti Prasertsith

COMPOSITOR

Cenveo Publisher Services

*“I close my eyes, and think of home. Another city goes by in the night. Ain’t it funny
how it is, you never miss it ‘til it has gone away.
And my heart is lying there and will be ‘til my dying day.”*

*We would like to dedicate this book to the greatest
band on earth, Iron Maiden.*

—Cormac & Duncan

This page intentionally left blank

Contents

Foreword by Christos Karamanolis xvii

About the Author xix

About the Technical Reviewers xxi

Acknowledgments xxiii

We Want to Hear from You! xxv

1 Introduction to VSAN 1

Software-Defined Datacenter 1

Software-Defined Storage 2

Hyper-Convergence/Server SAN Solutions 3

Introducing Virtual SAN 4

What Is Virtual SAN? 6

What Does VSAN Look Like to an Administrator? 9

Summary 12

2 VSAN Prerequisites and Requirements for Deployment 13

VMware vSphere 13

ESXi 14

Cache and Capacity Devices 14

ESXi Boot Considerations 15

VSAN Requirements 15

VMware Hardware Compatibility Guide 16

VSAN Ready Nodes 16

Storage Controllers 17

Capacity Tier Devices 19

Cache Tier Devices 21

Network Requirements 22

Network Interface Cards 22

Supported Virtual Switch Types 22

Layer 2 or Layer 3 23

VMkernel Network 23

VSAN Network Traffic 24

Jumbo Frames 24

NIC Teaming 25

Network I/O Control 25

VSAN Stretched Cluster 25

VSAN 2-Node Remote Office/Branch Office (ROBO) 26

Firewall Ports 26

Summary 26

3 VSAN Installation and Configuration 29

VSAN Networking	29
VMkernel Network for VSAN	30
VSAN Network Configuration: VMware Standard Switch	31
VSAN Network Configuration: vSphere Distributed Switch	32
Step 1: Create the Distributed Switch	32
Step 2: Create a Distributed Port Group	33
Step 3: Build VMkernel Ports	34
Possible Network Configuration Issues	38
Network I/O Control Configuration Example	40
Design Considerations: Distributed Switch and Network I/O Control	42
Scenario 1: Redundant 10 GbE Switch Without “Link Aggregation” Capability	43
Scenario 2: Redundant 10 GbE Switch with Link Aggregation Capability	45
Creating a VSAN Cluster	48
vSphere HA	49
vSphere HA Communication Network	49
vSphere HA Heartbeat Datastores	50
vSphere HA Admission Control	50
vSphere HA Isolation Response	51
vSphere HA Component Protection	51
The Role of Disk Groups	51
Disk Group Maximums	52
Why Configure Multiple Disk Groups in VSAN?	52
Cache Device to Capacity Device Sizing Ratio	53
Automatically Add Disks to VSAN Disk Groups	54
Manually Adding Disks to a VSAN Disk Group	55
Disk Group Creation Example	55
VSAN Datastore Properties	58
Summary	59

4 VM Storage Policies on VSAN 61

Introducing Storage Policy-Based Management in a VSAN Environment	62
Number of Failures to Tolerate	65
Failure Tolerance Method	66
Number of Disk Stripes Per Object	69
IOPS Limit for Object	70
Flash Read Cache Reservation	71
Object Space Reservation	71
Force Provisioning	71
Disable Object Checksum	73

VASA Vendor Provider	73
An Introduction to VASA	73
Storage Providers	74
VSAN Storage Providers: Highly Available	75
Changing VM Storage Policy On-the-Fly	75
Objects, Components, and Witnesses	80
VM Storage Policies	80
Enabling VM Storage Policies	81
Creating VM Storage Policies	81
Assigning a VM Storage Policy During VM Provisioning	81
Summary	82
5 Architectural Details	83
Distributed RAID	83
Objects and Components	86
Component Limits	87
Virtual Machine Storage Objects	88
Namespace	89
Virtual Machine Swap	90
VMDKs and Deltas	90
Witnesses and Replicas	90
Object Layout	91
VSAN Software Components	94
Component Management	95
Data Paths for Objects	95
Object Ownership	96
Placement and Migration for Objects	96
Cluster Monitoring, Membership, and Directory Services	97
Host Roles (Master, Slave, Agent)	97
Reliable Datagram Transport	98
On-Disk Formats	98
Cache Devices	99
Capacity Devices	99
VSAN I/O Flow	100
Caching Algorithms	100
The Role of the Cache Layer	100
Anatomy of a VSAN Read on Hybrid VSAN	102
Anatomy of a VSAN Read on All-Flash VSAN	103
Anatomy of a VSAN Write on Hybrid VSAN	103
Anatomy of a VSAN Write on All-Flash VSAN	104

Retiring Writes to Capacity Tier on Hybrid VSAN	105
Deduplication and Compression	105
Data Locality	107
Data Locality in VSAN Stretched Clusters	108
Storage Policy-Based Management	109
VSAN Capabilities	109
Number of Failures to Tolerate Policy Setting	110
Best Practice for Number of Failures to Tolerate	112
Stripe Width Policy Setting	113
RAID-0 Used When No Striping Specified in the Policy	117
Stripe Width Maximum	119
Stripe Width Configuration Error	120
Stripe Width Chunk Size	121
Stripe Width Best Practice	122
Flash Read Cache Reservation Policy Setting	122
Object Space Reservation Policy Setting	122
VM Home Namespace Revisited	123
VM Swap Revisited	123
How to Examine the VM Swap Storage Object	124
Delta Disk / Snapshot Caveat	126
Verifying How Much Space Is Actually Consumed	126
Force Provisioning Policy Setting	127
Witnesses and Replicas: Failure Scenarios	127
Data Integrity Through Checksum	130
Recovery from Failure	131
Problematic Device Handling	134
What About Stretching VSAN?	134
Summary	135

6 VM Storage Policies and Virtual Machine Provisioning 137

Policy Setting: Number of Failures to Tolerate = 1	137
Policy Setting: Failures to Tolerate = 1, Stripe Width = 2	144
Policy Setting: Failures to Tolerate = 2, Stripe Width = 2	148
Policy Setting: Failures to Tolerate = 1, Object Space Reservation = 50%	152
Policy Setting: Failures to Tolerate = 1, Object Space Reservation = 100%	155
Policy Setting: RAID-5	157
Policy Setting: RAID-6	158
Policy Setting: RAID-5/6 and Stripe Width = 2	159
Default Policy	160
Summary	164

7 Management and Maintenance 165

Health Check 165
Health Check Tests 165
Proactive Health Checks 167
Performance Service 168
Host Management 169
Adding Hosts to the Cluster 169
Removing Hosts from the Cluster 170
ESXCLI VSAN Cluster Commands 171
Maintenance Mode 172
Default Maintenance Mode/Decommission Mode 175
Recommended Maintenance Mode Option for Updates and Patching 175
Disk Management 177
Adding a Disk Group 177
Removing a Disk Group 178
Adding Disks to the Disk Group 179
Removing Disks from the Disk Group 180
Wiping a Disk 182
Blinking the LED on a Disk 183
ESXCLI VSAN Disk Commands 184
Failure Scenarios 185
Capacity Device Failure 185
Cache Device Failure 186
Host Failure 187
Network Partition 188
Disk Full Scenario 193
Thin Provisioning Considerations 194
vCenter Management 195
vCenter Server Failure Scenario 196
Running vCenter Server on VSAN 196
Bootstrapping vCenter Server 197
Summary 199

8 Stretched Cluster 201

What is a Stretched Cluster? 201
Requirements and Constraints 203
Networking and Latency Requirements 205
New Concepts in VSAN Stretched Cluster 206
Configuration of a Stretched Cluster 208

Failure Scenarios	216
Summary	224
9 Designing a VSAN Cluster	225
Ready Node Profiles	225
Sizing Constraints	227
Cache to Capacity Ratio	228
Designing for Performance	229
Impact of the Disk Controller	231
VSAN Performance Capabilities	235
Design and Sizing Tools	236
Scenario 1: Server Virtualization—Hybrid	237
Determining Your Host Configuration	238
Scenario 2—Server Virtualization—All-flash	241
Summary	244
10 Troubleshooting, Monitoring, and Performance	245
Health Check	246
Ask VMware	246
Health Check Categories	247
Proactive Health Checks	253
ESXCLI	256
esxcli vsan datastore	256
esxcli vsan network	257
esxcli vsan storage	258
esxcli vsan cluster	262
esxcli vsan faultdomain	263
esxcli vsan maintenancemode	264
esxcli vsan policy	264
esxcli vsan trace	267
Additional Non-ESXCLI Commands for Troubleshooting VSAN	268
Ruby vSphere Console	275
VSAN Commands	276
SPBM Commands	300
Troubleshooting VSAN on the ESXi	303
Log Files	304
VSAN Traces	304
VSAN VMkernel Modules and Drivers	305
Performance Monitoring	305
Introducing the Performance Service	305

ESXTOP Performance Counters for VSAN	308
vSphere Web Client Performance Counters for VSAN	309
VSAN Observer	310
Sample VSAN Observer Use Case	316
Summary	318

Index 319

This page intentionally left blank

Foreword by Christos Karamanolis

We are in the midst of a major evolution of hardware, especially storage technologies and networking. With NAND-based Flash, a single device can deliver 1M I/O operations a second. That's more than what a high-end disk array could deliver just a few years ago. With 10 Gbps networks becoming commodity, the industry is already eyeing 40 or even 100 Gbps Ethernet. Efficient protocols like RDMA and NVMe reduce the cycles CPUs need to consume to harness the capabilities of emerging storage and networking technologies. At the same time, CPU densities are doubling every two years.

It was within this technology forecast that I started the Virtual SAN project together with a handful of brilliant engineers back in 2010. We recognized that commodity servers have enough oomph to take over services like storage and networking traditionally offered by dedicated hardware appliances. We argued that software running on general-purpose (commodity) hardware is the way to go. Simply put, the technology evolution allows for the democratization of enterprise storage and networking. A software-defined architecture allows the customer to ride the wave of technology trends and offers enormous benefits—lower costs, flexibility with hardware choices, and, if it comes with the right tools, a unified way to manage one's infrastructure: compute, storage, network.

The operational complexity and scalability challenges of storage management have been amongst the biggest pain points of virtualization adopters. So, we focused on Software-Defined Storage at the same time that a little-known Stanford graduate, named Martin Casado, was making the case for Software-Defined Networks.

Fast forward to March 2016: VMware's Virtual SAN is in its fourth release (version 6.2), the most important release of the product since its debut in March 2014. With a panoply of data integrity, availability and space efficiency features, Virtual SAN stands tall against the most sophisticated storage products in the industry. By combining commodity hardware economics and a flexible, highly efficient software stack, VMware is making all-flash storage affordable and applicable to a broad range of use cases. With thousands of enterprise customers using VSAN for business-critical use cases that range from Hospitals to popular E-commerce platforms, and from oil rigs to aircraft carriers, Virtual SAN has the credentials to be deemed battle tested. And we proudly carry the scars that come with that!

However, if we focus entirely on the (impressive) storage features of Virtual SAN, we are missing an important point: the fundamental shift in the management paradigm that the product introduces. In the current instantiation of the technology, Virtual SAN is designed to be an ideal storage platform for hyper-converged infrastructures (HCI). The key benefit of HCI is its simple operational model that allows customers to manage the entire IT infrastructure with a single set of tools. Virtual SAN is blazing the path for VMware products and for the industry at large.

Since version 6.1, Virtual SAN has been offering a “Health Service” that monitors, probes and assists with remediation of hardware, configuration, and performance issues. With version 6.2, the management toolkit is extended with a “Performance Service”, which provides advanced performance monitoring and analysis. Using these tools and data, the user has end-to-end visibility into the state of their infrastructure and the consumption of resources by different VMs and workloads. For example, one can do performance troubleshooting, pinpoint the root cause of any issues (whether compute, network or storage) and decide on remediation actions, while using a single tool set.

All management features of Virtual SAN including enablement, configuration, upgrades, as well as the health and performance services, are built using a scalable distributed architecture. They are fully supported through APIs, which are extensions of the very popular vSphere API. In the same spirit, all VSAN and HCI management features are natively integrated with the vSphere Web Client.

