

Storage Management in Data Centers

Volker Herminghaus • Albrecht Scriba
Authors

Storage Management in Data Centers

Understanding, Exploiting, Tuning,
and Troubleshooting
Veritas Storage Foundation

 Springer

Volker Herminghaus
Nieder-Olm
Germany
v.herminghaus@anykey-dcs.de

Dr. Albrecht Scriba
Mainz
Germany
albrecht@albrecht-scriba.de

ISBN 978-3-540-85022-9

e-ISBN 978-3-540-85023-6

DOI 10.1007/978-3-540-85023-6

Library of Congress Control Number: 2009921159

© Springer-Verlag Berlin Heidelberg 2009

This work is subject to copyright. All rights are reserved, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilm or in any other way, and storage in data banks. Duplication of this publication or parts thereof is permitted only under the provisions of the German Copyright Law of September 9, 1965, in its current version, and permissions for use must always be obtained from Springer-Verlag. Violations are liable for prosecution under the German Copyright Law.

The use of general descriptive names, registered names, trademarks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

Cover design: KuenkelLopka GmbH Heidelberg,

Printed on acid-free paper

springer.com

For my wife and children, who taught me what really counts.

Volker Herminghaus

Dedicated to my wife and children.

Dr. Albrecht Scriba

PREFACE

This book is designed to meet the needs of UNIX architects and administrators working in data centers. While it will be useful for the computer science student or the newcomer who has been attracted to volume management by Symantec's release of a free version of its Volume Manager software its focus is on the data center. Most data center applications nowadays handle amounts of data that had been unconceivable at the time when the most commonly used storage media - the hard disk - was developed. As a consequence, the design of the hard disk simply cannot match the requirements posed by current applications. Its physical attributes and limits need to be overcome by additional layers of hardware or software. These layers, if properly designed and thoughtfully applied, convert a set of physical disks to a supply of storage space whose properties better match application requirements. Instead of physical disks with their physical limitations, logical entities known as volumes are now commonly used. These volumes can be fault tolerant, accelerated to the limits, replicated to remote locations, and made almost infinitely large. Additionally, volumes can even be reshaped and their features changed while they are in use, enabling the data center administrator to adapt to changing requirements without suffering an application downtime.

The technical term for this software or hardware layer is "volume management".

Veritas Volume Manager® is the most widely used software for volume management. It is used in data centers all over the world and has proven to be stable and deliver high performance under most circumstances. While there are other volume management software products on the market (e.g. AIX LVM, Sun microsystems' SunVM, or several Linux LVMs), most of them suffer from one or more limitations that hamper their widespread deployment. They are either limited to the manufacturer's operating system or they have less to offer than the Veritas product. In most cases, both is true at the same time. This has led to Veritas Volume Manager, or VxVM in short, being the most widely deployed product on the market, which in turn led to most administrators learning at least the basic skills required for its administration.

However, mastering the basic skills is something quite different from fully understanding a product and making full use of the available power. In data center operations it is imperative that the operators know precisely how things are supposed to work, rather than apply the skills of "experimental computer science". Even today's personal computers are too complex for any kind of experimental approach to solving a problem or finding a solution. This is much more true in data centers, where the motto must be: "If you do not know it, then learn it or leave it, but don't fumble it".

In my time as both an independent data center consultant and an independent trainer for the Veritas product suite I have tried to educate people enough so that they would at least realize what is possible if they could harness VxVM to its full extent. Staying in close

contact with my clients, it dawned on me that what they need is a written guide they can rely on when they actually try some of the more advanced features. If you are responsible for a mission critical application then the last thing you want is to incur a downtime. And with only some diffuse background knowledge and elementary skills left over from the last VxVM training, most of you would rather stick to established procedures than try something new.

