



*AIX 5L System  
Administration II: Problem  
Determination*

(Course Code AU16)

**Student Notebook**

ERC 12.0

IBM Certified Course Material

## Trademarks

The reader should recognize that the following terms, which appear in the content of this training document, are official trademarks of IBM or other companies:

IBM® is a registered trademark of International Business Machines Corporation.

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

AIX	AIX 5L	Micro-Partitioning
MVS	OS/2	POWER
POWER4	POWER5	POWER Gt1
POWER Gt3	PS/2	pSeries
Redbooks	RS/6000	SP

Java and all Java-based trademarks are trademarks of Sun Microsystems, Inc. in the United States, other countries, or both.

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

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

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

ALERTS is a registered trademark of Alphablox Corporation in the United States, other countries, or both.

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

## December 2004 Edition

The information contained in this document has not been submitted to any formal IBM test and is distributed on an "as is" basis without any warranty either express or implied. The use of this information or the implementation of any of these techniques is a customer responsibility and depends on the customer's ability to evaluate and integrate them into the customer's operational environment. While each item may have been reviewed by IBM for accuracy in a specific situation, there is no guarantee that the same or similar results will result elsewhere. Customers attempting to adapt these techniques to their own environments do so at their own risk.

© Copyright International Business Machines Corporation 1997, 2004. All rights reserved.

**This document may not be reproduced in whole or in part without the prior written permission of IBM.**

Note to U.S. Government Users — Documentation related to restricted rights — Use, duplication or disclosure is subject to restrictions set forth in GSA ADP Schedule Contract with IBM Corp.

# Contents

<b>Trademarks</b> .....	<b>xiii</b>
<b>Course Description</b> .....	<b>xv</b>
<b>Agenda</b> .....	<b>xvii</b>
<b>Unit 1. Problem Determination Introduction</b> .....	<b>1-1</b>
Unit Objectives .....	1-2
1.1 Problem Determination Introduction .....	1-3
Role of Problem Determination .....	1-4
Before Problems Occur .....	1-5
Before Problems Occur: A Few Good Commands .....	1-6
Problem Determination Techniques .....	1-7
Identify the Problem .....	1-8
Define the Problem (1 of 2) .....	1-9
Define the Problem (2 of 2) .....	1-10
Collect System Data .....	1-11
Problem Determination Tools .....	1-12
Resolve the Problem .....	1-13
Obtaining Software Fixes and Microcode Updates .....	1-14
Software Update Management Assistant (SUMA) .....	1-15
SUMA Modules .....	1-17
SUMA Examples (1 of 2) .....	1-19
SUMA Examples (2 of 2) .....	1-21
Relevant Documentation .....	1-23
1.2 pSeries Product Family .....	1-25
IBM @Server pSeries Product Family .....	1-26
AIX 5L 5.2 and 5.3 Logical Partition (LPAR) Support .....	1-28
Advance POWER Virtualization Feature (POWER5) .....	1-30
Virtual Ethernet (AIX 5.3 and POWER5) .....	1-32
Checkpoint Questions .....	1-33
Exercise 1 .....	1-34
Unit Summary .....	1-35
<b>Unit 2. The Object Data Manager (ODM)</b> .....	<b>2-1</b>
Unit Objectives .....	2-2
2.1 Introduction to the ODM .....	2-3
What Is the ODM? .....	2-4
Data Managed by the ODM .....	2-5
ODM Components .....	2-6
ODM Database Files .....	2-7
Device Configuration Summary .....	2-8
Configuration Manager .....	2-9
Location and Contents of ODM Repositories .....	2-10

How ODM Classes Act Together .....	2-12
Data Not Managed by the ODM .....	2-13
Let's Review: Device Configuration and the ODM .....	2-14
ODM Commands .....	2-15
Changing Attribute Values .....	2-17
Changing Attribute Values Using odmchange .....	2-19
2.2 ODM Database Files .....	2-21
Software Vital Product Data .....	2-22
Software States You Should Know About .....	2-24
Predefined Devices (PdDv) .....	2-26
Predefined Attributes (PdAt) .....	2-29
Customized Devices (CuDv) .....	2-31
Customized Attributes (CuAt) .....	2-33
Additional Device Object Classes .....	2-34
Next Step .....	2-36
Checkpoint .....	2-37
Unit Summary .....	2-38
<b>Unit 3. System Initialization Part I .....</b>	<b>3-1</b>
Unit Objectives .....	3-2
3.1 System Startup Process .....	3-3
How Does An AIX System Boot? .....	3-4
Loading of a Boot Image .....	3-6
Content of Boot Logical Volume (hd5) .....	3-7
How to Fix a Corrupted BLV .....	3-8
Working with Boot Lists (PCI) .....	3-10
Working with Boot Lists - SMS .....	3-12
System Management Services .....	3-15
Service Processors and Boot Failures .....	3-16
Let's Review .....	3-18
3.2 Solving Boot Problems .....	3-19
Accessing a System That Will Not Boot .....	3-20
Booting in Maintenance Mode .....	3-22
Working in Maintenance Mode .....	3-23
Boot Problem References .....	3-25
Firmware Checkpoints and Error Codes .....	3-26
Flashing 888 .....	3-27
Understanding the 103 Message .....	3-28
Location Codes: Model 150 .....	3-29
SCSI Addressing .....	3-31
Problem Summary Form .....	3-33
Getting Firmware Updates from Internet .....	3-34
Next Step .....	3-35
Checkpoint .....	3-36
Unit Summary .....	3-37

<b>Unit 4. System Initialization Part II</b> .....	<b>4-1</b>
Unit Objectives .....	4-2
4.1 AIX Initialization Part 1 .....	4-3
System Software Initialization - Overview .....	4-4
rc.boot 1 .....	4-6
rc.boot 2 (Part 1) .....	4-7
rc.boot 2 (Part 2) .....	4-9
rc.boot 3 (Part 1) .....	4-11
rc.boot 3 (Part 2) .....	4-13
rc.boot Summary .....	4-14
Let's Review: Review rc.boot 1 .....	4-15
Let's Review: Review rc.boot 2 .....	4-16
Let's Review: Review rc.boot 3 .....	4-17
4.2 AIX Initialization Part 2 .....	4-19
Configuration Manager .....	4-20
Config_Rules Object Class .....	4-22
Output of cfgmgr in the Boot Log Using alog .....	4-24
/etc/inittab File .....	4-25
System Hang Detection .....	4-27
Configuring shdaemon .....	4-29
Resource Monitoring and Control .....	4-31
RMC Conditions Property Screen: General Tab .....	4-33
RMC Conditions Property Screen: Monitored Resources Tab .....	4-34
RMC Actions Property Screen: General Tab .....	4-35
RMC Actions Property Screen: When in Effect Tab .....	4-36
/etc/inittab: Entries You Should Know About .....	4-37
Boot Problem Management .....	4-39
Next Step .....	4-42
Checkpoint .....	4-43
Unit Summary .....	4-44
<b>Unit 5. Disk Management Theory</b> .....	<b>5-1</b>
Unit Objectives .....	5-2
5.1 Basic LVM Tasks .....	5-3
LVM Terms .....	5-4
Volume Group Limits .....	5-6
Scalable Volume Groups - AIX 5.3 .....	5-8
Configuration Limits for Volume Groups .....	5-10
Mirroring .....	5-12
Striping .....	5-13
Mirroring and Striping with RAID .....	5-15
RAID Levels You Should Know About .....	5-17
Let's Review: Basic LVM Tasks .....	5-19
Review Activity: Basic LVM Tasks .....	5-20
5.2 LVM Data Representation .....	5-23
LVM Identifiers .....	5-24
LVM Data on Disk Control Blocks .....	5-25

LVM Data in the Operating System .....	5-27
Contents of the VGDA .....	5-28
VGDA Example .....	5-30
The Logical Volume Control Block (LVCB) .....	5-32
How LVM Interacts with ODM and VGDA .....	5-33
ODM Entries for Physical Volumes (1 of 3) .....	5-34
ODM Entries for Physical Volumes (2 of 3) .....	5-35
ODM Entries for Physical Volumes (3 of 3) .....	5-36
ODM Entries for Volume Groups (1 of 2) .....	5-37
ODM Entries for Volume Groups (2 of 2) .....	5-38
ODM Entries for Logical Volumes (1 of 2) .....	5-39
ODM Entries for Logical Volumes (2 of 2) .....	5-40
ODM-Related LVM Problems .....	5-41
Fixing ODM Problems (1 of 2) .....	5-42
Fixing ODM Problems (2 of 2) .....	5-44
Next Step .....	5-45
<b>5.3 Mirroring and Quorum .....</b>	<b>5-47</b>
Mirroring .....	5-48
Stale Partitions .....	5-49
Creating Mirrored LVs (smit mklv) .....	5-51
Scheduling Policies: Sequential .....	5-52
Scheduling Policies: Parallel .....	5-53
Mirror Write Consistency (MWC) .....	5-55
Adding Mirrors to Existing LVs (mklvcopy) .....	5-57
Mirroring rootvg .....	5-58
Mirroring Volume Groups (mirrorvg) .....	5-60
VGDA Count .....	5-61
Quorum .....	5-62
Nonquorum Volume Groups .....	5-63
Forced Varyon (varyonvg -f) .....	5-64
Physical Volume States .....	5-65
Summary Quorum .....	5-67
Next Step... .....	5-68
Checkpoint .....	5-69
Unit Summary .....	5-70
<b>Unit 6. Disk Management Procedures .....</b>	<b>6-1</b>
Unit Objectives .....	6-2
<b>6.1 Disk Replacement Techniques .....</b>	<b>6-3</b>
Disk Replacement: Starting Point .....	6-4
Procedure 1: Disk Mirrored .....	6-6
Procedure 2: Disk Still Working .....	6-8
Procedure 2: Special Steps for rootvg .....	6-10
Procedure 3: Total Disk Failure .....	6-12
Procedure 4: Total rootvg Failure .....	6-14
Procedure 5: Total non-rootvg Failure .....	6-16
Frequent Disk Replacement Errors (1 of 4) .....	6-18

Frequent Disk Replacement Errors (2 of 4) .....	6-19
Frequent Disk Replacement Errors (3 of 4) .....	6-20
Frequent Disk Replacement Errors (4 of 4) .....	6-21
6.2 Export and Import .....	6-23
Exporting a Volume Group .....	6-24
Importing a Volume Group .....	6-26
importvg and Existing Logical Volumes .....	6-28
importvg and Existing Filesystems (1 of 2) .....	6-29
importvg and Existing Filesystems (2 of 2) .....	6-30
importvg -L (1 of 2) .....	6-31
importvg -L (2 of 2) .....	6-32
Next Step .....	6-33
Checkpoint .....	6-34
Unit Summary .....	6-35
<b>Unit 7. Saving and Restoring Volume Groups and Online JFS/JFS2 Backups. .</b>	<b>7-1</b>
Unit Objectives .....	7-2
7.1 Saving and Restoring the rootvg .....	7-3
Creating a System Backup: mksysb .....	7-4
mksysb Tape Images .....	7-5
CD or DVD mksysb .....	7-7
Required Hardware and Software for Backup CDs and DVDs .....	7-8
The mkcd Command .....	7-9
Verifying a System Backup After mksysb Completion (1 of 2) .....	7-11
Verifying a System Backup After mksysb Completion (2 of 2) .....	7-12
mksysb Control File: bosinst.data .....	7-14
Restoring a mksysb (1 of 2) .....	7-16
Restoring a mksysb (2 of 2) .....	7-17
Cloning Systems Using mksysb Tapes .....	7-19
Changing the Partition Size in rootvg .....	7-21
Reducing a File System in rootvg .....	7-23
Let's Review: Working with mksysb Images .....	7-25
7.2 Alternate Disk Installation .....	7-27
Alternate Disk Installation .....	7-28
Alternate mksysb Disk Installation (1 of 2) .....	7-29
Alternate mksysb Disk Installation (2 of 2) .....	7-31
Alternate Disk rootvg Cloning (1 of 2) .....	7-33
Alternate Disk rootvg Cloning (2 of 2) .....	7-34
Removing an Alternate Disk Installation .....	7-35
Let's Review: Alternate Disk Installation .....	7-37
7.3 Saving and Restoring non-rootvg Volume Groups .....	7-39
Saving a non-rootvg Volume Group .....	7-40
savevg/restvg Control File: vgname.data .....	7-41
Restoring a non-rootvg Volume Group .....	7-42
7.4 Online JFS and JFS2 Backup; JFS2 Snapshot; VG Snapshot .....	7-43
Online jfs and jfs2 Backup .....	7-44
Splitting the Mirror .....	7-45

Reintegrate a Mirror Backup Copy	7-46
JFS2 Snapshot Image	7-47
Creation of a JFS2 Snapshot	7-48
Using a JFS2 Snapshot	7-49
Snapshot Support for Mirrored VGs	7-50
Snapshot VG Commands	7-51
Next Step	7-52
Checkpoint	7-53
Unit Summary	7-54
<b>Unit 8. Error Log and syslogd</b>	<b>8-1</b>
Unit Objectives	8-2
8.1 Working With Error Log	8-3
Error Logging Components	8-4
Generating an Error Report via smit	8-6
The errpt Command	8-8
A Summary Report (errpt)	8-10
A Detailed Report (errpt -a)	8-11
Types Of Disk Errors	8-13
LVM Error Log Entries	8-15
Maintaining the Error Log	8-16
Activity: Working with the Error Log	8-17
8.2 Error Notification and syslogd	8-19
Error Notification Methods	8-20
Self-made Error Notification	8-22
ODM-based Error Notification: errnotify	8-23
syslogd Daemon	8-26
syslogd Configuration Examples	8-27
Redirecting syslog Messages to Error Log	8-30
Directing Error Log Messages to syslogd	8-31
Next Step	8-32
Checkpoint	8-33
Unit Summary	8-34
<b>Unit 9. Diagnostics</b>	<b>9-1</b>
Unit Objectives	9-2
9.1 Diagnostics	9-3
When Do I Need Diagnostics?	9-4
The diag Command	9-5
Working with diag (1 of 2)	9-6
Working with diag (2 of 2)	9-8
What Happens If a Device Is Busy?	9-9
Diagnostic Modes (1 of 2)	9-10
Diagnostic Modes (2 of 2)	9-12
diag: Using Task Selection	9-13
Diagnostic Log	9-14
PCI: Using SMS for Diagnostics	9-15



Activity: Diagnostics .....	9-17
Checkpoint .....	9-21
Unit Summary .....	9-22
<b>Unit 10. The AIX System Dump Facility .....</b>	<b>10-1</b>
Unit Objectives .....	10-2
10.1 Working with System Dumps .....	10-3
How a System Dump Is Invoked .....	10-4
When a Dump Occurs .....	10-5
The sysdumpdev Command .....	10-6
Dedicated Dump Device (1 of 2) .....	10-9
Dedicated Dump Device (2 of 2) .....	10-10
The sysdumpdev Command .....	10-11
dumpcheck Utility .....	10-12
Methods of Starting a Dump .....	10-14
Start a Dump from a TTY .....	10-16
Generating Dumps with smit .....	10-17
Dump-related LED Codes .....	10-18
Copying System Dump .....	10-20
Automatically Reboot After a Crash .....	10-22
Sending a Dump to IBM .....	10-23
Use kdb to Analyze a Dump .....	10-25
Next Step .....	10-27
Checkpoint .....	10-28
Unit Summary .....	10-29
<b>Unit 11. Performance and Workload Management .....</b>	<b>11-1</b>
Unit Objectives .....	11-2
11.1 Basic Performance Analysis and Workload Management .....	11-3
Performance Problems .....	11-4
Understand the Workload .....	11-5
Critical Resource: The Four Bottlenecks .....	11-7
Identify CPU-Intensive Programs: ps aux .....	11-9
Identify High-Priority Processes: ps -elf .....	11-11
Basic Performance Analysis .....	11-12
Monitoring CPU Usage: sar -u .....	11-13
Simultaneous Multi-Threading (SMT) .....	11-16
Monitoring Memory Usage: vmstat .....	11-18
Monitoring Disk I/O: iostat .....	11-20
topas .....	11-22
topas, vmstat, and iostat Enhancements for Micro-Partitioning (AIX 5.3) .....	11-23
AIX Performance Tools .....	11-25
AIX Tools: tprof .....	11-26
AIX Tools: svmon .....	11-28
AIX Tools: filemon .....	11-29
There Is Always a Next Bottleneck! .....	11-31
Workload Management Techniques (1 of 3) .....	11-32

Workload Management Techniques (2 of 3) .....	11-33
Workload Management Techniques (3 of 3) .....	11-34
Next Step .....	11-36
11.2 Performance Diagnostic Tool (PDT) .....	11-37
Performance Diagnostic Tool (PDT) .....	11-38
Enabling PDT .....	11-40
cron Control of PDT Components .....	11-42
PDT Files .....	11-43
Customizing PDT: Changing Thresholds .....	11-45
Customizing PDT: Specific Monitors .....	11-48
PDT Report Example (Part 1) .....	11-49
PDT Report Example (Part 2) .....	11-51
Next Step .....	11-53
Checkpoint .....	11-54
Unit Summary .....	11-55
<b>Unit 12. Security .....</b>	<b>12-1</b>
Unit Objectives .....	12-2
12.1 Authentication and Access Control Lists (ACLs) .....	12-3
Protecting Your System .....	12-4
How Do You Set Up Your PATH? .....	12-6
Trojan Horse: An Easy Example (1 of 3) .....	12-7
Trojan Horse: An Easy Example (2 of 3) .....	12-8
Trojan Horse: An Easy Example (3 of 3) .....	12-9
login.cfg: login prompts .....	12-10
login.cfg: Restricted Shell .....	12-12
Customized Authentication .....	12-13
Authentication Methods (1 of 2) .....	12-14
Authentication Methods (2 of 2) .....	12-15
Two-Key Authentication .....	12-16
Base Permissions .....	12-17
Extended Permissions: Access Control Lists .....	12-18
ACL Commands .....	12-19
ACL Keywords: permit and specify .....	12-21
ACL Keywords: deny .....	12-22
JFS2 Extended Attributes Version 2 (AIX 5.3) .....	12-23
Next Step .....	12-24
12.2 The Trusted Computing Base (TCB) .....	12-25
The Trusted Computing Base (TCB) .....	12-26
TCB Components .....	12-27
Checking the Trusted Computing Base .....	12-28
The sysck.cfg File .....	12-29
tcbck: Checking Mode Examples .....	12-31
tcbck: Checking Mode Options .....	12-32
tcbck: Update Mode Examples .....	12-34
chtcb: Marking Files As Trusted .....	12-35
tcbck: Effective Usage .....	12-36

---

Trusted Communication Path .....	12-37
Trusted Communication Path: Trojan Horse .....	12-38
Trusted Communication Path Elements .....	12-39
Using the Secure Attention Key (SAK) .....	12-40
Configuring the Secure Attention Key .....	12-41
chtcb: Changing the TCB Attribute .....	12-42
Checkpoint (1 of 2) .....	12-43
Checkpoint (2 of 2) .....	12-44
Unit Summary .....	12-45
Challenge LAB .....	12-46
<b>Appendix A. Command Summary .....</b>	<b>A-1</b>
<b>Appendix B. Checkpoint Solutions .....</b>	<b>B-1</b>
<b>Appendix C. RS/6000 Three-Digit Display Values .....</b>	<b>C-1</b>
<b>Appendix D. PCI Firmware Checkpoints and Error Codes .....</b>	<b>D-1</b>
<b>Appendix E. Location Codes .....</b>	<b>E-1</b>
<b>Appendix F. Challenge Exercise .....</b>	<b>F-1</b>
<b>Appendix G. Auditing Security Related Events .....</b>	<b>G-1</b>
<b>Glossary .....</b>	<b>X-1</b>



# Trademarks

The reader should recognize that the following terms, which appear in the content of this training document, are official trademarks of IBM or other companies:

IBM® is a registered trademark of International Business Machines Corporation.

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

AIX®	AIX 5L™	Micro-Partitioning™
MVS™	OS/2®	POWER™
POWER4™	POWER5™	POWER Gt1™
POWER Gt3™	PS/2®	pSeries®
Redbooks™	RS/6000®	SP™

Java and all Java-based trademarks are trademarks of Sun Microsystems, Inc. in the United States, other countries, or both.

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

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

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

ALERTS is a registered trademark of Alphablox Corporation in the United States, other countries, or both.

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



# Course Description

## AIX 5L System Administration II: Problem Determination

**Duration: 5 days**

### Purpose

The purpose of this course is to add to the system administrator's skills in determining the cause of a problem and carrying out the appropriate steps to fix the problem. Also, there is heavy emphasis on customizing the system.

### Audience

This course is targeted for system administrators with at least three months' experience in AIX and with other relevant education.

### Prerequisites

- Be familiar with the basic tools and commands in AIX. These include vi, SMIT, the Web-based documentation, and other commonly used commands, such as grep, find, mail, chmod, and ls
- Perform basic file manipulation and navigation of the file system
- Define basic file system and LVM terminology
- Carry out basic system installation activities including basic setup of printers, disks, terminals, users, and software
- Create and kill processes, prioritize them, and change their environment via profiles

### Objectives

On completion of this course, students should be able to:

- Perform problem determination and analyze the problem by performing the relevant steps, such as running diagnostics, analyzing the error logs, and carrying out dumps on the system.

## Contents

- Problem Determination Introduction
- The ODM
- System Initialization
- Disk Management Theory
- Disk Management Procedures
- Saving and Restoring Volume Groups
- Error Log and syslogd
- Diagnostics
- The AIX System Dump Facility
- Performance and Workload Management
- Security (Auditing, Authentication and ACLs, TCB)



# Agenda

## Day 1

Welcome

Unit 1

Problem Determination Introduction

Exercise 1 - Problem Determination Introduction

Unit 2

The ODM, Topic 1

The ODM, Topic 2

Exercise 2 - The Object Data Manager (ODM)

Unit 3

System Initialization Part I, Topic 1

System Initialization Part I, Topic 2

Exercise 3 - System Initialization Part 1

## Day 2

Unit 4

System Initialization Part II, Topic 1

System Initialization Part II, Topic 2

Exercise 4 - System Initialization Part 2

Unit 5

Disk Management Theory, Topic 1

Disk Management Theory, Topic 2

Exercise 5 - Fixing LVM-Related ODM Problems

Disk Management Theory, Topic 3

Exercise 6 - Mirroring rootvg

## Day 3

Unit 6

Disk Management Procedures, Topic 1

Disk Management Procedures, Topic 2

Exercise 7 - Exporting and Importing Volume Groups

Unit 7

Saving and Restoring Volume Groups, Topic 1

Saving and Restoring Volume Groups, Topic 2

Saving and Restoring Volume Groups, Topic 3

Saving and Restoring Volume Groups, Topic 4

Exercise 8 - Saving and Restoring a User Volume Group

Unit 8

Error Log and syslogd, Topic 1

Error Log and syslogd, Topic 2  
Exercise 9 - Working with syslogd and errnotify

## Day 4

Unit 9  
Diagnostics  
Unit 10  
The AIX System Dump Facility  
Exercise 10 - System Dump  
Unit 11  
Performance and Workload Management, Topic 1  
Exercise 11 - Basic Performance Commands  
Performance and Workload Management, Topic 2  
Exercise 12 - PDT

## Day 5

Unit 12  
Authentication  
Exercise 13 - Authentication and Access Control Lists  
Trusted Computing Base

# Unit 1. Problem Determination Introduction

## What This Unit Is About

This unit introduces the problem determination process and gives an overview of what will be covered in the course.

## What You Should Be Able to Do

After completing this unit you should be able to:

- Understand the process of resolving system problems
- Describe the four primary techniques for start to finish troubleshooting
- Know how to find the appropriate documentation

## How You Will Check Your Progress

Accountability:

- Checkpoint questions
- Lab Exercise

## References

SG24-5496     *Problem Solving and Troubleshooting in AIX 5L*

# Unit Objectives

---

After completing this unit, you should be able to:

- Understand the role of problem determination
- Provide methods for describing a problem and collecting the necessary information about the problem in order to take the best corrective course of action

© Copyright IBM Corporation 2004

Figure 1-1. Unit Objectives

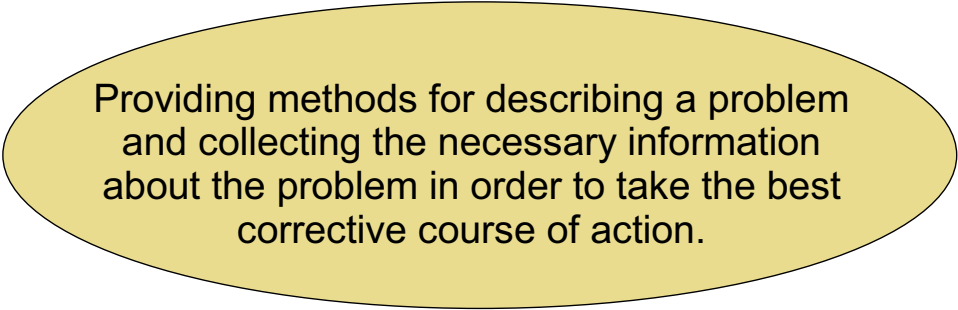
AU1612.0

## **Notes:**

## 1.1 Problem Determination Introduction

# Role of Problem Determination

---



Providing methods for describing a problem and collecting the necessary information about the problem in order to take the best corrective course of action.

© Copyright IBM Corporation 2004

Figure 1-2. Role of Problem Determination

AU1612.0

## **Notes:**

This course introduces problem determination and troubleshooting on the IBM p-Series and RS/6000 platforms running AIX 5L Version 5.2.

A problem can manifest itself in many ways, and very often the root cause might not be immediately obvious to system administrators and other support personnel. Once the problem and its cause are identified, the administrator should be able to identify the appropriate course of action to take.

The units will describe some common problems that can occur with AIX systems and will offer approaches to be taken to resolve them.

---

## Before Problems Occur

---

- Effective problem determination starts with a good understanding of the system and its components.
- The more information you have about the normal operation of a system, the better.
  - System configuration
  - Operating system level
  - Applications installed
  - Baseline performance
  - Installation, configuration, and service manuals

© Copyright IBM Corporation 2004

Figure 1-3. Before Problems Occur

AU1612.0

### **Notes:**

It's a good idea, whenever you approach a new system, to learn as much as you can about that system.

It is also critical to document both the logical and physical device information so that it is available when troubleshooting is necessary.

For example, look up information about the following:

- Machine architecture (model, cpu type)
- Physical volumes (type and size of disks)
- Volume groups (names, JBOD (just a bunch of disks) or RAID)
- Logical volumes (mirrored or not, which VG, type)
- Filesystems (which VG, what applications)
- Memory (size) and paging spaces (how many, location)

## Before Problems Occur: A Few Good Commands

---

- **lspv** - lists physical volumes, PVID, VG membership
- **lscfg** - provides information of system components
- **prtconf** - displays system configuration information
- **lsvg** - lists the volume groups
- **lsps** - displays information about paging spaces
- **lsfs** - give file system information
- **lsdev** - provides device information
- **getconf** - displays system configuration information
- **bootinfo** - displays system configuration information (unsupported)
- **snap** - collects system data

© Copyright IBM Corporation 2003

Figure 1-4. Before Problems Occur: A Few Good Commands

AU1612.0

### **Notes:**

This list provides a starting point for gathering documentation about the system.

There are many other commands as well.

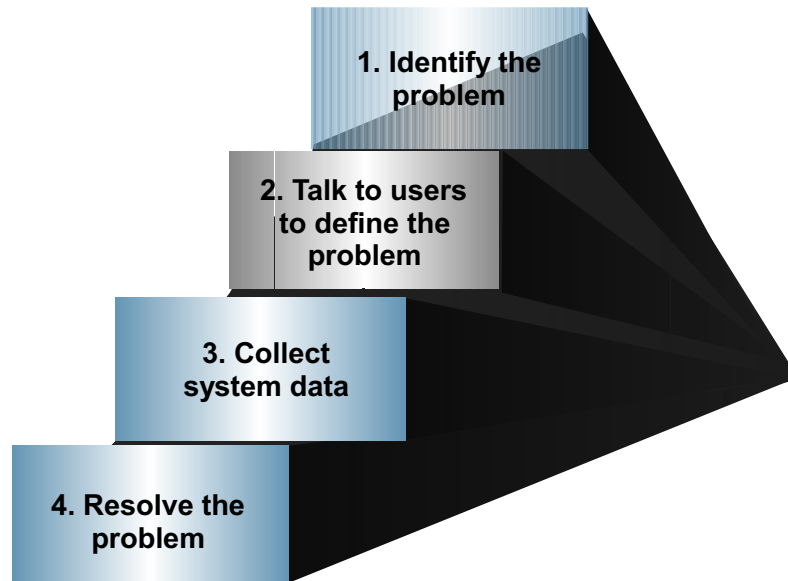
Be sure to check the man pages or the Commands Reference for correct syntax and option flags to be used to provide more specific information.



---

# Problem Determination Techniques

---



© Copyright IBM Corporation 2004

Figure 1-5. Problem Determination Techniques

AU1612.0

## **Notes:**

The “start-to-finish” method for resolving problems consists primarily of the four major components--identify the problem, talk to users, collect system data, and fix the problem.

## Identify the Problem

---

A clear definition of the problem:

- Gives clues as to the cause of the problem
- Aids in the choice of troubleshooting methods to apply

© Copyright IBM Corporation 2004

Figure 1-6. Identify the Problem

AU1612.0

### **Notes:**

The first step in problem resolution is to find out what the problem is. It is important to understand exactly what the users of the system perceive the problem to be.

---

## Define the Problem (1 of 2)

---

Understand what the users\* of the system perceive the problem to be.



\* **users** = data entry staff, programmers, system administrators, technical support personnel, management, application developers, operations staff, network users, etc.

© Copyright IBM Corporation 2004

Figure 1-7. Define the Problem (1 of 2)

AU1612.0

### **Notes:**

A problem can be identified by just about anyone who has use of or a need to interact with the system. If a problem is reported to you, it may be necessary to get details from the reporting user and then query others on the system for additional details or for a clear picture of what happened.

## Define the Problem (2 of 2)

---

- Ask questions:
  - What is the problem?
  - What is the system doing (or NOT doing)?
  - How did you first notice the problem?
  - When did it happen?
  - Have any changes been made recently?

"Keep 'em talking until the picture is clear!"



© Copyright IBM Corporation 2004

---

Figure 1-8. Define the Problem (2 of 2)

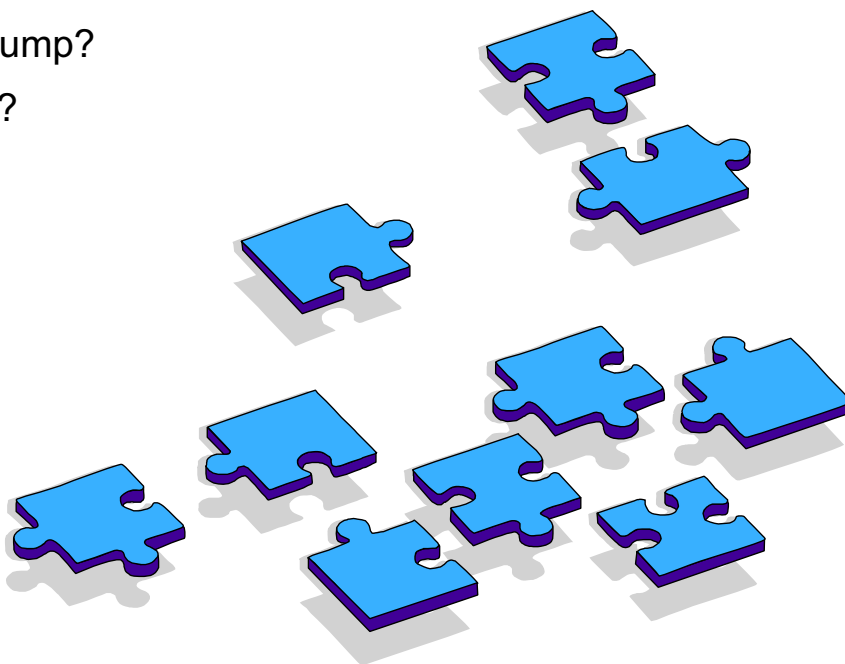
AU1612.0

### **Notes:**

Ask as many questions as you need to in order to get the entire history of the problem.

## Collect System Data

- How is the machine configured?
- What errors are being produced?
- What is the state of the OS?
- Is there a system dump?
- What log files exist?



© Copyright IBM Corporation 2004

Figure 1-9. Collect System Data

AU1612.0

### **Notes:**

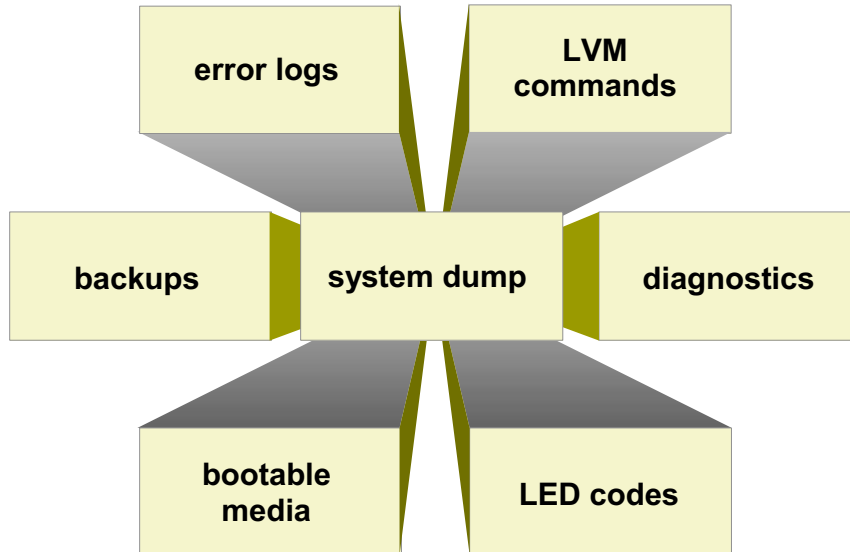
Some information about the system will have already been collected from the user during the process of defining the problem.

By using various commands, such as `lsdev`, `lspv`, `lsvg`, `lspp`, `lsattr` and others, you can gather further information about the system configuration.

If SMIT and the Web-based System Manager have been used, there will be system logs that could provide further information. The log files are normally contained in the home directory of the root user and are named `/smit.log` for SMIT and `/websm.log` for the Web-based System Manager, by default.

# Problem Determination Tools

---



© Copyright IBM Corporation 2004

Figure 1-10. Problem Determination Tools

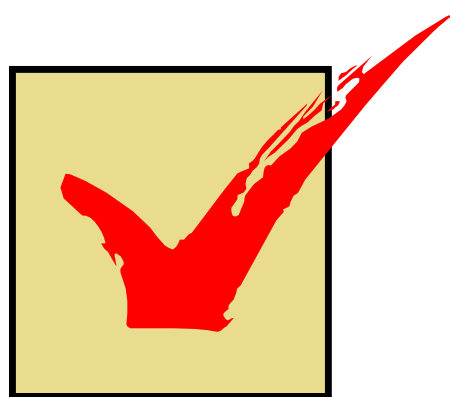
AU1612.0

## **Notes:**

## Resolve the Problem

---

- Use the information gathered.
- Use the tools available--commands documentation, downloadable fixes and updates.
- Contact IBM Support, if necessary.
- Keep a log of actions taken to correct the problem.



© Copyright IBM Corporation 2004

Figure 1-11. Resolve the Problem

AU1612.0

### **Notes:**

After all the information is gathered, select the procedure necessary to solve the problem. Keep a log of all actions you perform in trying to determine the cause of the problem, and any actions you perform to correct the problem.

The IBM e-server pSeries Information Center is a Web site that serves as a focal point for all information pertaining to pSeries and AIX. It provides a link to the entire pSeries library. A message database is available to search on error number, identifiers, LEDs and FAQs, how-to's, a troubleshooting guide, and more.

The URL is:

[http://publib16.boulder.ibm.com/pseries/en\\_US/infocenter/base](http://publib16.boulder.ibm.com/pseries/en_US/infocenter/base)

## Obtaining Software Fixes and Microcode Updates

---

Software fixes for AIX and hardware microcode updates are available on the Internet from the following URL:

<http://techsupport.services.ibm.com/server/fixes>

Access the Web site and register as a user



© Copyright IBM Corporation 2003

---

Figure 1-12. Obtaining Software Fixes and Microcode Updates

AU1612.0

### **Notes:**

Once you have determined the nature of your problem, you should try searching the Web site to see if you are experiencing known problems for which a fix has already been made available.