Virtual SAN is natively integrated with ESXi, the best hypervisor in the industry, to offer unmatched HCI efficiency and performance. The product works “out of the box” for vSphere customers. Moreover, the scale-out nature of its control plane opens the path for new use cases and applications in private and public clouds. We are just at the beginning of an exciting journey.

Much (digital) ink has been spilled describing Virtual SAN, its unique architecture, and the operational benefits it brings to customers. Not a day goes by without an article been published, without a friend or a foe expressing their opinion on the product and what it means for the customers. This book, however, is different. It is the most comprehensive and authoritative reference for Virtual SAN, the product and the technology behind it.

Duncan and Cormac have been on board since the early days of the project. They have contributed invaluable input during the design phases based on their personal experience and on customer feedback. Not only they are intimately familiar with the architecture of the product, but they also have extensive operational experience in the lab and with customer deployments around the world. As such, they approach the subject through the eyes of the IT professional.

I hope you find the book as informative as I have.

Christos Karamanolis
VMware Fellow. CTO, Storage and Availability Business Unit.

About the Authors

Cormac Hogan is a Senior Staff Engineer in the Office of the CTO in the Storage and Availability business unit at VMware. Cormac was one of the first VMware employees at the EMEA headquarters in Cork, Ireland, back in 2005, and has previously held roles in VMware's Technical Marketing, Integration Engineering and Support organizations. Cormac has written a number of storage-related white papers and has given numerous presentations on storage best practices and new features. Cormac is the owner of CormacHogan.com, a blog site dedicated to storage and virtualization.

He can be followed on twitter @CormacJHogan.

Duncan Epping is a Chief Technologist working for VMware in the Office of CTO of the Storage and Availability business unit. Duncan is responsible for ensuring VMware's future innovations align with essential customer needs, translating customer problems to opportunities, and function as the global lead evangelist for Storage and Availability. Duncan specializes in Software Defined Storage, hyper-converged infrastructures and business continuity/disaster recovery solutions. He has four patents pending and one granted on the topic of availability, storage and resource management. Duncan is the author/owner of VMware Virtualization blog Yellow-Bricks.com and has various books on the topic of VMware including the “vSphere Clustering Deepdive” series.

He can be followed on twitter @DuncanYB.

This page intentionally left blank

About the Technical Reviewers

Christian Dickmann is a Virtual SAN Architect and Sr. Staff Engineer in the Storage and Availability Business Unit at VMware. Besides system architecture, he currently focuses on management functionality for ease of use, troubleshooting, monitoring, installation, and lifecycle management. Since joining VMware in 2007, he worked on extensibility of the networking stack as well as built the vSphere R&D dev/test cloud from the ground up. The latter experience made him an avid user experience and customer advocate.

John Nicholson is a Senior Technical Marketing Manager in the Storage and Availability Business Unit. He focuses on delivering technical guidance around VMware Virtual SAN solutions. John previously worked for partners and customers in architecting and implementing enterprise storage and VDI solutions. He holds the VCP certification (VCP5-DT, VCP6-DCV) and is a member of the vExpert and Veeam Vanguard programs.

This page intentionally left blank

Acknowledgments

The authors of this book both work for VMware. The opinions expressed in the book are the authors' personal opinions and experience with the product. Statements made throughout the book do not necessarily reflect the views and opinions of VMware.

We would like to thank Christian Dickmann and John Nicholson for keeping us honest as our technical editors. Of course, we also want to thank the Virtual SAN engineering team for their help and patience. In particular, we want to call out a couple of individuals of the Storage and Availability BU, Christian Dickmann, Paudie O'Riordan and Christos Karamanolis, whose deep knowledge and understanding of VSAN, and storage technology in general, was leveraged throughout this book. We also want to thank William Lam for his contributions to the book.

Lastly, we want to thank our VMware management team (Yanbing Li, Christos Karamanolis) for supporting us on this and other projects.

Go VSAN!

Cormac Hogan and Duncan Epping

This page intentionally left blank

We Want to Hear from You!

As the reader of this book, *you* are our most important critic and commentator. We value your opinion and want to know what we're doing right, what we could do better, what areas you'd like to see us publish in, and any other words of wisdom you're willing to pass our way.

We welcome your comments. You can email or write us directly to let us know what you did or didn't like about this book—as well as what we can do to make our books better.

Please note that we cannot help you with technical problems related to the topic of this book.

When you write, please be sure to include this book's title and author as well as your name, email address, and phone number. We will carefully review your comments and share them with the author and editors who worked on the book.

Email: VMwarePress@vmware.com

Mail: VMware Press

ATTN: Reader Feedback

800 East 96th Street

Indianapolis, IN 46240 USA

Reader Services

Register your copy of *Essential Virtual SAN (VSAN)* at www.informit.com for convenient access to downloads, updates, and corrections as they become available. To start the registration process, go to informit.com/register and log in or create an account*. Enter the product ISBN, 9780133854992, and click Submit. Once the process is complete, you will find any available bonus content under Registered Products.

*Be sure to check the box that you would like to hear from us in order to receive exclusive discounts on future editions of this product.

Introduction

When talking about virtualization and the underlying infrastructure that it runs on, one component that always comes up in conversation is storage. The reason for this is fairly simple: In many environments, storage is a pain point. Although the storage landscape has changed with the introduction of flash technologies that mitigate many of the traditional storage issues, many organizations have not yet adopted these new architectures and are still running into the same challenges.

Storage challenges range from operational effort or complexity to performance problems or even availability constraints. The majority of these problems stem from the same fundamental problem: legacy architecture. The reason is that most storage platform architectures were developed long before virtualization existed, and virtualization changed the way these shared storage platforms were used.

In a way, you could say that virtualization forced the storage industry to look for new ways of building storage systems. Instead of having a single server connect to a single storage device (also known as a logical unit or LUN for short), virtualization typically entails having one (or many) physical server(s) running many virtual machines connecting to one or multiple storage devices. This did not only increase the load on these storage systems, it also changed the workload patterns and increased the total capacity required.

As you can imagine, for most storage administrators, this required a major shift in thinking. What should the size of my LUN be? What are my performance requirements, and how many spindles will that result in? What kind of data services are required on these LUNs, and where will virtual machines be stored? Not only did it require a major shift in thinking, but it also required working in tandem with other IT teams. Whereas in the past server admins and network and storage admins could all live in their own isolated worlds, they now needed to communicate and work together to ensure availability of the platform they were building. Whereas in the past a mistake, such as a misconfiguration or underprovisioning, would only impact a single server, it could now impact many virtual machines.

There was a fundamental shift in how we collectively thought about how to operate and architect IT infrastructures when virtualization was introduced. Now another collective shift is happening all over again. This time it is due to the introduction of software-defined networking and software-defined storage. But let's not let history repeat itself, and let's avoid the mistakes we all made when virtualization first arrived. Let's all have frank and open discussions with our fellow datacenter administrators as we all aim to revolutionize datacenter architecture and operations!

You, the Reader

This book is targeted at IT professionals who are involved in the care and feeding of a VMware vSphere environment. Ideally, you have been working with VMware vSphere for some time and perhaps you have attended an authorized course in vSphere, such as the “Install, Configure, and Manage” class. This book is not a starters guide, but there should be enough in the book for administrators and architects of all levels.

How to Use This Book

This book is split into ten chapters, as described here:

- **Chapter 1, “Introduction to VSAN”:** This chapter provides a high-level introduction to software-defined storage and VSAN.
- **Chapter 2, “VSAN Prerequisites and Requirements for Deployment”:** This chapter describes the requirements from a physical and virtual perspective to safely implement VSAN.
- **Chapter 3, “VSAN Installation and Configuration”:** This chapter goes over the steps needed to install and configure VSAN.
- **Chapter 4, “VM Storage Policies on VSAN”:** This chapter explains the concept of storage policy-based management.
- **Chapter 5, “Architectural Details”:** This chapter provides in-depth architectural details of VSAN.
- **Chapter 6, “VM Storage Policies and Virtual Machine Provisioning”:** This chapter describes how VM storage policies can be used to simplify VM deployment.
- **Chapter 7, “Management and Maintenance”:** This chapter describes the steps for most common management and maintenance tasks.
- **Chapter 8, “Stretched Clusters”:** This chapter covers the operational and architectural aspects and design decisions around the introduction of VSAN stretched clusters.
- **Chapter 9, “Designing a VSAN Cluster”:** This chapter provides various examples around designing a VSAN cluster, including sizing exercises.
- **Chapter 10, “Troubleshooting, Monitoring, and Performance”:** This chapter covers the various (command line) tools available to troubleshoot and monitor VSAN.

This page intentionally left blank

Chapter 4

VM Storage Policies on VSAN

In vSphere 5.0, VMware introduced a feature called profile-driven storage. Profile-driven storage is a feature that allows vSphere administrators to easily select the correct datastore on which to deploy virtual machines (VMs). The selection of the datastore is based on the capabilities of that datastore, or to be more specific, the underlying capabilities of the storage array that have been assigned to this datastore. Examples of the capabilities are RAID level, thin provisioning, deduplication, encryption, replication, etc. The capabilities are completely dependent on the storage array.

Throughout the life cycle of the VM, profile-driven storage allows the administrator to check whether its underlying storage is still *compatible*. In other words, does the datastore on which the VM resides still have the correct capabilities for this VM? The reason why this is useful is because if the VM is migrated to a different datastore for whatever reason, the administrator can ensure that it has moved to a datastore that continues to meet its requirements. If the VM is migrated to a datastore without paying attention to the capabilities of the destination storage, the administrator can still check the compliance of the VM storage from the vSphere client at any time and take corrective actions if it no longer resides on a datastore that meets its storage requirements (in other words, move it back to a compliant datastore).

However, VM storage policies and storage policy-based management (SPBM) have taken this a step further. In the previous paragraph, we described a sort of storage quality of service driven by the storage. All VMs residing on the same datastore would inherit the capabilities of the datastore. With VSAN, the storage quality of service no longer resides with the datastore; instead, it resides with the VM and is enforced by the VM storage policy associated with the VM and the VM disks (VMDKs). Once the policy is pushed down to the storage layer, in this case VSAN, the underlying storage is then responsible for creating storage for the VM that meets the requirements placed in the policy.

Introducing Storage Policy-Based Management in a VSAN Environment

VSAN leverages this approach to VM deployment, using an updated method called storage policy-based management (SPBM). All VMs deployed to a VSAN datastore must use a VM storage policy, although if one is not specifically created, a default one that is associated with the datastore is assigned to the VM. The VM storage policy contains one or more VSAN capabilities. This chapter will describe the VSAN capabilities. After the VSAN cluster has been configured and the VSAN datastore has been created, VSAN surfaces up a set of capabilities to the vCenter Server. These capabilities, which are surfaced by the vSphere APIs for Storage Awareness (VASA) storage provider (more on this shortly) when the cluster is configured successfully, are used to set the availability, capacity, and performance policies on a per-VM (and per-VMDK) basis when that VM is deployed on the VSAN datastore.

As previously mentioned, this differs significantly from the previous VM storage profile mechanism that we had in vSphere in the past. With the VM storage profile feature, the capabilities were associated with datastores, and were used for VM placement decisions. Now, through SPBM, administrators create a policy defining the storage requirements for the VM, and this policy is pushed out to the storage, which in turn instantiates per-VM (and per-VMDK) storage for virtual machines. In vSphere 6.0, VMware introduced Virtual Volumes (VVols). Storage policy-based management for VMs using VVols is very similar to storage policy-based management for VMs deployed on VSAN. In other words, administrators no longer need to carve up logical unit numbers (LUNs) or volumes for virtual machine storage. Instead, the underlying storage infrastructure instantiates the virtual machine storage based on the contents of the policy. What we have now with SPBM is a mechanism whereby we can specify the requirements of the VM, and the VMDKs. These requirements are then used to create a policy. This policy is then sent to the storage layer [in the case of VVols, this is a SAN or network-attached storage (NAS) storage array] asking it to build a storage object for this VM that meets these policy requirements. In fact, a VM can have multiple policies associated with it, different policies for different VMDKs.