My first attempt at writing down what I knew was a training companion book called "Veritas Storage Foundation" published by Springer in 2006 (ISBN: 3-540-34610-4) and endorsed by Switzerland's biggest Symantec partner, Infonika SQL AG (www.sql.ch) as their official training material. This book had been written together with Albrecht Scriba, one of the most respected Veritas trainers in Europe. It covered Veritas Volume Manager (VxVM) and Veritas Cluster Server (VCS) and was received very well by the administrators. However, its drawback was that it was written in German, our native tongue, so its distribution was severely limited by the language barrier. Having been approached numerous times by international colleagues I decided to take the next step and write a new, English book that concentrates on VxVM and the Veritas File System (VxFS), again with Albrecht acting as the co-author for some of the toughest chapters. It is not a training companion like the first one but uses a more classical approach. There are many walkthroughs to make you understand what you can do, how you can do it and what exactly is going on inside VxVM and VxFS so you can understand it and repeat it step by step on your own systems. It also holds a large section on troubleshooting that points out how problems can be found and fixed.

So here is your guide that helps you understand - in detail - the principles and the problems of mass storage, volume management, and file systems and how to manage them. It also tries to correct some common misconceptions about storage and UNIX, and highlights the most limiting factors in today's data center environments: anachronistic thinking and the sluggish speed of light!

ABOUT THE AUTHORS

Volker Herminghaus

Born in the stone ages (1963) and raised in a family full of physicists, he studied mechanical engineering and computer science in Darmstadt, Germany. Took some really deep looks at the kernel of AT&T UNIX System V Release 3.2 for his thesis and has been claiming to know what he's talking about since then.

He started computing on a Commodore 64 and switched to Atari ST as it became available, then finally to NeXTSTEP. Deeply enamoured with its elegance and power he has since stuck to descendants of this operating system (MacOS X) for his own use. Professionally, he has been working on Solaris and other UNIX variants as a consultant since the early 1990s. He has just co-founded his second data center consulting company: the **anykey-dcs**.

Dr. Albrecht Scriba

Albrecht studied mathematics and religious science in Mainz until 1998 (including thesis and State doctorate). Being familiar with ancient languages like Aramaic as well as computer programming he hacked the Atari ST's Signum! program in assembler to optimize printing of those languages' fonts. Albrecht has long been working as a Consultant and as a Trainer for Veritas/Symantec. He left Symantec in 2008 and is now working for anykey-dcs, Symantec and other companies as a free lancer. Fallen in love with Unix since his very first encounter, his motto is: "Never type a command twice, write a program for it!"

ACKNOWLEDGEMENTS

A book like this is not created out of thin air. It takes a lot of work, energy, resources, and determination to keep going all the way to the publication. Many people have helped us finish this task and have thus contributed to the successful completion of this book. The first round of thanks goes to the contact persons at Springer: Hermann Engesser and Gabriele Fischer, for providing a perfect office into which to throw all the unstructured suggestions, ideas, questions, and draft versions of this book. But most of all, for saying "YES" before we could even finish our sentence asking if they would like to publish our second book. That immediate and unquestioning positive reply provided the motivation required to kick off the project.

The second round goes to our wives and kids, who always suffer most when fathers decide to dedicate the better part of two years to sitting down late at night hacking, experimenting, and writing.

The third round goes to all the people that gave us gems of background information on storage management: Of these, Ron Karr and particularly Oleg Kiselev, two of the inventors of Veritas Volume Manager and extremely smart people, provided the most insight into VxVM's design ideas and implementation as well as a broad overview about modern storage systems in general.

Sun Microsystems' benchmarking center in Langen, Germany, allowed us to access their powerful Sun servers tied to high-end storage in order to run a multitude of tests and simulations. Special thanks go to all involved at Sun for their efficient, competent and friendly support: Kirsten Prahst, Rüdiger Frenk (who also maintains the only complete Sun hardware museum in the world), and Peter Hausdorf.

A final round of thanks goes to the many hundreds of people who have participated in our trainings or been our consulting clients. They never failed to come up with new questions, setups or problems that kept our brains busy.