# Software Update Management Assistant (SUMA)

- Task-oriented utility which automates the retrieval of the following fix types:
  - Specific APAR
  - Specific PTF
  - Latest critical PTFs
  - Latest security PTFs
  - All latest PTFs
  - Specific fileset
  - Specific maintenance level
- Interfaces
  - SMIT (**smit suma fastpath**)
  - Command (**/usr/bin/suma**)
- Documentation
  - Man pages
  - Infocenter
  - AIX 5.3 Differences Guide



© Copyright IBM Corporation 2003

Figure 1-13. Software Update Management Assistant (SUMA)

AU1612.0

## Notes:

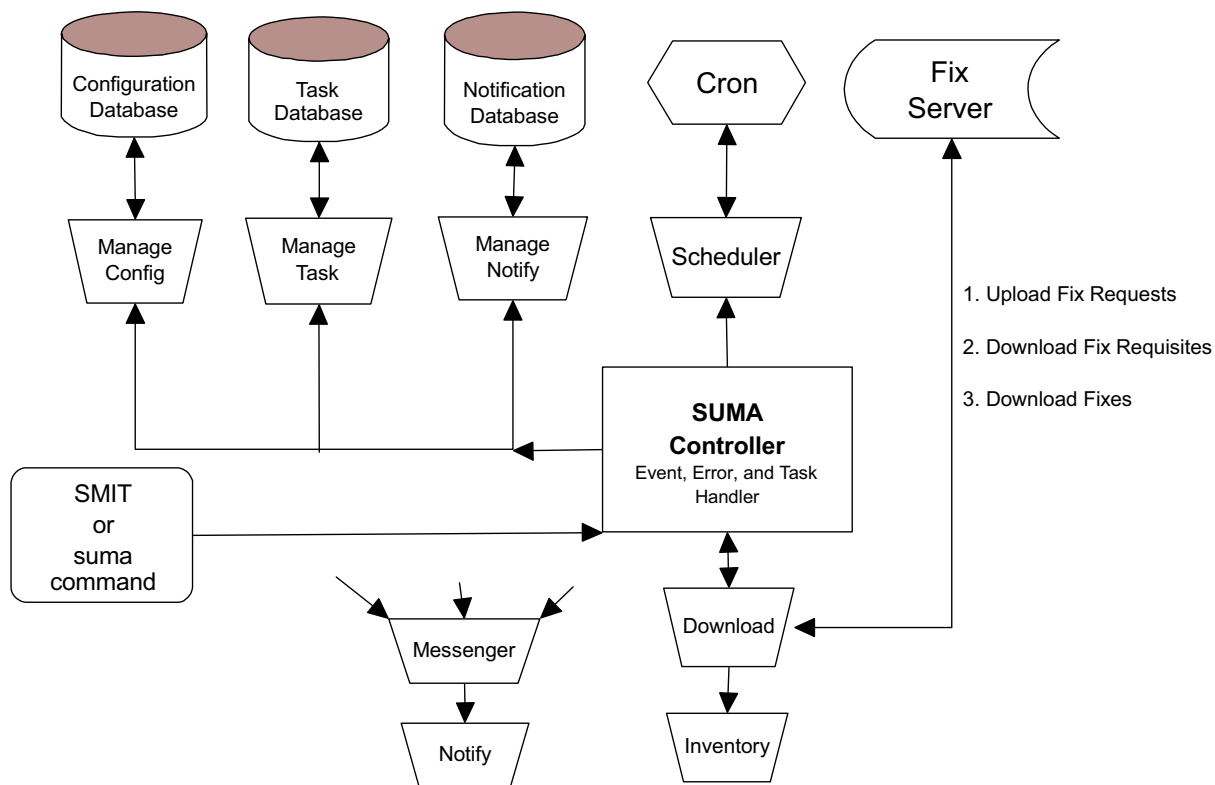
AIX 5L Version 5.3 introduces automatic download, scheduling and notification capabilities through the new Service Update Management Assistant (SUMA) tool. SUMA is fully integrated into the AIX Base Operating System and supports scheduled and unattended task-based download of Authorized Program Analysis Reports (APARs), Program Temporary Fixes (PTFs) and recommended maintenance levels (MLs). SUMA can also be configured to periodically check the availability of specific new fixes and entire maintenance levels, so that the time spent on such system administration tasks is cut significantly. The SUMA implementation allows for multiple concurrent downloads to optimize performance and has no dependency on any Web browser.

The Service Update Management Assistant will be available by default after any AIX 5L Version 5.3 operating system installation. All SUMA modules and the suma executable itself are contained in the bos.suma fileset. SUMA is implemented using the Perl programming language and therefore the Perl library extensions fileset perl.libext and the Perl runtime environment fileset perl.rte are prerequisites.

Highlights of this new feature include:

- Moves administrators away from the task of manually retrieving maintenance updates from the Web.
- Provides clients with flexible options.
- Schedule to run periodically. (For example, download the latest critical fixes weekly.)
- Can compare fixes needed against software inventory, fix repository, or a maintenance level.
- Receive e-mail notification after a fileset preview or download operation.
- Allows for ftp, http, or https transfers.
- Provides same requisite checking as the IBM fix distribution Web site.

# SUMA Modules



© Copyright IBM Corporation 2004

Figure 1-14. SUMA Modules

AU1612.0

## Notes:

The SUMA Controller utilizes certain SUMA modules to execute SUMA operations and functions.

### Download module

The download module provides functions related to network activities and is solely responsible for communicating with the IBM @Server pSeries support server. This communication manifests itself in two different transaction types. In the first a list of filesets is requested from the fix server based on the SUMA task data passed to download module. The second consists solely of downloading the requested files from the IBM @Server support server.

### Manage configuration module

The manage configuration module represents a utility class containing global configuration data and general-purpose methods. These methods allow for the validation of field names and field values since this information is predefined, meaning that there is a known set of

supported global configuration fields and their corresponding supported values. This module provides the interface to the global configuration database file.

#### Messenger module

The Messenger module provides messaging, logging, and notification capability. Messages will be logged (or displayed) when their specified verbosity level is not greater than the threshold defined by the SUMA global configuration. The log files themselves will be no larger than a known size (by default, 1 MB), as defined by the SUMA global configuration facility. When the maximum size is reached, a backup of the file will be created, and a new log file started, initially containing the last few lines of the previous file. Backup files are always created in the same directory as the current log file. Therefore, minimum free space for log files should keep this in mind. There are two log files which are located in the `/var/adm/ras/` directory. The log file `/var/adm/ras/suma.log` contains any messages that pertain to SUMA Controller operations. The other log file, `/var/adm/ras/suma_dl.log` tracks the download history of SUMA download operations and contains only entries of the form `DateStamp:FileName`. The download history file is appended when a new file is downloaded. The two logs are treated the same in respect to maximum size and creation/definition. The messenger module relies on contact information (e-mail addresses) from the notification database file which is managed by the notify module.

#### Notify module

The notify module manages the file which holds the contact information for SUMA event notifications. This database stores a list of email addresses for use by SMIT when populating the list of notification addresses as part of SUMA task configuration. Task module SUMA makes use of the task module to create, retrieve, view, modify, and delete SUMA tasks. All SUMA task related information is stored in a dedicated and private task database file.

#### Scheduler module

The scheduler module is responsible for handling scheduling of SUMA task execution and interacts with the AIX cron daemon and the files in `/var/spool/cron/crontabs` directory.

#### Inventory module

The inventory module returns the software inventory (installed or in a repository) of the local system (localhost) or a NIM client. It covers all software which is in the `installp`, `RPM`, or `ISMP` packaging format. If the system specified to the module is not local then the system must be a NIM client of the local system.

#### Utility and database modules

Other modules supply private utilities for SUMA code and utilities for handling the stanza-style SUMA databases. The Configuration, Task, and Notification Database are within the `/var/suma/data` path.

## SUMA Examples (1 of 2)

To immediately execute a task that will preview downloading any critical fixes that have become available and are not already installed on your system:

```
# suma -x -a RqType=Critical -a Action=Preview
```

To create and schedule a task that will download the latest fixes monthly (For example, on the 15th of every month at 2:30 AM):

```
# suma -s "30 2 15 * *" -a RqType=Latest \  
-a DisplayName="Critical fixes - 15th Monthly"  
Task ID 4 created.
```

To list the newly created SUMA task ID 4:

```
# suma -l 4
```

© Copyright IBM Corporation 2004

Figure 1-15. SUMA Examples (1 of 2)

AU1612.0

### Notes:

The first example will preview or pretend downloading all of the “Critical” fixes which are not already installed on the local machine. The output would show something like the following:

```
*****
```

```
Performing preview download.
```

```
*****
```

```
Download SKIPPED: Java131.adt.debug.1.3.1.13.bff  
Download SKIPPED: Java131.adt.includes.1.3.1.5.bff  
Download SKIPPED: Java131.ext.commapl.1.3.1.2.bff  
Download SKIPPED: Java131.ext.jaas.1.3.1.5.bff  
Download SKIPPED: Java131.ext.java3d.1.3.1.1.bff  
Download SKIPPED: Java131.ext.plugin.1.3.1.15.bff  
Download SKIPPED: Java131.ext.xml4j.1.3.1.1.bff  
Download SKIPPED: Java131.rte.bin.1.3.1.15.bff  
Download SUCCEEDED: /usr/sys/inst.images/installp/ppc/Java131.rte.bin.1.3.1.16.bff  
Download SUCCEEDED:
```

```
/usr/sys/inst.images/installp/ppc/Java131.rte.bin.1.3.1.2.bffDownload SUCCEEDED:  
/usr/sys/inst.images/installp/ppc/Java131.rte.lib.1.3.1.15.bff  
Download SUCCEEDED: /usr/sys/inst.images/installp/ppc/Java131.rte.lib.1.3.1.16.bff  
Download SUCCEEDED: /usr/sys/inst.images/installp/ppc/Java131.rte.lib.1.3.1.2.bff
```

```
.  
. .
```

```
Download SUCCEEDED: /usr/sys/inst.images/installp/ppc/xlsmp.rte.1.3.6.0.bff  
Download SUCCEEDED: /usr/sys/inst.images/installp/ppc/xlsmp.rte.1.3.8.0.bff
```

Summary:

```
    257 downloaded  
     0 failed  
     8 skipped
```

To download the files, rerun the command without the attribute “Action=Preview”. This will download the update filesets in the /usr/sys/inst.images path if we haven’t changes the default location. Use `suma -D` to display the default configuration options.

The second example creates a new SUMA task and a crontab job. The `-s` flag’s parameter value is in crontab time format. All saved SUMA tasks get a “Task ID” number. These tasks can be listed with `suma -l`.

---

## SUMA Examples (2 of 2)

---

To create and schedule a task that will check monthly (for example, on the 15th of every month at 2:30 AM) for all the latest new updates, and download any that are not already in the /tmp/latest repository, type the following:

```
suma -s "30 2 15 * *" -a RqType=Latest \  
-a DLTarget=/tmp/latest -a FilterDir=/tmp/latest
```

© Copyright IBM Corporation 2004

Figure 1-16. SUMA Examples (2 of 2)

AU1612.0

### **Notes:**

The `sum -D` shows configuration options as in the following output:

```
# suma -D  
  DisplayName=  
  Action=Download  
  RqType=Security  
  RqName=  
  RqLevel=  
  PreCoreqs=y  
  Ifreqs=y  
  Supersedes=n  
  ResolvePE=IfAvailable  
  Repeats=y  
  DLTarget=/usr/sys/inst.images  
  NotifyEmail=root  
  FilterDir=/usr/sys/inst.images
```

```
FilterML=  
FilterSysFile=localhost  
MaxDLSize=-1  
Extend=y  
MaxFSSize=-1
```

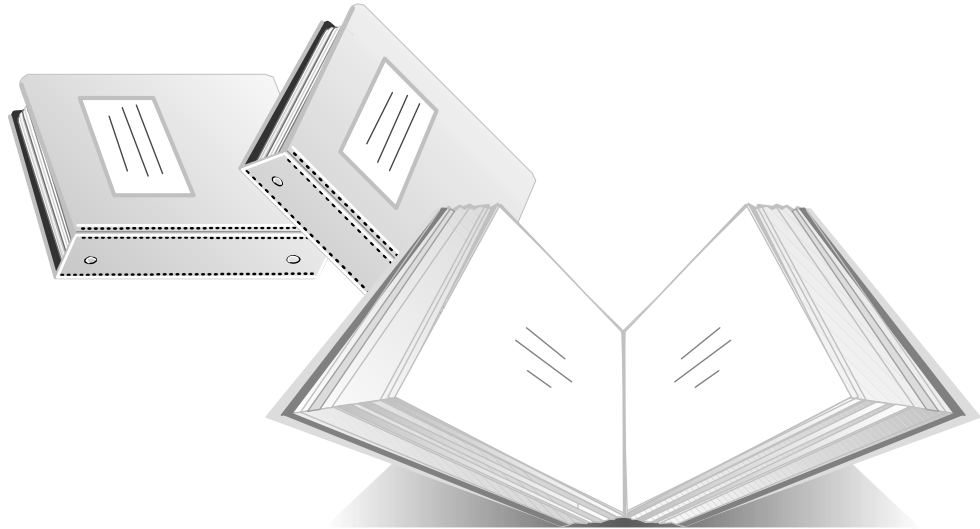
When running or creating a suma task, you can override the default settings. In the example above, we are overriding the “DLTarget” and “FilterDir” attribute values. This example is good for only downloading what you don’t already have in a directory which is being used as a repository for fixes.



## Relevant Documentation

---

- AIX Operating System Publications
- *pSeries and RS/6000 System Installation and Service Guides*
- IBM Redbooks
- Information Center documents
  - [http://publib16.boulder.ibm.com/pseries/en\\_US/infocenter/base/](http://publib16.boulder.ibm.com/pseries/en_US/infocenter/base/)



© Copyright IBM Corporation 2004

Figure 1-17. Relevant Documentation

AU1612.0

### **Notes:**

Most AIX software and hardware documentation can be viewed online at the IBM Web site:  
<http://www-1.ibm.com/servers/eserver/pseries/library>.

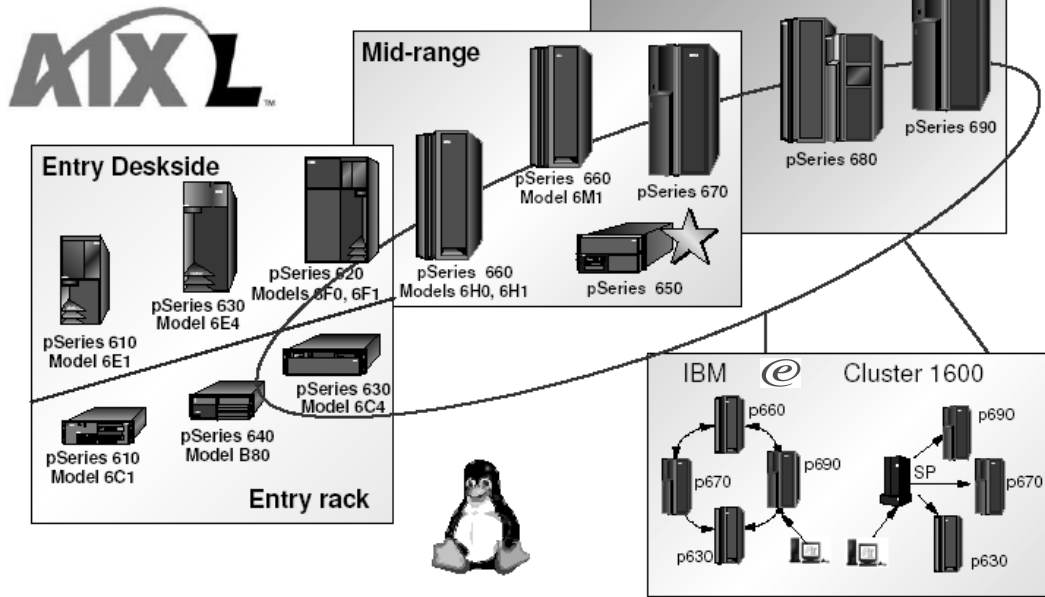
Redbooks can be viewed, downloaded or ordered from the Redbooks Web site:  
<http://www.ibm.com/redbooks>



## 1.2 pSeries Product Family

# IBM @ Server pSeries Product Family

**1,500,000 + Systems**  
**150,000 + Customers**



**Also includes eServer POWER5 Systems:**

- p5-520**
- p5-550**
- p5-570**
- p5-590**
- p5-595**

© Copyright IBM Corporation 2004

Figure 1-18. IBM @Server pSeries Product Family

AU1612.0

## Notes:

AIX 5L Version 5.2 and above exclusively supports PCI architecture machines. There is a minimum hardware requirement of 128 MB of RAM and 2.2 GB of disk space.

World-class UNIX and Linux implementations from IBM pSeries are the result of leading-edge IBM technologies. Through high-performance and flexibility between AIX and Linux operating environments, IBM pSeries delivers reliable, cost-effective solutions for commercial and technical computing applications in the entry, mid-range and high-end UNIX segments.

pSeries solutions offer the flexibility and availability to handle your most mission-critical and data-intensive applications. pSeries solutions also deliver the performance and application versatility necessary to meet the dynamic requirements of today's e-infrastructure environments.

IBM Cluster 1600 lets customers consolidate hundreds of applications and manage from a single point of control. IBM clustering hardware and software provide the building blocks, with availability, scalability, security and single-point-of-management control, to satisfy these needs.

Interconnecting two or more computers into a single, unified computing resource offers a set of system-wide, shared resources that cooperate to provide flexibility, adaptability and increased availability for services essential to customers, business partners, suppliers, and employees.

# AIX 5L 5.2 and 5.3 Logical Partition (LPAR) Support

*Improved throughput and resource utilization through increased workload management flexibility*

**Dynamic LPAR**  
*Add or remove processors, adapters and memory without requiring a reboot*

- AIX enablement for 32 partitions

**Dynamic Reconfiguration APIs**  
*Applications and middleware can automatically adjust to changes in hardware resources.*

**Dynamic Capacity Upgrade On Demand**  
*Customer's can activate additional processors without having to reboot.*

**Hot Sparring w/CUoD**  
*Dynamic substitution of failed processors with spare, unlicensed processors*

- 1-32 processors per partition
- Single Adapter I/O allocation
- Memory in 256MB increments
- Hardware enforced isolation

© Copyright IBM Corporation 2004

Figure 1-19. AIX 5L 5.2 and 5.3 Logical Partition (LPAR) Support

AU1612.0

## Notes:

Put the four bullet items and their detail, which are located on the right side of the page, in the “notes” section. Also, add the following:

Logical partitioning is a server design feature that provides more end-user flexibility by making it possible to run multiple, independent operating system images concurrently on a single server.

Dynamic Logical Partitioning (DLPAR) increases the flexibility of partitioned systems by enabling administrators to add, remove, or move system resources such as memory, PCI Adapters, and CPU between partitions without the need to reboot each partition. This allows a systems administrator to assign resources where they are needed most, now dynamically, without having to reboot a partition after it is modified. In addition, system administrators can adjust to changing hardware requirements within an LPAR environment, without impacting systems availability.

Dynamic CuOD enables a customer to order and install systems with additional processors and keep those resources in reserve until they are required as future applications workloads dictate. To enable the additional resources, the system administrator can

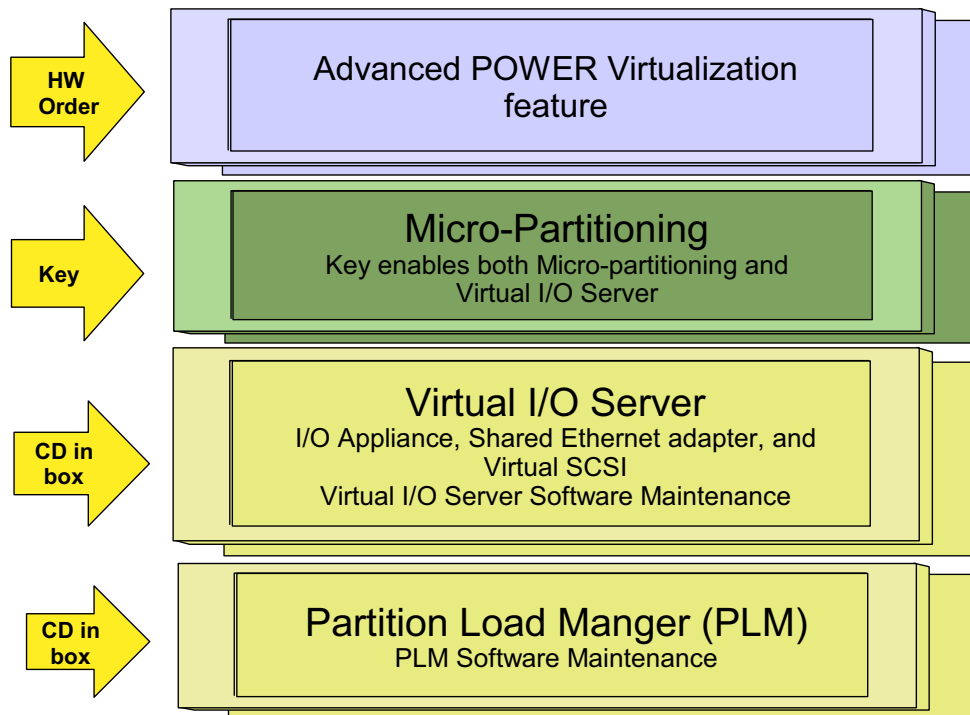
dynamically turn on the resources and then use dynamic LPAR services to assign those resource to one or more partitions without having to bring the system down. In addition, Dynamic CPU Guard is an important solution that can automatically and dynamically remove failing processors from a system image before they can cause a system failure. If spare processors are available on the systems, they can automatically replace the failing processors.

Most POWER4 pSeries servers implement LPAR. The LPAR information available on the links provided on this web site apply to all these servers unless otherwise noted. The introduction of logical partitioning (LPAR) technology to IBM pSeries™ systems has greatly expanded the options for deploying applications and workloads onto server hardware. IBM is adding to that LPAR capability with the introduction of dynamic LPAR (DLPAR), in which partition resources can be moved from one partition to another without requiring a reboot of the system or affected partitions.

Logical partitioning is intended to address a number of pervasive requirements, including:

- Server consolidation: The ability to consolidate a set of disparate workloads and applications onto a smaller number of hardware platforms, in order to reduce total cost of ownership (administrative and physical planning overhead).
- Production and test environments: The ability to have an environment to test and migrate software releases or applications, which runs on exactly the same platform as the production environment to ensure compatibility, but does not cause any exposure to the production environment.
- Data and workload isolation: The ability to support a set of disparate applications and data on the same server, while maintaining very strong isolation of resource utilization and data access.
- Scalability balancing: The ability to create resource configurations appropriate to the scaling characteristics of a particular application, without being limited by hardware upgrade granularities.
- Flexible configuration: The ability to change configurations easily to adapt to changing workload patterns and capacity requirements especially enhanced by the DLPAR feature

# Advance POWER Virtualization Feature (POWER5)



© Copyright IBM Corporation 2003

Figure 1-20. Advance POWER Virtualization Feature (POWER5)

AU1612.0

## Notes:

The Advanced POWER Virtualization feature is a combination of hardware enablement for Micro-partitions and software that supports the virtual I/O environment on POWER5 systems.

The Advanced POWER Virtualization optional feature includes:

Firmware enablement for Micro-partitions

Micro-partitioning is a mainframe-inspired technology that is based on two major advances in the area of server virtualization. Physical processors and I/O devices have been virtualized, enabling these resources to be shared by multiple partitions. There are several advantages associated with this technology, including finer grained resource allocations, more partitions, and higher resource utilization.

The virtualization of processors requires a new partitioning model, since it is fundamentally different from the partitioning model used on POWER4 processor-based servers, where whole processors are assigned to partitions. These processors are owned by the partition and are not easily shared with other partitions. They may be assigned through manual



dynamic logical partitioning (LPAR) procedures. In the new micro-partitioning model, physical processors are abstracted into virtual processors, which are assigned to partitions. These virtual processor objects cannot be shared, but the underlying physical processors are shared, since they are used to actualize virtual processors at the platform level. This sharing is the primary feature of this new partitioning model, and it happens automatically.

Installation image for the Virtual I/O Server software, which supports:

- Shared Ethernet Adapter

- Virtual SCSI server

The Virtual I/O Server provides the Virtual SCSI (VSCSI) Target and Shared Ethernet adapter virtual I/O function to client partitions. This is accomplished by assigning physical devices to the Virtual I/O Server partition, then configuring virtual adapters on the clients to allow communication between the client and the Virtual I/O Server. All aspects of Virtual I/O server administration are accomplished through a special command line interface.

### Partition Load Manager

PLM for AIX 5L is a resource manager that provides automated CPU and memory resource management across DLPAR capable logical partitions running AIX 5L V5.2 or AIX 5L V5.3. PLM allocates resources to partitions on-demand, within the constraints of a user-defined policy. It assigns resources from partitions with low usage to partitions with a higher demand, improving the overall resource utilization of the system. PLM works with both dedicated and shared processor environment partitions. The only restriction is that all partitions in a group must be of the same type. In dedicated LPARs, it will work by adding or removing real processors. In shared processor LPARs, it will work by adding or removing processing units from the capacity entitlement.

## Virtual Ethernet (AIX 5.3 and POWER5)

---

- Enables inter-partition communication.
  - In-memory point to point connections
- Physical network adapters are not needed.
- Similar to high-bandwidth Ethernet connections.
- Supports multiple protocols (IPv4, IPv6, and ICMP).
- No Advanced POWER Virtualization feature required.
  - POWER5 Systems
  - AIX 5L V5.3 or appropriate Linux level
  - Hardware management console (HMC)

© Copyright IBM Corporation 2004

Figure 1-21. Virtual Ethernet (AIX 5.3 and POWER5)

AU1612.0

### **Notes:**

The Virtual Ethernet enables inter-partition communication without the need for physical network adapters in each partition. The Virtual Ethernet allows the administrator to define in-memory point to point connections between partitions. These connections exhibit similar characteristics, as high bandwidth Ethernet connections supports multiple protocols (IPv4, IPv6, and ICMP). Virtual Ethernet requires a POWER5 system with either AIX 5L V5.3 or the appropriate level of Linux and a Hardware Management Console (HMC) to define the Virtual Ethernet devices. Virtual Ethernet does not require the purchase of any additional features or software, such as the Advanced Virtualization Feature.

Virtual Ethernet is also called Virtual LAN or even VLAN, which can be confusing, because these terms are also used in network topology topics. But the Virtual Ethernet, which uses virtual devices, has nothing to do with the VLAN known from Network-Topology, which divides a LAN in further Sub-LANs.

## Checkpoint Questions

---

- What are the four major problem determination steps?
- Who should provide information about the problems?
- **T/F:** If there is a problem with the software, it is necessary to get the next release of the product to resolve the problem.
- **T/F:** Documentation can be viewed or downloaded from the IBM Web site.

© Copyright IBM Corporation 2004

Figure 1-22. Checkpoint Questions

AU1612.0

### **Notes:**

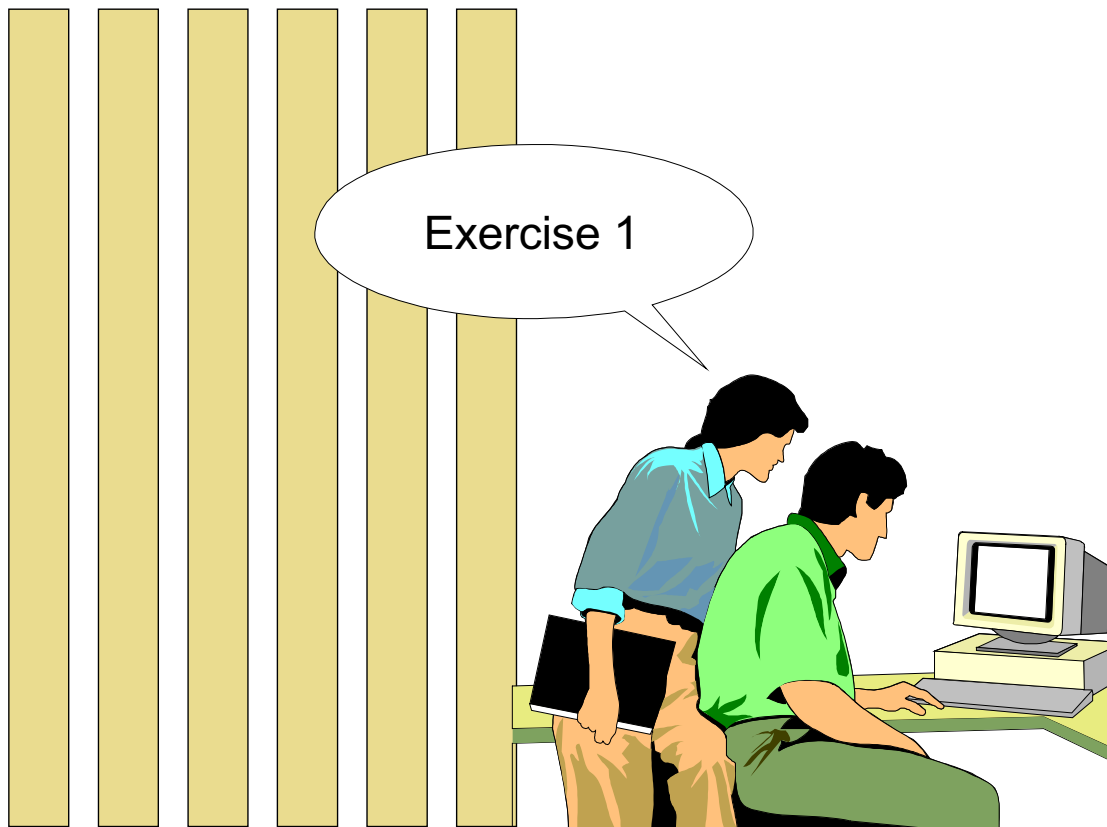


Figure 1-23. Exercise 1

AU1612.0

**Notes:**

## Unit Summary

---

Having completed this unit, you should be able to:

- Understand the role of problem determination
- Provide methods for describing a problem and collecting the necessary information about the problem in order to take the best corrective course of action

© Copyright IBM Corporation 2004

---

Figure 1-24. Unit Summary

AU1612.0

### **Notes:**



# Unit 2. The Object Data Manager (ODM)

## What This Unit Is About

This unit describes the structure of the ODM. It shows the use of the ODM command line interface and describes the role of ODM in device configuration. Also, the meaning of the most important ODM files is defined.

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Define the structure of the ODM
- Work with the ODM command line interface
- Define the meaning of the most important ODM files

## How You Will Check Your Progress

Accountability:

- Checkpoint questions
- Lab exercise

## References

Online	<i>AIX Commands Reference</i>
Online	<i>General Programming Concepts</i>
Online	<i>Technical Reference: Kernel and Subsystems</i>

# Unit Objectives

---

After completing this unit, students should be able to:

- Define the structure of the ODM
- Work with the ODM command line interface
- Describe the role of ODM in device configuration
- Define the meaning of the most important ODM files

© Copyright IBM Corporation 2004

Figure 2-1. Unit Objectives

AU1612.0

## **Notes:**

The ODM is a very important component of AIX and is one major difference to other UNIX systems. The structure of ODM database files is described in this unit, and how you can work with ODM files using the ODM command line interface.

From the administrator's point of view it is very important that you are able to understand the role of ODM during device configuration, which is another major point in this unit.



## 2.1 Introduction to the ODM

## What Is the ODM?

---

- The Object Data Manager (ODM) is a database intended for storing system information.
- Physical and logical device information is stored and maintained as objects with associated characteristics.

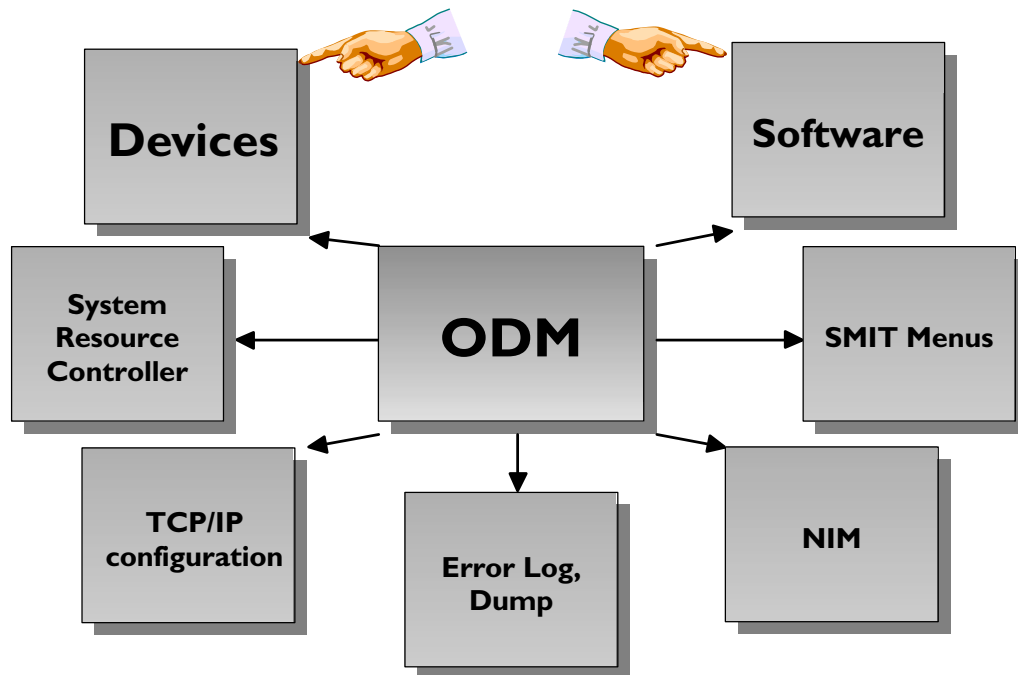
© Copyright IBM Corporation 2004

Figure 2-2. What Is the ODM?

AU1612.0

### **Notes:**

## Data Managed by the ODM



© Copyright IBM Corporation 2004

Figure 2-3. Data Managed by the ODM

AU1612.0

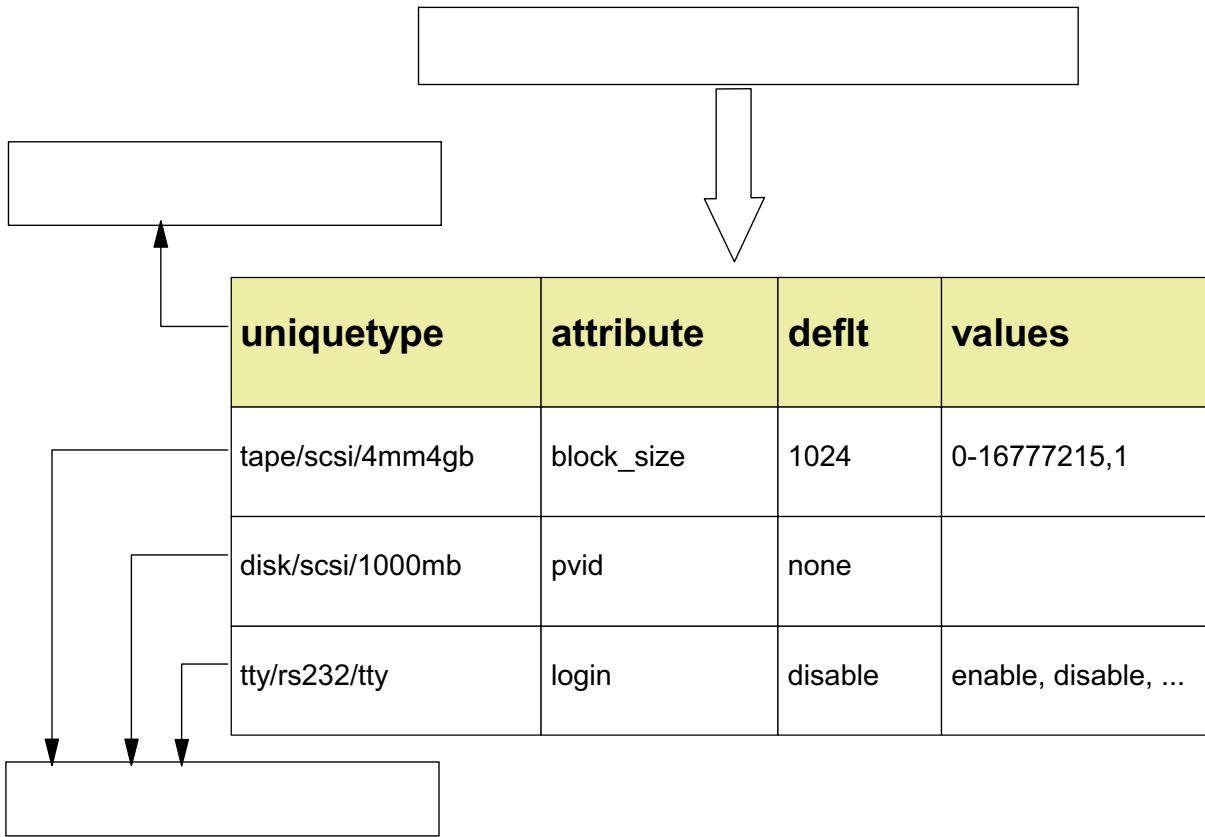
### Notes:

The ODM manages the following system data:

- Device configuration data
- Software Vital Product Data (SWVPD)
- System Resource Controller Data (SRC)
- TCP/IP configuration data
- Error Log and Dump information
- NIM (Network Installation Manager) information
- SMIT menus and commands

Our **main emphasis** in this unit is on **devices** and ODM files that are used to store **vital software product data**. During the course many other ODM classes are described.

# ODM Components



© Copyright IBM Corporation 2004

Figure 2-4. ODM Components

AU1612.0

**Notes:**

This page identifies the basic components of ODM. Your instructor will complete this page. Please complete the picture during the lesson.

For safety reasons the ODM data is stored in **binary** format. To work with ODM files you must use the ODM command line interface. It is not possible to update ODM files with an editor.

## ODM Database Files

<b>Predefined device information</b>	PdDv, PdAt, PdCn
<b>Customized device information</b>	CuDv, CuAt, CuDep, CuDvDr, CuVPD, Config_Rules
<b>Software vital product data</b>	history, inventory, lpp, product
SMIT menus	sm_menu_opt, sm_name_hdr, sm_cmd_hdr, sm_cmd_opt
Error log, alog and dump information	SWservAt
System Resource Controller	SRCsubsys, SRCsubsvr, ...
Network Installation Manager (NIM)	nim_attr, nim_object, nim_pdatr

© Copyright IBM Corporation 2004

Figure 2-5. ODM Database Files

AU1612.0

### Notes:

This list summarizes the major ODM files in AIX. In this unit we concentrate on ODM classes that are used to store device information and software product data.

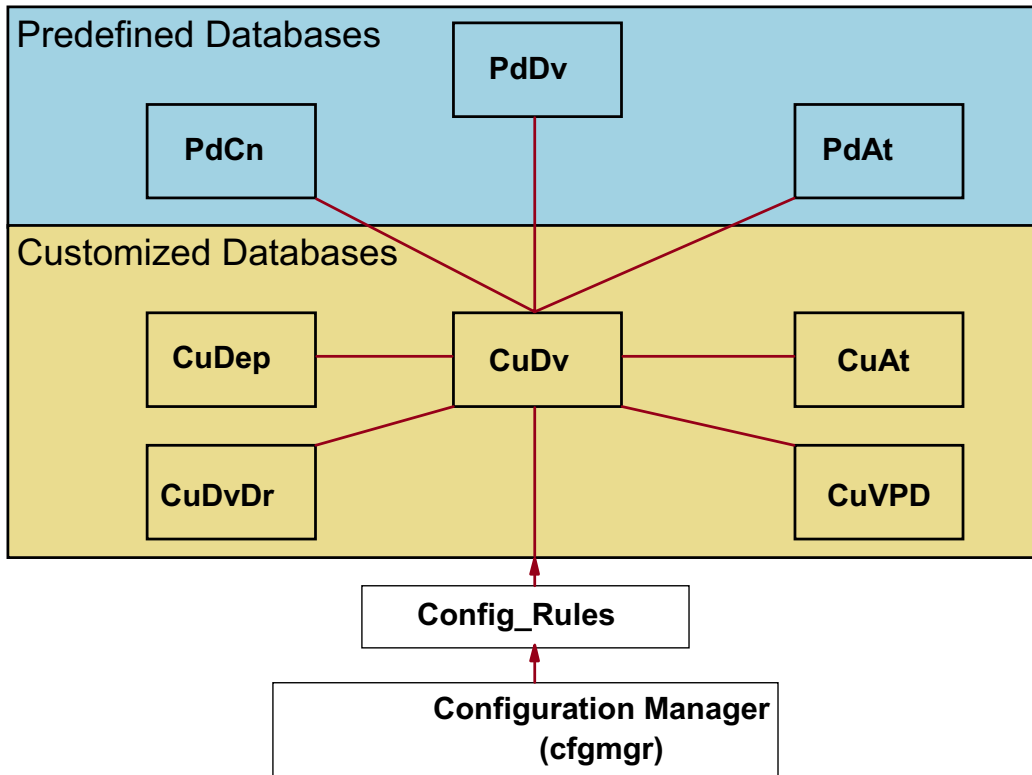
At this point you see ODM classes that contain predefined device configuration and others that contain customized device configuration. What is the difference between both?