By way of explaining capabilities, policies, and profiles, capabilities are what the underlying storage is capable of providing by way of availability, performance, and reliability. These capabilities are visible in vCenter Server. The capabilities are then used to create a VM storage policy (or just policy for short). A policy may contain one or more capabilities, and these capabilities reflect the requirements of your VM or application running in a VM. Previous versions of vSphere used the term *profiles*, but these are now known as *policies*.

Deploying VMs on a VSAN datastore is very different from previous approaches in vSphere. In the past, an administrator would present a LUN or volume to a group of ESXi hosts and in the case of block storage partition, format, and build a VMFS file system to create a datastore for storing VM files. In the case of network-attached storage (NAS), a network file system (NFS) volume is mounted to the ESXi host, and once again a VM

is created on the datastore. There is no way to specify a RAID-0 stripe width for these VMDKs, nor is there any way to specify a RAID-1 replica for the VMDK.

In the case of VSAN (and now VVols), the approach to deploying VMs is quite different. Consideration must be given to the availability, performance, and reliability factors of the application running in the VM. Based on these requirements, an appropriate VM storage policy must be created and associated with the VM during deployment.

There were five capabilities in the initial release of VSAN, as illustrated in Figure 4.1.

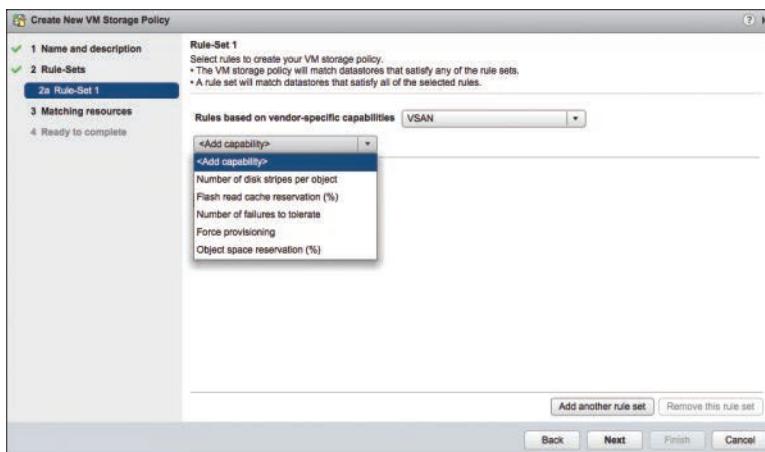


Figure 4.1 VSAN capabilities that can be used for VM storage policies

In VSAN 6.2, the number of capabilities is increased to support a number of new features. These features include the ability to implement RAID-5 and RAID-6 configurations for virtual machine objects deployed on an all-flash VSAN configuration, alongside the existing RAID-0 and RAID-1 configurations. With RAID-5 and RAID-6, it now allows VMs to tolerate one or two failures, but it means that the amount of space consumed is much less than a RAID-1 configuration to tolerate a similar amount of failures. There is also a new policy for software checksum. Checksum is enabled by default, but it can be disabled through policies if an administrator wishes to disable it. The last capability relates to quality of service and provides the ability to limit the number of input/output operations per second (IOPS) for a particular object.

You can select the capabilities when a VM storage policy is created. Note that certain capabilities are applicable to hybrid VSAN configurations (e.g., flash read cache reservation), while other capabilities are applicable to all-flash VSAN configurations (e.g., failure tolerance method set to performance).

VM storage policies are essential in VSAN deployments because they define how a VM is deployed on a VSAN datastore. Using VM storage policies, you can define the capabilities that can provide the number of VMDK RAID-0 stripe components or the number of

RAID-1 mirror copies of a VMDK. If an administrator desires a VM to tolerate one failure but does not want to consume as much capacity as a RAID-1 mirror, a RAID-5 configuration can be used. This requires a minimum of four hosts in the cluster and implements a distributed parity mechanism across the storage of all four hosts. If this configuration would be implemented with RAID-1, the amount of capacity consumed would be 200% the size of the VMDK. If this is implemented with RAID-5, the amount of capacity consumed would be 133% the size of the VMDK.

Similarly, if an administrator desires a VM to tolerate two failures using a RAID-1 mirroring configuration, there would need to be three copies of the VMDK, meaning the amount of capacity consumed would be 300% the size of the VMDK. With a RAID-6 implementation, a double parity is implemented, which is also distributed across all the hosts. For RAID-6, there must be a minimum of six hosts in the cluster. RAID-6 also allows a VM to tolerate two failures, but only consumes capacity equivalent to 150% the size of the VMDK.

Figure 4.2 shows the new policies introduced in VSAN 6.2.

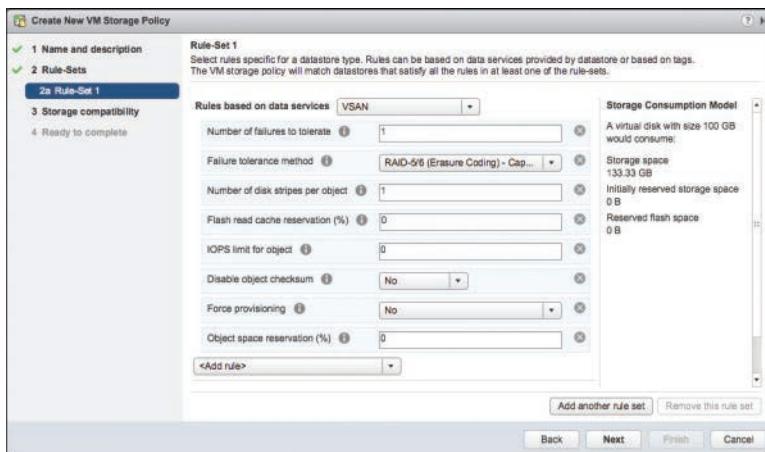


Figure 4.2 New VSAN capabilities

The sections that follow highlight where you should use these capabilities when creating a VM storage policy and when to tune these values to something other than the default. Remember that a VM storage policy will contain one or more capabilities.

In the initial release of VSAN, five capabilities were available for selection to be part of the VM storage policy. In VSAN 6.2, as previously highlighted, additional policies were introduced. As an administrator, you can decide which of these capabilities can be added to the policy, but this is, of course, dependent on the requirements of your VM. For example, what performance and availability requirements does the VM have? The capabilities are as follows:

- Number of failures to tolerate
- Number of disk stripes per object

- Failure tolerance method
- IOPS limit for object
- Disable object checksum
- Flash read cache reservation (hybrid configurations only)
- Object space reservation
- Force provisioning

The sections that follow describe the VSAN capabilities in detail.

Number of Failures to Tolerate

In this section, *number of failures to tolerate* is described having *failure tolerance method* set to its default value that is *Performance*. Later on we will describe a different scenario when *failure tolerance method* is set to *Capacity*.

This capability sets a requirement on the storage object to tolerate at least n number of failures in the cluster. This is the number of concurrent host, network, or disk failures that may occur in the cluster and still ensure the availability of the object. When the *failure tolerance method* is set to its default value of *RAID-1*, the VM's storage objects are mirrored; however, the mirroring is done across ESXi hosts, as shown in Figure 4.3.

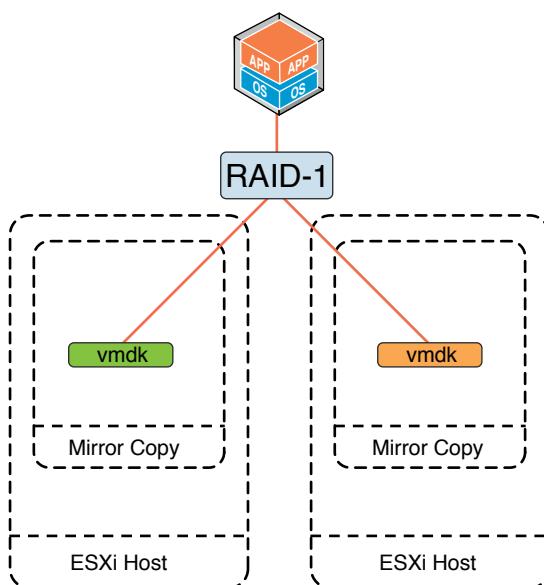


Figure 4.3 Number of failures to tolerate results in a RAID-1 configuration

When this capability is set to a value of n , it specifies that the VSAN configuration must contain at least $n + 1$ replicas (copies of the data); this also implies that there are $2n + 1$ hosts in the cluster.

Note that this requirement will create a configuration for the VM objects that may also contain an additional number of witness components being instantiated to ensure that the VM remains available even in the presence of up to *number of failures to tolerate* concurrent failures (see Table 4.1). Witnesses provide a quorum when failures occur in the cluster or a decision has to be made when a split-brain situation arises. These witnesses will be discussed in much greater detail later in the book, but suffice it to say that witness components play an integral part in maintaining VM availability during failures and maintenance tasks.

One aspect worth noting is that any disk failure on a single host is treated as a “failure” for this metric (although multiple disk failures on the same host are also treated as a single host failure). Therefore, the VM may not persist (remain accessible) if there is a disk failure on host A and a host failure of host B when number of failures to tolerate is set to 1.

Table 4.1 Witness and hosts required to meet number of failures to tolerate requirement

Number of Failures to Tolerate	Mirror Copies/Replicas	Witness Objects	Minimum Number of ESXi Hosts
0	1	0	1
1	2	1	3
2	3	2	5
3	4	3	7

Table 4.1 is true if the number of disk objects to stripe is set to 1. The behavior is subtly different if there is a stripe width greater than 1. The number of disk stripes per object will be discussed in more detail shortly.

If no policy is chosen when a VM is deployed, the default policy associated with the VSAN datastore is chosen, which in turn sets the number of failures to tolerate to 1. When a new policy is created, the default value of number of failures to tolerate is also 1. This means that even if this capability is not explicitly specified in the policy, it is implied.

Failure Tolerance Method

This is a new capability introduced in VSAN 6.2 and is how administrators can choose either RAID-1 or RAID-5/6 configuration for their virtual machine objects. The *failure tolerance method* is used in conjunction with *number of failures to tolerate*. The purpose of this

setting is to allow administrators to choose between performance and capacity. If performance is the absolute end goal for administrators, then RAID-1 (which is still the default) is the tolerance method that should be used. If administrators do not need maximum performance and are more concerned with capacity usage, then RAID-5/6 is the tolerance method that should be used. The easiest way to explain the behavior is to display the various policy settings and the resulting object configuration, as shown in Table 4.2.

As can be seen from Table 4.2, when the *failure tolerance method* RAID5/6 is selected, either RAID-5 or RAID-6 is implemented depending on the number of failures that you wish to tolerate (although it only supports a maximum setting of two for *number of failures to tolerate*). If performance is still the desired capability, then the traditional RAID-1 configuration is implemented, with the understanding that this uses mirror copies of the objects, and thus consumes significantly more space.

One might ask why RAID-5/6 less performing than RAID-1. The reason lies in I/O amplification. In steady state, where there are no failures in the cluster, there is no read amplification when using RAID-5/6 versus RAID-1. However, there is write amplification. This is because the parity component needs to be updated every time there is a write to the associated data components. So in the case of RAID-5, we need to read the component that is going to be updated with additional write data, read the current parity, merge the new write data with the current data, write this back, calculate the new parity value, and write this back also. In essence, a single write operation can amplify into two reads and two

Table 4.2 Object configuration when number of failures to tolerate and failure tolerance method are set

Number of Failures to Tolerate	Failure Tolerance Method	Object Configuration	Number of ESXi Hosts Required
0	RAID5/6 (Erasure Coding)	RAID-0	1
0	RAID-1 (mirroring)	RAID-0	1
1	RAID5/6 (Erasure Coding)	RAID-5	4
1	RAID-1 (mirroring)	RAID-1	3
2	RAID5/6 (Erasure Coding)	RAID-6	6
2	RAID-1 (mirroring)	RAID-1	5
3	RAID5/6 (Erasure Coding)	N/A	N/A
3	RAID-1 (mirroring)	RAID-1	7

writes. With RAID-6, which has double parity, a single write can amplify into three reads and three writes.

And indeed, when there is a failure of some component in the RAID-5 and RAID-6 objects, and data needs to be determined using parity, then the I/O amplification is even higher. These are the considerations an administrator needs to evaluate when deciding on the *failure tolerance method*.