TABLE OF CONTENTS

Preface	VII
About the authors	IX
Acknowledgements	XI
1 DISK AND STORAGE SYSTEM BASICS	1
1.1 Overview	1
1.1.1 Storage Hardware Situation and Outlook.	1
1.1.2 Physical Limits	3
1.1.3 Trying to Fix the Problems - and Failing!	7
1.1.4 SAN-Attached Hard Disks.	10
1.1.5 Storage Arrays and LUNs	10
1.1.6 Common Problems	15
1.1.7 Physical Disks vs. LUNs.	17
1.2 Disk Addressing and Layout.	19
1.3 Paths and path redundancy.	23
1.4 The Trouble with Networked Disk Access	30
1.4.1 Summary	36
2 EXPLORING VXVM	39
2.1 Getting Started	39
2.1.1 Hello, Volume!.	40
2.1.2 vxdisksetup: Turning Disks into VM Disks.	40
2.1.3 Disk Groups: Putting VM Disks into Virtual Storage Boxes	42
2.2 The Hard Way: a Low-level Walkthrough	45
2.2.1 Subdisks: Extents for Persistent Backing Store	45
2.2.2 Plexes: Mapping Virtual Extents to Physical Extents	46
2.2.3 Volumes: Virtual Partitions for Any Purpose	48
2.2.4 Volume Start: Prepare for Takeoff	52
2.3 The Easy Way: vxassist	53
2.3.1 Summary	53
3 INCORPORATING DISKS INTO VXVM	55
3.1 Solaris Disk Handling	56
3.1.1 Getting a New Disk into Solaris	56
3.1.2 You Don't Format with "format"	57
3.1.3 Finding New Disks in VxVM	57
3.1.4 What if My New Disk is Not Found?	59
3.1.5 Leaving Physics Behind – Welcome to VxVM!	61

Table Of Contents

3.2	VxVM disk handling	62
3.2.1	VxVM Disk Formats	62
3.2.2	cdsdisk and sliced	63
3.2.3	How to Mix CDS and Sliced Disks in a Disk Group?	66
3.2.4	Other Disk Formats	66
3.2.5	Encapsulation Overview – Integrating Legacy Data.	67
3.2.6	Summary	69
4	DISK GROUPS	71
4.1	Overview	71
4.1.1	What is a Disk Group?	71
4.2	Simple Disk Group Operations	74
4.3	Advanced Disk Group Operations	80
4.3.1	Options for Importing or Exporting a DG	81
4.3.2	Disk Group Operations for Off-Host Processing	83
4.3.3	Miscellaneous Disk Group Operations	85
4.3.4	Summary	87
4.4	Disk Group Implementation Details	89
4.4.1	Major and Minor Numbers for Volumes and Partitions	97
5	VOLUMES	99
5.1	Overview	99
5.1.1	What is a Volume?	99
5.2	Simple Volume Operations	101
5.2.1	Creating, Using and Displaying a Volume.	101
5.2.2	Useful vxprint Flags Explained	103
5.2.3	Starting and Stopping Volumes	105
5.3	Volume Layouts and RAID Levels	106
5.3.1	Volume Features Supported by VxVM	106
5.4	Volume Maintenance	114
5.5	Tuning vxassist Behavior.	120
5.5.1	Storage Attributes – Specifying Allocation Strategies	120
5.5.2	Skipping Initial Mirror Synchronisation	126
5.5.3	Changing the Layout of a Volume	127
5.6	Methods of Synchronisation	130
5.6.1	Atomic Copy	131
5.6.2	Read-Writeback, Schrödinger's Cat, and Quantum Physics	132
5.7	Volume Features in Detail	137
5.7.1	concat	137
5.7.2	stripe	137
5.7.3	mirror	139
5.7.4	RAID-4 and RAID-5.	142
5.7.5	mirror-concat	146