**Predefined** device information describes all **supported** devices. **Customized** device information describes all devices that are **actually attached** to the system.

It is very important that you understand the difference between both classifications.

The classes themselves are described in more detail in the next topic of this unit.

# Device Configuration Summary



© Copyright IBM Corporation 2004

Figure 2-6. Device Configuration Summary

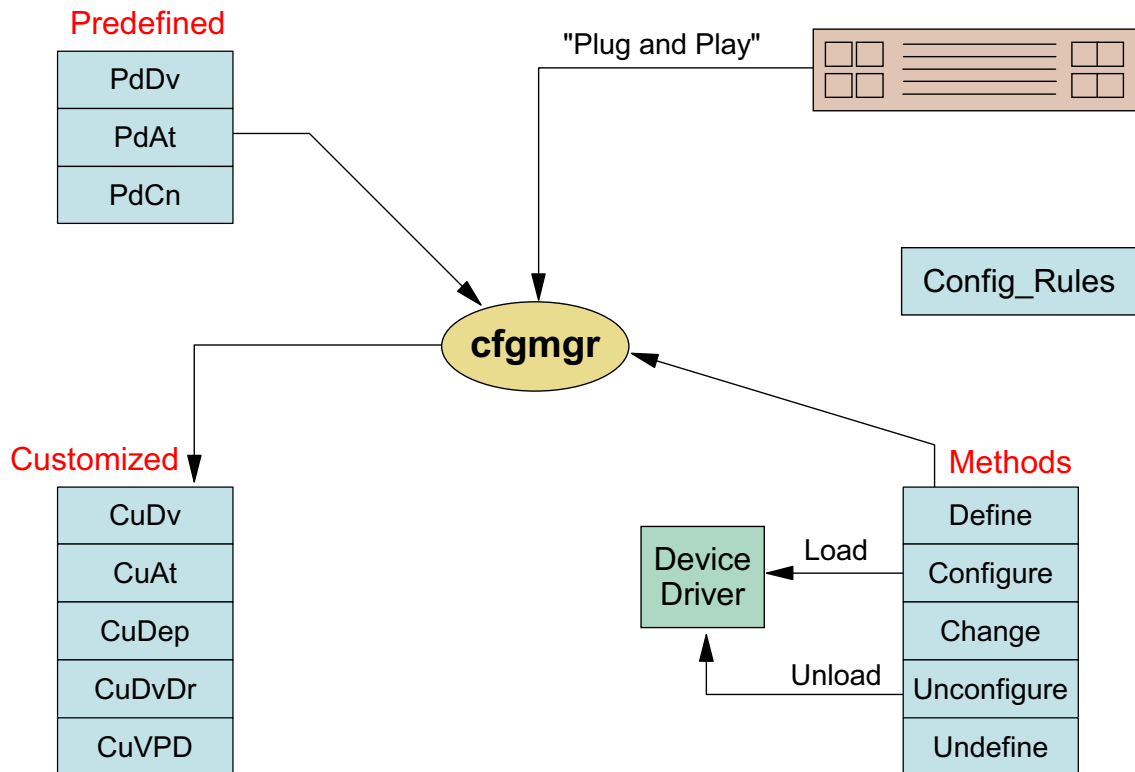
AU1612.0

## Notes:

This page shows the ODM object classes used during the configuration of a device.

When an AIX system boots, the **cfgmgr** is responsible for configuring devices. There is one ODM object class which the **cfgmgr** uses to determine the correct sequence when configuring devices: **Config\_Rules**

# Configuration Manager



© Copyright IBM Corporation 2004

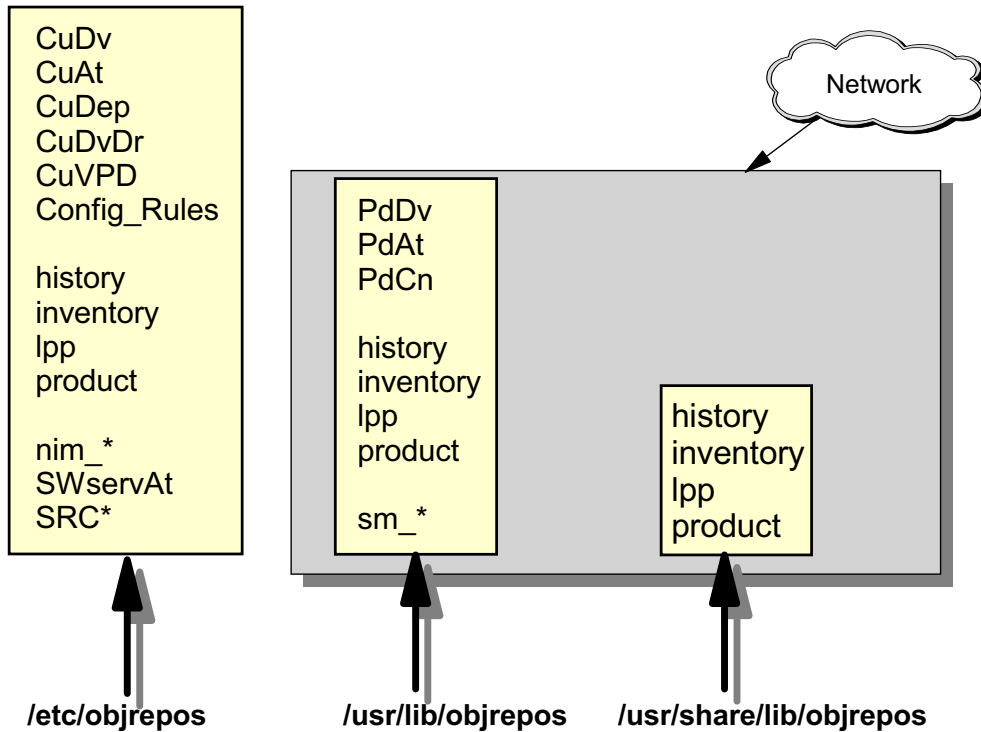
Figure 2-7. Configuration Manager

AU1612.0

## Notes:

Although **cfgmgr** gets credit for managing devices (adding, deleting, changing, and so forth) it is actually the **Config\_Rules** object class that does the work through various methods files.

# Location and Contents of ODM Repositories



© Copyright IBM Corporation 2004

Figure 2-8. Location and Contents of ODM Repositories

AU1612.0

## Notes:

To support diskless, dataless and other workstations, the ODM object classes are held in three repositories:

### /etc/objrepos

Contains the customized devices object classes and the four object classes used by the Software Vital Product Database (SWVPD) for the / (**root**) part of the installable software product. The root part of the software contains files that must be installed on the target system. To access information in the other directories this directory contains symbolic links to the predefined devices object classes. The links are needed because the **ODMDIR** variable points to only /etc/objrepos. It contains the part of the product that cannot be shared among machines. Each client must have its own copy. Most of this software requiring a separate copy for each machine is associated with the configuration of the machine or product.



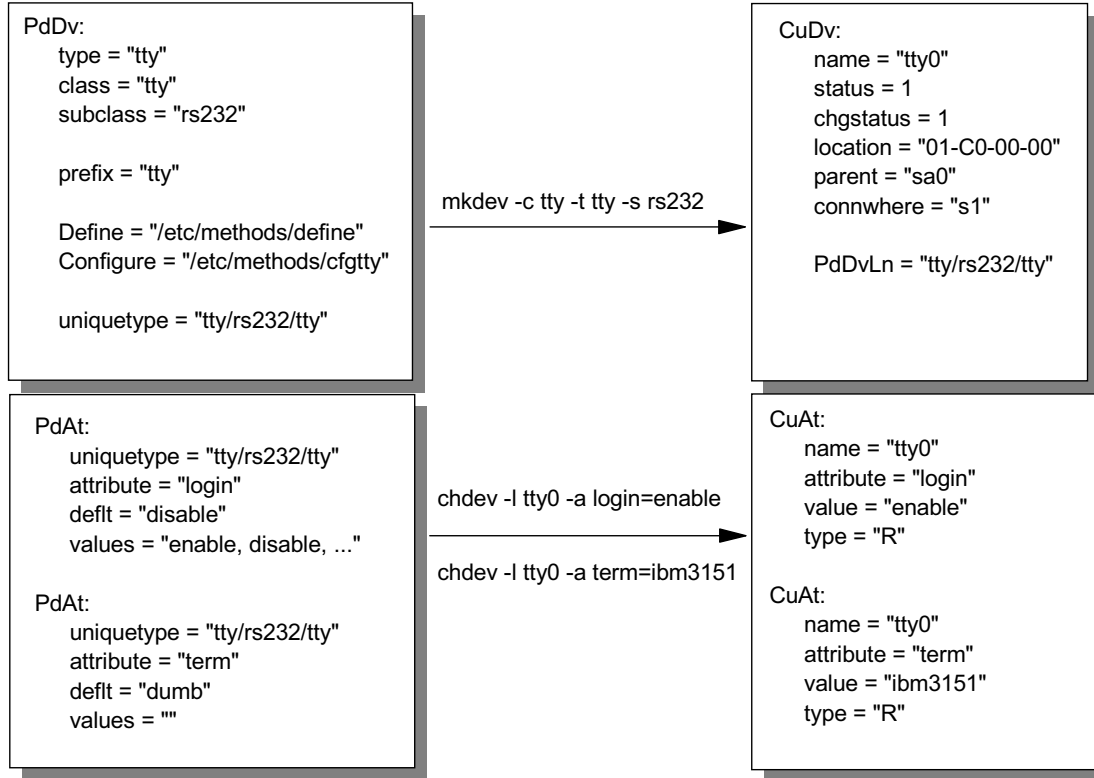
**/usr/lib/objrepos**

Contains the predefined devices object classes, SMIT menu object classes and the four object classes used by the SWVPD for the **/usr** part of the installable software product. The object classes in this repository can be shared across the network by **/usr** clients, dataless and diskless workstations. Software installed in the **/usr**-part can be shared among several machines with compatible hardware architectures.

**/usr/share/lib/objrepos**

Contains the four object classes used by the SWVPD for the **/usr/share** part of the installable software product. The **/usr/share** part of a software product contains files that are not hardware dependent. They can be shared among several machines, even if the machines have a different hardware architecture. An example for this are terminfo files that describe terminal capabilities. As terminfo is used on many UNIX systems, terminfo files are part of the **/usr/share**-part of a system product.

# How ODM Classes Act Together



© Copyright IBM Corporation 2004

Figure 2-9. How ODM Classes Act Together

AU1612.0

## Notes:

This visual summarizes how ODM classes act together.

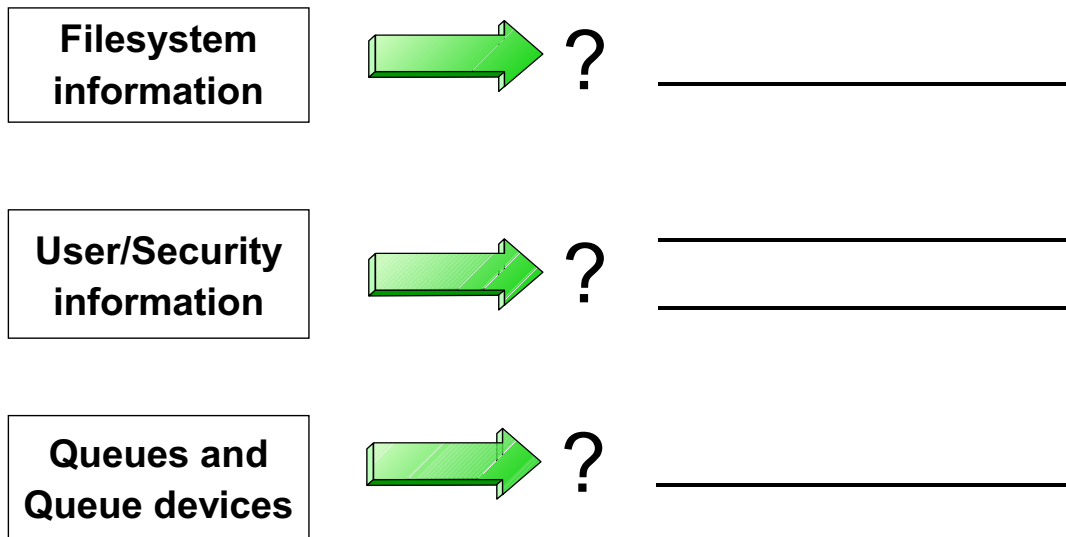
1. When a device is defined in AIX, the device must be defined in ODM class PdDv.
2. A device can be defined by either the **cfgmgr** (if the device is detectable), or by the **mkdev** command. Both commands use the **define method** to generate an instance in ODM class CuDv. The **configure method** is used to load a specific device driver and to generate an entry in the **/dev** directory.

Notice the link **PdDvLn** from CuDv back to PdDv.

3. At this point you only have default attribute values in PdAt, which means for a terminal you could not login (default is **disable**) and the terminal type is **dumb**. If you change the attributes, for example, login to **enable** and term to **ibm3151**, you get objects describing the nondefault values in CuAt.

## Data Not Managed by the ODM

---



© Copyright IBM Corporation 2004

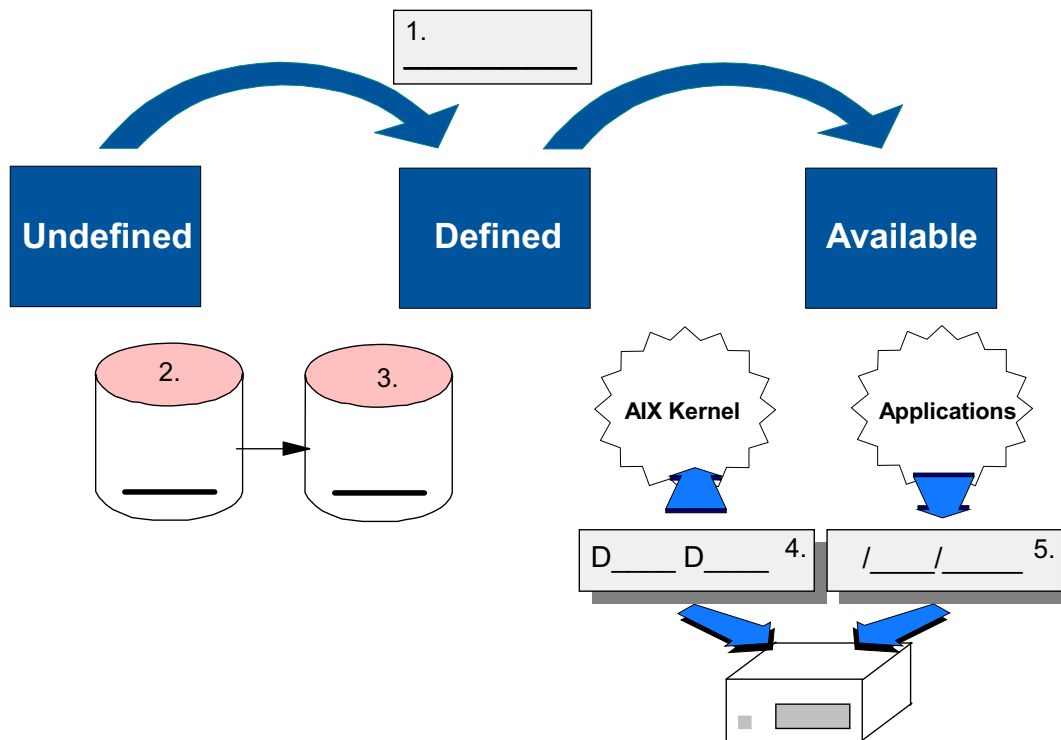
Figure 2-10. Data Not Managed by the ODM

AU1612.0

### **Notes:**

Your instructor will complete this page during the lesson.

# Let's Review: Device Configuration and the ODM



© Copyright IBM Corporation 2003

Figure 2-11. Let's Review: Device Configuration and the ODM

AU1612.0

## Notes:

Please answer the following questions. Please put the answers in the picture above. If you are unsure about a question, leave it out.

1. Which command configures devices in an AIX system? (Note: This is not an ODM command)?
2. Which ODM class contains all devices that your system supports?
3. Which ODM class contains all devices that are configured in your system?
4. Which programs are loaded into the AIX kernel that control access to the devices?
5. If you have a configured tape drive **rmt1**, which special file do applications access to work with this device?

# ODM Commands

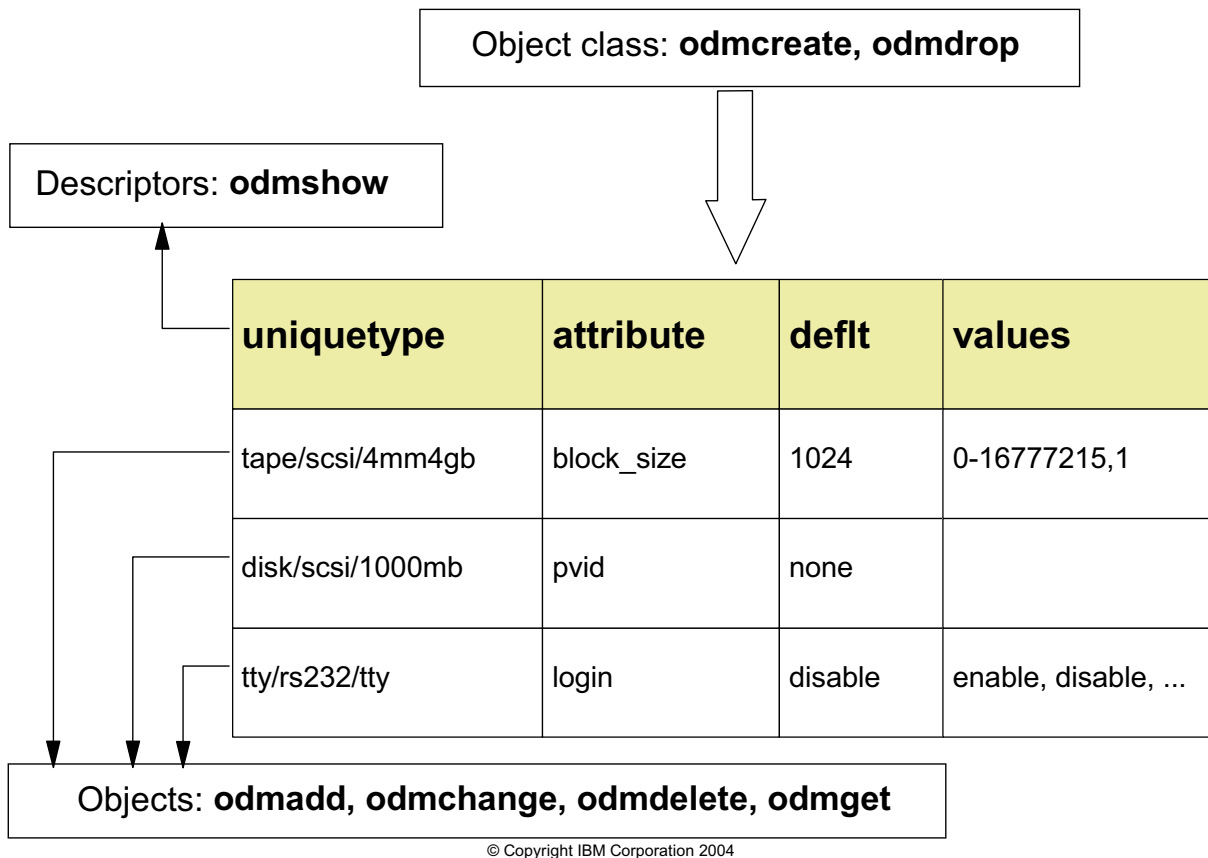


Figure 2-12. ODM Commands

AU1612.0

## Notes:

For each ODM component different commands are available:

1. You can create ODM classes using the **odmcreate** command. This command has the following syntax:

**odmcreate** *descriptor\_file.cre*

The file *descriptor\_file.cre* contains the class definition for the corresponding ODM class. Usually these files have the suffix **.cre**. Your exercise manual contains an optional part, that shows how to create self-defined ODM classes.

2. To delete an entire ODM class use the **odmdrop** command. This command has the following syntax:

**odmdrop -o** *object\_class\_name*

The name *object\_class\_name* is the name of the ODM class you want to remove. Be very careful with this command. It removes the complete class immediately.

3. To view the underlying layout of an object class use the **odmshow** command:

**odmshow** *object\_class\_name*

The picture shows an extraction from ODM class **PdAt**, where four descriptors are shown (uniquetype, attribute, deflt, and values).

4. Usually system administrators work with objects. The **odmget** command queries objects in classes (information just provided by the **odmshow** command). To add new objects use **odmadd**, to delete objects use **odmdelete** and to change objects use **odmchange**. Working on the object level is explained in more detail on the next pages.

All ODM commands use the **ODMDIR** environment variable, that is set in file **/etc/environment**. The default value of **ODMDIR** is **/etc/objrepos**.

## Changing Attribute Values

```
# odmget -q"uniquetype=tape/scsi/8mm and attribute=block_size" PdAt > file
# vi file
```

```
PdAt:
  uniquetype = "tape/scsi/8mm"
  attribute = "block_size"
  deflt = "1024"
  values = "0-245760,1"
  width = ""
  type = "R"
  generic = "DU"
  rep = "nr"
  nls_index = 6
```

← Modify deflt to 512

```
# odmdelete -o PdAt -q"uniquetype=tape/scsi/8mm and attribute=block_size"
# odmadd file
```

© Copyright IBM Corporation 2004

Figure 2-13. Changing Attribute Values

AU1612.0

### Notes:

The ODM objects are stored in a binary format; that means you need to work with the ODM commands to query or change any objects.

The **odmget** command in the example will pick all the records from the **PdAt** class, where **uniquetype** is equal to `tape/scsi/8mm` and **attribute** is equal to `block_size`. In this instance only one record should be matched. The information is redirected into a file which can be changed using an editor. In this example the default value for the attribute **block\_size** is changed to 512.

**Note:** Before the new value of 512 can be added into the ODM, the old object (which has the **block\_size** set to 1024) must be deleted, otherwise you would end up with two objects describing the same attribute in the database. The first object found will be used and can be quite confusing. This is why it is important to delete an entry before adding a replacement record.

The final operation is to add the file into the ODM.

As with any database you can perform queries for records matching certain criteria. The tests are on the values of the descriptors of the objects. A number of tests can be performed:

**Equality:** for example **uniquetype=tape/scsi/8mm** and **attribute=block\_size**

**Similarity:** for example **lpp\_name like bosext1.\***

Tests can be linked together using normal boolean operations. For example:

=	equal
!=	not equal
>	greater
>=	greater than or equal to
<	less than
<=	less than or equal to
<b>like</b>	similar to; finds path names in character string data

In addition to the \* wildcard, a ? can be used as a wildcard character.



## Changing Attribute Values Using odmchange

```
# odmget -q"uniquetype=tape/scsi/8mm and attribute=block_size" PdAt > file
```

```
# vi file
```

```
PdAt:
```

```
uniquetype = "tape/scsi/8mm"
```

```
attribute = "block_size"
```

```
deft = "1024" ←
```

```
values = "0-245760,1"
```

```
width = ""
```

```
type = "R"
```

```
generic = "DU"
```

```
rep = "nr"
```

```
nls_index = 6
```

Modify deflt to 512

```
# odmchange -o PdAt -q"uniquetype=tape/scsi/8mm and attribute=block_size" file
```

© Copyright IBM Corporation 2004

Figure 2-14. Changing Attribute Values Using odmchange

AU1612.0

### Notes:

The example shows how the **odmchange** command can be used instead of the **odmadd** and **odmdelete** steps (as in the previous example).



## 2.2 ODM Database Files

## Software Vital Product Data

<pre>lpp:   name = "bos.rte.printers"   state = 5   ver = 5   rel = 1   mod = 0   fix = 0   description = "Front End Printer Support"   lpp_id = 38</pre>	<pre>product:   lpp_name = "bos.rte.printers"   comp_id = "5765-C3403"   state = 5   ver = 5   rel = 1   mod = 0   fix = 0   ptf = ""   prereq = "*coreq bos.rte 5.1.0.0"   description = ""   supersedes = ""</pre>
<pre>inventory:   lpp_id = 38   file_type = 0   format = 1   loc0 = "/etc/qconfig"   loc1 = ""   loc2 = ""   size = 0   checksum = 0</pre>	<pre>history:   lpp_id = 38   ver = 5   rel = 1   mod = 0   fix = 0   ptf = ""   state = 1   time = 988820040   comment = ""</pre>

© Copyright IBM Corporation 2004

Figure 2-15. Software Vital Product Data

AU1612.0

### Notes:

Whenever installing a product or update in AIX, the **installp** command uses the ODM to maintain the software vital product database. The following information is part of this database:

- The name of the software product (for example, bos.rte.printers)
- The version, release and modification level of the software product (for example, 5.2.0)
- The fix level, which contains a summary of fixes implemented in a product
- Any PTFs (program temporary fix) that have been installed on the system
- The state of the software product:
  - Available (state = 1)
  - Applying (state = 2)
  - Applied (state = 3)
  - Committing (state = 4)
  - Committed (state = 5)
  - Rejecting (state = 6)
  - Broken (state = 7)

The Software Vital Product Data is stored in the following ODM classes:

- lpp** The lpp object class contains information about the installed software products, including the current software product state and description.
- inventory** The inventory object class contains information about the files associated with a software product.
- product** The product object class contains product information about the installation and updates of software products and their prerequisites.
- history** The history object class contains historical information about the installation and updates of software products.

Let's introduce the software states you should know about.

# Software States You Should Know About

<b>Applied</b>	<ul style="list-style-type: none"> <li>• Only possible for PTFs or Updates</li> <li>• Previous version stored in <code>/usr/lpp/Package_Name</code></li> <li>• Rejecting update recovers to saved version</li> <li>• Committing update deletes previous version</li> </ul>
<b>Committed</b>	<ul style="list-style-type: none"> <li>• Removing committed software is possible</li> <li>• No return to previous version</li> </ul>
<b>Applying, Committing, Rejecting, Deinstalling</b>	<p>If installation was not successful:</p> <ol style="list-style-type: none"> <li>installp -C</li> <li>smit maintain_software</li> </ol>
<b>Broken</b>	<ul style="list-style-type: none"> <li>• Cleanup failed</li> <li>• Remove software and reinstall</li> </ul>

© Copyright IBM Corporation 2004

Figure 2-16. Software States You Should Know About

AU1612.0

## Notes:

The AIX software vital product database uses software states that describe the status information of an install or update package:

1. When installing a PTF (program temporary fix) or update package, you can install the software into an **applied** state. Software in an applied state contains the newly installed version (which is active) and a backup of the old version (which is inactive). This gives you the opportunity to test the new software. If it works as expected, you can **commit** the software which will remove the old version. If it doesn't work as planned, you can **reject** the software which will remove the new software and reactivate the old version. Install packages cannot be **applied**. These will always be **committed**.
2. Once a product is committed, if you would like to return to the old version, you must remove the current version and reinstall the old version.
3. If an installation does not complete successfully, for example, if the power fails during the install, you may find software states like **applying**, **committing**, **rejecting**, or **deinstalling**. To recover from this failure, execute the command **installp -C** or use the

smit fastpath **smit maintain\_software**. Select *Clean Up After Failed or Interrupted Installation* when working in smit.

4. After a cleanup of a failed installation, you might detect a **broken** software status. In this case the only way to recover from this failure is to remove and reinstall the software package.

# Predefined Devices (PdDv)

```

PdDv:
  type = "8mm"
  class = "tape"
  subclass = "scsi"

  prefix = "rmt"
  ...
  base = 0
  ...
  detectable = 1
  ...
  led = 2418

  setno = 54
  msgno = 2
  catalog = "devices.cat"

  DvDr = "tape"

  Define = "/etc/methods/define"
  Configure = "/etc/methods/cfgsctape"
  Change = "/etc/methods/chggen"
  Unconfigure = "/etc/methods/ucfgdevice"
  Undefine = "/etc/methods/undefine"
  Start = ""
  Stop = ""
  ...
  uniquetype = "tape/scsi/8mm"

```

© Copyright IBM Corporation 2004

Figure 2-17. Predefined Devices (PdDv)

AU1612.0

## Notes:

The Predefined Devices (PdDv) object class contains entries for all devices supported by the system. A device that is not part of this ODM class could not be configured on an AIX system.

The attributes you should know about are:

<b>type</b>	Specifies the product name or model number (for example 8 mm (tape)).
<b>class</b>	Specifies the functional class name. A functional class is a group of device instances sharing the same high-level function. For example, tape is a functional class name representing all tape devices.
<b>subclass</b>	Device classes are grouped into subclasses. The subclass <b>scsi</b> specifies all tape devices that may be attached to an SCSI system.



---

<b>prefix</b>	Specifies the Assigned Prefix in the customized database, which is used to derive the device instance name and /dev name. For example, <b>rmt</b> is the prefix name assigned to tape devices. Names of tape devices would then look like rmt0, rmt1, or rmt2.
<b>base</b>	<p>This descriptor specifies whether a device is a base device or not. A base device is any device that forms part of a minimal base system. During system boot, a minimal base system is configured to permit access to the root volume group and hence to the root file system. This minimal base system can include, for example, the standard I/O diskette adapter and a SCSI hard drive. The device shown in the picture is not a base device.</p> <p>This flag is also used by the <b>bosboot</b> and <b>savebase</b> command, which are introduced in the next unit.</p>
<b>detectable</b>	Specifies whether the device instance is detectable or undetectable. A device whose presence and type can be determined by the <b>cfgmgr</b> once it is actually powered on and attached to the system, is said to be detectable. A value of 1 means that the device is detectable, and a value of 0 that it is not (for example, a printer or tty).
<b>led</b>	Indicates the value displayed on the LEDs when the configure method begins to run. The value stored is decimal, the value shown on the LEDs is hexadecimal (2418 is 972 in hex).
<b>setno, msgno</b>	Each device has a specific description (for example, 4.0 GB 8 mm Tape Drive) that is shown when the device attributes are listed by the <b>lsdev</b> command. These two descriptors are used to lookup the description in a message catalog.
<b>catalog</b>	Identifies the file name of the NLS (national language support) catalog. The <b>LANG</b> variable on a system controls which catalog file is used to show a message. For example, if LANG is set to en_US, the catalog file /usr/lib/nls/msg/en_US/devices.cat is used. If LANG is de_DE, catalog /usr/lib/nls/msg/de_DE/devices.cat is used.
<b>DvDr</b>	Identifies the name of the device driver associated with the device (for example, tape). Usually, device drivers are stored in directory <b>/usr/lib/drivers</b> . Device drivers are loaded into the AIX kernel when a device is made <b>available</b> .
<b>Define</b>	Names the define method associated with the device type. This program is called when a device is brought into the <b>defined</b> state.

<b>Configure</b>	Names the configure method associated with the device type. This program is called when a device is brought into the <b>available</b> state.
<b>Change</b>	Names the change method associated with the device type. This program is called when a device attribute is changed via the <b>chdev</b> command.
<b>Unconfigure</b>	Names the unconfigure method associated with the device type. This program is called when a device is unconfigured by <b>rmdev -l</b> .
<b>Undefine</b>	Names the undefine method associated with the device type. This program is called when a device is undefined by <b>rmdev -l -d</b> .
<b>Start, Stop</b>	Few devices support a stopped state (only logical devices). A stopped state means that the device driver is loaded, but no application can access the device. These two attributes name the methods to start or stop a device.
<b>uniquetype</b>	A key that is referenced by other object classes. Objects use this descriptor as pointer back to the device description in PdDv. The key is a concatenation of the class, subclass and type values.

## Predefined Attributes (PdAt)

```

PdAt:
  uniquetype = "tape/scsi/8mm"
  attribute = "block_size"
  deflt = "1024"
  values = "0-245760,1"
  ...

PdAt:
  uniquetype = "disk/scsi/1000mb"
  attribute = "pvid"
  deflt = "none"
  values = ""
  ...

PdAt:
  uniquetype = "tty/rs232/tty"
  attribute = "term"
  deflt = "dumb"
  values = ""
  ...

```

© Copyright IBM Corporation 2004

Figure 2-18. Predefined Attributes (PdAt)

AU1612.0

### Notes:

The Predefined Attribute object class contains an entry for each existing attribute for each device represented in the PdDv object class. An attribute is any device-dependent information, such as interrupt levels, bus I/O address ranges, baud rates, parity settings or block sizes. The extract out of PdAt shows three attributes (block size, physical volume identifier and terminal name) and their default values.

The meanings of the key fields shown on the visual are as follows:

- uniquetype** This descriptor is used as a pointer back to the device defined in the PdDv object class.
- attribute** Identifies the name of the attribute. This is the name that can be passed to the **mkdev** or **chdev** commands. For example to change the default name of **dumb** to **ibm3151** for a terminal name, you can issue:

```
# chdev -l tty0 -a term=ibm3151
```

<b>deflt</b>	Identifies the default value for an attribute. Nondefault values are stored in <b>CuAt</b> .
<b>values</b>	Identifies the possible values that can be associated with the attribute name. For example, allowed values for the <code>block_size</code> attribute range from 0 to 245760, with an increment of 1.

## Customized Devices (CuDv)

```

CuDv:
  name = "rmt0"
  status = 1
  chgstatus = 2
  ddins = "tape"
  location = "04-C0-00-1,0"
  parent = "scsi0"
  connwhere = "1,0"
  PdDvLn = "tape/scsi/8mm"

CuDv:
  name = "tty0"
  status = 1
  chgstatus = 1
  ddins = ""
  location = "01-C0-00-00"
  parent = "sa0"
  connwhere = "S1"
  PdDvLn = "tty/rs232/tty"

```

© Copyright IBM Corporation 2004

Figure 2-19. Customized Devices (CuDv)

AU1612.0

### Notes:

The Customized Devices (CuDv) object class contains entries for all device instances defined in the system. As the name implies, a defined device object is an object that a define method has created in the CuDv object class. A defined device object may or may not have a corresponding actual device attached to the system.

CuDv object class contains objects that provide device and connection information for each device. Each device is distinguished by a unique logical name. The customized database is updated twice, during system bootup and at run time, to define new devices, remove undefined devices and update the information for a device that has changed.

The key descriptors in CuDv are:

**name**            A customized device object for a device instance is assigned a unique logical name to distinguish the device from other devices. The visual shows two devices, a tape device **rmt0** and a tty, **tty0**.

<b>status</b>	Identifies the current status of the device instance. Possible values are: <ul style="list-style-type: none"><li>• status = 0: Defined</li><li>• status = 1: Available</li><li>• status = 2: Stopped</li></ul>
<b>chgstatus</b>	This flag tells whether the device instance has been altered since the last system boot. The diagnostics facility uses this flag to validate system configuration. The flag can take these values: <ul style="list-style-type: none"><li>• chgstatus = 0: New device</li><li>• chgstatus = 1: Don't care</li><li>• chgstatus = 2: Same</li><li>• chgstatus = 3: Device is missing</li></ul>
<b>ddins</b>	This descriptor typically contains the same value as the Device Driver Name descriptor in the Predefined Devices (PdDv) object class. It specifies the name of the device driver that is loaded into the AIX kernel.
<b>location</b>	Identifies the physical location of a device. The location code is a path from the system unit through the adapter to the device. In case of a hardware problem, the location code is used by technical support to identify a failing device. In many RS/6000 systems the location codes are labeled in the hardware, to facilitate the finding of devices.
<b>parent</b>	Identifies the logical name of the parent device. For example, the parent device of <b>rmt0</b> is <b>scsi0</b> .
<b>connwhere</b>	Identifies the specific location on the parent device where the device is connected. For example, the device <b>rmt0</b> uses the SCSI address <b>1,0</b> .
<b>PdDvLn</b>	Provides a link to the device instance's predefined information through the unique type descriptor in the PdDv object class.

## Customized Attributes (CuAt)

```
CuAt:
  name = "tty0"
  attribute = "login"
  value = "enable"
  ...
CuAt:
  name = "hdisk0"
  attribute = "pvid"
  value = "0016203392072a540000000000000000"
  ...
```

© Copyright IBM Corporation 2004

Figure 2-20. Customized Attributes (CuAt)

AU1612.0

### Notes:

The Customized Attribute object class contains customized device-specific attribute information.

Devices represented in the Customized Devices (CuDv) object class have attributes found in the Predefined Attribute (PdAt) object class and the CuAt object class. There is an entry in the CuAt object class for attributes that take **customized** values. Attributes taking the default value are found in the PdAt object class. Each entry describes the current value of the attribute.

These objects out of the CuAt object class show two attributes that take customized values. The attribute **login** has been changed to **enable**. The attribute **pvid** shows the physical volume identifier that has been assigned to disk hdisk0.

# Additional Device Object Classes

<p><b>PdCn:</b>            uniquetype = "adapter/pci/sym875"            connkey = "scsi"            connwhere = "1,0"</p> <p><b>PdCn:</b>            uniquetype = "adapter/pci/sym875"            connkey = "scsi"            connwhere = "2,0"</p>	<p><b>CuDvDr:</b>            resource = "devno"            value1 = "22"            value2 = "0"            value3 = "rmt0"</p> <p><b>CuDvDr:</b>            resource = "devno"            value1 = "22"            value2 = "1"            value3 = "rmt0.1"</p>
<p><b>CuDep:</b>            name = "rootvg"            dependency = "hd6"</p> <p><b>CuDep:</b>            name = "datavg"            dependency = "lv01"</p>	<p><b>CuVPD:</b>            name = "rmt0"            vpd = "*MFEXABYTE            PN21F8842"</p>

© Copyright IBM Corporation 2004

Figure 2-21. Additional Device Object Classes

AU1612.0

## Notes:

- PdCn** The Predefined Connection (PdCn) object class contains connection information for adapters (or sometimes called intermediate devices). This object class also includes predefined dependency information. For each connection location, there are one or more objects describing the subclasses of devices that can be connected.
- The example objects show that at the given locations all devices belonging to subclass SCSI could be attached.
- CuDep** The Customized Dependency (CuDep) object class describes device instances that depend on other device instances. This object class describes the dependence links between logical devices and physical devices as well as dependence links between logical devices, exclusively. Physical dependencies of one device on another device are recorded in the Customized Device (CuDep) object class.
- The example object show the dependencies between logical volumes and the volume groups they belong to.