One item to keep in mind is that even though *failure tolerance method* set to RAID5/6 consumes less capacity, it does require more hosts than the traditional RAID-1 approach and is only supported (in this release) on an all-flash VSAN configuration. When using RAID-1, the rule is that to tolerate n failures, there must be $2n + 1$ hosts for the mirrors/replicas and witness. So to tolerate one failure, there must be three hosts; to tolerate two failures, there must be five hosts; or to tolerate three failures, there must be seven hosts in the cluster, all contributing storage to the VSAN datastore. With *failure tolerance method* set to RAID5/6, four hosts are needed to tolerate one failure and six hosts are needed to tolerate two failures, even though less space is consumed on each host. Figure 4.4 shows an example of a RAID-5 configuration for an object, deployed across four hosts with a distributed parity.

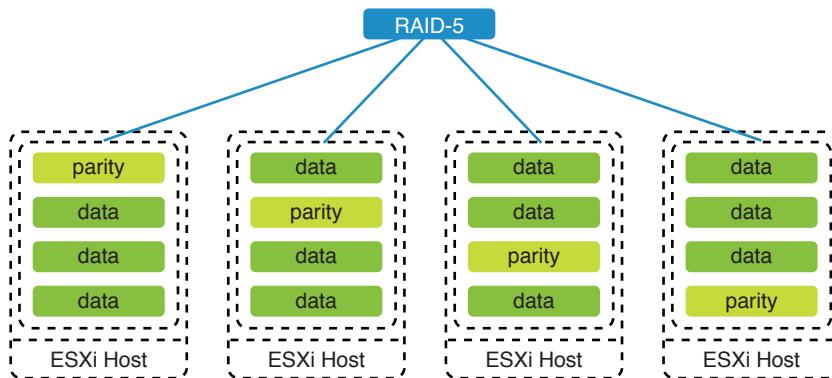


Figure 4.4 RAID-5 configuration, a result of *failure tolerance method* RAID5/6 and number of failures to tolerate set to 1

The RAID-5 or RAID-6 configurations also work with *number of disk stripes per object*. If stripe width is also specified as part of the policy along with *failure tolerance method* set to RAID5/6 each of the components on each host is striped in a RAID-0 configuration, and these are in turn placed in either a RAID-5 or-6 configuration.

One final note is in relation to having a number of failures to tolerate setting of zero or three. If you deploy a VM with this policy setting, which includes a *failure tolerance method*

RAID5/6 setting, the VM provisioning wizard will display a warning stating that this policy setting is only effective when the number of failures to tolerate is set to either one or two. You can still proceed with the deployment, but the object is deployed as a single RAID-0 object.

Number of Disk Stripes Per Object

This capability defines the number of physical disks across which each replica of a storage object (e.g., VMDK) is striped. When *failure tolerance method* is set to performance, this policy setting can be considered in the context of a RAID-0 configuration on each RAID-1 mirror/replica where I/O traverses a number of physical disk spindles. When *failure tolerance method* is set to capacity, each component of the RAID-5 or RAID-6 stripe may also be configured as a RAID-0 stripe. Typically, when the number of disk stripes per object is defined, the number of failures to tolerate is also defined. Figure 4.5 shows what a combination of these two capabilities could result in, once again assuming that the new VSAN 6.2 policy setting of *failure tolerance method* is set to its default value *RAID-1*.

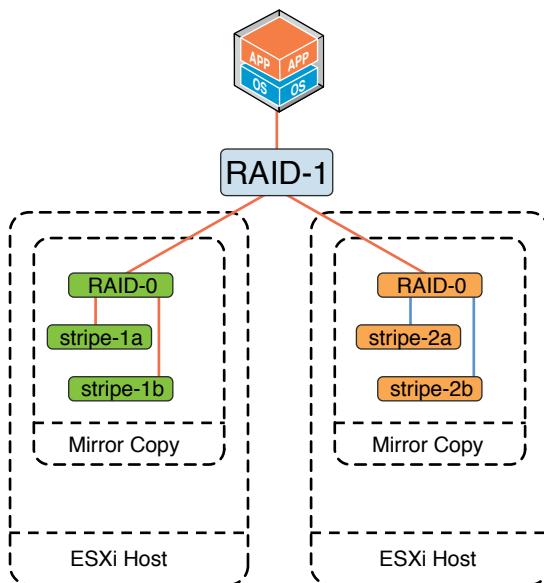


Figure 4.5 Storage object configuration when stripe width set is to 2 and failures to tolerate is set to 1 and replication method optimizes for is not set

To understand the impact of stripe width, let's examine it first in the context of write operations and then in the context of read operations.

Because all writes go to the cache device write buffer, the value of an increased stripe width may or may not improve performance. This is because there is no guarantee that the new stripe will use a different cache device; the new stripe may be placed on a capacity device in the same disk group, and thus the new stripe will use the same cache device. If the new stripe is placed in a different disk group, either on the same host or on a different host, and thus leverages a different cache device, performance might improve. However, you as the vSphere administrator have no control over this behavior. The only occasion where an increased stripe width could definitely add value is when there is a large amount of data to destage from the cache tier to the capacity tier. In this case, having a stripe could improve destage performance.

From a read perspective, an increased stripe width will help when you are experiencing many read cache misses, but note that this is a consideration in hybrid configurations only. All-flash VSAN considerations do not have a read cache. Consider the example of a VM deployed on a hybrid VSAN consuming 2,000 read operations per second and experiencing a hit rate of 90%. In this case, there are still 200 read operations that need to be serviced from magnetic disk in the capacity tier. If we make the assumption that a single magnetic disk can provide 150 input/output operations per second (IOPS), then it is obvious that it is not able to service all of those read operations, so an increase in stripe width would help on this occasion to meet the VM I/O requirements. In an all-flash VSAN, which is extremely read intensive, striping across multiple capacity flash devices can also improve performance.

In general, the default stripe width of 1 should meet most, if not all VM workloads. Stripe width is a capability that should change only when write destaging or read cache misses are identified as a performance constraint.

IOPS Limit for Object

IOPS limit for object is a new Quality of Service (QoS) capability introduced with VSAN 6.2. This allows administrators to ensure that an object, such as a VMDK, does not generate more than a predefined number of I/O operations per second. This is a great way of ensuring that a “noisy neighbor” virtual machine does not impact other virtual machine components in the same disk group by consuming more than its fair share of resources. By default, VSAN uses an I/O size of 32 KB as a base. **This means that a 64 KB I/O will therefore represent two I/O operations in the limits calculation.** I/Os that are less than or equal to 32 KB will be considered single I/O operations. For example, 2×4 KB I/Os are considered as two distinct I/Os. It should also be noted that both read and write IOPS are regarded as equivalent. Neither cache hit rate nor sequential I/O are taken into account. If the IOPS limit threshold is passed, the I/O is throttled back to bring the IOPS value back under the threshold. The default value for this capability is 0, meaning that there is no IOPS limit threshold and VMs can consume as many IOPS as they want, subject to available resources.

Flash Read Cache Reservation

This capability is applicable to hybrid VSAN configurations only. It is the amount of flash capacity reserved on the cache tier device as read cache for the storage object. It is specified as a percentage of the logical size of the storage object (i.e., VMDK). This is specified as a percentage value (%), with up to four decimal places. This fine granular unit size is needed so that administrators can express sub 1% units. Take the example of a 1 TB VMDK. If you limited the read cache reservation to 1% increments, this would mean cache reservations in increments of 10 GB, which in most cases is far too much for a single VM.

Note that you do not have to set a reservation to allow a storage object to use cache. All VMs equally share the read cache of cache devices. The reservation should be left unset (default) unless you are trying to solve a real performance problem and you believe dedicating read cache is the solution. If you add this capability to the VM storage policy and set it to a value 0 (zero), however, you will not have any read cache reserved to the VM that uses this policy. In the current version of VSAN, there is no proportional share mechanism for this resource when multiple VMs are consuming read cache, so every VM consuming read cache will share it equally.

Object Space Reservation

All objects deployed on VSAN are thinly provisioned. This means that no space is reserved at VM deployment time, but rather space is consumed as the VM uses storage. The object space reservation capability defines the percentage of the logical size of the VM storage object that may be reserved during initialization. The object space reservation is the amount of space to reserve specified as a percentage of the total object address space. This is a property used for specifying a thick provisioned storage object. If object space reservation is set to 100%, all of the storage capacity requirements of the VM storage are reserved up front (thick). This will be lazy zeroed thick (LZT) format and not eager zeroed thick (EZT). The difference between LZT and EZT is that EZT virtual disks are zeroed out at creation time; LZT virtual disks are zeroed out gradually at first write time.

One thing to bring to the readers' attention is the special case of using object space reservation when deduplication and compression are enabled on the VSAN cluster. When deduplication and compression are enabled, any objects that wish to use object space reservation in a policy must have it set to either 0% (no space reservation) or 100% (fully reserved). Values between 1% and 99% are not allowed. Any existing objects that have object space reservation between 1% and 99% will need to be reconfigured with 0% or 100% prior to enabling deduplication and compression on the cluster.

Force Provisioning

If the force provisioning parameter is set to a nonzero value, the object that has this setting in its policy will be provisioned even if the requirements specified in the VM storage

policy cannot be satisfied by the VSAN datastore. The VM will be shown as noncompliant in the VM summary tab and relevant VM storage policy views in the vSphere client. If there is not enough space in the cluster to satisfy the reservation requirements of at least one replica, however, the provisioning will fail even if force provisioning is turned on. When additional resources become available in the cluster, VSAN will bring this object to a compliant state.

One thing that might not be well understood regarding *force provisioning* is that if a policy cannot be met, it attempts a much simpler placement with requirements which reduces to *number of failures to tolerate* to 0, *number of disk stripes per object* to 1, and *flash read cache reservation* to 0 (on hybrid configurations). This means Virtual SAN will attempt to create an object with just a single copy of data. Any *object space reservation (OSR)* policy setting is still honored. Therefore there is no gradual reduction in capabilities as VSAN tries to find a placement for an object. For example, if policy contains *number of failures to tolerate* = 2, VSAN won't attempt an object placement using *number of failures to tolerate* = 1. Instead, it immediately looks to implement *number of failures to tolerate* = 0.

Similarly, if the requirement was *number of failures to tolerate* = 1, *number of disk stripes per object* = 4, but Virtual SAN doesn't have enough capacity devices to accommodate *number of disk stripes per object* = 4, then it will fall back to *number of failures to tolerate* = 0, *number of disk stripes per object* = 1, even though a policy of *number of failures to tolerate* = 1, *number of disk stripes per object* = 2 or *number of failures to tolerate* = 1, *number of disk stripes per object* = 3 may have succeeded.

Caution should be exercised if this policy setting is implemented. Since this allows VMs to be provisioned with no protection, it can lead to scenarios where VMs and data are at risk.

Administrators who use this option to force provision virtual machines need to be aware that although virtual machine objects may be provisioned with only one replica copy (perhaps due to lack of space), once additional resources become available in the cluster, VSAN may immediately consume these resources to try to satisfy the policy settings of virtual machines.

Some commonly used cases where force provisioning is used are (a) when boot-strapping a VSAN management cluster, starting with a single node that will host the vCenter Server, which is then used to configure a larger VSAN cluster, and (b) when allowing the provisioning of virtual machine/desktops when a cluster is under maintenance, such as a virtual desktop infrastructure (VDI) running on VSAN.

Remember that this parameter should be used only when *absolutely* needed and as an exception. When used by default, this could easily lead to scenarios where VMs, and all data associated with it, are at risk. Use with caution!

Disable Object Checksum

VSAN 6.2 introduced this new capability. This feature, which is enabled by default, is looking for data corruption (bit rot), and if found, automatically corrects it. Checksum is validated on the complete I/O path, which means that when writing data the checksum is calculated and automatically stored. Upon a read the checksum of the data is validated, and if there is a mismatch, the data is repaired. VSAN 6.2 also includes a scrubber mechanism. This mechanism is configured to run once a year (by default) to check all data on the VSAN datastore; however, this value can be changed by setting an advanced host setting. We recommend leaving this configured to the default value of once a year. In some cases you may desire to disable checksums completely. The reason for this could be performance, although the overhead is negligible and most customers prefer data integrity over a 1% to 3% performance increase. However in some cases, this performance increase may be desired. Another reason for disabling checksums is in the situation where the application already provides a checksum mechanism, or the workload does not require checksum. If that is the case, then checksums can be disabled through the “disable object checksum capability,” which should be set to “Yes” to disable it.

That completes the capabilities overview. Let’s now look at some other aspects of the storage policy-based management mechanism.

VASA Vendor Provider

As part of the VSAN cluster creation step, each ESXi host has a VSAN storage provider registered with vCenter. This uses the vSphere APIs for Storage Awareness (VASA) to surface up the VSAN capabilities to the vCenter Server. The capabilities can then be used to create VM storage policies for the VMs deployed on the VSAN datastore. If you are familiar with VASA and have used it with traditional storage environments, you’ll find this functionality familiar; however, with traditional storage environments that leverage VASA, some configuration work needs to be done to add the storage provider for that particular storage. In the context of VSAN, a vSphere administrator does not need to worry about registering these; these are automatically registered when a VSAN cluster is created.

An Introduction to VASA

VASA allows storage vendors to publish the capabilities of their storage to vCenter Server, which in turn can display these capabilities in the vSphere Web Client. VASA may also provide information about storage health status, configuration info, capacity and thin provisioning info, and so on. VASA enables VMware to have an end-to-end story regarding storage. Traditionally, this enabled storage arrays to inform the VASA storage provider of capabilities, and then the storage provider informed vCenter Server,

so now users can see storage array capabilities from vSphere Web Client. Through VM storage policies, these storage capabilities are used in the vSphere Web Client to assist administrators in choosing the right storage in terms of space, performance, and service-level agreement (SLA) requirements. This was true for both traditional storage arrays, and now it is true for VSAN also. Prior to the release of virtual volumes (VVols), there was a notable difference in workflow when using VASA and VM storage policies when comparing traditional storage to VSAN. With traditional storage, VASA historically surfaced information about the datastore capabilities, and a vSphere administrator had to choose the appropriate storage on which to place the VM. With VSAN, and now VVols, you define the capabilities you want to have for your VM storage in a VM storage policy. This policy information is then pushed down to the storage layer, basically informing it that these are the requirements you have for storage. VASA will then tell you whether the underlying storage (e.g., VSAN) can meet these requirements, effectively communicating compliance information on a per-storage object basis. The major difference is that this functionality is now working in a bidirectional mode. Previously, VASA would just surface up capabilities. Now it not only surfaces up capabilities but also verifies whether a VM's storage requirements are being met based on the contents of the policy.

Storage Providers

Figure 4.6 illustrates an example of what the storage provider looks like. When a VSAN cluster is created, the VASA storage provider from every ESXi host in the cluster is registered to the vCenter Server. In a four-node VSAN cluster, the VASA VSAN storage provider configuration would look similar to this.

The screenshot shows the vSphere Web Client interface for managing storage providers. At the top, there is a navigation bar with tabs: Getting Started, Summary, Monitor, Manage, and Related Objects. Below this is a secondary navigation bar with tabs: Settings, Scheduled Tasks, Alarm Definitions, Tags, Permissions, Sessions, and Storage Providers. The Storage Providers tab is selected.

The main content area displays a table titled "Storage Providers". The table has columns: Storage Provider/Storage System, Status, Active/Standby, Priority, URL, Last Report Time, VASA API Version, and Certificate Expiry. There are four entries in the table:

- VASA Provider esxi-0-and-rainpole - Online, Active, 128, https://esxi-0-and-rainpole.com, 12/10/2015 3:..., 1.0, 1821 days
- VASA Provider esxi-0-and-rainpole - Online, Standby, 3, https://esxi-0-and-rainpole.com, —, 1.0, 1821 days
- VASA Provider esxi-0-pre-rainpole - Online, Standby, 128, https://esxi-0-pre.rainpole.com, —, 1.0, 1821 days (highlighted with a blue background)
- VASA Provider esxi-0-pre-rainpole - Online, Standby, 128, https://esxi-0-pre.rainpole.com, —, 1.0, 1821 days

Below the table, there is a section titled "Storage Provider Details" with tabs: General and General. The General tab is selected. It contains the following information:

Supported vendor ID#:	Provider name: VASA Provider esxi-0-pre.rainpole.com
Certificate info:	Provider status: Online
	Advertiser status: —
	Activation: Automatic
URL:	https://esxi-0-pre.rainpole.com/808/Discovery.xml
Provider version:	1.0
VASA API version:	1.0
Default namespace:	VSAN

Figure 4.6 VSAN storage providers, added when the VSAN cluster is created

You can always check the status of the storage providers by navigating in the Web Client to the vCenter Server inventory item, selecting the **Manage** tab and then the **Storage Providers** view. One VSAN provider should always be online. The other storage providers should be in standby mode. This is all done automatically by VSAN. There is typically no management of the VASA providers required by administrators.

In VSAN clusters that have more than eight ESXi hosts, and thus more than eight VASA storage providers, the list of storage providers is shortened to eight in the user interface (UI) for display purposes. The number of standby storage providers is still displayed correctly; you simply won't be able to interrogate them.

VSAN Storage Providers: Highly Available

You might ask why every ESXi host registers this storage provider. The reason for this is high availability. Should one ESXi host fail, another ESXi host in the cluster can take over the presentation of these VSAN capabilities. If you examine the storage providers shown in Figure 4.6, you will see that only one of the VSAN providers is online. The other storage providers from the other two ESXi hosts in this three-node cluster are in a standby state. Should the storage provider that is currently active go offline or fail for whatever reason (most likely because of a host failure), one of the standby providers will be promoted to active.

There is very little work that a vSphere administrator needs to do with storage providers to create a VSAN cluster. This is simply for your own reference. However, if you do run into a situation where the VSAN capabilities are not surfacing up in the VM storage policies section, it is worth visiting this part of the configuration and verifying that at least one of the storage providers is active. If you have no active storage providers, you will not discover any VSAN capabilities when trying to build a VM storage policy. At this point, as a troubleshooting step, you could consider doing a refresh of the storage providers by clicking on the refresh icon (orange circular arrows) in the storage provider screen.

What should be noted is that the VASA storage providers do not play any role in the data path for VSAN. If storage providers fail, this has no impact on VMs running on the VSAN datastore. The impact of not having a storage provider is lack of visibility into the underlying capabilities, so you will not be able to create new storage policies. However, already running VMs and policies are unaffected.

Changing VM Storage Policy On-the-Fly

Being able to change a VM storage policy on-the-fly is quite a unique aspect of VSAN. We will use an example to explain the concept of how you can change a VM storage policy

on-the-fly and how it changes the layout of a VM without impacting the application or the guest operating system running in the VM.

Consider the following scenario, briefly mentioned earlier in the context of stripe width. A vSphere administrator has deployed a VM on a hybrid VSAN configuration with the default VM storage policy, which is that the VM storage objects should have no disk striping and should tolerate one failure. The layout of the VM disk file would look something like Figure 4.7.

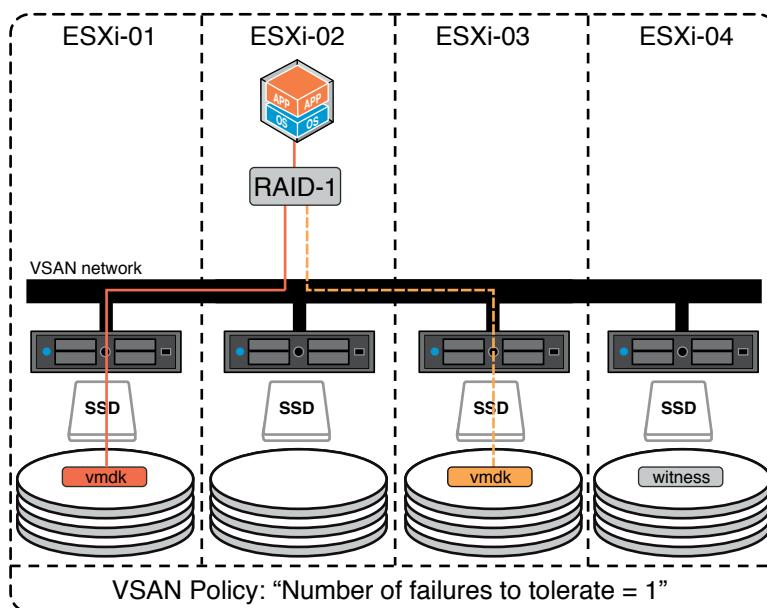


Figure 4.7 VSAN policy with the capability number of failures to tolerate = 1

The VM and its associated applications initially appeared to perform satisfactorily with a 100% cache hit rate; however, over time, an increasing number of VMs were added to the VSAN cluster. The vSphere administrator starts to notice that the VM deployed on VSAN is getting a 90% read cache hit rate. This implies that 10% of reads need to be serviced from magnetic disk/capacity tier. At peak time, this VM is doing 2,000 read operations per second. Therefore, there are 200 reads that need to be serviced from magnetic disk (the 10% of reads that are cache misses). The specifications on the magnetic disks imply that each disk can do 150 IOPS, meaning that a single disk cannot service these additional 200 IOPS. To meet the I/O requirements of the VM, the vSphere administrator correctly decides to create a RAID-0 stripe across two disks.

On VSAN, the vSphere administrator has two options to address this.

The first option is to simply modify the VM storage policy currently associated with the VM and add a stripe width requirement to the policy; however, this would change the storage layout of all the other VMs using this same policy.

Another approach is to create a brand-new policy that is identical to the previous policy but has an additional capability for stripe width. This new policy can then be attached to the VM (and VMDKs) suffering from cache misses. Once the new policy is associated with the VM, the administrator can synchronize the new/updated policy with the VM. This can be done immediately, or can be deferred to a maintenance window if necessary. If it is deferred, the VM is shown as noncompliant with its new policy. When the policy change is implemented, VSAN takes care of changing the underlying VM storage layout required to meet the new policy, *while the VM is still running* without the loss of any failure protection. It does this by mirroring the new storage objects with the additional components (in this case additional RAID-0 stripe width) to the original storage objects.

As seen, the workflow to change the VM storage policy can be done in two ways; either the original current VM storage policy can be edited to include the new capability of a stripe width = 2 or a new VM storage policy can be created that contains the failures to tolerate = 1 and stripe width = 2. The latter is probably more desirable because you may have other VMs using the original policy, and editing that policy will affect all VMs using it. When the new policy is created, this can be associated with the VM and the storage objects in a number of places in the vSphere Web Client. In fact, policies can be changed at the granularity of individual VM storage objects (e.g., VMDK) if necessary.

After making the change, the new components reflecting the new configuration (e.g., a RAID-0 stripe) will enter a state of reconfiguring. This will temporarily build out additional replicas or components, in addition to keeping the original replicas/components, so additional space will be needed on the VSAN datastore to accommodate this on-the-fly change. When the new replicas or components are ready and the configuration is completed, the original replicas/components are discarded.

Note that not all policy changes require the creation of new replicas or components. For example, adding an IOPS limit, or reducing the number of failures to tolerate, or reducing space reservation does not require this. However, in many cases, policy changes will trigger the creation of new replicas or components.

Your VM storage objects may now reflect the changes in the Web Client, for example, a RAID-0 stripe as well as a RAID-1 replica configuration, as shown in Figure 4.8.