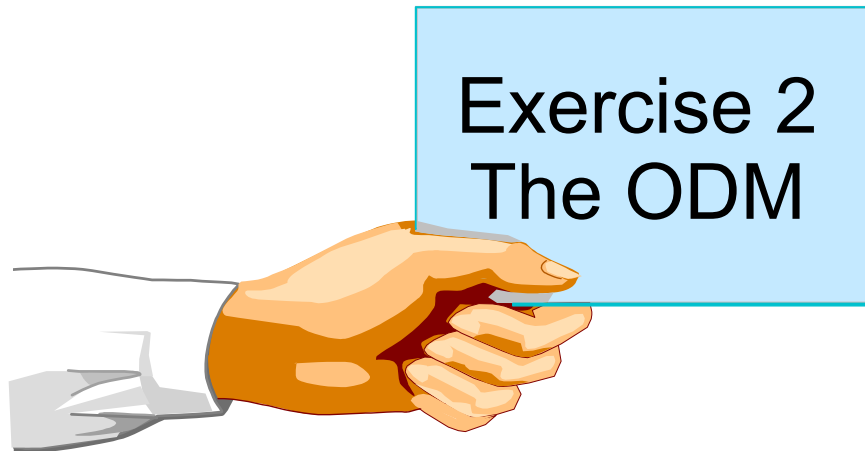
**CuDvDr** The Customized Device Driver (CuDvDr) object class is used to create the entries in the **/dev** directory. These special files are used from applications to access a device driver that is part of the AIX kernel. The attribute **value1** is called the **major number** and is a unique key for a device driver. The attribute **value2** specifies a certain operating mode of a device driver.

The example objects reflect the device driver for tape rmt0. The major number 22 specifies the driver in the kernel, the minor numbers 0 and 1 specify two different operating modes. The operating mode **0** specifies a *rewind on close* for the tape drive, the operating mode **1** specifies *no rewind on close* for a tape drive.

**CuVPD** The Customized Vital Product Data (CuVPD) object class contains vital product data (manufacturer of device, engineering level, part number, and so forth) that is useful for technical support. When an error occurs with a specific device the vital product data is shown in the error log.

## Next Step

---



© Copyright IBM Corporation 2004

Figure 2-22. Next Step

AU1612.0

### **Notes:**

At the end of the exercise you should be able to:

- Define the meaning of the most important ODM files
- Work with the ODM command line interface
- Describe how ODM classes are used from device configuration commands

An optional part provides how to create self-defined ODM classes, which is very interesting for AIX system programmers.

---

## Checkpoint

---

1. In which ODM class do you find the physical volume IDs of your disks?

---

2. What is the difference between state **defined** and **available**?

---

---

---

---

---

© Copyright IBM Corporation 2004

Figure 2-23. Checkpoint

AU1612.0

### **Notes:**

## Unit Summary

---

- The ODM is made from object classes, which are broken into individual objects and descriptors.
- AIX offers a command line interface to work with the ODM files.
- The device information is held in the customized and the predefined databases (Cu\*, Pd\*).

© Copyright IBM Corporation 2004

Figure 2-24. Unit Summary

AU1612.0

### **Notes:**

## Unit 3. System Initialization Part I

### What This Unit Is About

This unit describes the boot process to loading the boot logical volume. It provides the content of the boot logical volume and how it can be re-created if it's corrupted.

The meaning of the LED codes is described and how they can be analyzed to fix boot problems.

### What You Should Be Able to Do

After completing this unit, you should be able to:

- Describe the boot process to loading the boot logical volume
- Describe the contents of the boot logical volume
- Interpret LED codes displayed during system boot and at system halt
- Re-create the boot logical volume on a system which is failing to boot
- Describe the features of a service processor

### How You Will Check Your Progress

Accountability:

- Activity
- Checkpoint questions
- Lab exercise

### References

Online	System Management Concepts: Operating System and Devices
Online	System Management Guide: Operating System and Devices
Online	<a href="http://publib16.boulder.ibm.com/pseries/en-US/infocenter/base/aix52.htm">http://publib16.boulder.ibm.com/pseries/en-US/infocenter/base/aix52.htm</a>
SA38-0541	<i>RS/6000 7025 F50 Series Service Guide</i>
SA38-0547	<i>RS/6000 7026 Model H50 Service Guide</i>
SA38-0512	<i>RS/6000 7043 43P Series Service Guide</i>
SA38-0554	<i>RS/6000 7043 Model 260 Service Guide</i>
SA38-0548	<i>Enterprise Servers S70 and S7A Service Guide</i>

## Unit Objectives

---

After completing this unit, students should be able to:

- Describe the **boot process** to loading the **boot logical volume**
- Describe the **contents** of the **boot logical volume**
- Interpret **LED codes** displayed during boot and at **system halt**
- **Re-create the boot logical volume** on a system which is failing to boot

© Copyright IBM Corporation 2004

Figure 3-1. Unit Objectives

AU1612.0

### **Notes:**

Boot problems are the most frequent errors that occur. Hardware and software problems might cause a system to stop during the boot process.

This unit describes the boot process of loading the boot logical volume and provides the knowledge a system administrator needs to have to analyze the boot problem.

## 3.1 System Startup Process

# How Does An AIX System Boot?

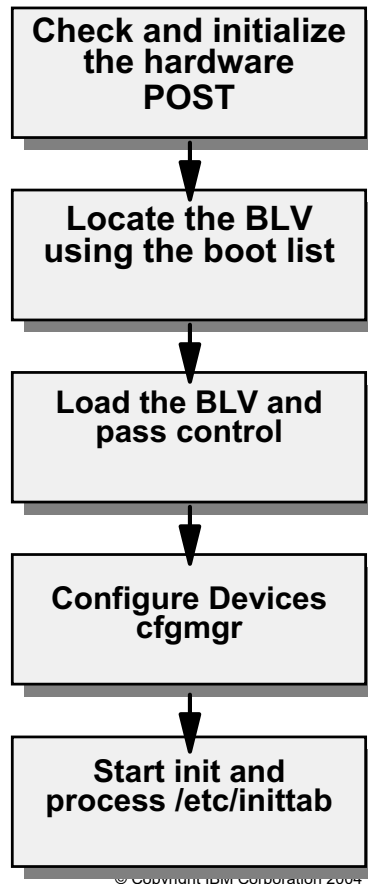


Figure 3-2. How Does An AIX System Boot?

AU1612.0

## Notes:

This is the basic overview of the boot process.

After powering on a machine the hardware is checked and initialized. This phase is called the POST (Power-On Self Test). The goal of the POST is to verify the functionality of the hardware.

After the POST is complete, a boot logical volume (BLV or boot image) is located from the boot list and is loaded into memory. During a normal boot, the location of the BLV is usually a hard drive. Besides hard drives, the BLV could be loaded from tape or CD-ROM. This is the case when booting into maintenance or service mode. If working with NIM (network install manager), the BLV is loaded via the network.

To use an alternate boot location you must invoke the appropriate boot list by depressing function keys during the boot process. There is more information on boot lists, later in the unit.

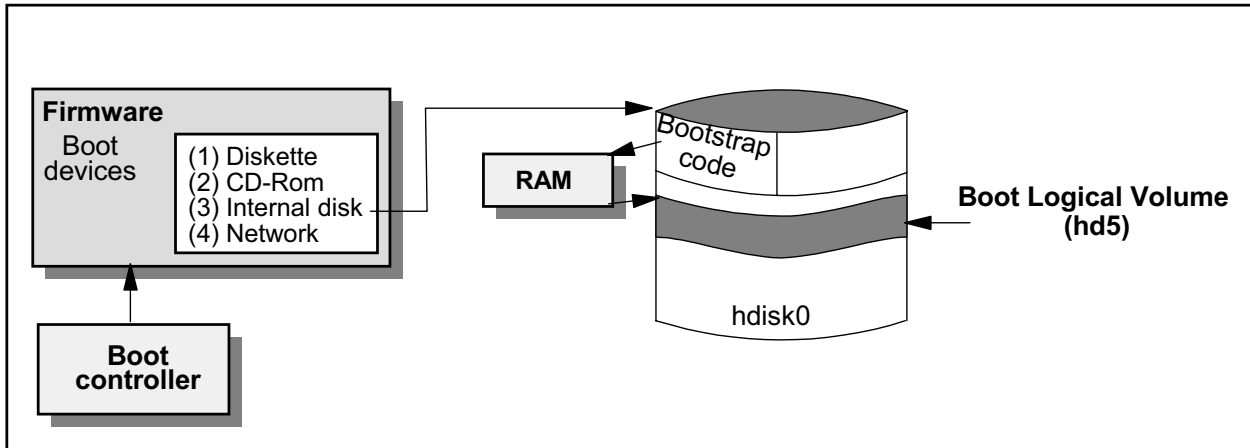


Passing control to the boot logical volume means that one component of the boot logical volume, the AIX kernel, gets control over the boot process. The components of the BLV are discussed later in the unit.

All devices are configured during the boot processes. This is done in different phases by the `cfgmgr`.

At the end, the `init` process is started and processes the `/etc/inittab` file.

# Loading of a Boot Image



© Copyright IBM Corporation 2004

Figure 3-3. Loading of a Boot Image

AU1612.0

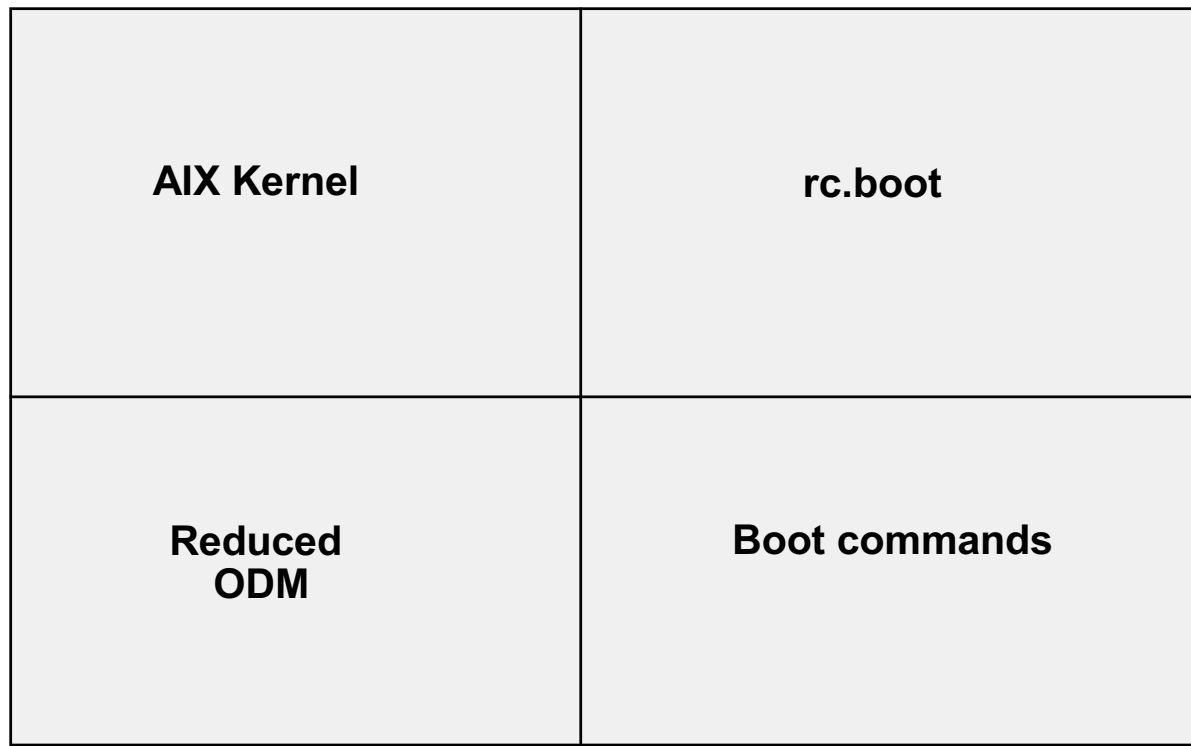
## Notes:

This picture shows how the boot logical volume is found during the AIX boot process. Machines use one or more boot lists to identify a boot device. The boot list is part of the firmware.

RS/6000s can manage several different operating systems. The hardware is not bound to the software. The first 512 bytes contain a bootstrap code that is loaded into RAM during the boot process. This part is sometimes referred to as System ROS (Read Only Storage). The bootstrap code gets control. The task of this code is to start up the operating system - in some technical manuals this second part is called the Software ROS. In the case of AIX, the boot image is loaded.

To save disk space, the boot logical volume is compressed on the disk (therefore it's called a boot image). During the boot process the boot logical volume is uncompressed and the AIX kernel gets boot control.

## Content of Boot Logical Volume (hd5)



© Copyright IBM Corporation 2004

Figure 3-4. Content of Boot Logical Volume (hd5)

AU1612.0

### Notes:

This picture shows the components of the boot logical volume.

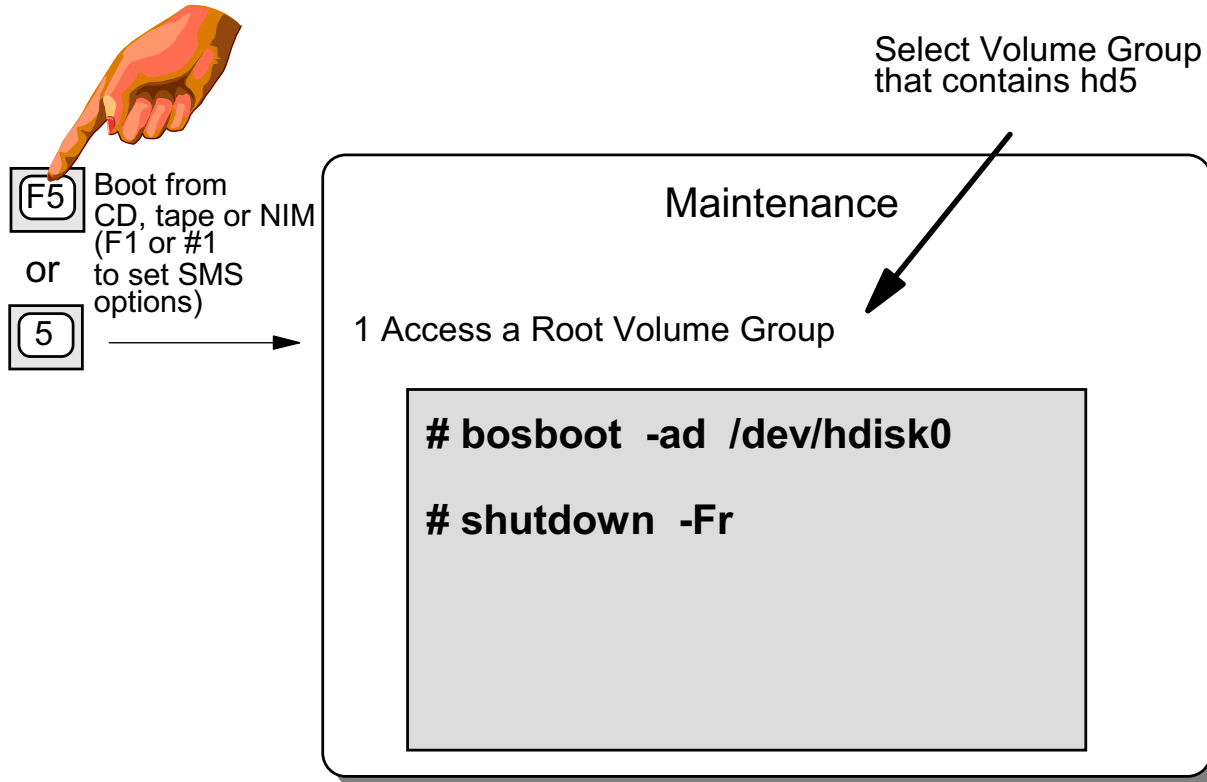
The AIX kernel is the core of the operating system and provides basic services like process, memory and device management. The AIX kernel is always loaded from the boot logical volume. There is a copy of the AIX kernel in the **hd4** file system (under the name **/unix**), but this program has no role in system initialization. Never remove **/unix**, because it's used for rebuilding the kernel in the boot logical volume.

The boot commands are programs that are called during the boot process. Examples are **bootinfo**, **cfgmgr** and more.

The boot logical volume contains a reduced copy of the ODM. During the boot process many devices are configured before **hd4** is available. For these devices the corresponding ODM files must be stored in the boot logical volume.

After starting the kernel, the boot script **rc.boot** gets control over the boot process. This is explained in the System Initialization Part II Unit.

# How to Fix a Corrupted BLV



© Copyright IBM Corporation 2004

Figure 3-5. How to Fix a Corrupted BLV

AU1612.0

## Notes:

If a boot logical volume is corrupted (for example, bad blocks on a disk might cause a corrupted BLV), a machine will not boot.

To fix this situation, you must boot your machine in **maintenance mode**, from a CD or tape. If NIM has been set up for a machine, you can also boot the machine from a NIM master in maintenance mode. By the way, that's what you would do on an SP node if an SP node does not boot.

The boot lists are set using the **bootlist** command or the System Management Services (SMS) program. Some machines support a normal and service boot list. If your model supports this, you will use a function key during bootup to select the appropriate list. Normally, pressing F5 when you hear the first tones during bootup, will force the machine to use the firmware default bootlist which lists media devices first. So it will check for a bootable CD or Tape before looking for a disk to boot. More on this later.

Be careful to use the correct AIX installation CD to boot your machine. You can't boot an AIX 5.2 installed machine with an AIX 5200-01 installation CD as well as

AIX 5.1 installed machine with an AIX 5100-03 installation CD (you must match the Version, Release and maintenance level).

After booting from CD, tape or NIM an **Installation and Maintenance Menu** is shown and you can startup the maintenance mode. We will cover this later in this unit. After accessing the rootvg, you can repair the boot logical volume with the **bosboot** command. You need to specify the corresponding disk device, for example **hdisk0**:

### **bosboot -ad /dev/hdisk0**

It is important that you do a proper shutdown. All changes need to be written from memory to disk.

The **bosboot** command requires that the boot logical volume **hd5** exists. If you ever need to re-create the BLV from scratch - maybe it had been deleted by mistake or the LVCB of hd5 has been damaged - the following steps should be followed:

1. Boot your machine in maintenance mode (from CD or tape (F5 or 5) or use (F1 or 1) to access the Systems Management Services (SMS) to select boot device).
2. Remove the old hd5 logical volume.  
**# rmlv hd5**
3. Clear the boot record at the beginning of the disk.  
**# chpv -c hdisk0**
4. Create a new hd5 logical volume: one physical partition in size, must be in rootvg and outer edge as intrapolicy. Specify boot as logical volume type.  
**# mklv -y hd5 -t boot -a e rootvg 1**
5. Run the bosboot command as described on the foil.  
**# bosboot -ad /dev/hdisk0**
6. Check the actual bootlist.  
**# bootlist -m normal -o**
7. Write data immediately to disk.  
**# sync**  
**# sync**
8. Shutdown and reboot the system.  
**# shutdown -Fr**

By using the internal command **ipl\_varyon -i** you can check the state of the boot record.

# Working with Boot Lists

## Normal Mode

```
# bootlist -m normal hdisk0 hdisk1
# bootlist -m normal -o
hdisk0
hdisk1
```

## Service Mode

```
# bootlist -m service -o
fd0
cd0
hdisk0
tok0
```

```
# diag
```

```
TASK SELECTION LIST
SCSI Bus Analyzer
Download Microcode
Display or Change Bootlist
Periodic Diagnostics
```



© Copyright IBM Corporation 2004

Figure 3-6. Working with Boot Lists (PCI)

AU1612.0

## Notes:

You can use the command **bootlist** or **diag** from the command line to change or display the boot lists. You can also use the **System Management Services (SMS)** programs. **SMS** is covered on the next page.

### 1. **bootlist** command

The **bootlist** command is the easiest way to change the boot list. The first example shows how to change the boot list for a normal boot. In this example, we boot either from hdisk0 or hdisk1. To query the boot list, you can use the option **-o** which was introduced in AIX 4.2.

The next example shows how a service boot list can be set.

### 2. **diag** command

The **diag** command is part of the package **bos.rte.diag** which allows diagnostic tasks. One part of these diagnostic tasks allows for displaying and changing boot lists. Working with the **diag** command is covered later in the course.

The custom boot list is the normal boot list set via the **bootlist** command, the **diag** command or the **SMS programs**. The normal boot list is used during a normal boot. The default boot list is called when F5 or F6 is pressed during the boot sequence.

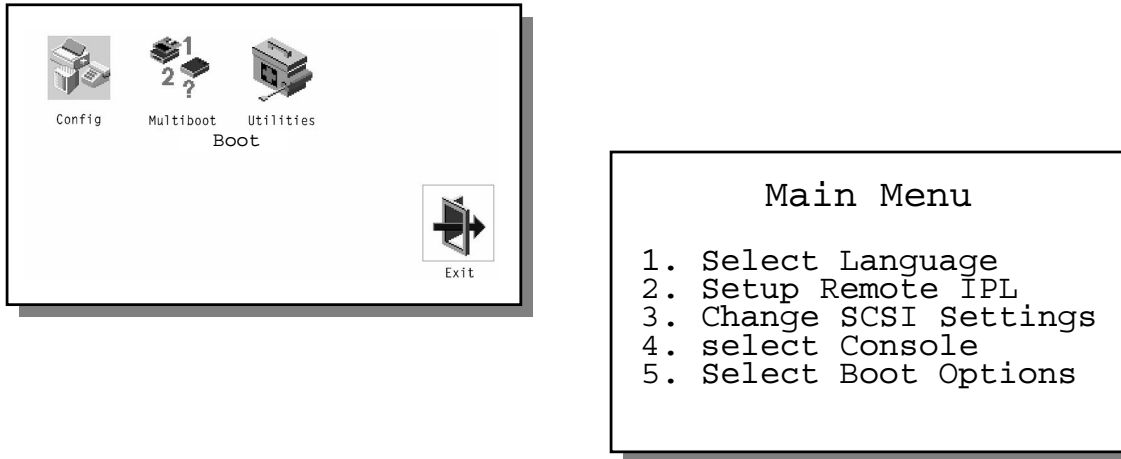
Other machines, in addition to the default boot list and the custom boot list, allow for a customized service boot list. This is set using mode service with the **bootlist** command. The default boot list is called when F5 is pressed during boot. The service boot list is called when F6 is pressed during boot.

You may find variations on the different models of RS/6000s. Refer to the *User's Guide* for your specific model ([www.rs6000.ibm.com/resource/hardware\\_docs/#index6](http://www.rs6000.ibm.com/resource/hardware_docs/#index6)).

# Working with Boot Lists - SMS

1. Reboot or power on the system.
2. Press **F1** or #1 when tone is heard.
3. Select **Boot** Options.

## System Management Services



© Copyright IBM Corporation 2004

Figure 3-7. Working with Boot Lists - SMS

AU1612.0

### Notes:

You can also change the boot list with the **System Management Services**. The SMS programs are integrated into the hardware (they reside on ROM).

The picture shows how to start the **System Management Services** in graphic mode seen on older systems as well as the ascii menus seen on newer systems. After power-on you need to press **F1** to start up the graphic version of the **System Management Services**. You must press this key when the tone is heard and before the fifth of five icons appear.

If your model does not have a graphic adapter, you need to set up an ASCII terminal on the S1 port. In this case a text version of the **System Management Services** will be started on your terminal.

Newer systems (with graphical or ascii console) use the number 1 key and this should be depressed before the fifth text icon appears (it shows the text for the icon; that is, memory, speaker, and so forth)

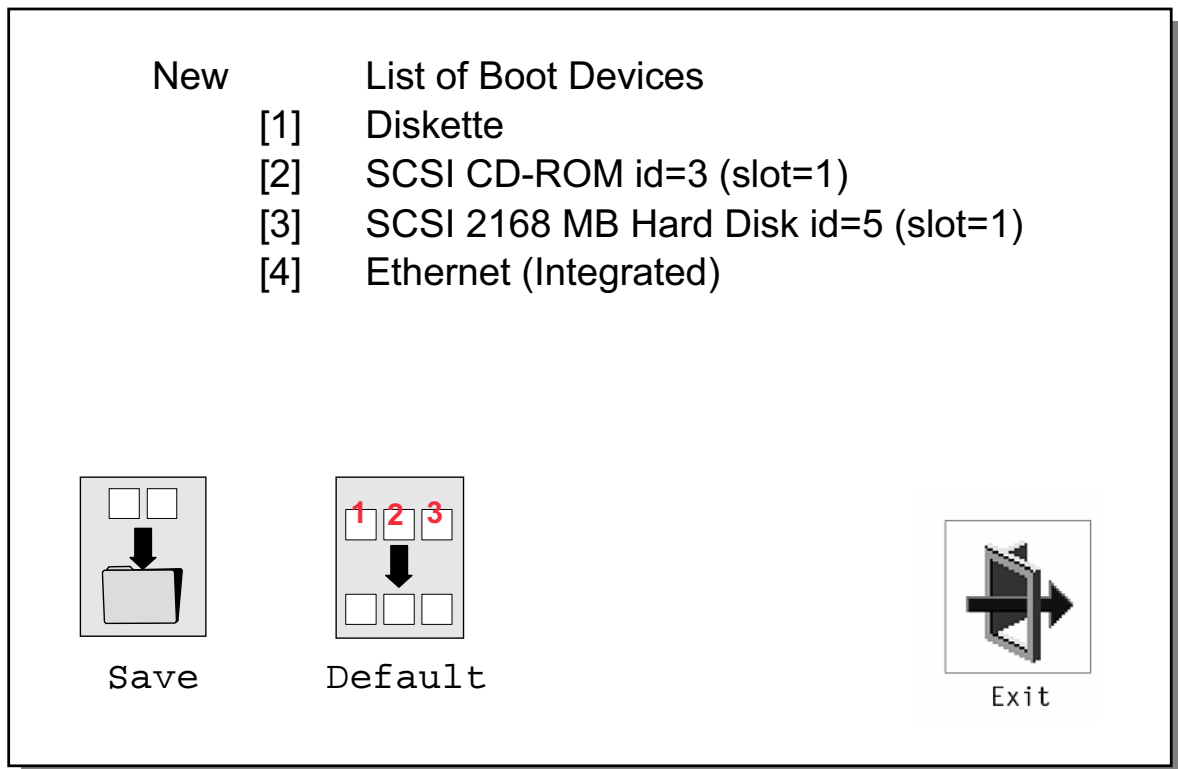


In the **System Management Service** menu, select Boot or Multiboot or Select Boot Options (model-dependent) to work with the boot list. The look of the menu differs on the various models and firmware levels.

All new RS/6000 models use the following key allocation standard:

1. **F1 or 1 on ASCII terminal:** Start System Management Services.
2. **F5 or 5 on ASCII terminal:** Boot diagnostics from disk, use default boot list.
3. **F6 or 6 on ASCII terminal:** Boot diagnostics from disk, use custom service boot list.

# System Management Services



© Copyright IBM Corporation 2004

Figure 3-8. System Management Services

AU1612.0

## Notes:

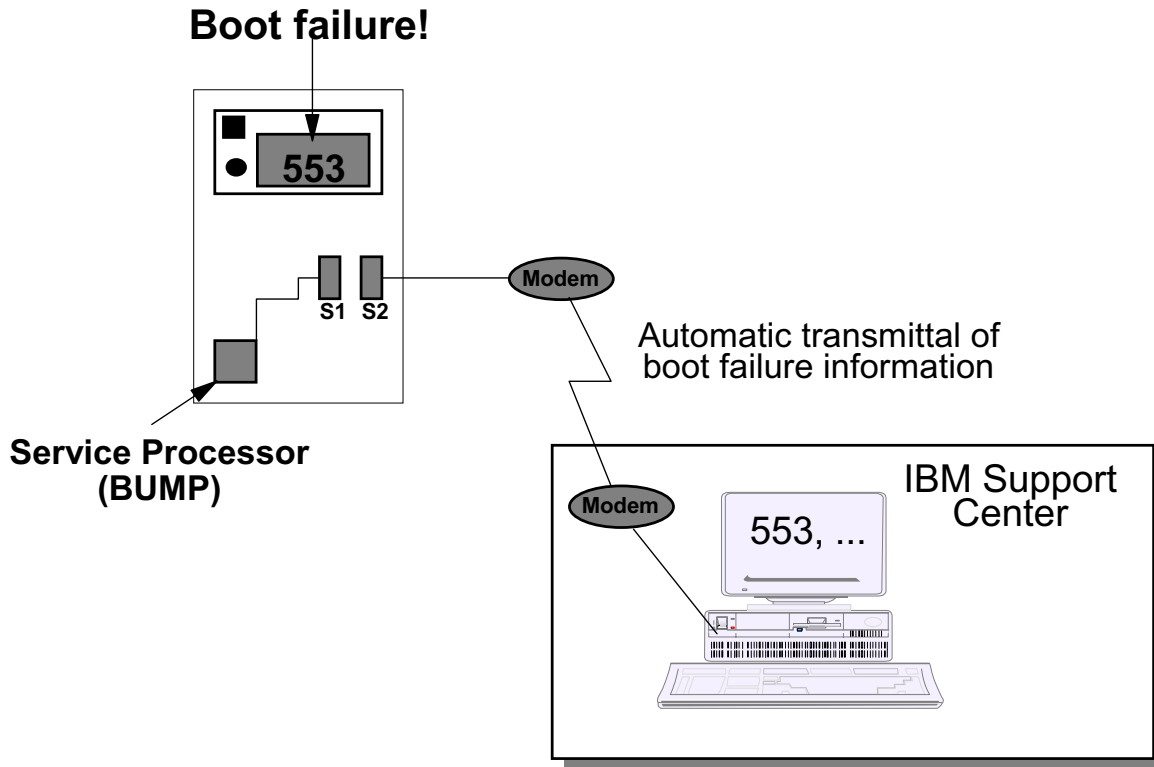
RS/6000s support up to **five boot devices**. Some models only support four. A **default boot list** is stored with the following sequence:

1. Diskette drive
2. CD-ROM
3. Internal disk
4. Communication adapter (like Ethernet or token-ring)

To set a new boot sequence, type the sequence number in the **new** column. Be sure to **save** your changes before exiting.

Only SCSI disks containing a boot record are shown.

# Service Processors and Boot Failures



© Copyright IBM Corporation 2004

Figure 3-9. Service Processors and Boot Failures

AU1612.0

## Notes:

IBM's family of SMP servers includes a service processor. This processor allows actions to occur even when the regular processors are down.

The SMP servers can be set up to automatically call an IBM support center (or any other site) in case of a boot failure. An automatic transmittal of boot failure information takes place. This information includes LED codes and service request numbers, that describe the cause of the boot failure.

If the data is sent to an IBM Service Center, the information is extracted and placed in a problem record. IBM Service personnel will call the customer to find out if assistance is requested.

A valid service contract is a prerequisite for this dial-out feature of the service processor.

Other features of the service processor are:

- Console mirroring to make actions performed by a remote technician visible and controllable by the customer.
- Remote as well as local control of the system (power-on/off, diagnostics, reconfiguration, maintenance).
- Run-time hardware and operating system surveillance. If, for example, a CPU fails, the service processor would detect this, reboot itself automatically and run without the failed CPU.
- Timed power-on and power-off, reboot on crash, reboot on power loss.

# Let's Review

---



© Copyright IBM Corporation 2004

Figure 3-10. Let's Review

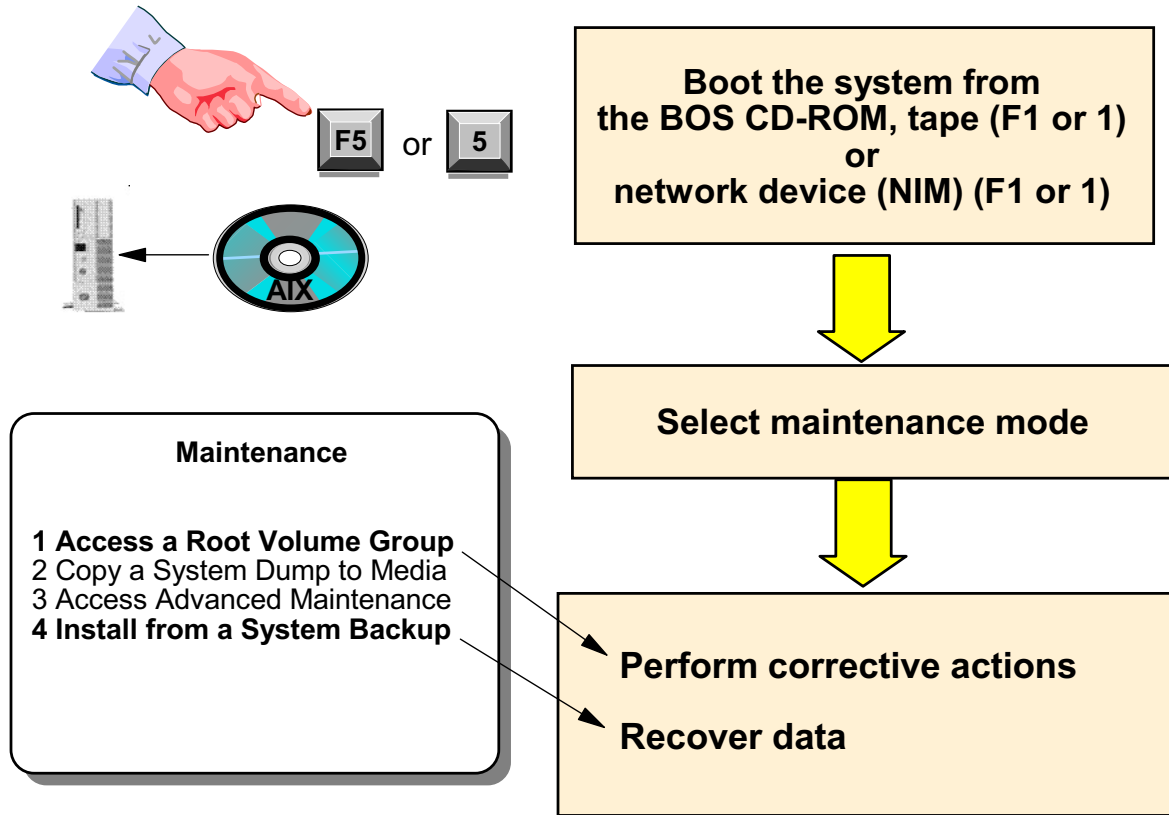
AU1612.0

## ***Let's Review***

1. T/F: You must have AIX loaded on your RS/6000 to use the System Management Services Programs.
2. Your RS/6000 is currently powered off. AIX is installed on hdisk1 but the boot list is set to boot from hdisk0. How can you fix the problem and make the machine boot from hdisk1?
3. Your machine is booted and you are sitting at the # prompt. What is the command that will display the boot list? How could you change the boot list?
4. What command is used to fix the boot logical volume?
5. What script controls the boot sequence?

## 3.2 Solving Boot Problems

# Accessing a System That Will Not Boot



© Copyright IBM Corporation 2004

Figure 3-11. Accessing a System That Will Not Boot

AU1612.0

## Notes:

Before discussing LED/LCD codes that are shown during the boot process we want to identify how a system can be accessed that will not boot. The maintenance mode can be started from an AIX CD, an AIX bootable tape (like an mksysb) or a network device, that has been prepared on a NIM master. The devices that contain the boot media must be stored in the boot lists.

To boot into maintenance modes:

- Newer PCI systems support the **bootlist** command and booting from a **mksysb** tape, but the tape device is by default not part of the boot sequence.
- Verify your boot list, but do not forget that some machines do not have a service boot list. Check that your boot device is part of the boot list:

```
# bootlist -m normal -o
```

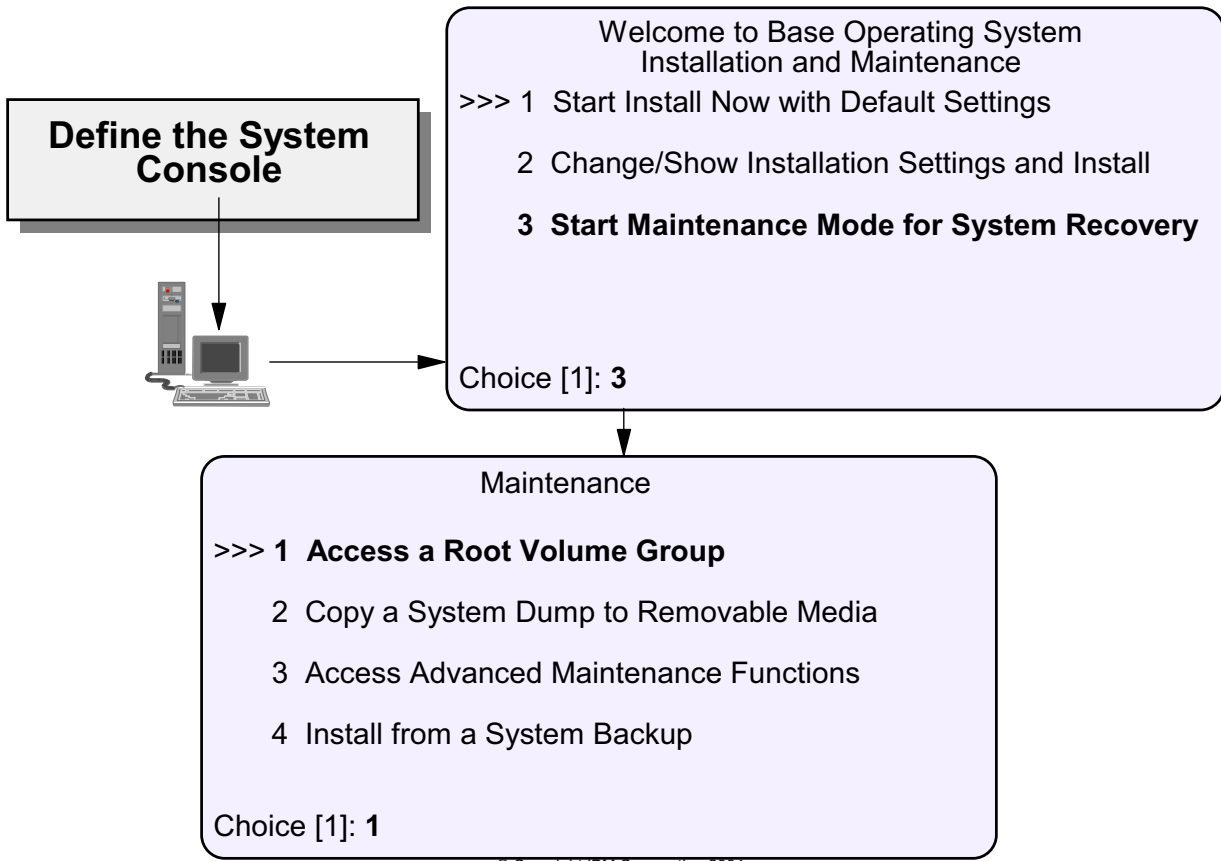
- If you want to boot from your internal tape device you need to change the boot list because the tape device by default is not part of the boot list. For example:

```
# bootlist -m normal cd0 rmt0 hdisk0
```



- Insert the boot media (either tape or CD) into the drive.
- Power on the system. The system begins booting from the installation media. After several minutes, **c31** is displayed in the LED/LCD panel. After a few minutes you will see the **Installation and Maintenance** menu.