Compare this to the tasks you may have to perform on many traditional storage arrays to achieve this. It would involve, at the very least, the following:

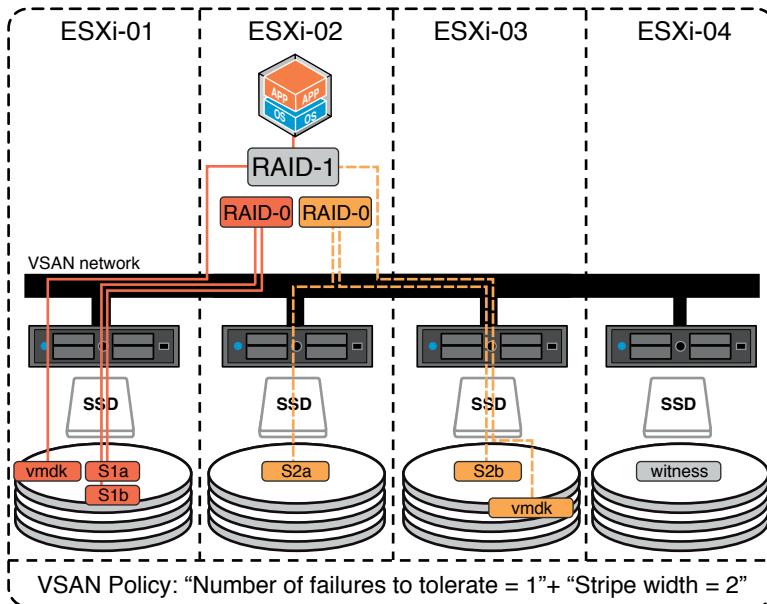


Figure 4.8 VSAN RAID-0 and RAID-1 configuration

- The migration of VMs from the original datastore.
- The decommissioning of the said LUN/volume.
- The creation of a new LUN with the new storage requirements (different RAID level).
- Possibly the reformatting of the LUN with VMFS in the case of block storage.
- Finally, you have to migrate your VMs back to the new datastore.

In the case of VSAN, after the new storage replicas or components have been created and synchronized, the older storage replicas and/or components will be automatically removed. Note that VSAN is capable of striping across disks, disk groups, and hosts when required, as depicted in Figure 4.8, where stripes S1a and S1b are located on the same host but stripes S2a and S2b are located on different hosts. It should also be noted that VSAN can create the new replicas or components without the need to move any data between hosts; in many cases the new components can be instantiated on the same storage on the same host.

We have not shown that there are, of course, additional witness components that could be created with such a change to the configuration. For a VM to continue to access all its components, a full replica copy of the data must be available and more than 50% of the components (votes) of that object must also be available in the cluster. Therefore, changes to the VM storage policy could result in additional witness components being created, or indeed, in the case of introducing a policy with less requirements, there could be fewer witnesses.

You can actually see the configuration changes taking place in the vSphere UI during this process. Select the VM that is being changed, click its **manage** tab, and then choose the **VM storage policies** view, as shown in Figure 4.9. Although this view does not show all the VM storage objects, it does display the VM home namespace, and the VMDKs are visible.

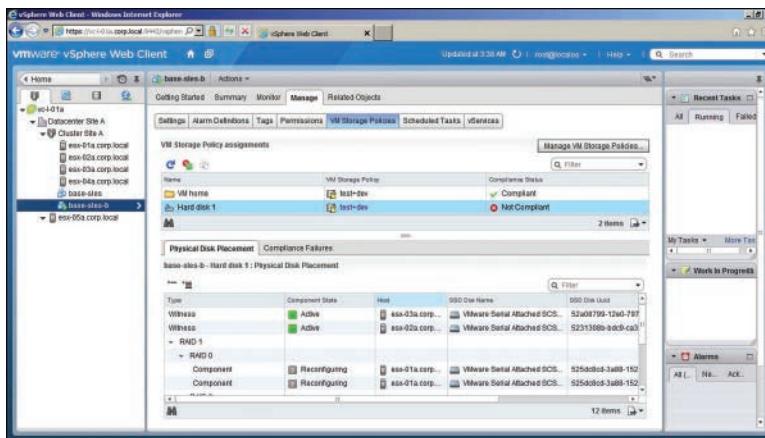


Figure 4.9 VM Storage Policy view in the vSphere client showing component reconfiguration

In VSAN 6.0, there is also a way to examine all resyncing components. Select the VSAN cluster object in the vCenter inventory, then select monitor, Virtual SAN, and finally “resyncing components” in the menu. This will display all components that are currently resyncing/rebuilding. Figure 4.10 shows the resyncing dashboard view, albeit without any resyncing activity taking place.

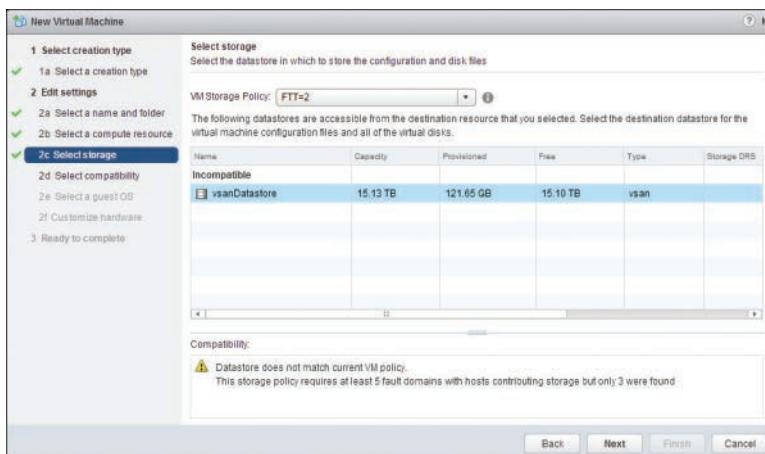


Figure 4.10 Resyncing activity as seen from the vSphere Web Client

Objects, Components, and Witnesses

A number of new concepts have been introduced in this chapter so far, including some new terminology. Chapter 5, “Architectural Details,” covers in greater detail objects, components, and indeed witness disks, as well as which VM storage objects are impacted by a particular capability in the VM storage policy. For the moment, it is enough to understand that on VSAN, a VM is no longer represented by a set of files but rather a set of storage objects. There are five types of storage objects:

- VM home namespace
- VMDKs
- VM swap
- Snapshot delta disks
- Snapshot memory

Although the vSphere Web Client displays only the VM home namespace and the VMDKs (hard disks) in the VM > monitor > policies > physical disk placement, snapshot deltas and VM swap can be viewed in the cluster > monitor > Virtual SAN > virtual disks view. We will also show ways of looking at detailed views of all the storage objects, namely delta and VM swap, in Chapter 10, “Troubleshooting, Monitoring, and Performance,” when we look at various monitoring tools available to VSAN.

VM Storage Policies

VM storage policies work in an identical fashion to storage profiles introduced in vSphere 5.0, insofar as you simply build a policy containing your VM provisioning requirements. There is a major difference in how storage policies work when compared to the original storage profiles feature. With storage profiles, you simply used the requirements in the policy to select an appropriate datastore when provisioning the VM. The storage policies not only select the appropriate datastore, but also inform the underlying storage layer that there are also certain availability and performance requirements associated with this VM. So while the VSAN datastore may be the destination datastore when the VM is provisioned with a VM storage policy, settings within the policy will stipulate additional requirements. For example, it may state that this VM has a requirement for a number of replica copies of the VM files for availability, a stripe width and read cache requirement for high performance, and a thin provisioning requirement.

VM storage policies are held inside VSAN, as well as being stored in the vCenter inventory database. Every object stores its policy inside its own metadata. This means that

vCenter is not required for VM storage policy enforcement. So if for some reason the vCenter Server is unavailable, policies can continue to be enforced.

Enabling VM Storage Policies

In the initial release of VSAN, VM storage policies could be enabled or disabled via the UI. This option is not available in later releases. However, VM storage policies are automatically enabled on a cluster when VSAN is enabled on the cluster. Although VM storage policies are normally only available with certain vSphere editions, a VSAN license will also provide this feature.

Creating VM Storage Policies

vSphere administrators have the ability to create multiple policies. As already mentioned, a number of VSAN capabilities are surfaced up by VASA related to availability and performance, and it is at this point that the administrator must decide what the requirements are for the applications running inside of the VMs from a performance and availability perspective. For example, how many component failures (hosts, network, and disk drives) does the administrator require this VM to tolerate and continue to function? Also, is the application running in this VM demanding from an IOPS perspective? If so, an adequate read cache should be provided as a possible requirement so that the performance requirement is met. Other considerations include whether the VM should be thinly provisioned or thickly provisioned, if RAID-5 or RAID-6 configurations are desired to save storage space, if checksum should be disabled, or if an IOPS limit is required for a particular VM to avoid a “noisy neighbor” situation.

Another point to note is that since vSphere 5.5, policies also support the use of tags for provisioning. Therefore, instead of using VSAN datastore capabilities for the creation of requirements within a VM storage policy, tag-based policies may also be created. The use of tag-based policies is outside the scope of this book, but further information may be found in the generic vSphere storage documentation.

Assigning a VM Storage Policy During VM Provisioning

The assignment of a VM storage policy is done during the VM provisioning. At the point where the vSphere administrator must select a destination datastore, the appropriate policy is selected from the drop-down menu of the available VM storage policies. The datastores are then separated into compatible and incompatible datastores, allowing the vSphere administrator to make the appropriate and correct choice for VM placement.

This matching of datastores does not necessarily mean that the datastore will meet the requirements in the VM storage policy. What it means is that the datastore understands

the set of requirements placed in the policy. It may still fail to provision this VM if there are not enough resources available to meet the requirements placed in the policy. However, if a policy cannot be met, the compatibility section in the lower part of the screen displays a warning that states why a policy may not be met.

This three-node cluster example shows a policy that contains a *number of failures to tolerate* = 2. A three-node cluster cannot meet this policy, but when the policy was originally created, the VSAN datastore shows up as a matching resource as it understood the contents of the policy. However, on trying to use this policy when deploying a VM, the VSAN datastore shows up as noncompliant, as Figure 4.11 demonstrates.

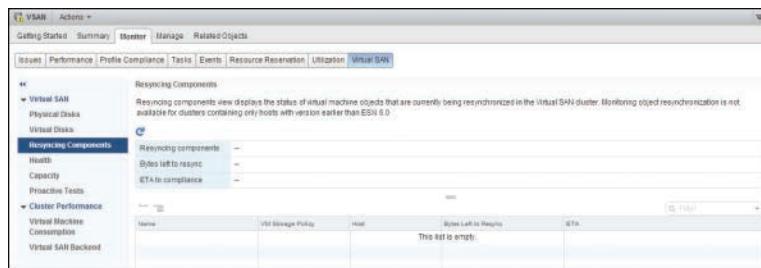


Figure 4.11 The VSAN datastore is shown as noncompliant when a policy cannot be met

This is an important point to keep in mind: Just because VSAN tells you it is compatible with a particular policy when that policy is originally created, this in no way implies that it can deploy a VM that uses the policy.

Summary

You may have used VM storage profiles in the past. VM storage policies differ significantly. Although we continue to use VASA, the vSphere APIs for Storage Awareness, VM storage policies have allowed us to switch the storage characteristics away from datastores and to the VMs. VMs, or more specifically the applications running in VMs, can now specify their own requirements using a policy that contains underlying storage capabilities around performance, reliability, and availability.

Index

A

absent components, 185
administrators, views of VSAN, 9–11
admission control, vSphere HA, 50–51
agent role (hosts), 97–98
all-flash VSAN
 read cache in, 101
 VSAN read on, anatomy of, 103
`apply_license_to_cluster` command, 297
Ask VMware link, 246

B

boot considerations, ESXi, 15
bootstrapping, vCenter server, 197–199

C

cache devices, 14–15, 21–22, 99
 failure, 186–187

`cacheReservation` command option, 266
cache tier destaging, 117–118
caching algorithms, I/O flow, 100
capacity devices, 14–15, 19–20, 99
 failure, 185–186
`check_limits` command, 297–298
`check_state` command, 292–293
checksum
 data integrity through, 130–131
 object, disabling, 73
CLI (command-line interface), 171
CLOM (cluster level object manager), 96–97, 249
`cluster_info` command, 282–284
cluster level object manager (CLOM), 96–97
cluster monitoring, membership, and directory services (CMMDS), 97, 187
clusters, 10
 adding hosts to, 169–170
 creation of, 48–49
 designing (*see* design considerations, clusters)

- clusters (*Continued*)
- health, 248–249
 - removing hosts from, 170–171
 - stretched (*see* stretched clusters)
 - vSphere HA (*see* vSphere HA (high availability) cluster)
- `cluster_set_default_policy` command, 284–285
- CMMDS (cluster monitoring, membership, and directory services), 97, 187, 263
- `cmmds_find` command, 289–290
- `cmmds-tool` command, 269–272
- command-line interface (CLI), 171
- commands. *see* specific commands
- communication network, vSphere HA, 49–50
- component protection, vSphere HA, 51
- components, 86–87
- absent, 185
 - degraded, 185
 - limits, 87–88
 - management, 95
 - replicas, 90–91
 - software (*see* software components)
 - witness, 87, 90–91
- compression, deduplication and, 105–107
- configuration
- IGMP snooping querier, 39
 - networks
 - issues, 38–40
 - VDS, 32–37
 - VSS, 31–32
 - NIOC example, 40–42
 - stretched clusters, 202–203, 208–216
- congestion test, disks, 251
-
- D**
- data health checks, 250
- data integrity, through checksum, 130–131
- data locality, I/O flow, 107–108
- in VSAN stretched clusters, 108–109
- data paths, for objects, 95–96
- datastores, usage warning, 195
- decommission mode, 175
- deduplication, compression and, 105–107
- default maintenance mode, 175
- default policy, 160–164
- hard disk 1 layout, 163
 - number of failures to tolerate = 1, 162
 - storage, 93–94
- degraded components, 185
- Dell R730XD server, 238–241
- delta disks, 87, 90, 126
- design considerations
- clusters
 - cache to capacity ratio, 228–229
 - ready node profiles, 225–226
 - server virtualization (*see* server virtualization)
 - sizing constraints, 227–228
 - tools, design and sizing, 236
 - performance, 229–231
 - disk controllers, 231–235
- VDS and NIOC, 42–43
- redundant 10GbE switch with link aggregation, 45–48
 - redundant 10GbE switch without link aggregation, 43–45
- DHCP (dynamic host configuration protocol), 36
- `disable_vsan_on_cluster` command, 279–280
- disabling
- IGMP (internet group management protocol) snooping, 39
 - object checksum capability, 73

disk controllers
 queue depth, impact of, 231–235
 RAID-0, 18–19
 disk failure, 133, 185
 disk full scenario, 193–194
 disk groups
 adding, 177–178
 disks to, 54–55, 179–180
 creation, example, 55–58
 maximums, 52
 multiple, configuration, 52–53
 removing, 178–179
 disks from, 180–181
 role of, 51–52
 sizing ratio, cache device to capacity
 device, 53–54
 disk management, 177
 adding
 disk groups, 177–178
 disks to disk groups, 179–180
 removing
 disk groups, 178–179
 disks from, 180–181
 disk_object_info command, 287–288
 disks
 adding to disk groups, 177–178
 automatically, 54–55
 manually, 55
 blinking LED on, 183
 removing from disk groups, 180–181
 wiping, 182–183
 esxcli vsan disk commands, 184
 disks_info command, 284
 disk_stats command, 294–295
 distributed object manager (DOM), 95–96,
 100, 207
 distributed port group, creating, 33–34
 distributed RAID, 83–85

Distributed Resource Scheduler (DRS), 195,
 204
 DOM (distributed object manager), 95–96,
 100, 207
 DRS (Distributed Resource Scheduler), 195,
 204

E

enable_vsan_on_cluster command, 279–280
 Ensure Accessibility option (maintenance
 mode), 173
 enter_maintenance_mode command, 296
 ESXCLI, 245, 256
 esxcli vsan cluster command, 258–259
 esxcli vsan datastore namespace, 256–257
 esxcli vsan faultdomain command,
 263–264
 esxcli vsan maintenance mode
 command, 264
 esxcli vsan network namespace, 257–258
 esxcli vsan policy command,
 264–267
 esxcli vsan storage namespace, 258–262
 esxcli vsan trace command, 267–268
 esxcli network diag ping command,
 258
 esxcli network ip connection list
 command, 258
 esxcli network ip neighbor list
 command, 258
 esxcli vsan cluster command, 171–172,
 258–259
 esxcli vsan disk command, 184
 esxcli vsan faultdomain command,
 263–264
 esxcli vsan maintenance mode
 command, 264
 esxcli vsan network namespace, 257–258
 esxcli vsan policy command, 264–267

esxcli vsan storage, 258–262
esxcli vsan trace command, 267–268
ESXi, 14
boot considerations, 15
firewall ports, 26
NIC teaming in, 25
troubleshooting VSAN on, 303
log files, 304
VMkernel modules and drivers, 305
VSAN traces, 304–305
ESXTOP, 245
esxtop command, 308–309
ESXTOP performance counters, 308–309
explicit failover order, 44

F

failures
cache devices, 186–187
capacity devices, 185–186
hosts, 187–188
number to tolerate (*see* number of failures to tolerate)
recovery from, 131–134
write failures, 185–186
failure scenarios, 127–130, 185
cache device failure, 186–187
capacity device failure, 185–186
disk full scenario, 193–194
host failure, 187–188
network partition, 188–193
stretched clusters, 216–223
network failure, 221–223
preferred site, 219–220
secondary site, 217–218
witness host failure, 220–221
vCenter server, 195

failures to tolerate = 1, object space reservation = 50 percent, 152–155
failures to tolerate = 1, object space reservation = 100 percent, 155–156
failures to tolerate = 1, stripe width = 2, 144–148
failures to tolerate = 2, stripe width = 2, 148–152
failure tolerance method, 66–69, 119–120
fault domains, stretched clusters *vs.*, 206–208
find option (cmmds-tool command), 270
firewall ports, 26
fix_renamed_vms command, 291
flash devices, 21–22
flash read cache reservation, 71
policy setting, 122
force provisioning, 71–72, 94, 127
forceProvisioning command option, 266
Full Data Migration option (maintenance mode), 173

G

get option (esxcli vsan cluster command), 262
global support services (GSS), 267
GSS (global support services), 267

H

HA (high availability) cluster
VSAN storage providers, 75
changing policies on-the-fly, 75–79
vSphere, 49–51
hard disk 1 layout (default policy), 163

HBA mode, 17
HCL (hardware compatibility list), 4
 health checks, 247–248
health check, VSAN, 165, 245, 246
 Ask VMware link, 246
 cluster health, 248–249
 data health, 250
 limits, 250
 network health, 249–250
 performance service, 168–169
 physical disk, 251
 proactive (*see* proactive health checks)
 stretched cluster, 252
 tests, 165–167
 Virtual SAN HCL health, 247–248
 VSAN performance service, 252
heartbeat datastores, vSphere HA, 50
`help` command, 276–279, 301
`host_consume_disks` command, 282
`host_evacuate_data` command, 282
`hostFailuresToTolerate` command
 option, 266
`host_info` command, 280–282
host management, 169
 adding hosts to clusters, 169–170
 esxcli vsan cluster commands, 171–172
 removing hosts from clusters, 170–171
hosts, 10
 adding to cluster, 169–170
 configuration, determination, 238–241
 failure, 132, 187–188
 management (*see* host management)
 removing from cluster, 170–171
 roles, 97–98
 witness, 206
`host_wipe_non_vsan_disks` command, 282
`host_wipe_vsan_disks` command, 282
hybrid VSAN

capacity tier on, 105
VSAN read on, anatomy of, 102–103
VSAN write on, anatomy of, 103–104

I

IEEE 802.3ad, 42
IGMP (internet group management protocol) snooping, 39
I/O flow, 100
 cache layer, role of, 100–102
 read cache, 100–101
 write cache, 101–102
caching algorithms, 100
data locality, 107–108
 in VSAN stretched clusters, 108–109
deduplication and compression, 105–107
striping, two hosts, 145
VSAN read
 on all-flash VSAN, 103
 on hybrid VSAN, 102–103
VSAN write
 on all-flash VSAN, 104–105
 on hybrid VSAN, 103–104
IOPS (input/output operations per second)
 limit threshold, 70
`iopsLimit` command option, 266
IP-Hash, 46
isolation response, vSphere HA, 51

J

JBOD mode, 17
jumbo frames, 24–25

L

LACP (link aggregation control protocol), 46
Lam, William, 197

latency issues, VSAN Observer sample use cases, 317
latency requirements, stretched clusters, 205
layer 3 (L3) (routed) network, 23
layer 2 (L2) (switched) network, 23
layout, objects, 91–93
 default storage policy, 93–94
LBT (load based teaming) mechanism, 45
license editions, stretched clusters, 203–204
limits, components, 87–88
limits health checks, 250
`lldpnetmap` command, 296
local log structured object manager (LSOM), 95, 101
log files, troubleshooting VSAN on ESXi, 304
LSOM (local log structured object manager), 95, 101
LUN (logical unit number), 5, 61

M

magnetic disks, 20
maintenance mode, 172–174
 updates and patching, 175–176
master role (hosts), 97–98
migration, for objects, 96–97
multicast heartbeats, 25
multicast performance test, 168
multicast performance tests, 254–255

N

namespaces, 89–90
`esxcli vsan datastore`, 256–257
`esxcli vsan network`, 257–258
`esxcli vsan storage`, 258–262
NAS (network-attached storage), 62

Network file system (NFS), 62
network health checks, 249–250
networking, 29–30
 configuration
 issues, 38–40
 VDS, 32–37
 VSS, 31–32
redundant 10GbE switch with link aggregation capability, 45–48
redundant 10GbE switch without link aggregation capability, 43–45
stretched clusters, 205
VMkernel network for VSAN, 30–31
network interface card (NIC), 7, 25
network I/O control. *see* NIOC (network I/O control)
network partition, 188–193
network policies, 2
network requirements, VSAN, 22–25
 jumbo frames, 24–25
 layer 2/layer 3 networks, 23
NIC (network interface cards), 22
NIC teaming, 25
NIOC (network I/O control), 25
switches, 22
traffic, 24
VMkernel network, 23, 30–31
NFS (Network file System), 62
NIC (network interface card), 7, 22
NIC teaming, 25
NIOC (network I/O control), 25
 configuration example, 40–42
 design considerations, 42–43
NL-SAS drives, 234
No Data Migration option (maintenance mode), 174
NSX virtualization platform, 2
number of disk stripes per object (SPBM), 69–70

number of failures to tolerate, 65–66, 84, 228
best practice for, 112–113
policy setting, 110–112
number of failures to tolerate = 1, 137–144
default policy, 162

O

`object_info` command, 285–289
`object_reconfigure` command, 285
objects, 86–87
 data paths for, 95–96
 delta disks, 87
 layout, 91–93
 default storage policy, 93–94
 namespace, 89–90
 ownership, 96
 placement and migration for, 96–97
 storage (*see* storage objects)
 VM swap, 90
object space reservations, 71, 152
 policy setting, 122–123
`obj_status_report` command, 291–292
on-disk formats, 98
 cache devices, 99
 capacity devices, 99
`osfs-ls` command, 268–269

P

`partedUtil` method, 182–183
pass-through mode, 17
patching (maintenance mode), 175–176
performance, 245
 cache tier destaging, 116–117
 data (VSAN Observer), 313–316
 designing for, 229–231
 disk controllers, 231–235

monitoring (*see* performance monitoring) and RAID caching, 19
read cache misses, 115–116
VSAN capabilities, 235
writes, 114
performance monitoring, 305
 ESXTOP performance counters, 308–309
performance service (*see* performance service)
VSAN observer (*see* VSAN observer)
vSphere Web Client performance counters, 309–310
performance service, 168–169, 305–306
 enabling, 306
 health checks, 252
 metrics, 307–308
 VSAN, 306–307
physical disk health check, 251
physical disk placement
 failures to tolerate = 1, stripe width = 2, 146–147
 failures to tolerate = 2, stripe width = 2, 151
 number of failures to tolerate = 1, 144
PIM (protocol-independent multicast), 205
placement, for objects, 96–97
policy settings
 failures to tolerate = 1, object space reservation = 50 percent, 152–155
 failures to tolerate = 1, object space reservation = 100 percent, 155–156
 failures to tolerate = 1, stripe width = 2, 144–148
 failures to tolerate = 2, stripe width = 2, 148–152
 number of failures to tolerate = 1, 137–144
RAID-5, 157–158
RAID-6, 158–159
RAID-5/6 and Stripe Width = 2, 159–160