# Booting in Maintenance Mode



© Copyright IBM Corporation 2004

Figure 3-12. Booting in Maintenance Mode

AU1612.0

## Notes:

When booting in maintenance mode you first have to identify the system console that will be used, for example your **lft** terminal or a tty that is attached to the S1 port.

After selecting the console the **Installation and Maintenance** menu is shown.

As we want to work in maintenance mode, we use selection **3** to start up the **Maintenance** menu.

From this point we access our **rootvg** to execute any system recovery steps that may be necessary.

## Working in Maintenance Mode

### Access a Root Volume Group

- 1) Volume Group 001620336e1bc8a3 contains these disks:  
hdisk0 2063 04-C0-00-4,0
- 2) Volume Group 001620333C9b1b8e contains these disks:  
hdisk1 2063 04-C0-00-5,0

Choice: 1

### Volume Group Information

Volume Group ID 001620336e1bc8a3 includes the following logical volumes:

hd6    hd5    hd8    hd4    hd2    hd9var    hd3

- 1) **Access this Volume Group and start a shell**
- 2) **Access this Volume Group and start a shell before mounting file systems**
- 99) Previous Menu

Choice [99]:

© Copyright IBM Corporation 2004

Figure 3-13. Working in Maintenance Mode

AU1612.0

### Notes:

When accessing the rootvg in maintenance mode, you need to select the volume group that is the rootvg. In the example two volume groups exist on the system. Note that only the volume group IDs are shown and not the names of the volume groups. Check with your system documentation that you select the correct disk. Do not rely too much on the physical volume name but more on the PVID, VGID or SCSI ID.

After selecting the volume group it will show the list of LVs contained in the VG. This is how you confirm you have selected rootvg. Two selections are then offered:

#### 1. Access this Volume Group and start a shell

When you choose this selection the **rootvg** will be activated (varyonvg command), and all file systems belonging to the **rootvg** will be mounted. A shell will be offered to you to execute any system recovery steps.

Typical scenarios where this selection must be chosen are:

- Changing a **forgotten root password**
- Re-creating the **boot logical volume**
- Changing a **corrupted boot list**

## 2. Access this Volume Group and start a shell before mounting file systems

When you choose this selection the **rootvg** will be activated, but the file system belonging to the **rootvg** will **not be mounted**.

A typical scenario where this selection is chosen is when a corrupted file system needs to be repaired by the **fsck** command. Repairing a corrupted file system is only possible if the file system is not mounted.

Another scenario might be a corrupted **hd8** transaction log. Any changes that take place in the superblock or i-nodes are stored in the log logical volume. When these changes are written to disk, the corresponding transaction logs are removed from the log logical volume.

A corrupted transaction log must be reinitialized by the **logform** command, which is only possible, when no file system is mounted. After initializing the log device, you need to do a file system repair for all file systems that use this transaction log. Beginning with AIX V5.1 you have explicitly to specify the filesystem type jfs or jfs2:

```
# logform -V jfs /dev/hd8
# fsck -y -V jfs /dev/hd1
# fsck -y -V jfs /dev/hd2
# fsck -y -V jfs /dev/hd3
# fsck -y -V jfs /dev/hd4
# fsck -y -V jfs /dev/hd9var
# fsck -y -V jfs /dev/hd10opt
# exit
```

Keep in mind that US keyboard layout is used but you can use the retrieve function by using `set -o emacs` or `set -o vi`.

## Boot Problem References

AIX Message Guide and Reference	Contains: ▶ AIX boot codes
AIX Problem Solving Guide and Reference	Contains: ▶ Problem Solving Procedures ▶ Problem Summary Form
RS/6000 Service Guides	Contains: ▶ PCI firmware checkpoints ▶ PCI error codes

© Copyright IBM Corporation 2004

Figure 3-14. Boot Problem References

AU1612.0

### Notes:

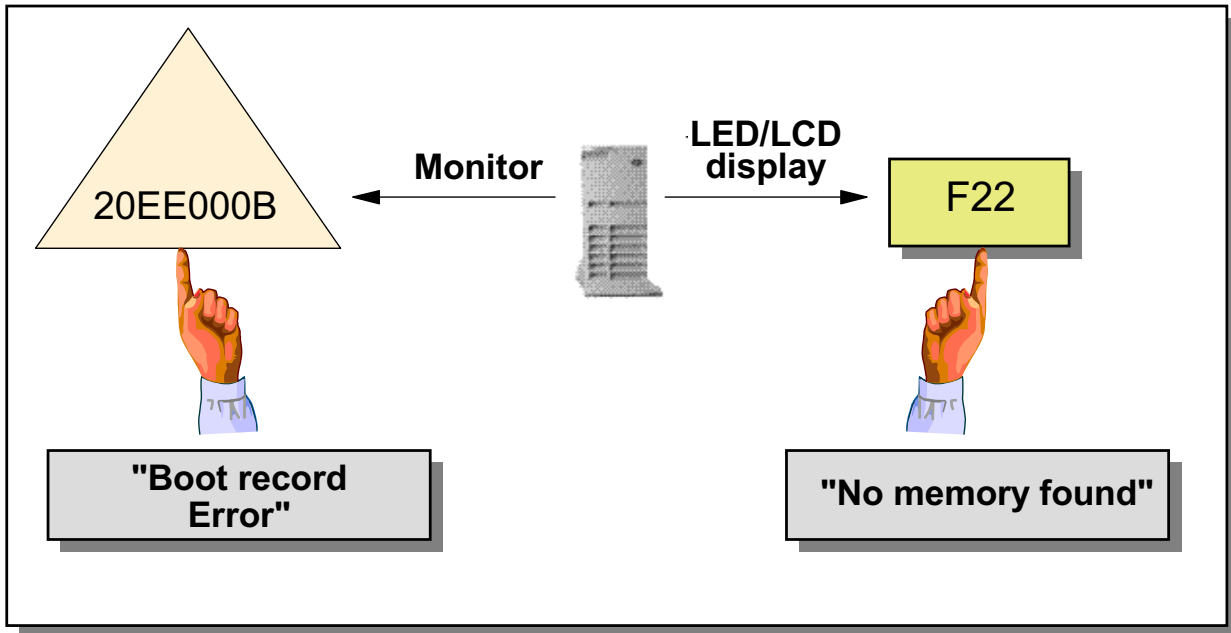
Whenever your machine does not boot and you are not sure what is causing the boot problem, look up the LED code in the *AIX Messages Guide and Reference*. It recommends actions that you should follow to fix the problem.

Many other problem solving procedures are described in the *AIX Problem Solving Guide and Reference*. These are manuals which an AIX administrator needs to resolve problems.

PCI **firmware checkpoints** and **error codes** are not explained in the *AIX Messages Guide and Reference*. Since they are hardware related, you need to look them up in your *RS/6000 Service Guide* that belongs to your PCI system.

All RS/6000 service guides are online at:  
[www.rs6000.ibm.com/resource/hardware\\_docs](http://www.rs6000.ibm.com/resource/hardware_docs).

# Firmware Checkpoints and Error Codes



- Explained in *RS/6000 Service Guide*
- Online available on [www-1.ibm.com/servers/eserver/pseries/library/hardware\\_docs](http://www-1.ibm.com/servers/eserver/pseries/library/hardware_docs)

© Copyright IBM Corporation 2004

Figure 3-15. Firmware Checkpoints and Error Codes

AU1612.0

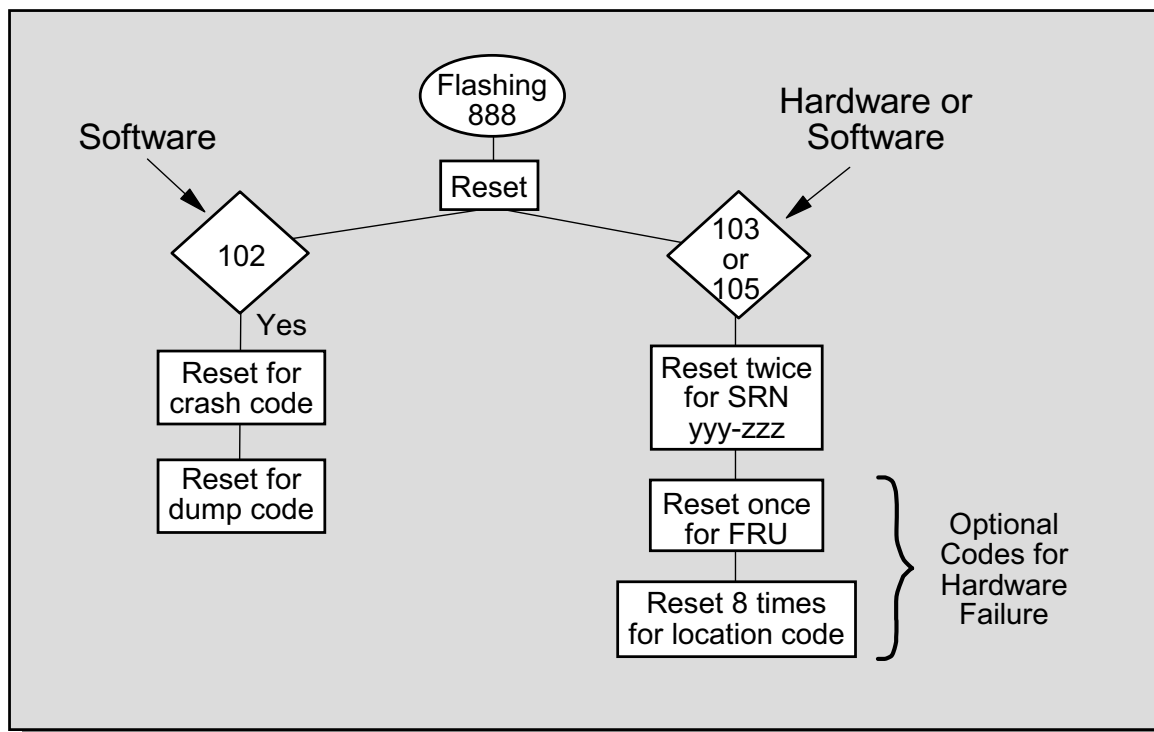
## Notes:

RS/6000s use the LED/LCD display to show the current boot status. These boot codes are called **firmware checkpoints**.

If errors are detected by the firmware during the boot process, an error code is shown on the monitor. For example, the error code 20EE000B indicates that a boot record error has occurred.

Firmware checkpoints and error codes are different on various models and they are not listed in the *AIX Messages Guide and Reference*. They are provided in the *RS/6000 Service Guides* of your model. The service guides are available online at: [http://www-1.ibm.com/servers/eserver/pseries/libraryhardware\\_docs](http://www-1.ibm.com/servers/eserver/pseries/libraryhardware_docs)

## Flashing 888



© Copyright IBM Corporation 2004

Figure 3-16. Flashing 888

AU1612.0

### Notes:

Another type of error you may encounter is a flashing 888.

A flashing 888 indicates that there is more information to be extracted from the system by pressing the reset button.

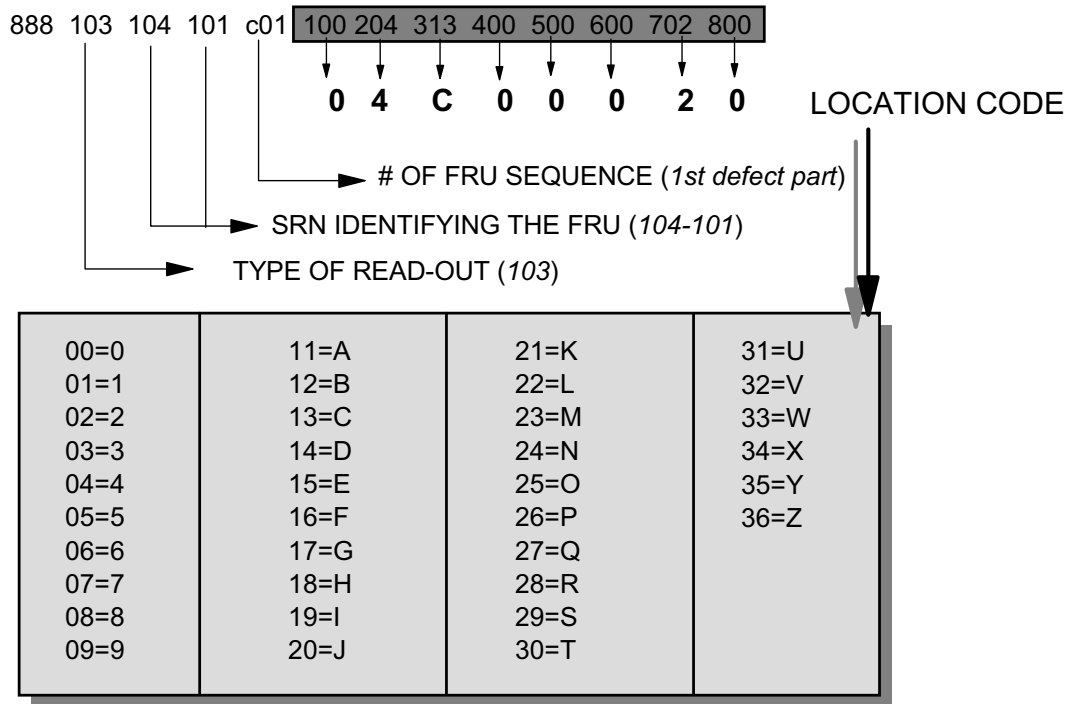
A **102** indicates that a dump has occurred - your AIX kernel crashed due to bad circumstances. By pressing the reset button the dump code can be obtained. We will cover more on dump in Unit 10 - The AIX Dump Facility.

A **103** may be hardware or software related. More frequent are hardware errors, but a corrupted boot logical volume may also lead to a flashing **888-103**.

If you press the reset button twice you get a **Service Request Number**, that may be used by IBM support to analyze the problem.

In case of a hardware failure, you get the sequence number of the **FRU** (Field Replaceable Unit) and a **location code**. The location code identifies the **physical location** of a device.

# Understanding the 103 Message



**FRU** = Field Replaceable Unit

**SRN** = Service Request Number

© Copyright IBM Corporation 2004

Figure 3-17. Understanding the 103 Message

AU1612.0

## Notes:

This picture shows an example 888 sequence.

- 103 determines that the error may be hardware or software related.
- 104-101 provides the **Service Request Number** for technical support. This number together with other system related data is used to analyze the problem.
- c01 identifies the first defect part. More than one part could be described in a 888 sequence.
- The next eight identifiers describe the **location code** of the defect part. These identifiers must be mapped with the shown table to identify the location code. In this example the location code is **04-C0-00-2,0**, which means that the SCSI device with address 2,0 on the built-in SCSI controller causes the flashing 888.



## Location Codes: Model 150

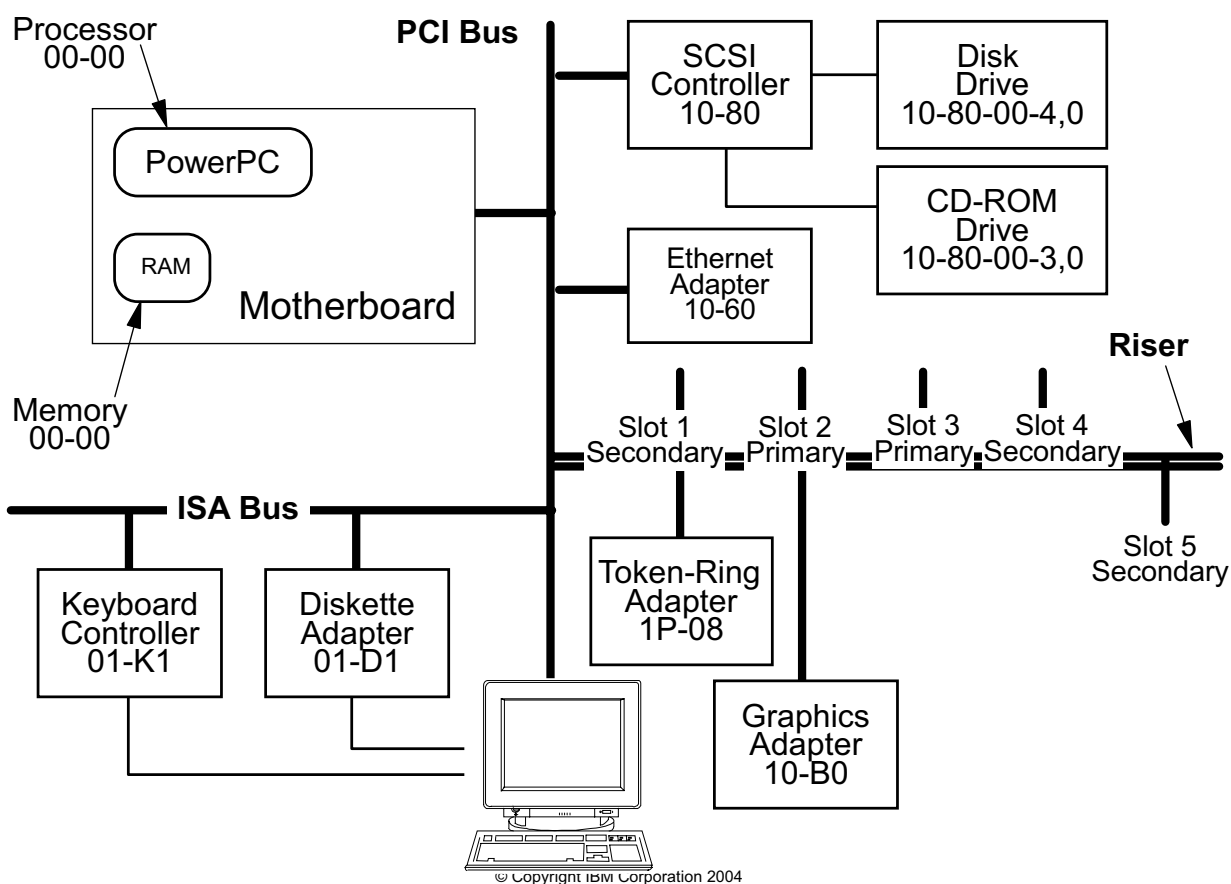


Figure 3-18. Location Codes: Model 150

AU1612.0

### Notes:

The location codes vary among PCI systems. The 43P Model 150 has a different addressing scheme than the 44P Model 270, for example. The same concept is still here - providing information about where the device is attached. The information on this page pertains only to the Model 150.

The processor bus still contains the processor and memory (addresses start with 00). The integrated ISA devices still start with 01, but the follow-on codes differ from the Model 140. You can see examples in the picture. For instance, the keyboard adapter is 01-K1 and the diskette adapter is 01-D1. On the Model 150, the integrated PCI device addresses start with 10. You can see the SCSI controller has an address of 10-80 and the Ethernet adapter has an address of 10-60.

Attached to the PCI bus is a riser card that has slots for the pluggable PCI cards. There are five slots on this card. Slots 1, 4, and 5 are on a secondary bus (addresses start with 1P), while, slots 2 and 3 are on the primary bus (as we have already seen start with 10). Here are the valid address ranges for those slots:

**1P-08 to 1P-0f** Slot 1

**10-b0 to 10-b7** Slot 2

**10-90 to 10-97** Slot 3

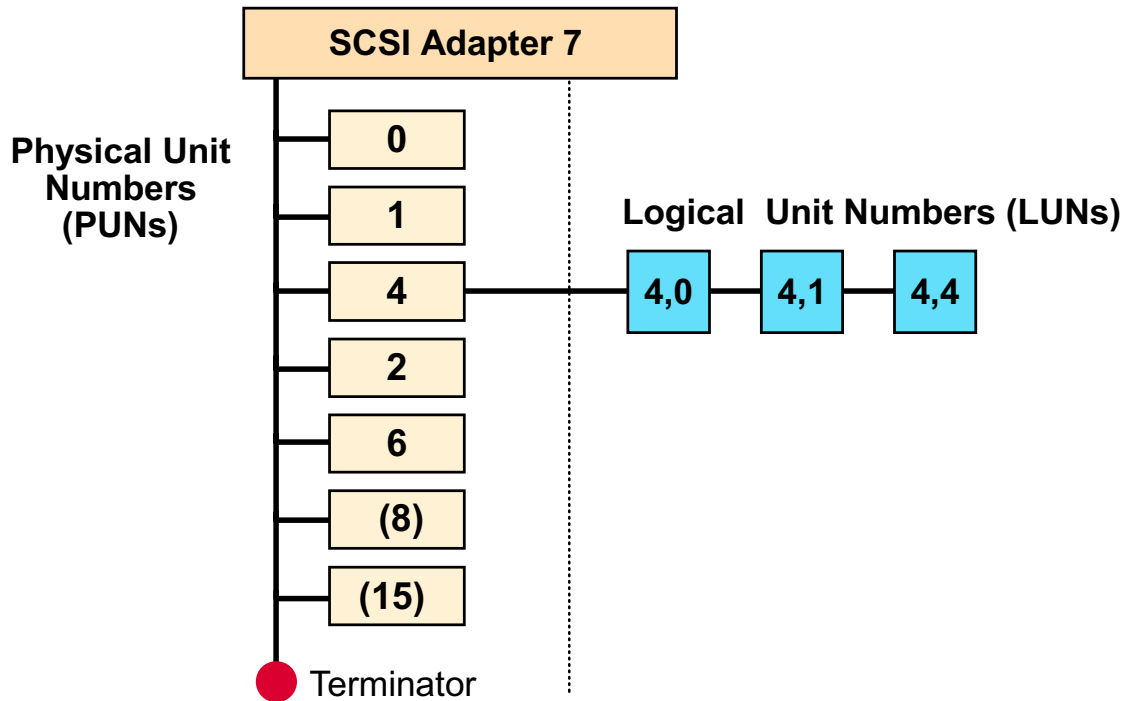
**1P-18 to 1P-1f** Slot 4

**1P-10 to 1P-17** Slot 5

In our example, the token-ring card is plugged into slot 1 (part of the secondary bus) and is assigned the address 1P-08. The graphics card is in Slot 2 (on the primary PCI bus) and is assigned the address 10-b0. The system will ensure there is a unique pair of numbers for each device.

For specifics on your type of machine, you should refer to the *RS/6000 User's Guide* for your model.

# SCSI Addressing



- Both ends, internal and external, of SCSI bus must be terminated

© Copyright IBM Corporation 2004

Figure 3-19. SCSI Addressing

AU1612.0

## Notes:

SCSI devices must use a unique SCSI address that has to be set on the SCSI device. It is very important that each device on an SCSI bus have a unique SCSI ID. To find out which addresses are already used, use the **lsdev** command:

```
# lsdev -Cs scsi -H
name      status      location      description
hdisk0    Available   04-C0-00-4,0  16 Bit SCSI Disk Drive
hdisk1    Available   04-C0-00-5,0  16 Bit SCSI Disk Drive
hdisk2    Available   04-C0-00-11,0 16 Bit SCSI Disk Drive
rmt0      Available   04-C0-00-2,0  2.3GB 8mm Tape Drive
          |
          SCSI address
```

The SCSI address consists of a physical unit number and a logical unit number. The physical unit number identifies a SCSI device, for example hdisk0 or rmt0. Some SCSI devices, for example, CD changers where more than one CD could be inserted, use logical

unit numbers. In this case, the logical unit number reflects the first, second, and so forth, CD in the drive.

Today most internal SCSI devices are self-terminating. However, a terminator resistor pack has to be attached to the device at the end of the daisy-chain externally. The SCSI adapter broadcasts to all devices attached to the SCSI system; each device reads the broadcast to determine whether the data is for them and, if so, reads the data. The data, however, continues down the SCSI bus. If no terminator is present the data will bounce back up the SCSI bus, and the receiving device will read the data again.

On an AIX system the lack of a terminator will in most cases not cause a problem. However, when it does, it is usually a serious problem, such as a system crash, or a system that does not boot.

Typically, SCSI controllers support up to seven devices, with SCSI addresses 0 through 6. If the SCSI controller supports **wide SCSI**, it supports up to 15 devices per SCSI bus, with addresses ranging from 0 through 15, excluding 7.

Never use the address 7 as SCSI address. This address is used by the adapter itself.

# Problem Summary Form

<b>Background Information</b>	
1. Record the Current Date and Time	_____
2. Record the System Date and Time (if available)	_____
3. Record the Symptom	_____
4. Record the Service Request Number (SRN)	_____
5. Record the Three-Digit Display Codes (if available)	__-__-__
6. Record the Location Codes:	
• First FRU	__-__-__
• Second FRU	__-__-__
• Third FRU	__-__-__
• Fourth FRU	__-__-__
<b>Problem Description</b>	
<b>Data Captured</b>	
(Describe data captured, such as system dumps, core dumps, error IDs error logs, or messages that needs to be examined by your service organization)	
(After completing this form, copy it and keep it on hand for future problem solving reference.)	

© Copyright IBM Corporation 2004

Figure 3-20. Problem Summary Form

AU1612.0

## Notes:

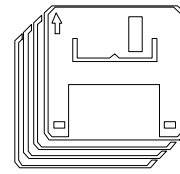
For every problem that comes up on your AIX system, not only boot problems, fill out the **Problem Summary Form**.

This information is used by IBM Support to analyze your problem.

# Getting Firmware Updates from Internet

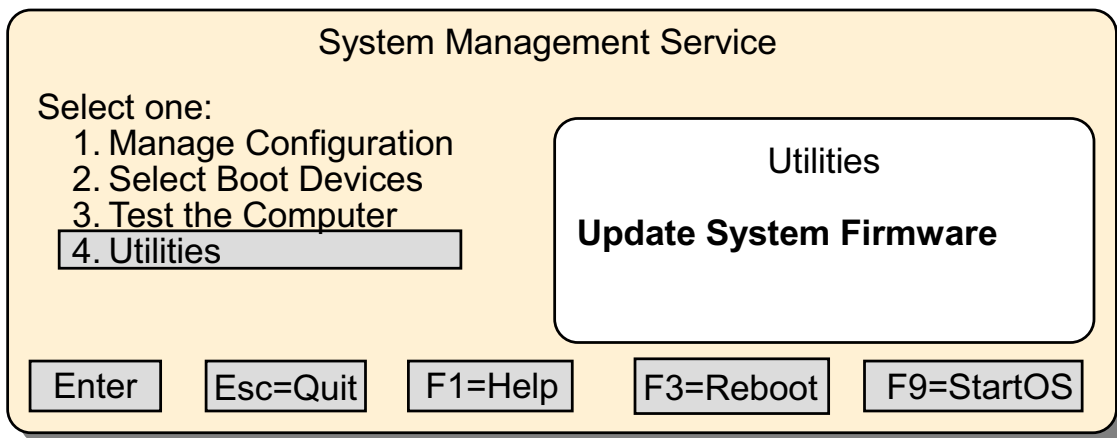
## 1. Get firmware update from IBM

<http://www.rs6000.ibm.com/support/micro>



Firmware-  
Update-  
Diskette

## 2. Update firmware via System Management Services



© Copyright IBM Corporation 2004

Figure 3-21. Getting Firmware Updates from Internet

AU1612.0

### Notes:

If you ever need a firmware update for your PCI model, for example, you want to install new hardware that requires a higher firmware level, download a **firmware update diskette** from the Internet. Use URL **<http://www.rs6000.ibm.com/support/micro>** to download the firmware update. After downloading the package follow the instructions in the **README** that comes with the package to create the diskette.

To install the new firmware level, start the **System Management Services** and select **Utilities**. From there, select **Update System Firmware**.

This will install a new firmware level on your PCI model.

This shows the ASCII interface of the SMS programs. If you are using the graphical interface, you would select Utilities followed by Update.

## Next Step

---



© Copyright IBM Corporation 2004

Figure 3-22. Next Step

AU1612.0

### **Notes:**

At the end of the exercise, you should be able to:

- Boot a machine in maintenance mode
- Repair a corrupted boot logical volume
- Alter boot lists on different RS/6000 hardware models

# Checkpoint

---

1. During the AIX boot process, the AIX kernel is loaded from the root file system. True or False?

---

2. Which RS/6000 models do not have a bootlist for the service mode?

---

3. How do you boot an AIX machine in maintenance mode?

---

---

4. Your machine keeps rebooting and repeating the POST. What can be the reason for this?

---

---

© Copyright IBM Corporation 2004

Figure 3-23. Checkpoint

AU1612.0

## Notes:



---

## Unit Summary

---

- During the boot process a **boot logical volume is loaded into memory**.
- Boot devices and sequences can be updated via the **bootlist**-command and the **diag**-command.
- The boot logical volume contains an **AIX kernel**, an **ODM** and a boot script **rc.boot** that controls the AIX boot process.
- The boot logical volume can be re-created using the **bosboot** command.
- LED codes produced during the boot process can be used to **diagnose boot problems**. PCIs additionally use **visual boot signals**.

© Copyright IBM Corporation 2004

Figure 3-24. Unit Summary

AU1612.0

### **Notes:**



# Unit 4. System Initialization Part II

## What This Unit Is About

This unit describes the final stages of the boot process and outlines how devices are configured for the system.

Common boot errors are described and how they can be analyzed to fix boot problems.

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Identify the steps in system initialization from loading the boot image to boot completion
- Identify how devices are configured during the boot process
- Analyze and solve boot problems

## How You Will Check Your Progress

Accountability:

- Checkpoint questions
- Lab exercise

# Unit Objectives

---

After completing this unit, students should be able to:

- Identify the steps in system initialization from **loading the boot image** to **boot completion**
- Identify **how devices** are **configured** during the **boot process**
- **Analyze and solve boot problems**

© Copyright IBM Corporation 2004

Figure 4-1. Unit Objectives

AU1612.0

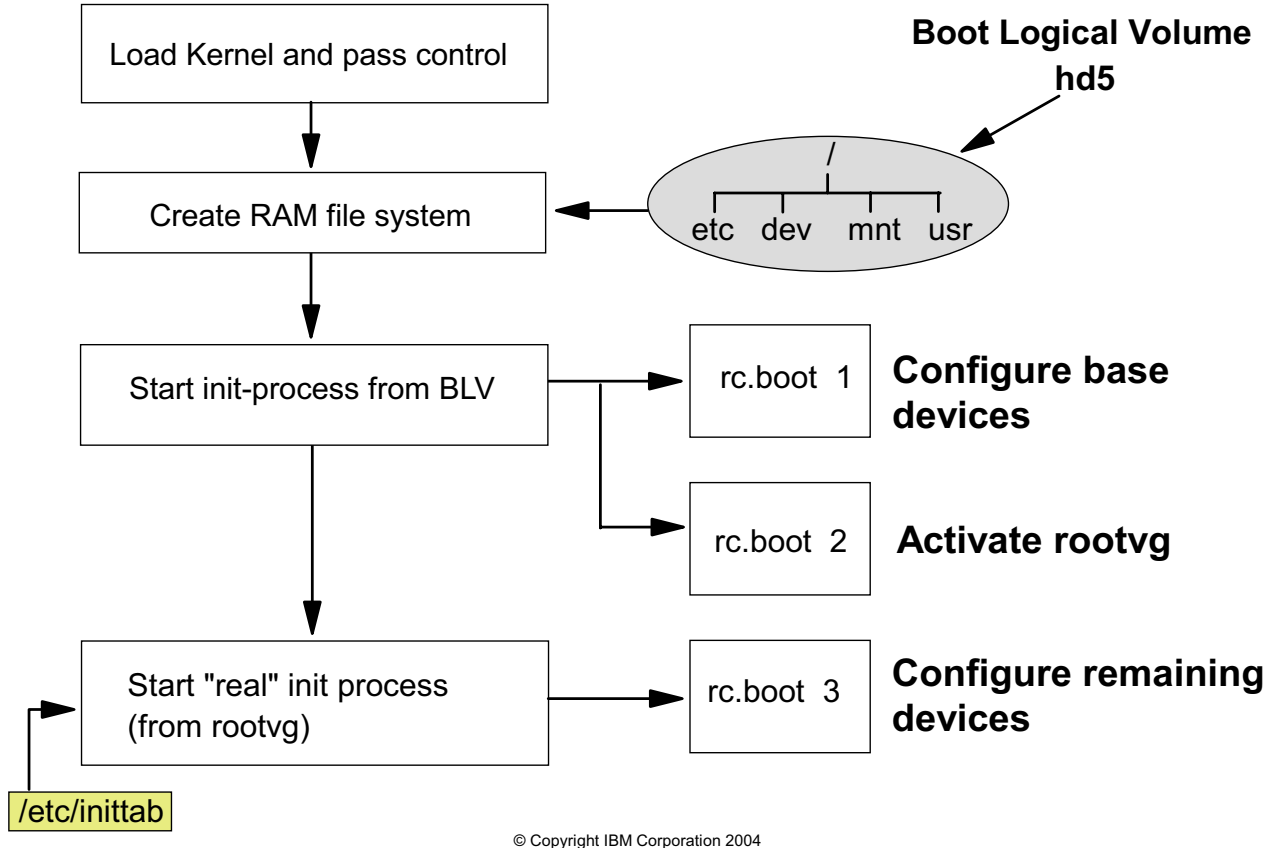
## **Notes:**

There are many reasons for boot failures. The hardware might be damaged or, due to user errors, the operating system might not be able to complete the boot process.

A good knowledge of the AIX boot process is a prerequisite for all AIX system administrators.

## 4.1 AIX Initialization Part 1

# System Software Initialization - Overview



© Copyright IBM Corporation 2004

Figure 4-2. System Software Initialization - Overview

AU1612.0

## Notes:

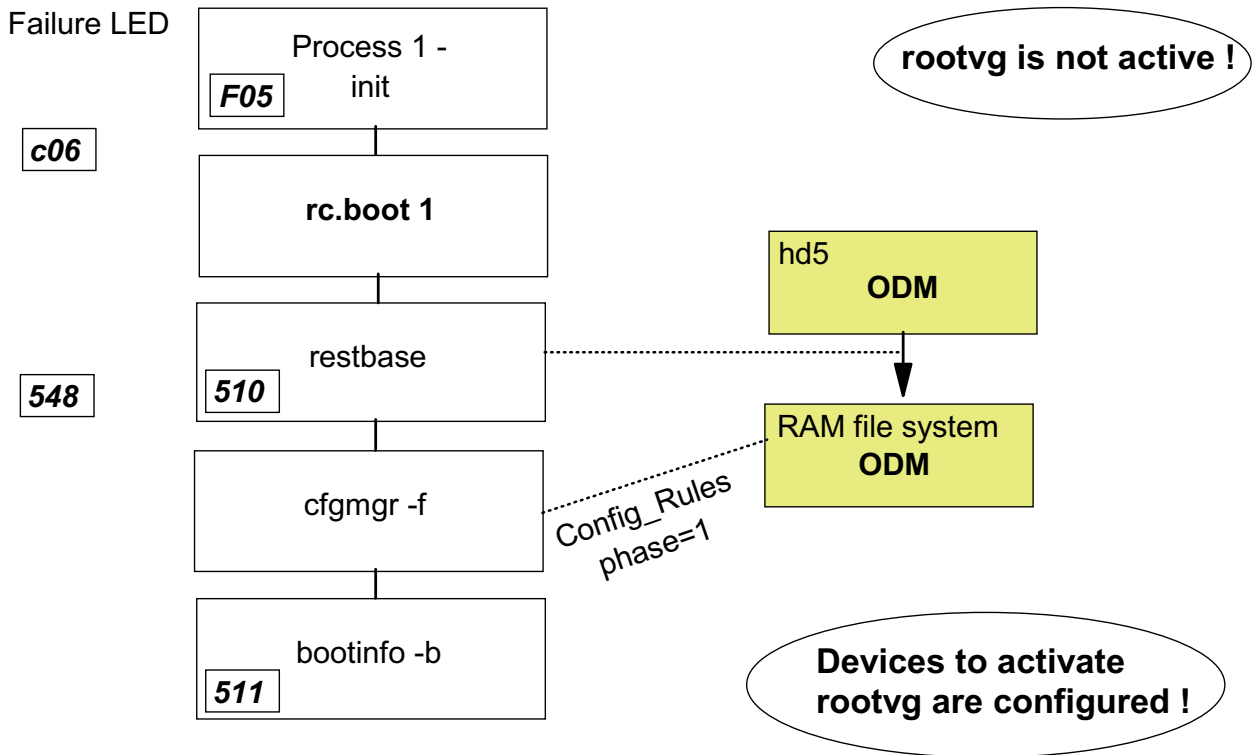
This page provides the boot sequence after loading the AIX kernel from the boot logical volume.

The AIX kernel gets control and executes the following steps:

- The kernel creates a RAM file system by using the components from the boot logical volume. At this stage the rootvg is not available, so the kernel needs to work with the boot logical volume. You can consider this RAM file system as a small AIX operating system.
- The kernel starts the **init** process which was loaded out of the boot logical volume (not from the root file system). This **init** process executes a boot script **rc.boot**.
- **rc.boot** controls the boot process. In the first phase (it is called by **init** with **rc.boot 1**), the base devices are configured. In the second phase (**rc.boot 2**), the rootvg is activated (or varied on).

- After activating the rootvg at the end of rc.boot 2, the kernel destroys the RAM file system and accesses the rootvg file systems from disks. The **init** from the boot logical volume is replaced by the **init** from the root file system **hd4**.
- This **init** processes the **/etc/inittab** file. Out of this file, **rc.boot** is called a third time (**rc.boot 3**) and all remaining devices are configured.

# rc.boot 1



© Copyright IBM Corporation 2004

Figure 4-3. rc.boot 1

AU1612.0

## Notes:

The **init** process started from the RAM file system executes the boot script **rc.boot 1**. If **init** fails for some reason (for example, a bad boot logical volume), **c06** is shown on the LED display. The following steps are executed when **rc.boot 1** is called:

- The **restbase** command is called which copies the ODM from the boot logical volume into the RAM file system. After this step an ODM is available in the RAM file system. The LED shows **510** if **restbase** completes successfully, otherwise LED **548** is shown.
- When **restbase** has completed successfully, the configuration manager **cfgmgr** is run with the option **-f** (first). **cfgmgr** reads the **Config\_Rules** class and executes all methods that are stored under **phase=1**. Phase 1 configuration methods results in the configuration of base devices into the system, so that the rootvg can be activated in the next **rc.boot** phase.
- Base devices are all devices that are necessary to access the rootvg. If the rootvg is stored on a hdisk0, all devices from the motherboard to the disk itself must be configured in order to be able to access the rootvg.
- At the end of **rc.boot 1** the system determines the last boot device by calling **bootinfo -b**. The LED shows **511**.



## rc.boot 2 (Part 1)

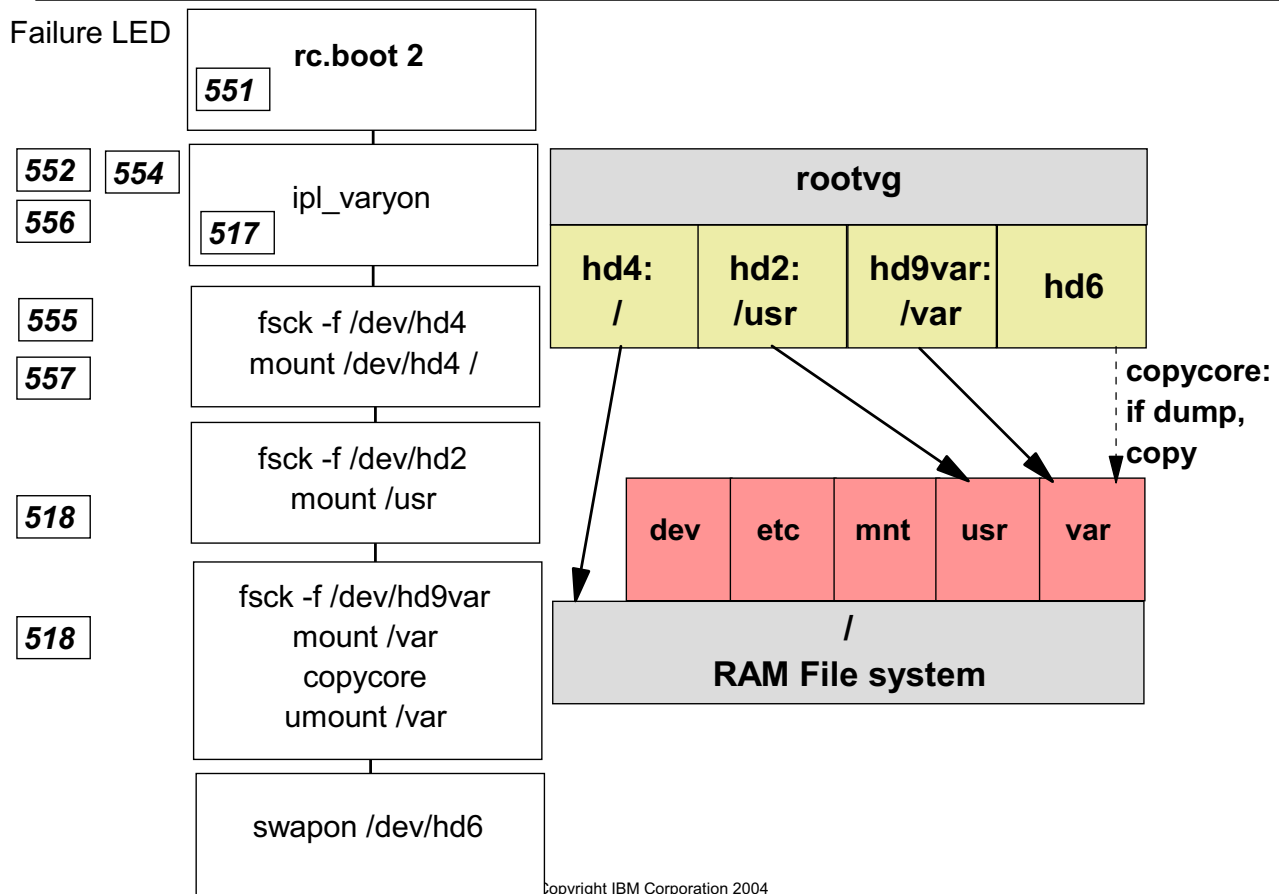


Figure 4-4. rc.boot 2 (Part 1)

AU1612.0

### Notes:

`rc.boot` is run for the second time and is passed to parameter 2. The LED shows **551**. The following steps take part in this boot phase:

- The `rootvg` is varied on with a special version of the `varyonvg` command designed to handle `rootvg`. If `ipl_varyon` completes successfully, **517** is shown on the LED, otherwise **552**, **554** or **556** are shown and the boot process stops.
- The root file system `hd4` is checked by `fsck`. The option `-f` means that the file system is checked only if it was mounted uncleanly during the last shutdown. This improves the boot performance. If the check fails, **555** is shown on the LED.
- Afterwards `/dev/hd4` is mounted directly onto the `root (/)` in the RAM file system. If the mount fails, for example, due to a **corrupted JFS log**, the LED shows **557** and the boot process stops.
- Next `/dev/hd2` is checked (again with option `-f`, that checks only if the file system wasn't unmounted cleanly) and mounted. If the mount fails, LED **518** is displayed and the boot stops.

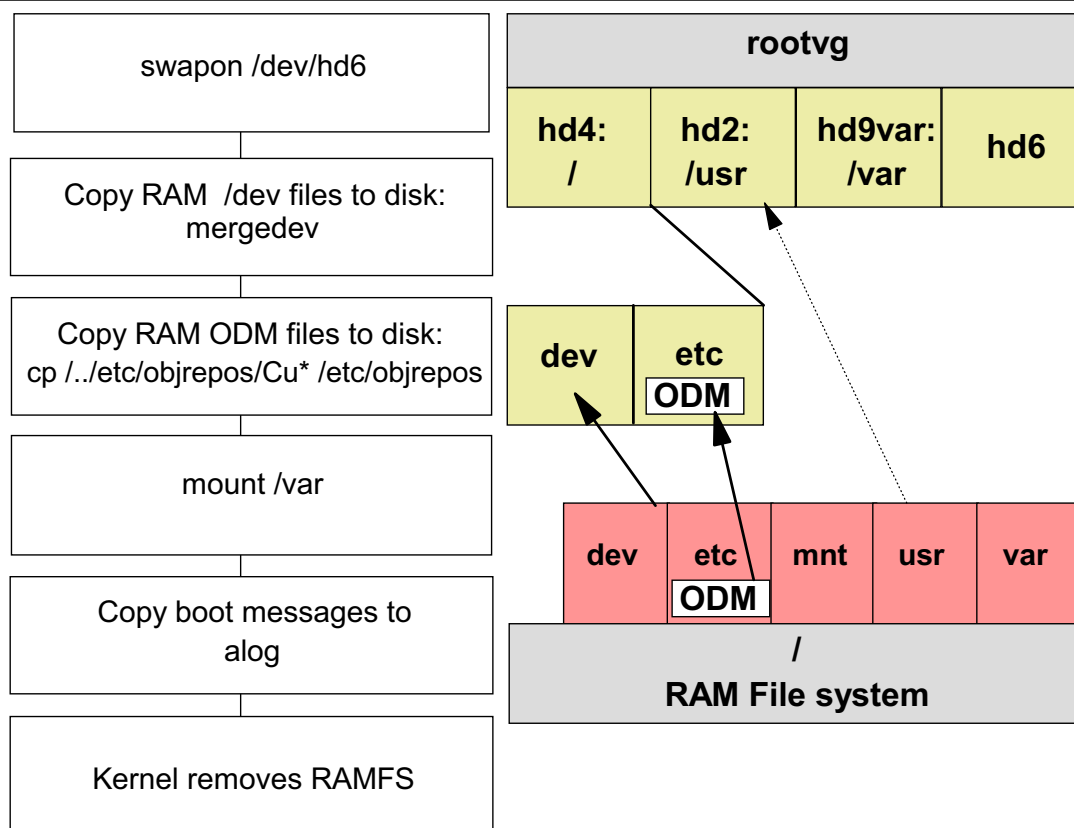
- Next the **/var** file system is checked and mounted. This is necessary at this stage, because the **copycore** command checks if a **dump** occurred. If a dump exists, it will be copied from the dump device **/dev/hd6** to the copy directory which is by default the directory **/var/adm/ras**. **/var** is unmounted afterwards.
- The primary paging space **/dev/hd6** is made available.

Once the disk-based root file system is mounted over the RAMFS, a special syntax is used in **rc.boot** to access the RAMFS files:

- RAMFS files are accessed using a prefix of **././**. For example to access the **fsck** command in the RAMFS (before the **/usr** file system is mounted) **rc.boot** uses **././usr/sbin/fsck**.
- Disk-based files are accessed using normal AIX file syntax. For example, to access the **fsck** command on the disk (after the **/usr** file system is mounted) **rc.boot** uses **/usr/sbin/fsck**.

**Note:** This syntax only works during the boot process. If you boot from the CD-ROM into maintenance mode and need to mount the root file system by hand, you will need to mount it over another directory, such as **/mnt**, or you will be unable to access the RAMFS files.

## rc.boot 2 (Part 2)



© Copyright IBM Corporation 2004

Figure 4-5. rc.boot 2 (Part 2)

AU1612.0

### Notes:

After the paging space `/dev/hd6` has been made available, the following tasks are executed in `rc.boot 2`:

- To understand the next step, remember two things:
  - a. `/dev/hd4` is mounted onto `root(/)` in the RAM file system.
  - b. In `rc.boot 1` the `cfgmgr` has been called and all base devices are configured. This configuration data has been written into the ODM of the RAM file system.
- Now `mergedev` is called and all `/dev` files from the RAM file system are copied to disk.
- All customized ODM files from the RAM file system ODM are copied to disk as well. At this stage both ODMs (in `hd5` and `hd4`) are in sync now.
- The `/var` file system (`hd9var`) is mounted.
- All messages during the boot process are copied into a special file. You must use the `alog` command to view this file:

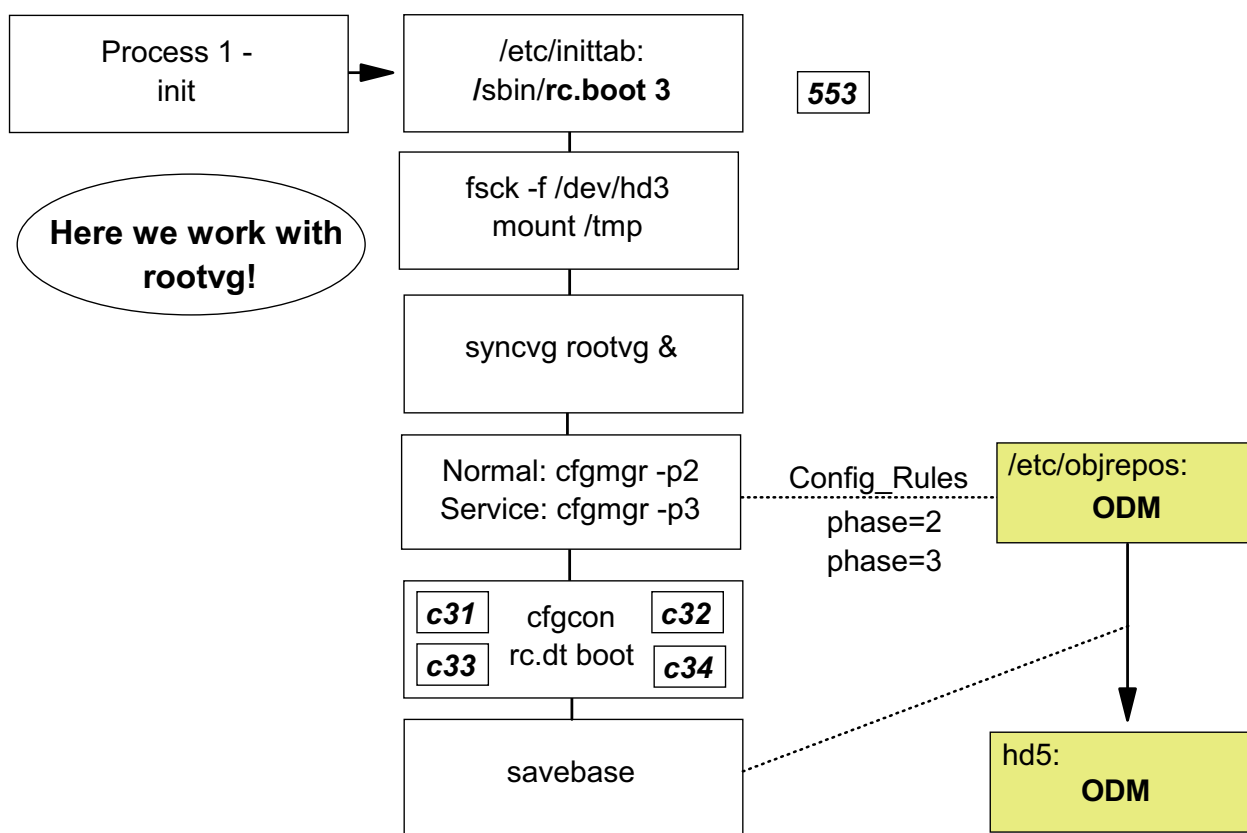
```
# alog -t boot -o
```

**As no console is available at this stage all boot information is collected in this file.**

**When rc.boot 2** is finished, the /, /usr and /var file systems in **rootvg** are active.

At this stage the AIX kernel removes the RAM file system (returns the memory to the free memory pool) and starts the **init** process from the / file system in **rootvg**.

## rc.boot 3 (Part 1)



© Copyright IBM Corporation 2004

Figure 4-6. rc.boot 3 (Part 1)

AU1612.0

### Notes:

At this boot stage, the **/etc/init** process is started. It reads the **/etc/inittab** file (LED displays 553) and executes the commands line by line. It runs **rc.boot** for the third time passing the argument 3, that indicates the last boot phase.

**rc.boot 3** executes the following tasks:

- The **/tmp** file system is checked and mounted.
- The rootvg is synchronized by **syncvg rootvg**. If rootvg contains any **stale partitions** (for example, a disk that is part of rootvg was not active), these partitions are updated and synchronized. **syncvg** is started as a background job.
- The configuration manager is called again. If the key switch is normal the **cfgmgr** is called with option **-p2** (phase 2). If the key switch is service (either the physical key switch of a microchannel or the logical key switch of a PCI model), the **cfgmgr** is called with option **-p3** (phase 3).

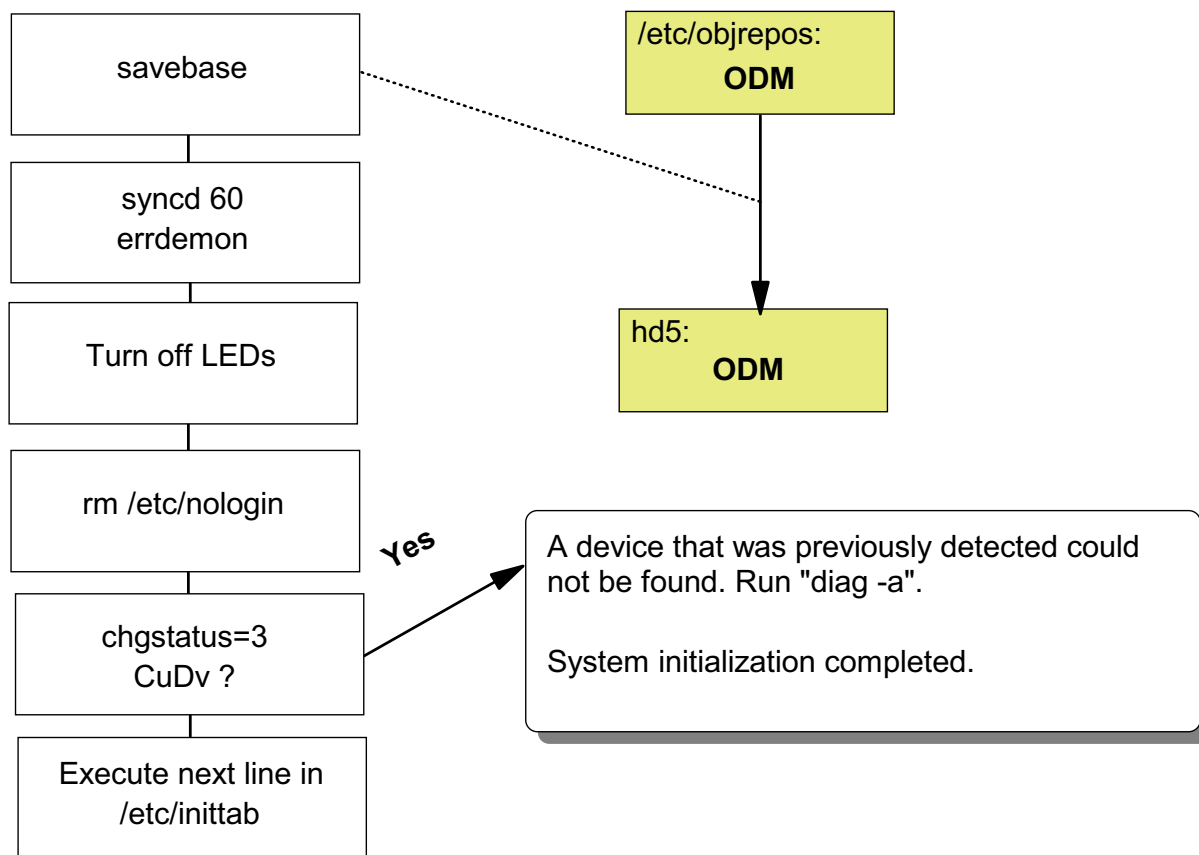
The configuration manager reads ODM class **Config\_Rules** and executes either all methods for **phase=2** or **phase=3**. All remaining devices that are not base devices are configured in this step.

- The console will be configured by **cfgcon**. The numbers **c31**, **c32**, **c33** or **c34** are displayed depending on the type of console:
  - **c31**: Console not yet configured. Provides instruction to select a console.
  - **c32**: Console is a **lft** terminal
  - **c33**: Console is a **tty**
  - **c34**: Console is a file on the disk

If CDE is specified in **/etc/inittab**, the CDE will be started and you get a graphical boot on the console.

- To synchronize the ODM in the boot logical volume with the ODM from the / file system, **savebase** is called.

## rc.boot 3 (Part 2)



© Copyright IBM Corporation 2004

Figure 4-7. rc.boot 3 (Part 2)

AU1612.0

### Notes:

After the ODMs have been synchronized again, the following steps take place:

- The **syncd** daemon is started. All data that is written to disk is first stored in a cache in memory before writing it to the disk. The **syncd** daemon writes the data from the cache each 60 seconds to the disk. Another daemon process, the **errdemon** daemon is started. This process allows errors triggered by applications or the kernel to be written to the error log.
- The LED display is turned off.
- If a file `/etc/nologin` exists, it will be removed. If a system administrator creates this file, a login to the AIX machine is not possible. During the boot process `/etc/nologin` will be removed.
- If devices exist that are flagged as **missing** in `CuDv` (`chgstatus=3`), a message is displayed on the console. For example, this could happen if external devices are not powered on during system boot.
- The last message **System initialization completed** is written to the console. **rc.boot 3** is finished. The **init** process executes the next command in `/etc/inittab`.

## rc.boot Summary

	Where From	Action	Phase Config_Rules
<b>rc.boot 1</b>	/dev/ram0	restbase cfgmgr -f	1
<b>rc.boot 2</b>	/dev/ram0	ipl_varyon rootvg Merge /dev Copy ODM	
<b>rc.boot 3</b>	rootvg	cfgmgr -p2 cfgmgr -p3 savebase	2-normal 3-service

© Copyright IBM Corporation 2004

Figure 4-8. rc.boot Summary

AU1612.0

### Notes:

This page summarizes the **rc.boot** script.

During **rc.boot 1** all base devices are configured. This is done by **cfgmgr -f** which executes all phase 1 methods from **Config\_Rules**.

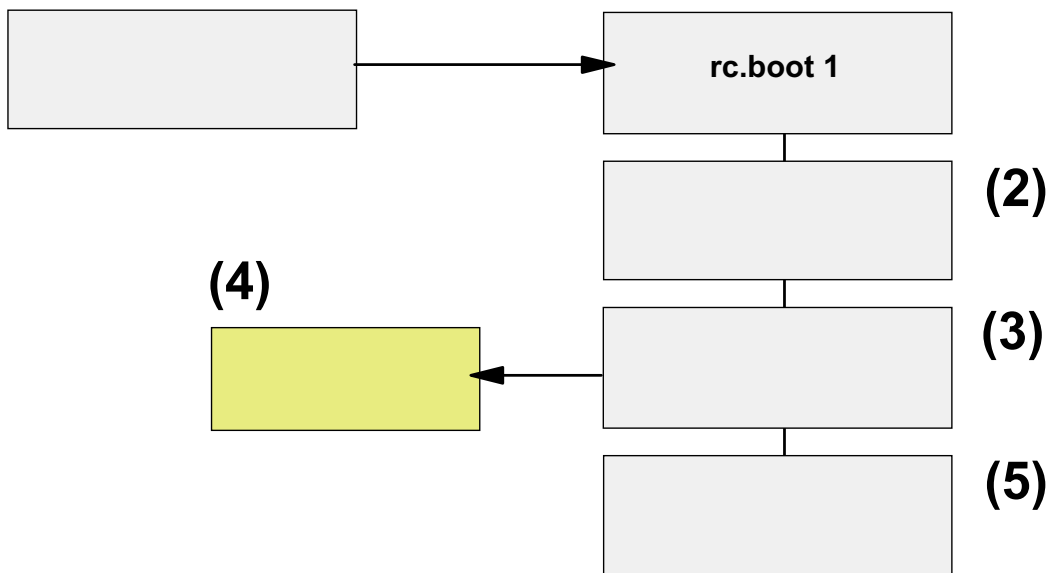
During **rc.boot 2** the rootvg is varied on. All **/dev** files and the customized ODM files from the RAM file system are merged to disk.

During **rc.boot 3** all remaining devices are configured by **cfgmgr -p**. The configuration manager reads the **Config\_Rules** class and executes the corresponding methods. To synchronize the ODMs, **savebase** is called that writes the ODM from the disk back to the boot logical volume.



## Let's Review: Review rc.boot 1

(1)



© Copyright IBM Corporation 2004

Figure 4-9. Let's Review: Review rc.boot 1

AU1612.0

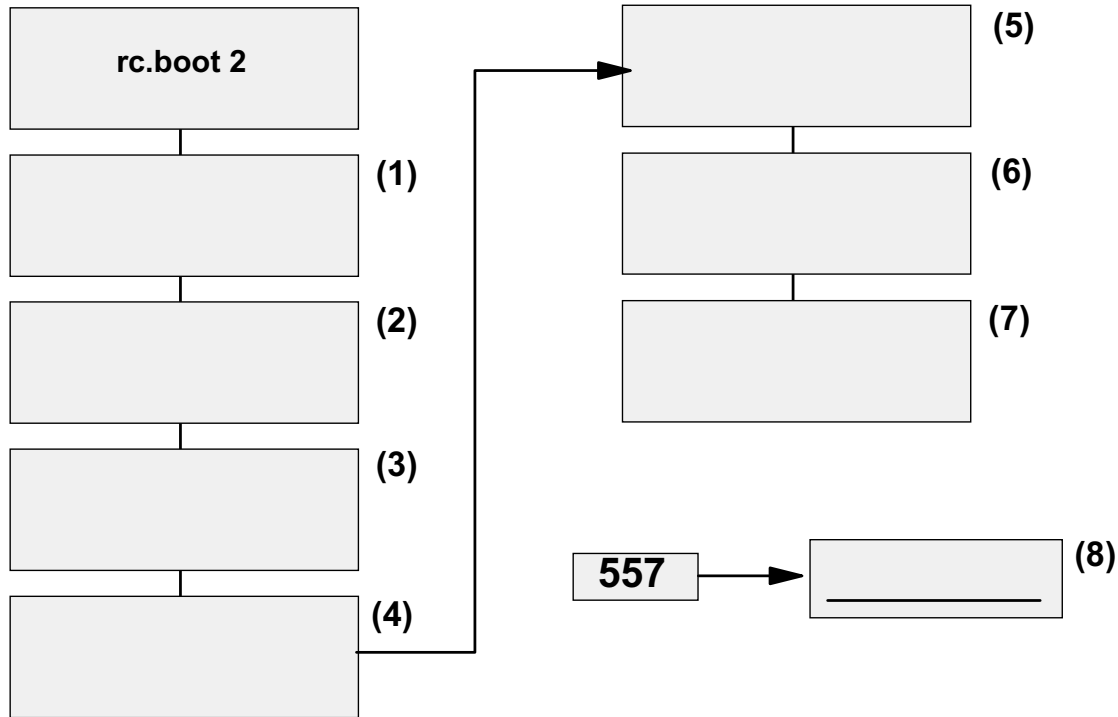
### Notes:

Please answer the following question and put the solutions into the picture above.

1. Who calls **rc.boot 1**? Is it:
  - /etc/init from hd4
  - /etc/init from hd5
2. Which command copies the ODM files from the boot logical volume into the RAM file system?
3. Which command triggers the execution of all phase 1 methods in Config\_Rules?
4. Which ODM files contains the devices that have been configured in **rc.boot 1**?
  - ODM files in hd4
  - ODM files in RAM file system
5. How can you determine the last boot device?

When you completed these questions, please go ahead with the review of **rc.boot 2**.

# Let's Review: Review rc.boot 2



© Copyright IBM Corporation 2004

Figure 4-10. Let's Review: Review rc.boot 2

AU1612.0

### Notes:

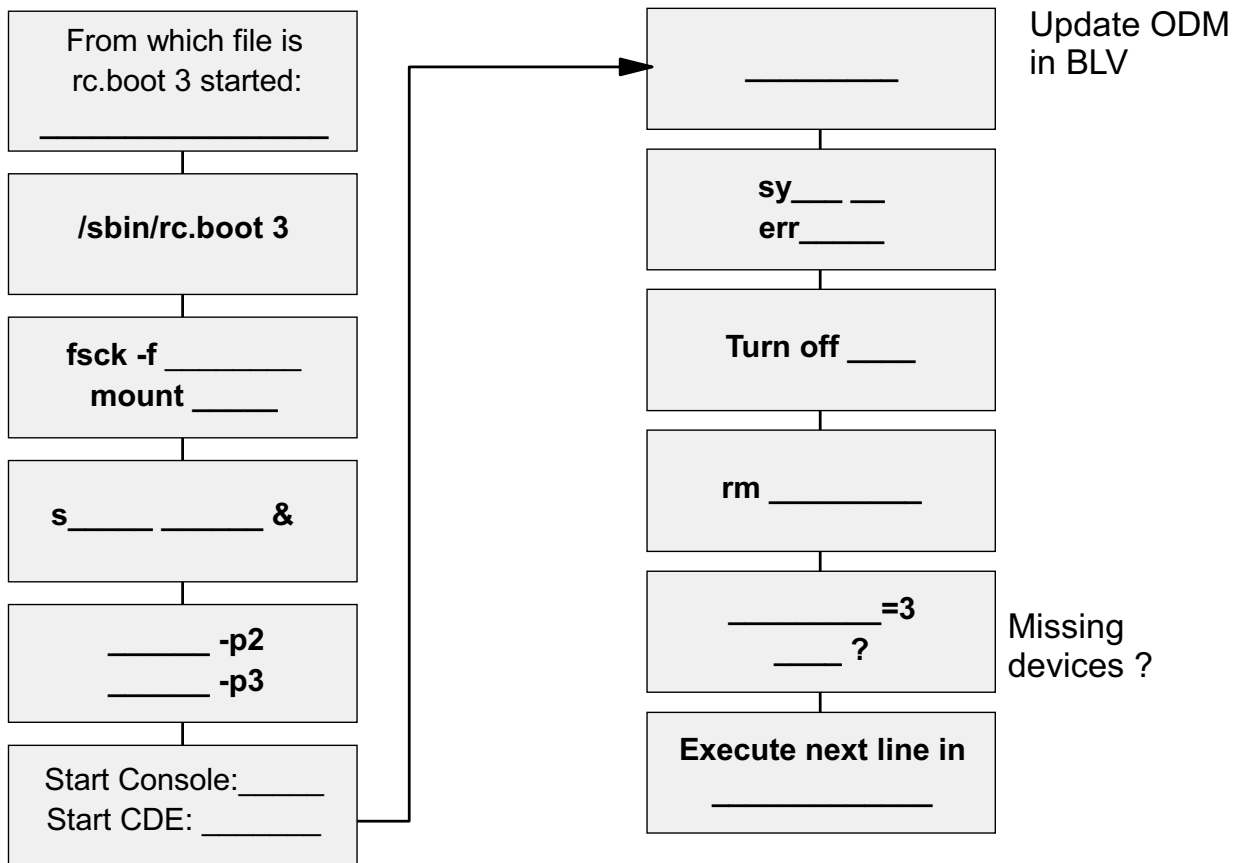
This page reviews **rc.boot 2**. Please order the following nine expressions in the correct sequence:

1. Turn on paging
2. Merge RAM /dev files
3. Copy boot messages to alog
4. Activate rootvg
5. Mount /var; copy dump; Unmount /var
6. Mount /dev/hd4 onto / in RAMFS
7. Copy RAM ODM files

Finally answer the following question. Put the answer in box 8:

Your system stops booting with an LED 557. Which command failed?

# Let's Review: Review rc.boot 3



© Copyright IBM Corporation 2004

Figure 4-11. Let's Review: Review rc.boot 3

AU1612.0

## Notes:

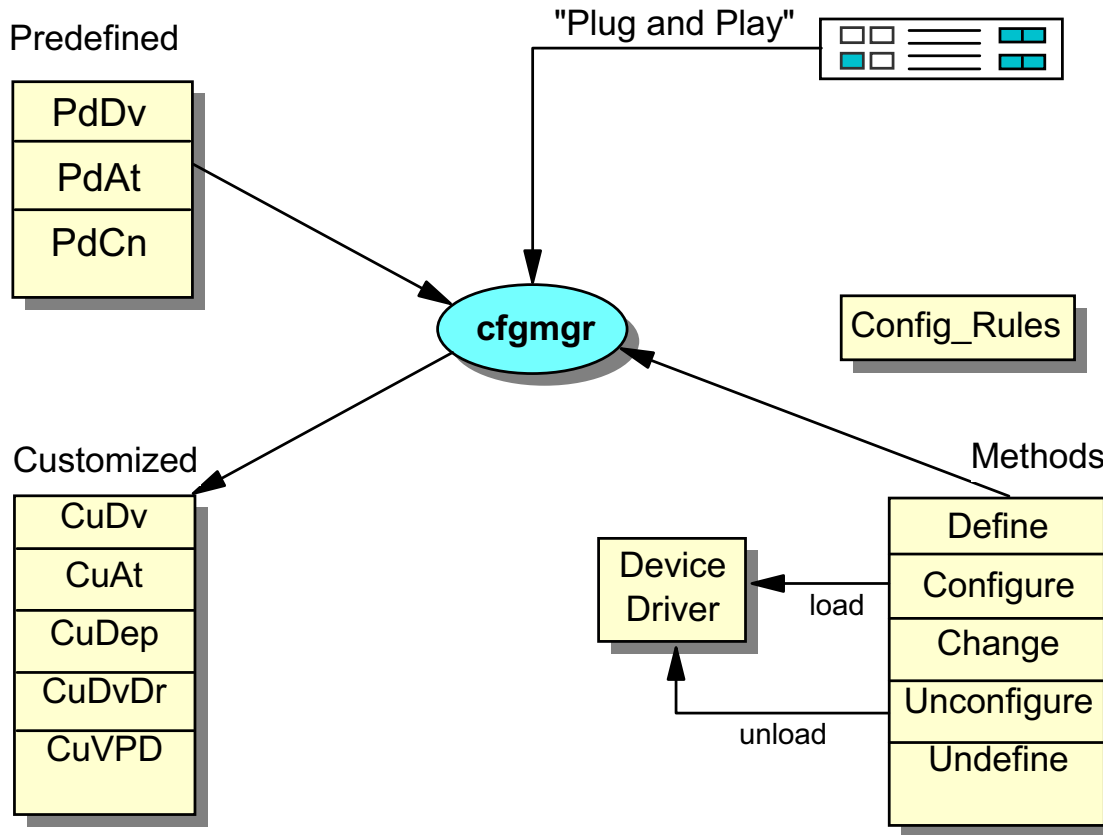
Please complete the missing information in the picture.

Your instructor will review the activity with you.



## 4.2 AIX Initialization Part 2

# Configuration Manager



© Copyright IBM Corporation 2004

Figure 4-12. Configuration Manager

AU1612.0

## Notes:

This page summarizes the tasks of the configuration manager in AIX.

During system boot the configuration manager is invoked to configure all devices detected as well as any device whose device information is stored in the configuration database. At run time, you can configure a specific device by directly invoking the **cfgmgr** command.

If you encounter problems during the configuration of a device, use **cfgmgr -v**. With this option **cfgmgr** shows the devices as they are configured.

Many devices are automatically detected by the configuration manager. For this to occur, device entries must exist in the predefined object classes. The configuration manager uses the methods from **PdDv** to manage the device state, for example, to bring a device into the defined or available state.

**cfgmgr** can be used to install new device support. If you invoke **cfgmgr** with the **-i** flag, the command attempts to install device software support for each newly detected device.

High-level device commands like **mkdev** invoke methods and allow the user to add, delete, show or change devices and their attributes.

When a device is defined through its `define` method, the information from the predefined database for that type of device is used to create the information describing the device specific instance. This device specific information is then stored in the customized database.

The process of configuring a device is often device-specific. The `configure` method for a kernel device must:

1. Load the device driver into the kernel.
2. Pass device-dependent information describing the device instance to the driver.
3. Create a special file for the device in the `/dev` directory.

Of course, many devices do not have device drivers, such as logical volumes or volume groups which are **pseudodevices**. For this type of device the configured state is not as meaningful. However, it still has a configuration method that simply marks the device as configured or performs more complex operations to determine if there are any devices attached to it.

The configuration process requires that a device be defined or configured before a device attached to it can be defined or configured. At system boot time, the configuration manager configures the system in a hierarchical fashion. First the motherboard is configured, then the buses, then the adapters that are attached, and finally the devices that are connected to the adapters. The configuration manager then configures any pseudodevices (volume groups, logical volumes, and so forth) that need to be configured.

# Config\_Rules Object Class

phase	seq	boot mask	rule	
1	1	0	/etc/methods/defsys	← <b>cfgmgr -f</b>
1	2	0	/usr/lib/methods/deflvm	
2	10	0	/etc/methods/defsys	← <b>cfgmgr -p2 (Normal boot)</b>
2	10	0	/usr/lib/methods/deflvm	
2	15	0	/etc/methods/ptynode	
2	20	0	/etc/methods/startlft	
3	10	0	/etc/methods/defsys	← <b>cfgmgr -p3 (Service boot)</b>
3	10	0	/usr/lib/methods/deflvm	
3	15	0	/etc/methods/ptynode	
3	20	0	/etc/methods/startlft	
3	25	0	/etc/methods/starttty	

© Copyright IBM Corporation 2004

Figure 4-13. Config\_Rules Object Class

AU1612.0

## Notes:

This page shows the ODM class **Config\_Rules** that is used by **cfgmgr** during the boot process. The attribute **phase** determines when the respective method is called:

- All methods with **phase=1** are executed when **cfgmgr -f** is called. The first method that is started is **/etc/methods/defsys**, which is responsible for the configuration of all base devices. The second method **/usr/lib/methods/deflvm** loads the logical volume device driver (LVDD) into the AIX kernel.

If you have devices that must be configured in **rc.boot 1**, that means before the rootvg is active, you need to place phase 1 configuration methods into **Config\_Rules**. A **bosboot** is required afterwards.

- All methods with **phase=2** are executed when **cfgmgr -p2** is called. This takes place in the third **rc.boot** phase, when the key switch is in normal position or for a normal boot on a PCI machine. The **seq** attribute controls the sequence of the execution: The lower the value, the higher the priority.



- All methods with **phase=3** are executed when **cfgmgr -p3** is called. This takes place in the third **rc.boot** phase, when the key switch is in service position, or a service boot has been issued on a PCI system.

Each configuration method has an associated **boot mask**. If the `boot_mask` is zero, the rule applies to all types of boot. If the `boot_mask` is non-zero, the rule then only applies to the boot type specified. For example, if `boot_mask = DISK_BOOT`, the rule would only be used for boots from disk versus `NETWORK_BOOT` which only applies when booting via the network.

## Output of cfgmgr in the Boot Log Using alog

```
# alog -t boot -o
-----
attempting to configure device 'sys0'
invoking /usr/lib/methods/cfgsys_rspc -l sys0
return code = 0
***** stdout *****
bus0
***** no stderr *****
-----
attempting to configure device 'bus0'
invoking /usr/lib/methods/cfgbus_pci bus0
return code = 0
***** stdout *****
bus1, scsi0
***** no stderr *****
-----
attempting to configure device 'bus1'
invoking /usr/lib/methods/cfgbus_isa bus1
return code = 0
***** stdout *****
fda0, ppa0, sa0, sioka0, kbd0
***** no stderr *****
```

Figure 4-14. Output of cfgmgr in the Boot Log Using alog

AU1612.0

### Notes:

Because no console is available during the boot phase, the boot messages are collected in a special file, which, by default, is **/var/adm/ras/bootlog**. As shown, you have to use the **alog** command to view the contents of this file.

To view the boot log, issue the command as shown, or use the **smit alog** fastpath.

If you get boot problems, it's always a good idea to check the boot alog file for potential boot error messages. All output from **cfgmgr** is shown in the boot log, as well as other information that is produced in the **rc.boot** script.