- ports
 firewall, 26
 VMkernel, building, 34–37
- preferred site, 207
 failure scenarios, 219–220
- proactive health checks, 167–168, 253–255
 multicast performance tests, 254–255
 storage performance tests, 255
- VM creation test, 253–254
- profile-driven storage, 61
- proof-of-concept (PoC), 167
- proportionalCapacity command
 option, 266
- protocol-independent multicast (PIM), 205
- provisioning
 force, 71–72, 94, 127
 thin, 194–195
-
- Q**
- queue depth, disk controllers, 231–235
- queuing layers, disk controllers, 231–232
-
- R**
- Raghuram, Raghu, 1
- RAID-0, 18–19, 87, 144
 stripe configuration, 151
- usage, when no striping specified in
 policy, 117
 test 1, 118
 test 2, 119
 test 3, 119
- RAID-1, 67, 84, 87, 144, 228
- RAID-5, 67–69, 84, 157–158, 228
- RAID-6, 67–69, 84–85, 158–159, 228
- RAID caching, 19
- RAID (redundant array of inexpensive
 disks), 17
- distributed, 83–85
- RDT (reliable datagram transport), 98
- read cache, I/O flow, 100–101
 in all-flash VSAN configurations, 101
- ready node profiles, clusters, 225–226
- ready nodes, VSAN, 16–17
- reapply_vsan_vmknic_config command,
 298
- reconfiguration, 133
- recover_spbm command, 298–300
- recovery, from failures, 131–134, 223
- redundant 10GbE switch with link
 aggregation capability, 45–48
- redundant 10GbE switch without link
 aggregation capability, 43–45
- reliable datagram transport (RDT), 98
- replicaPreference command option, 266
- replicas, 90–91
 failure scenarios, 127–130
- reservations
 flash read cache, 71, 122
 object space, 71, 122–123, 152
- resiliency, 11
- ROBO (remote office/branch office)
 VSAN 2-node, 26
- Ruby vSphere console (RVC), 245, 275–276
 SPBM commands, 300–303
- VSAN commands (*see* VSAN commands
 (RVC))
- rule-sets, 139
- RVC. *see* Ruby vSphere console (RVC)
-
- S**
- SATA drives, 234
- SDDC (software-defined datacenter), 1–2
- SDN (software-defined network), 1
- secondary site, 207
 failure scenarios, 217–218

-
- server virtualization
 - all-flash configuration, 241–244
 - hybrid configuration, 237–238
 - host configuration, determination, 238–241
 - service level objective (SLO), 11
 - sizing calculator, 236
 - sizing constraints, 227–228
 - SKU (stock keeping unit), 4
 - slave role (hosts), 97–98
 - SLO (service level objective), 11
 - snapshots, 126
 - software components, 94
 - cluster monitoring, membership, and directory services, 97
 - data paths for objects, 95–96
 - host roles, 97–98
 - management, 95
 - object ownership, 96
 - placement and migration for objects, 96–97
 - reliable datagram transport, 98
 - software-defined datacenter (SDDC), 1–2
 - software-defined network (SDN), 1
 - software-defined storage, 2
 - software-only storage solutions, 4
 - space consumption, verifying, 126–127
 - SPBM (storage policy-based management), 5, 62–65, 109
 - commands in RVC, 279, 300–303
 - disable object checksum capability, 73
 - failure tolerance method, 66–69
 - flash read cache reservation, 71
 - force provisioning, 71–72
 - IOPS limit for object, 70
 - number of disk stripes per object, 69–70
 - number of failures to tolerate, 65–66
 - object space reservations, 71
 - SSD (solid-state disks), 17
 - stock keeping unit (SKU), 4
 - storage controllers, 17–19
 - storage objects, 80
 - virtual machine, 88
 - VM swap, 124–126
 - storage performance tests, 168, 255
 - storage policies, 11
 - changing on-the-fly, 75–79
 - example, 11–12
 - VM storage policy (*see* VM storage policy)
 - storage policy-based management. *see* SPBM (storage policy-based management)
 - storage providers
 - VASA, 74–75
 - VSAN (highly available), 75
 - changing policies on-the-fly, 75–79
 - storage reports, 126–127
 - storage sizing, 239
 - storage traffic, 25
 - stretched clusters, 25–26, 201–223
 - configuration, 202–203, 208–216
 - procedural steps, 211–216
 - size of, 209
 - constraints, 203–204
 - data locality in, 108–109
 - defined, 201–203
 - failure scenarios, 216–223
 - network failure, 221–223
 - preferred site, 219–220
 - secondary site, 217–218
 - witness host failure, 220–221
 - vs. fault domains, 206–208
 - health checks, 252
 - latency requirements, 205
 - license editions, 203–204
 - networking, 205
 - products/features not supported by, 204

stretched clusters (*Continued*)

- requirements, 203–204

stripe width

- best practices, 122

- chunk size, 121

- configuration error, 120–121

- maximum, 119–120

policy setting

- cache tier destaging, 117–118

- read cache misses, 115–116

- writes, 114

`stripeWidth` command option, 266

Supermicro 2U 4-Node TwinPro2 server, 243

switches

- VDS (VMware vSphere Distributed Switches), 22, 29–30, 32–37

- VSS (VMware standard switches), 22, 29, 31–32

T

`tcpdump-uw` command, 258

`-t DISK` option (`cmmcmd_find` command), 290

`-t DISK_USAGE` option (`cmmcmd_find` command), 289–290

tests, VSAN health check, 165–167

thin provisioning, 194–195

tools, for design/sizing clusters, 236

troubleshooting

- VSAN, non-ESXCLI commands for

- `cmmcmd_tool`, 269–272

- `osfs-ls`, 268–269

- `vdq`, 272–275

VSAN on ESXi, 303

- log files, 304

- VMkernel modules and drivers, 305

- VSAN traces, 304–305

U

updates (maintenance mode), 175–176

`upgrade_status` command, 300

usage warnings (datastore), 195

- `-u UUID` option (`cmmcmd_find` command), 290

V

VASA (vSphere APIs for Storage Awareness)

- overview, 73–74

- storage providers, 74–75

- vendor providers, 73

vCenter server

- bootstrapping, 197–199

- failure scenario, 196

- management, 195–199

- running on VSAN, 196–197

VCG (VMware Compatibility Guide), 225

VCVA (VMware vCenter Virtual Appliance), 275

VDI (virtual desktop infrastructure), 117

`vdq` command, 272–275

VDS (VMware vSphere Distributed Switches), 22, 29–30

- creation, 32

- design considerations, 42–43

- distributed port group, creation, 33–34

- VMkernel ports, building, 34–37

- VSAN network configuration, 32–37

virtual desktop infrastructure (VDI), 117

virtualization layer, 1

Virtual Machine File System (VMFS), 99

virtual machine (VM), 2

- storage objects, 88

Virtual SAN 1.0, 8

Virtual SAN 6.0, 8

Virtual SAN 6.1, 8

Virtual SAN 6.2, 8
Virtual SAN HCL health, 247–248
Virtual SAN object health test, 250
Virtual SAN traffic, 36
virtual storage appliance (VSA), 4
Virtual Volumes (VVols), 62
VM. *see* virtual machine (VM)
VM creation test, 167–168, 253–254
VM disk (VMDK), 84, 89, 90
VMFS-L (VMFS local), 99
VMFS (Virtual Machine File System), 99
VM home namespace, 87, 89–90, 123, 148
VMkernel network, 23, 30–31
 adapter, 34
VMkernel ports, building, 34–37
`vm_object_info` command, 288–289
vMotion traffic, 40
`vm_perf_stats` command, 295
VM provisioning, VM storage policy
 assignment during, 81–82
VM storage policy, 5, 61–82, 138, 139
 assignment during VM provisioning,
 81–82
 creating, 81
 default, 93–94
 enabling, 81
 storage policy-based management
 (*see* SPBM (storage policy-based
 management))
 VSAN capabilities and, 63–65
VM swap, 90, 123–124
 descriptor file, 125
 storage objects, 124–126
VM (virtual machine), 2
VMware Compatibility Guide
 (VCG), 225
VMware hardware compatibility guide, 16
VMware standard switches. *see* VSS
 (VMware standard switches)

VMware vCenter Virtual Appliance
 (VCVA), 275
VMware VSAN sizing calculator, 236
VMware vSphere, 13
 ESXi, 14, 15
VSAN capabilities, 109–110
 data integrity through checksum,
 130–131
 delta disks, 126
 failure scenarios, 127–130
 flash read cache reservation policy
 setting, 122
 force provisioning, 127
 number of failures to tolerate
 best practice for, 112–113
 policy setting, 110–112
 object space reservation policy setting,
 122–123
 performance, 235
 problematic device handling, 134
 recovery from failure, 131–134
 snapshots, 126
 storage reports, 126–127
 stretching, 134–135
 stripe width, best practices, 122
 stripe width chunk size, 121
 stripe width configuration error, 120–121
 stripe width maximum, 119–120
 stripe width policy setting, 113
 cache tier destaging, 116–117
 read cache misses, 115–116
 writes, 114
 VM home namespace, 123
 VM swap, 123–126
VSAN commands (RVC)
 `apply_license_to_cluster`, 297
 `check_limits`, 297–298
 `check_state`, 292–293
 `cluster_info`, 282–284

- VSAN commands (RVC) (*Continued*)
- cluster_set_default_policy, 284–285
 - cmmms_find, 289–290
 - disable_vsan_on_cluster, 279–280
 - disk_object_info, 287–288
 - disks_info, 284
 - disk_stats, 294–295
 - enable_vsan_on_cluster, 279–280
 - enter_maintenance_mode, 296
 - fix_renamed_vms, 291
 - help, 276–279
 - host_consume_disks, 282
 - host_evacuate_data, 282
 - host_info, 280–282
 - host_wipe_non_vsan_disks, 282
 - host_wipe_vsan_disks, 282
 - lldpnetmap, 296
 - object_info, 285–289
 - object_reconfigure, 285
 - obj_status_report, 291–292
 - reapply_vsan_vmknic_config, 298
 - recover_spbm, 298–300
 - upgrade_status, 300
 - vm_object_info, 288–289
 - vm_perf_stats, 295
 - whatif_host_failures, 293–294
- VSAN datastore, 48
- properties, 58–59
- VSAN observer, 245, 310
- performance data, 313–316
 - requirements, 310–313
 - sample use cases, 316–318
- vsan.observer command, 311
- VSAN read, I/O flow
- on all-flash VSAN, 103
 - on hybrid VSAN, 102–103
- VSAN ready nodes, 16–17
- VSAN traces, troubleshooting VSAN on ESXi, 304–305
- VSAN traffic, 24, 40
- VSAN (virtual SAN), 4
- administrator’s view, 9–11
 - benefits of, 7
 - cache and capacity devices, 14–15
 - defined, 6–7
 - disk groups (*see* disk groups)
 - firewall ports, 26
 - networking, 29–30
 - overview, 4–5
 - requirements, 15–22
 - cache tier devices, 21–22
 - capacity tier devices, 19–20
 - hardware compatibility guide, 16
 - network (*see* network requirements, VSAN)
 - ready nodes, 16–17
 - storage controllers, 17–19
- software components (*see* software components, VSAN)
- stretched cluster, 25–26
- troubleshooting on ESXi, 303
- log files, 304
 - VMkernel modules and drivers, 305
 - VSAN traces, 304–305
- 2-node/ROBO, 26
- use cases, 7–8
- VMkernel network for, 30–31
- VSA (virtual storage appliance), 4
- vSphere APIs for Storage Awareness. *see* VASA (vSphere APIs for Storage Awareness)
- vSphere Distributed Switch. *see* VDS (VMware vSphere Distributed Switches)
- vSphere HA (high availability) cluster, 49
- admission control, 50–51

- communication network, 49–50
- component protection, 51
- heartbeat datastores, 50
- isolation response, 51
- vSphere Web Client performance counters, 309–310
- VSS (VMware standard switches), 22, 29
 - VSAN network configuration, 31–32
- wiping disks, 182–183
 - esxcli vsan disk commands, 184
- Witness appliance
 - deployment, 208
 - icon, 206
- witness component, 87, 90–91, 111
 - failure scenarios, 127–130
- write cache, I/O flow, 101–102
- write failures, 185–186
- writes
 - performance, 114
 - retiring to capacity tier, 105

W

- whatif_host_failures command, 293–294