The boot alog is created with a default size of 8192 bytes. If you want to increase the size of the boot log, for example to 64 KB, issue the following command:

```
# print "Resizing boot log" | alog -t boot -s 65536
```

## /etc/inittab File

```

init:2:initdefault:
brc::sysinit:/sbin/rc.boot 3 >/dev/console 2>&1 # Phase 3 of system boot
powerfail::powerfail:/etc/rc.powerfail 2>&1 | alog -tboot > /dev/console
rc:23456789:wait:/etc/rc 2>&1 | alog -tboot > /dev/console # Multi-User checks
fbcheck:23456789:wait:/usr/sbin/fbcheck 2>&1 | alog -tboot > /dev/console
srcmstr:23456789:respawn:/usr/sbin/srcmstr # System Resource Controller
rctcpip:23456789:wait:/etc/rc.tcpip > /dev/console 2>&1 # Start TCP/IP daemons
rcnfs:23456789:wait:/etc/rc.nfs > /dev/console 2>&1 # Start NFS Daemons
rchtcpd:23456789:wait:/etc/rc.httpd > /dev/console 2>&1 # Start HTTP daemon
cron:23456789:respawn:/usr/sbin/cron
piobe:2:wait:/usr/lib/lpd/pio/etc/pioint >/dev/null 2>&1 # pb cleanup
sqdaemon:23456789:wait:/usr/bin/startsrc -sqdaemon
writesrv:23456789:wait:/usr/bin/startsrc -swritesrv
uprintfd:23456789:respawn:/usr/sbin/uprintfd
shdaemon:2:off:/usr/sbin/shdaemon >/dev/console 2>&1
l2:2:wait:/etc/rc.d/rc 2
l2:3:wait:/etc/rc.d/rc 3
...
tty0:2:respawn:/usr/sbin/getty /dev/tty0
tty1:2:respawn:/usr/sbin/getty /dev/tty1
ctrmc:2:once:/usr/bin/startsrc -s ctrmc > /dev/console 2>&1
cons:0123456789:respawn:/usr/sbin/getty /dev/console

```

Do not use an editor to change /etc/inittab.  
Use **mkitab**, **chitab**, **rmitab** instead !

© Copyright IBM Corporation 2004

Figure 4-15. /etc/inittab File

AU1612.0

### Notes:

The **/etc/inittab** file supplies information for the **init** process. Before discussing the structure of this file, identify how the **rc.boot** script is executed out of the **inittab** file, to configure all remaining devices in the boot process.

Do not use an editor to change **/etc/inittab**. One small mistake in **/etc/inittab**, and your machine will not boot. Use instead the commands **mkitab**, **chitab** and **rmitab** to edit **/etc/inittab**.

Consider the following examples:

- To add a line to **inittab** use **mkitab**:  
**# mkitab "myid:2:once:/usr/local/bin/errlog.check"**
- Identify, in the sample **inittab**, the **tty1** line.  
 To change **inittab** so that **init** will ignore this line, issue the following command:  
**# chitab "tty1:2:off:/usr/sbin/getty /dev/tty1"**
- To remove the line **tty1** from **inittab** use the following command:

## # rmitab tty1

Besides these commands, the command **lsitab** views the **inittab** file:

```
# lsitab dt
dt:2:wait:/etc/rc.dt
```

If you issue **lsitab -a**, the complete **inittab** is shown.

The advantage of these commands is that they always guarantee a non-corrupted **inittab** file. If your machine stops booting with an LED **553**, this indicates a bad **inittab** file in most cases.

Another daemon (**shdaemon**) also started with **inittab**, called the system hang detection, provides a SMIT-configurable mechanism to detect certain types of system hangs and initiate the configured action. The **shdaemon** daemon uses a corresponding configuration program named **shconf**.

The system hang detection feature uses a **shdaemon** entry in the **/etc/inittab** file, as shown in the visual, with an action field that is set to off by default. Using the **shconf** command or SMIT (fastpath: **smit shd**), you can enable this daemon and configure the actions it takes when certain conditions are met. **shdaemon** is described in the next visual.

## System Hang Detection

- System hangs
  - High priority process
  - Other
- What does `shdaemon` do?
  - Monitors system's ability to run processes
  - Takes specified action if threshold is crossed
- Actions
  - Log Error in the Error Logging
  - Display a warning message on the console
  - Launch recovery login on a console
  - Launch a command
  - Automatically REBOOT system

© Copyright IBM Corporation 2004

Figure 4-16. System Hang Detection

AU1612.0

### Notes:

**shdaemon** can help recover from certain types of system hangs. For our purposes, we will divide system hangs into two types:

- High priority process

The system may appear to be hung if some applications have adjusted their process or thread priorities so high that regular processes are not scheduled. In this case, work is still being done, but only by the high priority processes. As currently implemented, `shdaemon` specifically addresses this type of hang.

- Other

Other types of hangs may be caused by a variety of problems (for example: system thrashing, kernel deadlock, kernel in tight loop, and so forth). In these cases, no (or very little) meaningful work will get done. `shdaemon` may help with some of these problems.

If enabled, **shdaemon** monitors the system to see if any process with a process priority number higher than a set threshold has been run during a set time-out period.

**Note:** Remember that a higher process priority number indicates a lower priority on the system.

In effect, **shdaemon** monitors to see if lower priority processes are being scheduled.

**shdaemon** runs at the highest priority (priority number = 0) so that it will always be able to get CPU time, even if a process is running at very high priority.

#### Actions

If lower priority processes are not being scheduled, shdaemon will perform the specified action. Each action can be individually enabled and has its own configurable priority and time-out values. There are five actions available:

- Log Error in the Error Logging
- Display a warning message on a console
- Launch a recovery login on a console
- Launch a command
- Automatically REBOOT system

# Configuring shdaemon

```
# shconf -E -l prio
sh_pp      enable      Enable Process Priority Problem

pp_errlog  enable      Log Error in the Error Logging
pp_eto     2           Detection Time-out
pp_eprio   60          Process Priority

pp_warning enable      Display a warning message on a console
pp_wto     2           Detection Time-out
pp_wprio   60          Process Priority
pp_wterm   /dev/console Terminal Device

pp_login   disable     Launch a recovering login on a console
pp_lto     2           Detection Time-out
pp_lprio   100        Process Priority
pp_lterm   /dev/console Terminal Device

pp_cmd     enable      Launch a command
pp_cto     5           Detection Time-out
pp_cprio   60          Process Priority
pp_cpath   /home/unhang     Script

pp_reboot  disable     Automatically REBOOT system
pp_rto     5           Detection Time-out
pp_rprio   39          Process Priority
```

© Copyright IBM Corporation 2004

Figure 4-17. Configuring shdaemon

AU1612.0

## Notes:

**shdaemon** configuration information is stored as attributes in the SWservAt ODM object class. Configuration changes take effect immediately and survive across reboots.

Use **shconf** (or **smit shd**) to configure or display the current configuration of shdaemon.

## Enabling shdaemon

At least two parameters must be modified to enable shdaemon:

- Enable priority monitoring (**sh\_pp**)
- Enable one or more actions (**pp\_errlog**, **pp\_warning**, and so forth)

When enabling shdaemon, shconf performs the following steps:

- Modifies the SWservAt parameters
- Starts **shdaemon**
- Modifies **/etc/inittab** so that shdaemon will be started on each system boot

## Action attributes

Each action has its own attributes, which set the priority and time-out thresholds and define the action to be taken.

### Example

In the example, **shdaemon** is enabled to monitor process priority (**sh\_pp=enable**), and the following actions are enabled:

- Log Error in the Error Logging (**pp\_log=enable**)

Every two minutes (**pp\_eto=2**), **shdaemon** will check to see if any process has been run with a process priority number greater than 60 (**pp\_eprio=60**). If not, **shdaemon** logs an error to the error log.

- Display a warning message on a console (**pp\_warning=enable**)

Every two minutes (**pp\_wto=2**), **shdaemon** will check to see if any process has been run with a process priority number greater than 60 (**pp\_wprio=60**). If not, **shdaemon** send a warning message to the console specified by **pp\_wterm**.

- Launch a command (**pp\_cmd=enable**)

Every five minutes (**pp\_cto=5**), **shdaemon** will check to see if any process has been run with a process priority number greater than 60 (**pp\_cprio=60**). If not, **shdaemon** runs the command specified by **pp\_cpath** (in this case, **/home/unhang**).



---

## Resource Monitoring and Control (RMC)

---

- Based on two concepts: conditions and responses
- Associates predefined responses with predefined conditions for monitoring system resources.
  - Example: Broadcast a message to the system administrator when the /tmp file system becomes 90% full.

© Copyright IBM Corporation 2004

Figure 4-18. Resource Monitoring and Control

AU1612.0

### **Notes:**

RMC is automatically installed and configured when AIX is installed.

A very good redbook describing this topic is:

SG24-6615 *A Practical Guide for Resource Monitoring and Control*

This redbook can be found under

<http://www.redbooks.ibm.com/redbooks/pdfs/sg246615.pdf>

RMC is started by an entry in /etc/inittab:

```
ctrmc:2:once:/usr/bin/startsrc -s ctrmc > /dev/console 2>&1
```

To provide a ready-to-use system, 84 conditions, 8 responses are predefined

- Use them as they are
- Customize them
- Use as templates to define your own

To monitor a condition, simply associate one or more responses with the condition.

A log file is maintained in /var/ct.

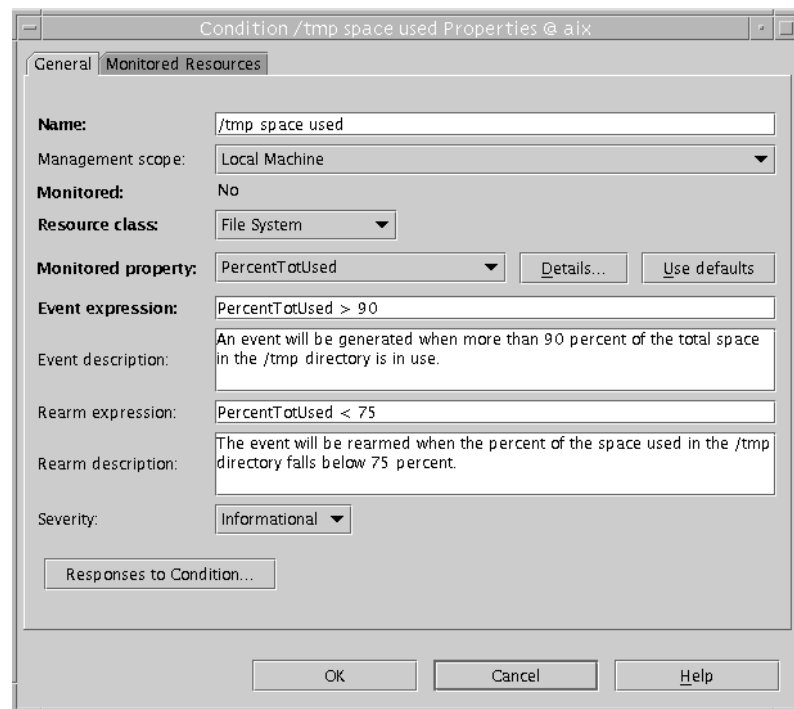
The following steps are provided to assist you in setting up an efficient monitoring system:

1. Review the predefined conditions of your interests. Use them as they are, customize them to fit your configurations, or use them as templates to create your own.
2. Review the predefined responses. Customize them to suit your environment and your working schedule. For example, the response "Critical notifications" is predefined with three actions:
  - a. Log events to /tmp/criticalEvents.
  - b. E-mail to root.
  - c. Broadcast message to all logged-in users any time when an event or a rearm event occurs.

You may modify the response, such as to log events to a different file any time when events occur, e-mail to you during non-working hours, and add a new action to page you only during working hours. With such a setup, different notification mechanisms can be automatically switched, based on your working schedule.

3. Reuse the responses for conditions. For example, you can customize the three severity responses, "Critical notifications," "Warning notifications," and "Informational notifications" to take actions in response to events of different severities, and associate the responses to the conditions of respective severities. With only three notification responses, you can be notified of all the events with respective notification mechanisms based on their urgencies.
4. Once the monitoring is set up, your system continues being monitored whether your Web-based System Manager session is running or not. To know the system status, you may bring up a Web-based System Manager session and view the Events plug-in, or simply use the `lsaudrec` command from the command line interface to view the audit log.

# RMC Conditions Property Screen: General Tab



© Copyright IBM Corporation 2003

Figure 4-19. RMC Conditions Property Screen: General Tab

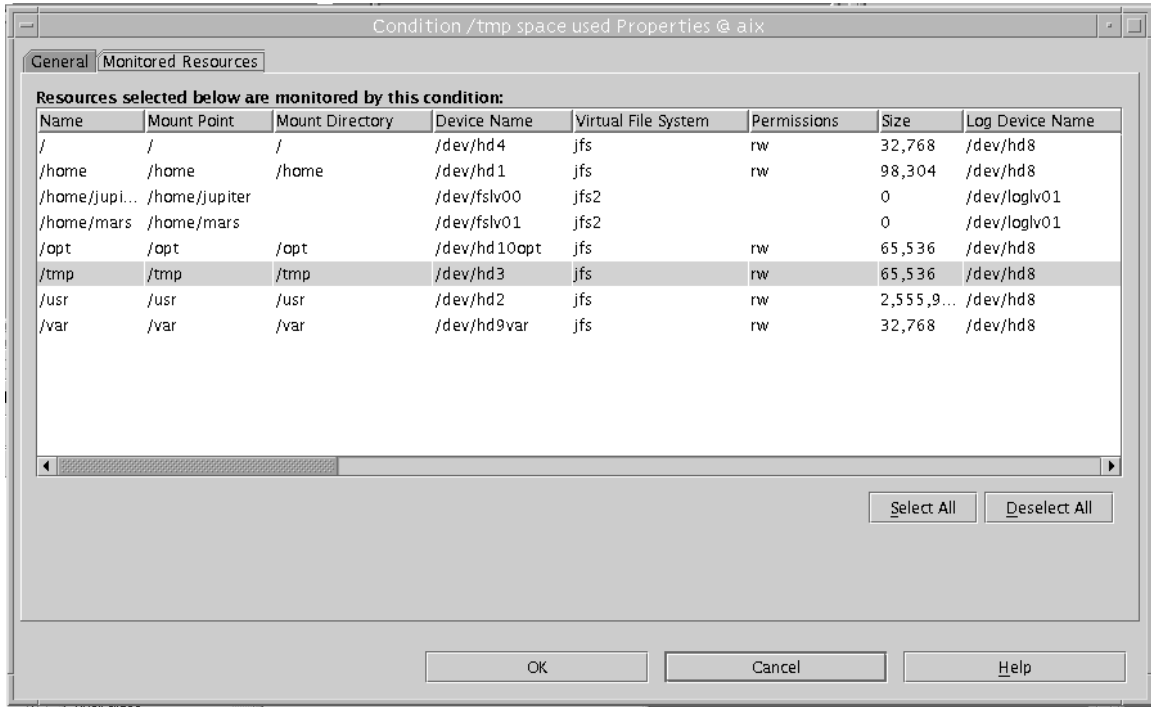
AU1612.0

## Notes:

A condition monitors a specific property, such as total percentage used, in a specific resource class, such as JFS.

Each condition contains an event expression to define an event and an optional re-arm event.

# RMC Conditions Property Screen: Monitored Resources Tab



© Copyright IBM Corporation 2003

Figure 4-20. RMC Conditions Property Screen: Monitored Resources Tab

AU1612.0

## Notes:

You can monitor the condition for one or more resources within the monitored property, such as /tmp, or /tmp and /var, or all of the file systems.

# RMC Actions Property Screen: General Tab

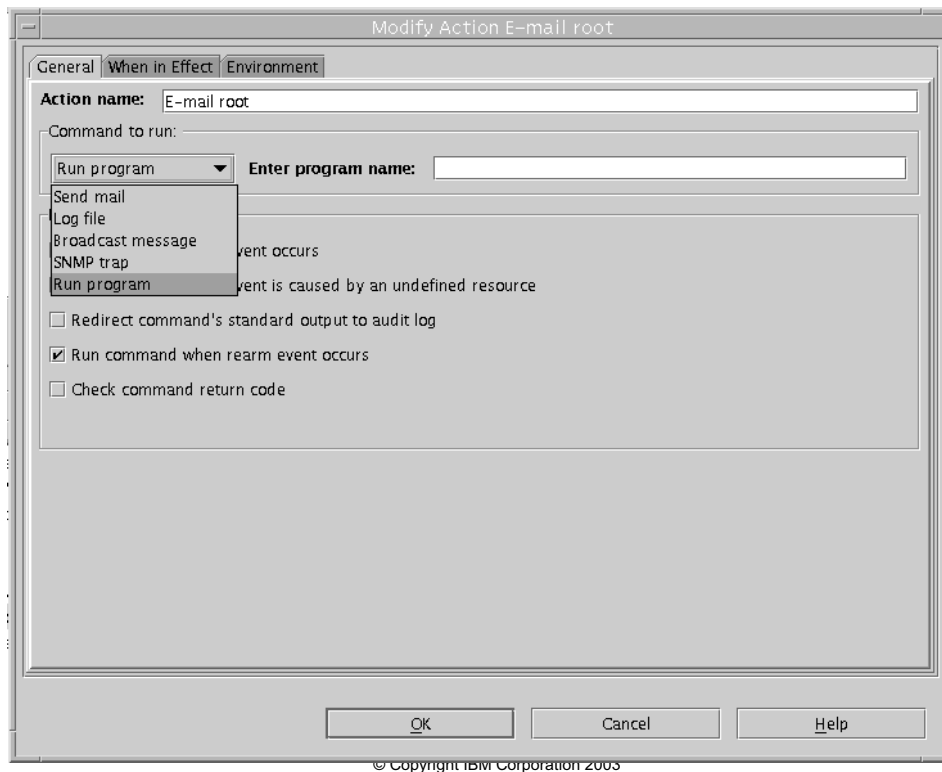


Figure 4-21. RMC Actions Property Screen: General Tab

AU1612.0

## Notes:

To define an action, you can choose one of the three predefined commands, Send Mail, Log an entry to a file, or Broadcast a message, or you can specify an arbitrary program or script of your own by using the Run option.

# RMC Actions Property Screen: When in Effect Tab

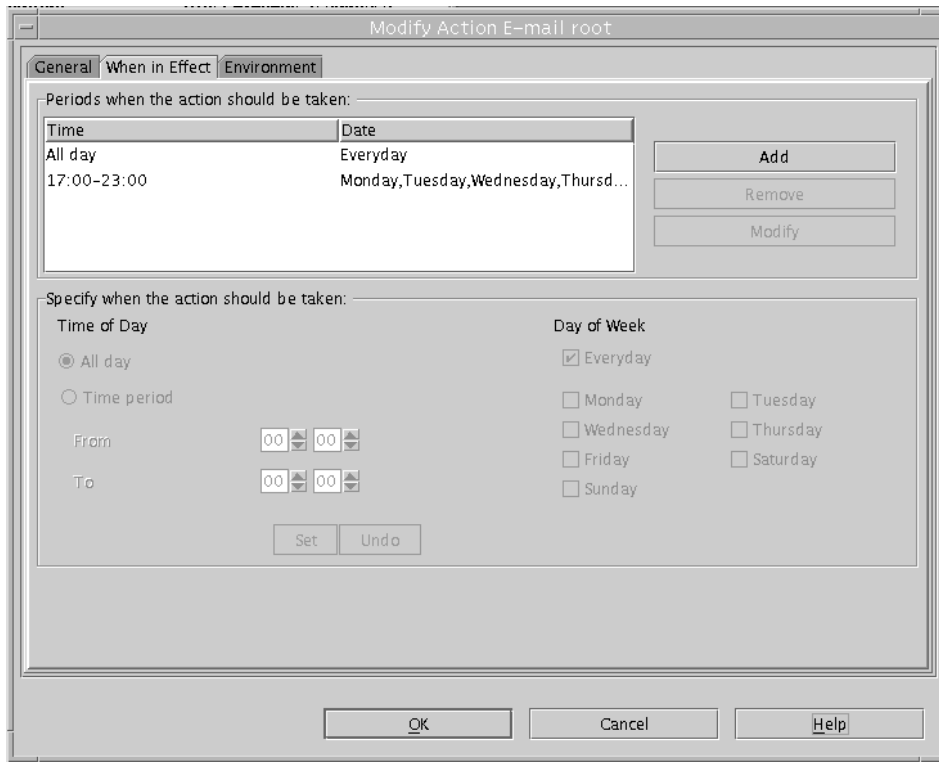


Figure 4-22. RMC Actions Property Screen: When in Effect Tab

AU1612.0

## Notes:

The action can be active for an event only, for a re-arm event only or for both.

You can also specify a time window in which the action is active, such as always, or only during on-shift on weekdays.

Once the monitoring is set up, the system continues to be monitored whether a WSM session is running or not.

## /etc/inittab: Entries You Should Know About

init:2:initdefault:	
brc::sysinit:/sbin/rc.boot 3	
rc:2:wait:/etc/rc	
fbcheck:2:wait:/usr/sbin/fbcheck	
srcmstr:2:respawn:/usr/sbin/srcmstr	
cron:2:respawn:/usr/sbin/cron	
rctcpip:2:wait:/etc/rc.tcpip rcnfs:2:wait:/etc/rc.nfs	
qdaemon:2:wait:/usr/bin/startsrc -sqdaemon	
dt:2:wait:/etc/rc.dt	
tty0:2:off:/usr/sbin/getty /dev/tty1	
myid:2:once:/usr/local/bin/errlog.check	

Figure 4-23. /etc/inittab: Entries You Should Know About

AU1612.0

### Notes:

Related to the shown **/etc/inittab**, please answer the following questions.

**Note:** Your instructor will complete the empty boxes in the visual after you have answered the questions.

1. Which process is started by the **init** process only one time? The **init** process does not wait for the initialization of this process.

---

2. Which process is involved in print activities on an AIX system?

---

3. Which line is ignored by the **init** process?

---

4. Which line determines that multiuser mode is the initial run level of the system?

---

5. Where is the System Resource Controller started?

---

6. Which line controls network processes?

---

---

7. Which component allows the execution of programs at a certain date or time?

---

8. Which line executes a file **/etc/firstboot** if it exists?

---

9. Which script controls starting of the CDE desktop?

---

10. Which line is executed in all run levels?

---

11. Which line takes care of varying on the volume groups, activating paging spaces and mounting file systems that are to be activated during boot?

---



# Boot Problem Management

Check:	LED:	User Action:
Bootlist wrong?	LED codes cycle	PowerOn, press F1, select Multi-Boot, select the correct bootdevice.
/etc/inittab? /etc/environment?	553	Access the rootvg. Check /etc/inittab (empty, missing or corrupt?). Check /etc/environment.
BLV or Boot record corrupt?	20EE000B	Access the rootvg. Re-create the BLV: # bosboot -ad /dev/hdiskx
JFS log corrupt?	551, 552, 554, 555, 556, 557	Access rootvg <b>before</b> mounting the rootvg file systems. Re-create the JFS log: # logform -V jfs /dev/hd8 Run fsck afterwards.
Superblock corrupt?	552, 554, 556	Run fsck against all rootvg-file systems. If fsck indicates errors (not an AIX file system), repair the superblock as described in the notes.
rootvg locked?	551	Access rootvg and unlock the rootvg: # chvg -u rootvg
ODM files missing?	523 - 534	ODM files are missing or inaccessible. Restore the missing files from a system backup.
Mount of /usr or /var failed?	518	Check /etc/filesystem. Check network (remote mount), file systems (fsck) and hardware.

Figure 4-24. Boot Problem Management

AU1612.0

## Notes:

This page shows some common boot errors that might happen during the AIX software boot process.

Some of the more common ones are shown above. Let's take a closer look.

### 1. Bootlist wrong?

If the bootlist is wrong the system cannot boot anymore. This is very easy to fix. Boot in SMS Menu by pressing F1, select Multi-Boot and select the correct boot device. Keep in mind that only harddisks with boot records are shown as selectable boot devices.

### 2. /etc/inittab corrupt? /etc/environment corrupt?

A LED of 553 mostly indicates a corrupted /etc/inittab file, but in some cases a bad /etc/environment may also lead to a 553. To fix this problem boot in maintenance mode and check both files. Consider using a **mksysb** to retrieve these files from a backup tape.

### 3. Boot logical volume or boot record corrupt?

The next thing to try if your machine does not boot, is to check the boot logical volume.

To fix a corrupted boot logical volume, boot in maintenance mode and use the **bosboot** command:

```
# bosboot -ad /dev/hdisk0
```

### 4. JFS log corrupt?

To fix a corrupted JFS log, boot in maintenance mode and access the rootvg but do not mount the file systems. In the maintenance shell issue the **logform** command and do a file system check for all file systems that use this JFS log. Keep in mind what filesystem type your rootvg had: jfs or jfs2:

```
# logform -V jfs /dev/hd8
# fsck -y -V jfs /dev/hd1
# fsck -y -V jfs /dev/hd2
# fsck -y -V jfs /dev/hd3
# fsck -y -V jfs /dev/hd4
# fsck -y -V jfs /dev/hd9var
# fsck -y -V jfs /dev/hd10opt
exit
```

The **logform** command initializes a new JFS transaction log and this may result in loss of data, because JFS transactions may be destroyed. But, your machine will boot afterwards, because the JFS log has been repaired.

### 5. Superblock corrupt?

Another thing you can try is to check the superblocks of your rootvg file systems. If you boot in maintenance mode and you get error messages like **Not an AIX file system** or **Not a recognized file system type** it is probably due to a corrupt superblock in the file system.

Each file system has two super blocks, one in logical block 1 and a copy in logical block 31. To copy the superblock from block 31 to block 1 for the root file system, issue the following command:

```
# dd count=1 bs=4k skip=31 seek=1 if=/dev/hd4 of=/dev/hd4
```

### 6. rootvg locked?

Many LVM commands place a lock into the ODM to prevent other commands working on the same time. If a lock remains in the ODM due to a crash of a command, this may lead to a hanging system.

To unlock the rootvg, boot in maintenance mode and access the rootvg with file systems. Issue the following command to unlock the rootvg:

```
# chvg -u rootvg
```

### 7. ODM files missing?

If you see LED codes in the range 523 to 534 ODM files are missing on your machine. Use a **mksysb** tape of the system to restore the missing files.

8. **Mount of /usr or /var failed?**

An LED of 518 indicates that the mount of the **/usr or /var** file system **failed**. If /usr is mounted from a network, check the network connection. If /usr or /var are locally mounted, use **fsck** to check the consistency of the file systems. If this does not help check the hardware (diag).

## Next Step

---



© Copyright IBM Corporation 2004

Figure 4-25. Next Step

AU1612.0

### **Notes:**

At the end of the exercise, you should be able to:

- Boot a machine in maintenance mode
- Repair a corrupted log logical volume
- Analyze and fix an unknown boot problem

---

## Checkpoint

---

1. From where is rc.boot 3 run?

---

2. Your system stops booting with LED 557. In which rc.boot phase does the system stop? What can be the reasons for this problem?

---

---

---

3. Which ODM file is used by the **cfgmgr** during boot to configure the devices in the correct sequence?

---

4. What does the line **init:2:initdefault:** in /etc/inittab mean?

---

---

© Copyright IBM Corporation 2004

Figure 4-26. Checkpoint

AU1612.0

### Notes:

## Unit Summary

---

- After the BLV is loaded into RAM, the **rc.boot** script is executed **three times** to configure the system
- During **rc.boot 1** devices to **varyon** the rootvg are configured
- During **rc.boot 2** the rootvg is varied on
- In **rc.boot 3** the remaining devices are configured. Processes defined in **/etc/inittab** file are initiated by the **init** process

© Copyright IBM Corporation 2004

Figure 4-27. Unit Summary

AU1612.0

### **Notes:**

---

# Unit 5. Disk Management Theory

## What This Unit Is About

This unit describes important concepts of the logical volume manager in AIX.

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Describe where the LVM information is stored
- Solve ODM-related LVM problems
- Set up mirroring according to different needs
- Explain the quorum mechanism
- Describe what physical volume states the LVM uses

## How You Will Check Your Progress

Accountability:

- Checkpoint questions
- Lab exercises

## References

Online      *Commands Reference*

Online      *System Management Guide: Operating System and Devices*

GG24-4484-00 *AIX Storage Management*

## Unit Objectives

---

After completing this unit, students should be able to:

- Describe where LVM information is kept
- Solve ODM-related LVM problems
- Set up Mirroring
- Explain the Quorum Mechanism
- Describe Physical Volume States

© Copyright IBM Corporation 2004

Figure 5-1. Unit Objectives

AU1612.0

### **Notes:**

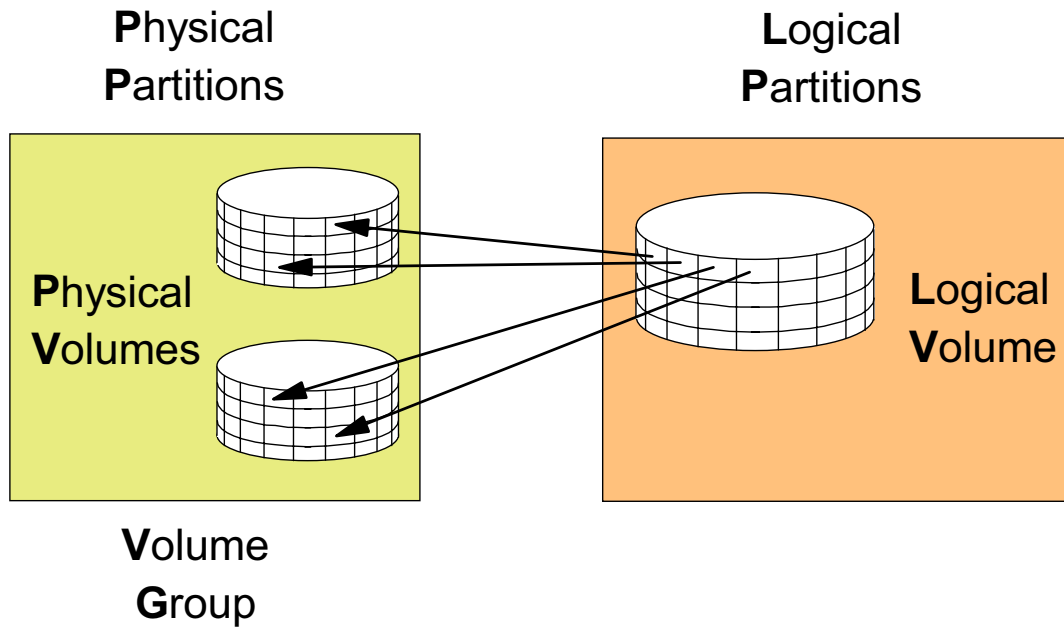
The LVM basic concepts are introduced in the basic system administration course.

We will review and extend your knowledge about LVM in this unit.



## 5.1 Basic LVM Tasks

# LVM Terms



© Copyright IBM Corporation 2004

Figure 5-2. LVM Terms

AU1612.0

## Notes:

Let's start with a review of basic LVM terms.

A **volume group** consists of one or more **physical volumes** that are divided into **physical partitions**. When a volume group is created, a physical partition size has to be specified. This partition size can range from 1 MB to 1024 MB. This physical partition size is the smallest allocation unit for the LVM. If it is not specified, the system will select the minimum size to create 1016 partitions.

The LVM provides **logical volumes**, that can be created, extended, moved and deleted at run time. Logical volumes may span several disks, which is one of the biggest advantages of the LVM.

Logical volumes contain the journaled file systems, paging spaces, journal logs, the boot logical volumes or nothing (when used as a raw logical volume).

Logical volumes are divided into **logical partitions** where each logical partition is associated with at least one physical partition.

Other features of LVM are **mirroring** and **striping**, which are discussed on the following pages.

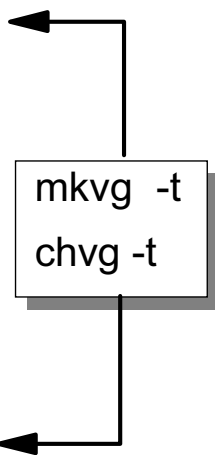
# Volume Group Limits

## • Normal Volume Groups (mkvg)

Number of disks:	Max. number of partitions/disk:
1	32512
2	16256
4	8128
8	4064
16	2032
<b>32</b>	<b>1016</b>

## • Big Volume Groups (mkvg -B or chvg -B)

Number of disks:	Max. number of partitions/disk:
1	130048
2	65024
4	32512
8	16256
16	8128
32	4064
64	2032
<b>128</b>	<b>1016</b>



© Copyright IBM Corporation 2004

Figure 5-3. Volume Group Limits

AU1612.0

### Notes:

Two different volume group types are available:

- **Normal volume groups:** When creating a volume group with **smit** or using the **mkvg** command, without specifying option **-B**, a normal volume group is created.  
The maximum number of logical volumes in a normal volume group is **256**.
- **Big volume groups:** This volume group type has been introduced with AIX 4.3.2. A big volume group must be created with **mkvg -B**.  
A big volume group cannot be imported into an AIX 4.3.1 or lower versions.  
The maximum number of logical volumes in a big volume group is **512**.

Volume groups are created with the **mkvg** command. Here are some examples:

1. Create a normal volume group **datavg**, that contains a disk **hdisk2**:

```
# mkvg -s 16 -t 2 -y datavg hdisk2
```

- The option **-s 16** specifies a partition size of **16 MB**.

- The option **-t 2** is a factor that must be multiplied by 1016. In this case the option indicates that the **maximum number of partitions** on a disk is 2032. That means that the volume group can have up to **16 disks**. Each disk must be less than 4064 megabytes (2032 \* 2).
- The option **-y** specifies the name of the volume group (datavg).

2. Create a big volume group **bigvg** with three disks:

```
# mkgv -B -t 16 -y bigvg hdisk2 hdisk3 hdisk4
```

- The option **-B** specifies that we are creating a **big** volume group.
- The option **-t 16** indicates that the **maximum number of partitions** on a disk is **16256**. That means that the volume group can have up to **8 disks**.
- The option **-y** specifies the name of the volume group.

Volume groups characteristics could be changed with the **chvg** command. For example, to change a normal volume group **datavg** into a big volume group, the following command must be executed:

```
# chvg -B datavg
```

## Scalable Volume Groups - AIX 5.3

---

- Supports 1024 disks per volume group.
- Supports 4096 logical volumes per volume group.
- Maximum number of PPs is VG instead of PV dependent.
- LV control information is kept in the VGDA.
- No need to set the maximum values at creation time; the initial settings can always be increased at a later date.

© Copyright IBM Corporation 2004

Figure 5-4. Scalable Volume Groups - AIX 5.3

AU1612.0

### **Notes:**

AIX 5L V5.3 takes the LVM scalability to the next higher level and offers a new scalable volume group (scalable VG) type. The scalable VG can accommodate a maximum of 1024 PVs and raises the limit for the number of LVs to 4096. The maximum number of PPs is no longer defined on a per disk basis, but applies to the entire VG. This opens up the prospect to configure VGs with a relatively small number of disks, but with fine grained storage allocation options through a large number of PPs that are small in size. The scalable VG can hold up to 2097152 (2048 K) PPs. Optimally, the size of a physical partition can also be configured for a scalable VG. As with the older VG types, the size is specified in units of megabytes and the size variable must be equal to a power of 2. The range of the PP size starts at 1 (1 MB) and goes up to 131072 (128 GB), which is more than two orders of magnitude above the 1024 (1 GB) maximum for AIX 5L V5.2. (The new maximum PP size provides an architectural support for 256 PB disks.)

Note that the maximum number of user definable LVs is given by the maximum number of LVs per VG minus 1, because one LV is reserved for system use. Consequently, system administrators can configure 255 LVs in normal VGs, 511 in big VGs, and 4095 in scalable VGs.

The LVCB contains meta data about a logical volume. For standard VGs, the LVCB resides in the first block of the user data within the LV. Big VGs keep additional LVCB information in the on disk VGDA. The LVCB structure on

the first LV user block and the LVCB structure within the VGDA are similar but not identical. (If a big VG was created with the -T option of the mkvg command, no LVCB will occupy the first block of the LV.) With scalable VGs, logical volume control information is no longer stored on the first user block of any LV. All relevant logical volume control information is kept in the VGDA as part of the LVCB information area and the LV entry area. So no precautions have to be met when using raw logical volumes because there is no longer a need to preserve the information held by the first 512 bytes of the logical device.

# Configuration Limits for Volume Groups

VG Type	Maximum PVs	Maximum LVs	Maximum PPs per VG	Maximum PP size
Normal VG	32	256	32512 (1016*32)	1 GB
Big VG	128	512	130048 (1016*128)	1 GB
Scalable VG	1024	4096	2097152	128 GB

© Copyright IBM Corporation 2004

Figure 5-5. Configuration Limits for Volume Groups

AU1612.0

## Notes:

To determine the type of a VG, use the `lsvg` command:

```
# lsvg data_svg
```

```
VOLUME GROUP:      mike_svg                VG IDENTIFIER:
000c91ad00004c00000000fd961161d9

VG STATE:          active                    PP SIZE:          16 megabyte(s)
VG PERMISSION:     read/write                TOTAL PPs:        1080 (17280 megabytes)
MAX LVs:           256                      FREE PPs:         1080 (17280 megabytes)
LVs:               0                          USED PPs:         0 (0 megabytes)
OPEN LVs:          0                          QUORUM:           2
TOTAL PVs:         1                          VG DESCRIPTORS:  2
STALE PVs:         0                          STALE PPs:        0
ACTIVE PVs:        1                          AUTO ON:          yes
MAX PPs per VG:    32512 MAX PVs:           1024
LTG size (Dynamic): 256 kilobyte(s)          AUTO SYNC:        no
```



---

HOT SPARE:                   no   BB POLICY:                   relocatable

The value MAX PVs should show which type the VG has. Scalable VGs will say 1024, big VGs will say 128, and original VGs will say 32 (if not changed with the `-t` factor). Additionally, the older VG types have one more line in the output:

...

MAX PPs per VG:   32512

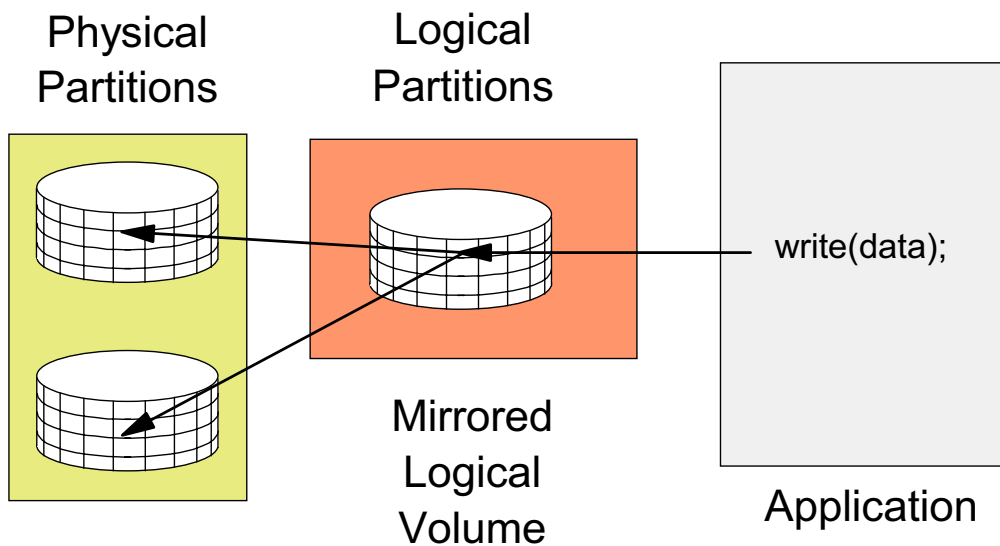
MAX PPs per PV:   1016                                   MAX PVs:       32

...

This lines shows that the VG cannot be a scalable VG, as it is not PP per PV dependent.

A volume group can be converted to a scalable VG using the **chvg -G <vg\_name>** command but the VG must be varied off.

# Mirroring



© Copyright IBM Corporation 2004

Figure 5-6. Mirroring

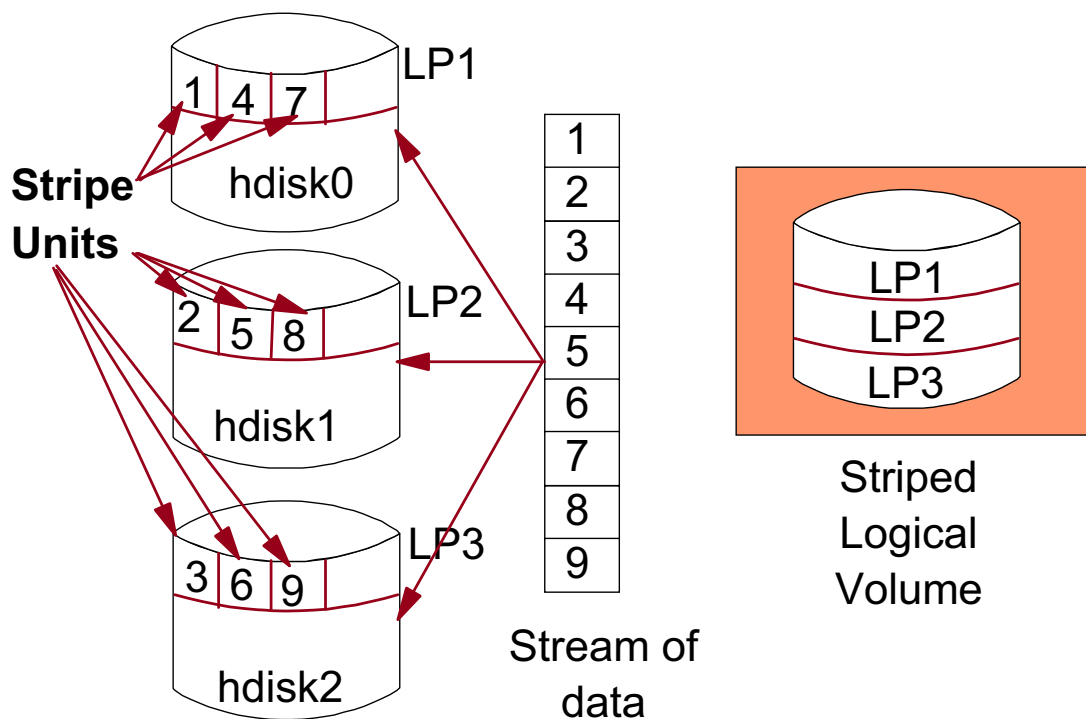
AU1612.0

## Notes:

Logical volumes can be **mirrored**, that means each logical partition gets more than one associated physical partition. The maximum ratio is 1:3; that means one logical partition has three associated physical partitions.

The picture shows a two-disk mirroring of a logical volume. An application writes data to the disk which is always handled by the LVM. The LVM recognizes that this partition is mirrored. The data will be written to both physical partitions. If one of the disks fails, there will be at least one good copy of the data.

# Striping



© Copyright IBM Corporation 2004

Figure 5-7. Striping

AU1612.0

## Notes:

Striping is an LVM feature where the partitions of the logical volume are spread across different disks. The number of disks involved is called **stripe width**.

Striping works by splitting write and read requests to a finer granularity, named **stripe size**. Strip sizes may vary from 4 KB to 128 KB. A single application write or read request is divided into parallel physical I/O requests. The LVM fits the pieces together by tricky buffer management.

Striping makes good sense, when the following conditions are true:

- The disks use separate adapters. Striping on the same adapter does not improve the performance very much.
- The disks are equal in size and speed.
- The disks contain striped logical volumes only.
- Accessing large sequential files. For writing or reading small files striping does not improve the performance.

AIX 5L V5.3 further enhances the LVM RAID implementation and provides striped columns support for logical volumes. This new feature allows you to extend a striped logical volume even if one of the physical volumes in the disk array became full.

In previous AIX releases, you could enlarge the size of a striped logical volume with the `extendlv` command, as long as enough physical partitions were available within the group of disks that define the RAID disk array. Rebuilding the entire LV was the only way to expand a striped logical volume beyond the hard limits imposed by the disk capacities. This workaround required you to back up and delete the striped LV and then to recreate the LV with a larger stripe width followed by a restore operation of the LV data.

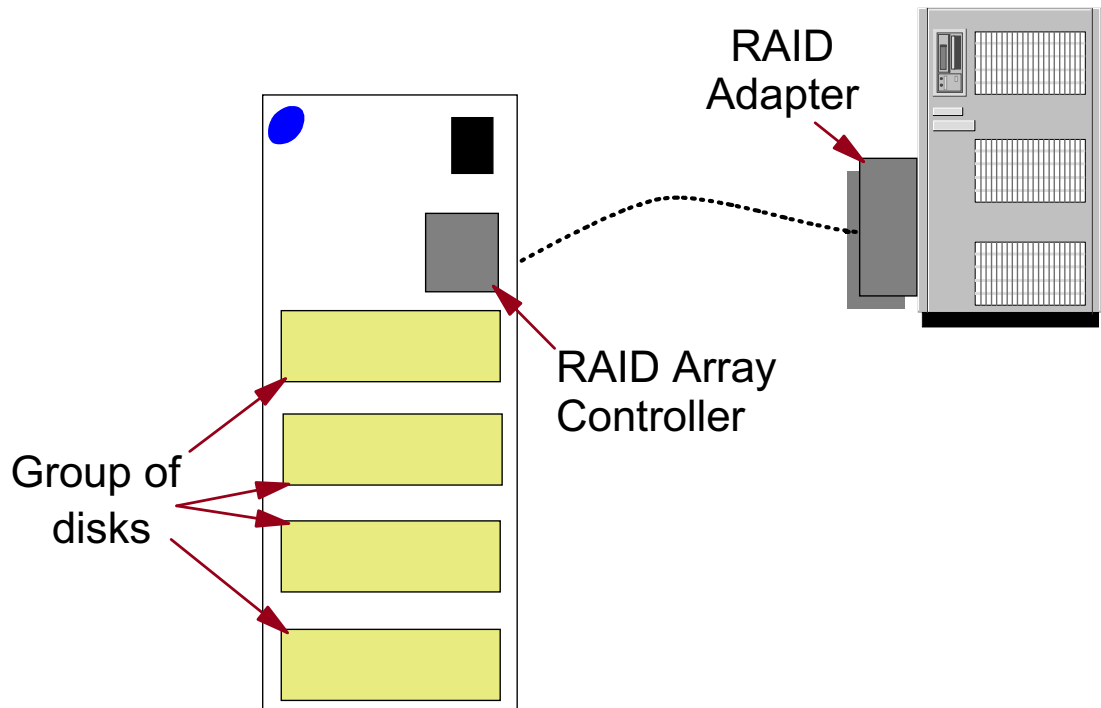
To overcome the disadvantages of this rather time-consuming procedure, AIX 5L V5.3 introduces the concept of striped columns for LVs.

In AIX 5L V5.3, the upper bound can be a multiple of the stripe width. One set of disks, as determined by the stripe width, can be considered as one striped column.

If you use the `extendlv` command to extend a striped logical volume beyond the physical limits of the first striped column, an entire new set of disks will be used to fulfill the allocation request for additional logical partitions. If you further expand the LV, more striped columns may get added as required and as long as you stay within the upper bound limit. The `-u` flag of the `chlv`, `extendlv`, and `mklvcopy` commands will now allow you to change the upper bound to be a multiple of the stripe width. The `extendlv -u` command can be used to change the upper bound and to extend the LV in a single operation.

# Mirroring and Striping with RAID

RAID = **R**edundant **A**rray of **I**ndependent **D**isks



© Copyright IBM Corporation 2004

Figure 5-8. Mirroring and Striping with RAID

AU1612.0

## Notes:

IBM offers storage subsystems (for example the model 7133) that allow mirroring and striping on a hardware level.

The term RAID stands for **Redundant Array of Independent Disks**. Disk arrays are groups of disks that work together to achieve higher data-transfer and I/O rates than those provided by single large drives. An array is a set of multiple disk drives plus an array controller that keeps track of how data is distributed across the drives.

By using multiple drives, the array can provide higher data-transfer rates and higher I/O rates when compared to a single large drive; this is achieved through the consequent ability to schedule reads and writes to the disks in parallel.

Arrays can also provide data redundancy so that no data is lost if a single physical disk in the array should fail. Depending on the RAID level, data is either mirrored or striped.

Striping involves splitting a data file into multiple blocks and writing a sequential set of blocks to each available drive in parallel, repeating this process until all blocks have been written.

Mirroring describes the situation where data written to one disk is also copied exactly to another disk, thereby providing a backup copy.

The most common RAID levels are **RAID 0, RAID 1 and RAID 5**. They are introduced on the next page.

## RAID Levels You Should Know About

RAID Level	Implementation	Explanation
0	Striping	Data is split into blocks. These blocks are written to or read from a series of disks in parallel. No data redundancy.
1	Mirroring	Data is split into blocks and duplicate copies are kept on separate disks. If any disk in the array fails, the mirrored data can be used.
5	Striping with parity drives	Data is split into blocks that are striped across the disks. For each block parity information is written that allows the reconstruction in case of a disk failure.

© Copyright IBM Corporation 2004

Figure 5-9. RAID Levels You Should Know About

AU1612.0

### Notes:

The most common RAID levels are **RAID 0**, **RAID 1** and **RAID 5**.

#### 1. RAID 0:

RAID 0 is known as disk striping. Conventionally, a file is written out to (or read from) a disk in blocks of data. With striping, the information is split into chunks (a fixed amount of data) and the chunks are written to (or read from) a series of disks in parallel.

RAID 0 is well suited for applications requiring fast read or write accesses. On the other hand, RAID 0 is only designed to increase performance, there is no data redundancy, so any disk failure will require reloading from backups.

Select RAID level 0 for applications that would benefit from the increased performance capabilities of this RAID level. Never use this level for critical applications that require high availability.

#### 2. RAID 1:

RAID 1 is known as disk mirroring. In this implementation, duplicate copies of each chunk of data are kept on separate disks, or more usually, each disk has a twin that

contains an exact replica (or mirror image) of the information. If any disk in the array fails, then the mirrored twin can take over.

Read performance can be enhanced as the disk with its actuator closest to the required data is always used, thereby minimizing seek times. The response time for writes can be somewhat slower than for a single disk, depending on the write policy; the writes can either be executed in parallel for speed, or serially for safety. This technique improves response time for read-mostly applications, and improves availability. The downside is you'll need twice as much disk space.

RAID 1 is most suited to applications that require high data availability, good read response times, and where cost is a secondary issue.

### 3. **RAID 5:**

RAID 5 can be considered as disk striping combined with a sort of mirroring. That means that data is split into blocks that are striped across the disks, but additionally parity information is written that allows recovery in the event of a disk failure.

Parity data is never stored on the same drive as the blocks that are protected. In the event of a disk failure, the information can be rebuilt by the using the parity information from the remaining drives.

Select RAID level 5 for applications that manipulate small amounts of data, such as transaction processing applications. This level is generally considered the best all-around RAID solution for commercial applications.

RAID algorithms can be implemented as part of the operating system's file system software, or as part of a disk device driver. AIX LVM supports the following RAID options:

RAID 0 Striping

RAID 1 Mirroring

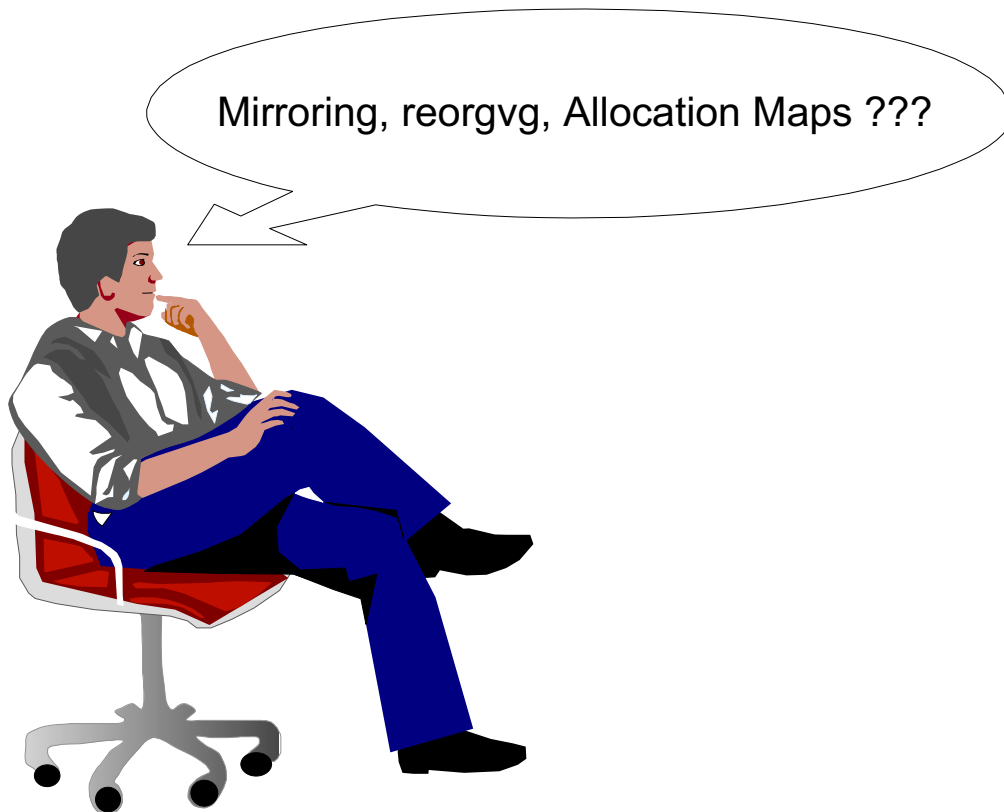
RAID 10 or 0+1 Mirroring and striping



---

## Let's Review: Basic LVM Tasks

---



© Copyright IBM Corporation 2004

Figure 5-10. Let's Review: Basic LVM Tasks

AU1612.0

### **Notes:**

On the next page you'll find a review activity where you will have to execute some basic LVM tasks.

The goal of this activity is to refresh important LVM terms.

# Review Activity: Basic LVM Tasks

Add a Logical Volume

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

<p>[TOP]</p> <p>Logical volume NAME</p> <p>VOLUME GROUP name</p> <p>Number of LOGICAL PARTITIONS</p> <p>PHYSICAL VOLUME names</p> <p>Logical Volume TYPE</p> <p>→ POSITION on physical volume</p> <p>RANGE of physical volumes</p> <p>MAXIMUM NUMBER of PHYSICAL VOLUMES to use for allocation</p> <p>→ Number of COPIES of each logical partition</p> <p>Mirror Write Consistency?</p> <p>Allocate each logical partition copy on a SEPARATE physical volume?</p> <p>...</p> <p>→ File containing ALLOCATION MAP</p>	<p>[Entry Fields]</p> <p>□</p> <p>rootvg</p> <p>□</p> <p>□</p> <p>□</p> <p>middle</p> <p>minimum</p> <p>□</p> <p>□</p> <p>active</p> <p>yes</p> <p>□</p>
---	--

© Copyright IBM Corporation 2004

Figure 5-11. Review Activity: Basic LVM Tasks

AU1612.0

## Notes:

In this activity you will execute basic LVM tasks. Do the following tasks without your instructor.

Only one person per machine can execute these commands.

- Using **smit mklv**, create a **mirrored logical volume** with the name **mirrorlv**. Make it two logical partitions in size.

Use **lslv -m** to identify the physical partitions that have been assigned to your logical partitions.

LP	PP1	PV1	PP2	PV2
0001				
0002				

Finally, remove the logical volume **mirrorlv**.

- Use **smit mklv** to create an unmirrored logical volume **lvtmp1** with a size of one partition. Choose an intraphysical policy where free partitions exist.

---

Use **lspv -p** to check where the partitions of **lvtmp1** reside.

- Using **smit chlv** change the intraphysical policy to another disk region. Have the partitions been moved to another region?

If not, use the **reorgvg** command. Use the man pages to identify how to reorganize a logical volume.

**Note: Do not reorganize the complete rootvg, because this takes too much time!**

Write down the command you used:

---

Using **lspv -p** check where the partitions of **lvtmp1** reside now.

Finally remove the logical volume **lvtmp1**.

- Find two free partitions on a disk. Write down the partition numbers:
- 

Create a logical volume **lvtmp2** that uses an **allocation map**. The logical volume should have a size of two partitions and should use the two partitions you identified before. Here is an example for an allocation map:

```
hdisk1:1-2
```

After creating the logical volume, check where the partitions reside.

Finally remove **lvtmp2**.

- What is the maximum number of disks in a volume group that would be created by the following command?

```
# mkvg -B -t 4 -y homevg hdisk11 hdisk99
```

---

## Review Activity Hints

1. Use these values with **smit mklv**:

```
Logical Volume NAME           mirrorlv
Number of LOGICAL PARTITIONS   2
Number of COPIES of each logical
partition                       2
```

[Allocate each logical partition copy  
on a SEPARATE physical volume]\*\*

\*\*You may need to set this to **no** if you only have one physical volume in your volume group.

To see the partitions:

**lslv -m mirrorlv**

To remove the logical volume:

**rmlv mirrorlv**

2. Use these values with **smit mklv**:

```
Logical Volume NAME           lvtmp1
Number of LOGICAL PARTITIONS   1
POSITION on physical volume     ***
```

\*\*\*Select a region that is available. You determined this with **lspv -p hdisk0**.

To check the position of **lvtmp1**:

**lspv -p hdiskX**

3. Change the value of POSITION on physical volume. Use **smit chlv**.

Did the partitions move?

**lspv -p hdiskX**

Reorganize the logical volume:

**reorgvg rootvg lvtmp1**

4. To create the logical volume using an allocation map:

Create the map file:

**vi /tmp/lvtmp2map**

Add the free partitions that you identified into the allocation file. For example,

**hdisk0:22-23**

Next, use **smit mklv** and modify the screen to use your map file:

File containing ALLOCATION MAP /tmp/lvtmp2map

5. Maximum amount of disks: **32**

## 5.2 LVM Data Representation

# LVM Identifiers

Goal: Unique worldwide identifiers for

- Hard disks
- Volume Groups (including logical volumes)

```
# lsvg rootvg
VOLUME GROUP:rootvg VG IDENTIFIER:00008371c98a229d4c0000000000000e
                                                                    (32 Bytes long)

# lspv
hdisk0      00008371b5969c35      rootvg
                                                                    (32 Bytes long)

# lslv hd4
LOGICAL VOLUME:  hd4      VOLUME GROUP: rootvg
LV IDENTIFIER: 00008371c98a229d4c0000000000000e.4
                                                                    (VGID.Minor Number)

# uname -m
000083714C00
```

© Copyright IBM Corporation 2004

Figure 5-12. LVM Identifiers

AU1612.0

## Notes:

The LVM uses identifiers for disks, volume groups, and logical volumes. As volume groups could be exported and imported between systems, these identifiers must be unique worldwide.

The volume groups identifiers (VGID) have a length of 32 bytes.

Hard disk identifiers have a length of 32 bytes, but currently the last 16 bytes are unused and are all set to 0 in the ODM.

If you ever have to manually update the disk identifiers in the ODM, do not forget to add 16 zeros to the physical volume ID.

The logical volume identifiers consist of the volume group identifier, a period and the minor number of the logical volume.

All identifiers are based on the CPU ID of the creating host and a timestamp.

# LVM Data on Disk Control Blocks

## Volume Group Descriptor Area (VGDA)

- Most important data structure of LVM
- Global to the volume group (same on each disk)
- One or two copies per disk

## Volume Group Status Area (VGSA)

- Tracks the state of mirrored copies
- One or two copies per disk

## Logical Volume Control Blocks (LVCB)

- First 512 bytes of each logical volume
- Contains LV attributes (Policies, Number of copies)
- Should not be overwritten by applications using raw devices!

© Copyright IBM Corporation 2004

Figure 5-13. LVM Data on Disk Control Blocks

AU1612.0

### Notes:

The LVM uses three different disk control blocks.

1. The **Volume Group Descriptor Area (VGDA)** is the most important data structure of the LVM. It is kept redundant on each disk that is contained in a volume group. Each disk contains the complete allocation information of the entire volume group.
2. The **Volume Group Status Area (VGSA)** is always present, but is only used when mirroring has been setup. It tracks the state of the mirrored copies, that means whether the copies are synchronized or **stale**.
3. The **Logical Volume Control Blocks (LVCB)** resides at the first 512 bytes of each logical volume. If raw devices are used (for example, many database systems use raw logical volumes), be careful that these programs do not destroy the LVCB.

The VGSA for scalable VGs consists of three areas: PV missing area (PVMA), MWC dirty bit area (MWC\_DBA), and PP status area (PPSA).

PV missing area

PVMA tracks if any of the disks are missing

MWC dirty bit area

MWC\_DBA holds the status for each LV if passive mirror write consistence is used

PP status area

PPSA logs any stale PPs

The overall size reserved for the VGSA is independent of the configuration parameters of the scalable VG and stays constant. However, the size of the contained PPSA changes proportional to the configured maximum number of PPs.

The LVCB contains metadata about a logical volume. For standard VGs the LVCB resides in the first block of the user data within the LV. Big VGs keep additional LVCB information in the ondisk VGDA. With scalable VGs logical volume control information is no longer stored on the first user block of any LV. All relevant logical volume control information is kept in the VGDA as part of the LVCB information area and the LV entry area. So, no precautions have to be met when using raw logical volumes because there is no longer a need to preserve the information held by the first 512 bytes of the logical device.



# LVM Data in the Operating System

## Object Data Manager (ODM)

- Physical volumes, volume groups and logical volumes are represented as devices (Customized devices)
- CuDv, CuAt, CuDvDr, CuDep

## AIX Files

- /etc/vg/vgVGID      Handle to the VGDA copy in memory
- /dev/hdiskX        Special file for a disk
- /dev/VGname        Special file for administrative access to a VG
- /dev/LVname        Special file for a logical volume
- /etc/filesystems    Used by the mount command to associate LV name, JFS log and mount point

© Copyright IBM Corporation 2004

Figure 5-14. LVM Data in the Operating System

AU1612.0

### Notes:

Physical volumes, volume groups, and logical volumes are handled as devices in AIX. Every physical volume, volume group, and logical volume is defined in the customized object classes in the ODM.

Additionally, many AIX files contain LVM-related data.

The VGDA is always stored by the kernel in memory to increase performance. This technique is called a memory-mapped file. The handle is always a file in the /etc/vg directory. This filename always reflects the volume group identifier.

## Contents of the VGDA

<b>Header Time Stamp</b>	- Updated when VG is changed
<b>Physical Volume List</b>	- PVIDs only (no PV names) - VGDA count and PV state
<b>Logical Volume List</b>	- LVIDs and LV names - Number of copies
<b>Physical Partition Map</b>	- Maps LPs to PPs
<b>Trailer Time Stamp</b>	- Must contain same value as header time stamp

© Copyright IBM Corporation 2004

Figure 5-15. Contents of the VGDA

AU1612.0

### Notes:

This table shows the contents of the VGDA.

The time stamps are used to check if a VGDA is valid. If the system crashes while changing the VGDA the time stamps will differ. The next time when the volume group is varied on, this VGDA is marked as invalid. The latest intact VGDA will then be used to overwrite the other VGDA's in the volume group.

The VGDA contains the physical volume list. Note that no disk names are stored, only the unique disk identifiers are used. For each disk the number of VGDA's on the disk and the physical volume state is stored. We talk about physical volume states later in this unit.

The VGDA contains the logical volumes that are part of the volume group. It stores the LV identifiers and the corresponding logical volume names. Additionally, the number of copies is stored for each LV.

The most important data structure is the physical partition map. It maps each logical partition to a physical partition. The size of the physical partition map is determined at volume group creation time (depending on the number of disks that can be in the volume group, specified by **mkvg -d**). This size is a hard limit when trying to extend the volume group.

# VGDA Example

```
# lqueryvg -p hdisk1 -At
```

```
Max LVs:          256
PP Size:          24      → 1:

Free PPs:         56
LV count:         3      → 2:
PV count:         2      → 3:

Total VGDA:       3      → 4:

MAX PPs per       1016
MAX PVs:          32

Logical:          00008371387fa8bb0000ce0001390000.1   lv_01  1
                  00008371387fa8bb0000ce0001390000.2   lv_02  1
                  00008371387fa8bb0000ce0001390000.3   lv_03  1
```

```
Physical:         00008371b5969c35      2      0
                  00008371b7866c77      1      0

6:                7:
```

© Copyright IBM Corporation 2004

Figure 5-16. VGDA Example

AU1612.0

## Notes:

The command **lqueryvg** is a low-level command that shows an extract from the VGDA on a disk, for example **hdisk1**. As you notice, the visual is not complete. Use the following unordered expressions and try to put each expression to the corresponding number in the picture.

- VGDA count on disk
- 3 VGDA in VG
- 3 LVs in VG
- PP size = 16 MB
- Quorum check on
- LVIDs (VGID.minor\_number)
- 2 PVs in VG
- PVIDs

The **lqueryvg** on newer AIX versions show more information. An AIX 5.3 might show the following:

Max LVs: 256  
PP Size: 26  
Free PPs: 464  
LV count: 11  
PV count: 1  
Total VGDA: 2  
Conc Allowed: 0  
MAX PPs per PV 1016  
MAX PVs: 32  
Conc Autovaryo 0  
Varied on Conc 0  
Logical: 00096baa00004c00000000ffdb801c14.1 hd5 1  
00096baa00004c00000000ffdb801c14.2 hd6 1  
00096baa00004c00000000ffdb801c14.3 hd8 1  
00096baa00004c00000000ffdb801c14.4 hd4 1  
00096baa00004c00000000ffdb801c14.5 hd2 1  
00096baa00004c00000000ffdb801c14.6 hd9var 1  
00096baa00004c00000000ffdb801c14.7 hd3 1  
00096baa00004c00000000ffdb801c14.8 hd1 1  
00096baa00004c00000000ffdb801c14.9 hd10opt 1  
00096baa00004c00000000ffdb801c14.10 loglv00 1  
00096baa00004c00000000ffdb801c14.11 fslv00 1  
Physical: 00096baa1ec9fa18 2 0  
Total PPs: 542  
LTG size: 128  
HOT SPARE: 0  
AUTO SYNC: 0  
VG PERMISSION: 0  
SNAPSHOT VG: 0  
IS\_PRIMARY VG: 0  
PSNFSTPP: 4352  
VARYON MODE: 0  
VG Type: 0  
Max PPs: 32512

# The Logical Volume Control Block (LVCB)

```
# getlvcb -AT hd2

AIX LVCB
intrapolicy = c
copies = 1
interpolicy = m
lvid = 0009301300004c00000000e63a42b585.5
lvname = hd2
label = /usr
machine id = 010193100
number lps = 103
relocatable = y
strict = y
stripe width = 0
stripe size in exponent = 0
type = jfs
upperbound = 32
fs = log=/dev/hd8:mount=automatic:type=bootfs:vol=/usr:free=false
time created = Mon Jan 19 14:20:27 2003
time modified = Fri Feb 14 10:18:46 2003
```

© Copyright IBM Corporation 2004

Figure 5-17. The Logical Volume Control Block (LVCB)

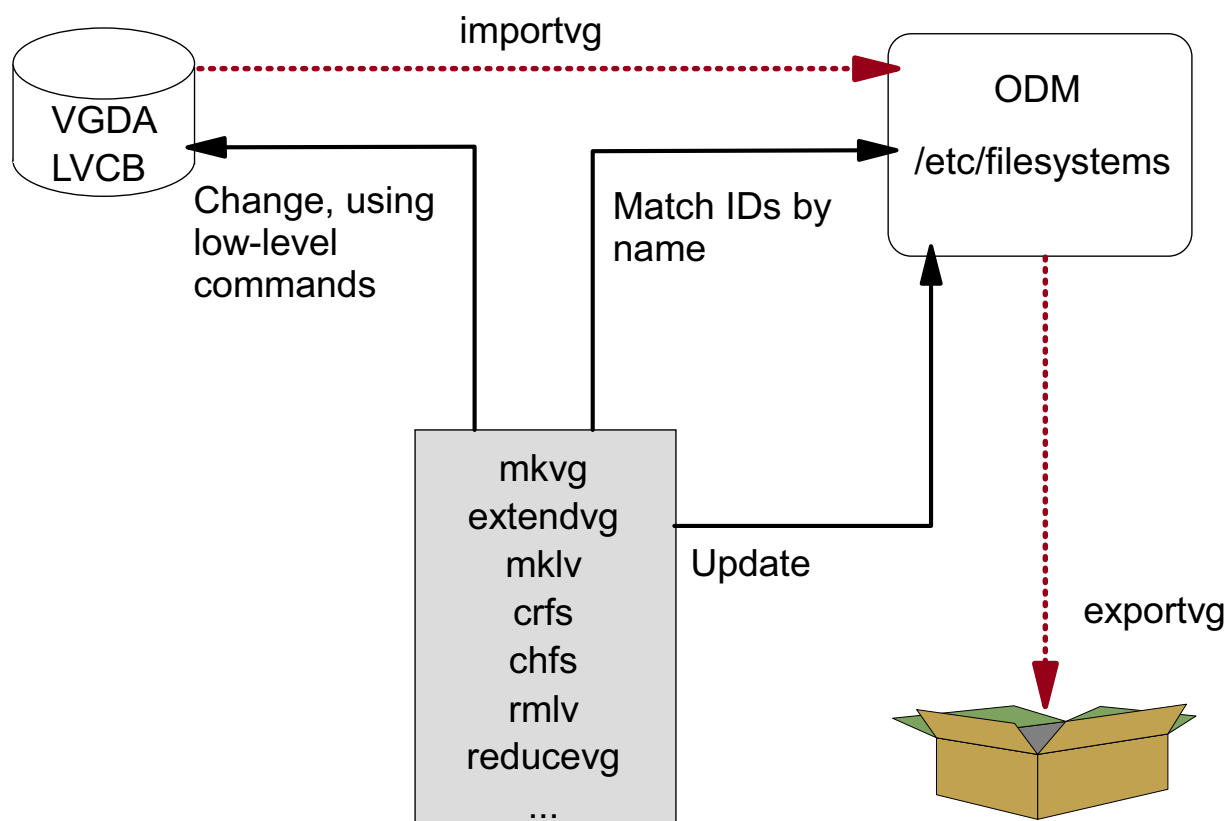
AU1612.0

## Notes:

The LVCB stores attributes of a logical volume. The command **getlvcb** queries an LVCB, for example the logical volume **hd2**. For example:

- Intrapolicy (c = Center)
- Number of copies (1 = No mirroring)
- Interpolicy (m = Minimum)
- LVID
- LV name (hd2)
- Number of logical partitions (103)
- Can the partitions be reorganized? (relocatable = y)
- Each mirror copy on a separate disk (strict = y)
- Number of disks involved in striping (stripe width)
- Stripe size
- Logical volume type (type = jfs)
- JFS file system information
- Creation and last update time

## How LVM Interacts with ODM and VGDA



© Copyright IBM Corporation 2004

Figure 5-18. How LVM Interacts with ODM and VGDA

AU1612.0

### Notes:

Most of the LVM commands that are used when working with volume groups, physical or logical volumes are high-level commands. These high-level commands (like **mkvg**, **extendvg**, **mklv**) are implemented as shell scripts and use names to reference a certain LVM object. To match a name, for example `rootvg` or `hdisk0`, to an identifier the ODM is consulted.

The high-level commands call intermediate or low-level commands that query or change the disk control blocks VGDA or LVCB. Additionally, the ODM has to be updated; for example, to add a new logical volume. The high-level commands contain signal handlers to clean up the configuration if the program is stopped abnormally. If a system crashes, or if high-level commands are stopped by **kill -9**, the system can end up in a situation where the VGDA/LVCB and the ODM are not in sync. The same situation may occur when low-level commands are used incorrectly.

This page shows two very important commands that are explained in detail later. The command **importvg** imports a complete new volume group based on a VGDA and LVCB on a disk. The command **exportvg** removes a complete volume group from the ODM.

## ODM Entries for Physical Volumes (1 of 3)

```
# odmget -q "name like hdisk?" CuDv
```

```
CuDv:
```

```
name = "hdisk0"  
status = 1  
chgstatus = 2  
ddins = "scdisk"  
location = "04-C0-00-2,0"  
parent = "scsi0"  
connwhere = "2,0"  
PdDvLn = "disk/scsi/scsd"
```

```
CuDv:
```

```
name = "hdisk1"  
status = 1  
chgstatus = 2  
ddins = "scdisk"  
location = "04-C0-00-3,0"  
parent = "scsi0"  
connwhere = "3,0"  
PdDvLn = "disk/scsi/scsd"
```

© Copyright IBM Corporation 2004

Figure 5-19. ODM Entries for Physical Volumes (1 of 3)

AU1612.0

### Notes:

All physical volumes are stored in **CuDv**.

Remember the most important attributes:

- status = 1 means the disk is available
- chgstatus = 2 means the status has not changed since last reboot
- location specifies the location code of the device
- parent specifies the parent device



---

## ODM Entries for Physical Volumes (2 of 3)

---

```
# odmget -q "name=hdisk0 and attribute=pvid" CuAt
CuAt:
    name = "hdisk0"
    attribute = "pvid"
    value = "250000010700040b000c0d0000000000"
    type = "R"
    generic = "D"
    rep = "s"
    nls_index = 2
```

© Copyright IBM Corporation 2004

Figure 5-20. ODM Entries for Physical Volumes (2 of 3)

AU1612.0

### **Notes:**

The disk's most important attribute is the PVID.

The PVID has a length of 32 bytes, where the last 16 bytes are set to zeros in the ODM. Whenever you must manually update a PVID in the ODM you must specify the complete 32-byte PVID of the disk.

Other attributes (for example, SCSI command queue depth, timeout values) may occur in **CuAt**.

## ODM Entries for Physical Volumes (3 of 3)

```
# odmget -q "value3 like hdisk?" CuDvDr
```

```
CuDvDr:
  resource = "devno"
  value1 = "22"
  value2 = "1"
  value3 = "hdisk0"
```

```
CuDvDr:
  resource = "devno"
  value1 = "22"
  value2 = "2"
  value3 = "hdisk1"
```

```
# ls -l /dev/hdisk*
```

```
brw----- 1 root system 22, 1 08 Jan 06:56 /dev/hdisk0
brw----- 1 root system 22, 2 08 Jan 07:12 /dev/hdisk1
```

© Copyright IBM Corporation 2004

Figure 5-21. ODM Entries for Physical Volumes (3 of 3)

AU1612.0

### Notes:

The ODM class **CuDvDr** is used to store the major and minor numbers of the devices. Applications or system programs use the special files to access a certain device.

## ODM Entries for Volume Groups (1 of 2)

```
# odmget -q "name=rootvg" CuDv
CuDv:
    name = "rootvg"
    status = 0
    chgstatus = 1
    ddins = ""
    location = ""
    parent = ""
    connwhere = ""
    PdDvLn = "logical_volume/vgsubclass/vgtype"

# odmget -q "name=rootvg" CuAt
CuAt:
    name = "rootvg"
    attribute = "vgserial_id"
    value = "0009301300004c00000000e63a42b585"
    type = "R"
    generic = "D"
    rep = "n"
    nls_index = 637                                (continues on next page)
```

Figure 5-22. ODM Entries for Volume Groups (1 of 2)

AU1612.0

### Notes:

The existence of a volume group is stored in **CuDv**, that means all volume groups must have an object in this class.

One of the most important pieces of information is the VGID, which is stored in **CuAt**.

All disks that belong to a volume group are stored in **CuAt**. That's shown on the next page.

## ODM Entries for Volume Groups (2 of 2)

```
# odmget -q "name=rootvg" CuAt
...

CuAt:
    name = "rootvg"
    attribute = "timestamp"
    value = "3ec3cb943749cbc3"
    type = "R"
    generic = "DU"
    rep = "s"
    nls_index = 0

CuAt:
    name = "rootvg"
    attribute = "pv"
    value = "00008371d11226670000000000000000"
    type = "R"
    generic = ""
    rep = "sl"
    nls_index = 0
```

© Copyright IBM Corporation 2004

Figure 5-23. ODM Entries for Volume Groups (2 of 2)

AU1612.0

### Notes:

All disks that belong to a volume group are stored in CuAt.

Remember that the PVID is a 32-number field, where the last 16 numbers are set to zeros.

## ODM Entries for Logical Volumes (1 of 2)

```
# odmget -q "name=hd2" CuDv
```

```
CuDv:
```

```
name = "hd2"
status = 0
chgstatus = 1
ddins = ""
location = ""
parent = "rootvg"
connwhere = ""
PdDvLn = "logical_volume/lvsubclass/lvtype"
```

```
# odmget -q "name=hd2" CuAt
```

```
CuAt:
```

```
name = "hd2"
attribute = "lvserial_id" (intra, stripe_width, size, label ...)
value = "0009301300004c00000000e63a42b585.5"
type = "R"
generic = "D"
rep = "n"
nls_index = 648
```

© Copyright IBM Corporation 2004

Figure 5-24. ODM Entries for Logical Volumes (1 of 2)

AU1612.0

### Notes:

All logical volumes are stored in the object class **CuDv**.

Attributes of a logical volume, for example its **LVID**, are stored in the object class **CuAt**. Other attributes that belong to a logical volume are the intra-policy, stripe\_width or the size.

## ODM Entries for Logical Volumes (2 of 2)

```
# odmget -q "value3=hd2" CuDvDr
CuDvDr:
    resource = "devno"
    value1 = "10"
    value2 = "5"
    value3 = "hd2"

# ls -l /dev/hd2
brw----- 1 root system 10, 5 08 Jan 06:56 /dev/hd2

# odmget -q "dependency=hd2" CuDep
CuDep:
    name = "rootvg"
    dependency = "hd2"
```

© Copyright IBM Corporation 2004

Figure 5-25. ODM Entries for Logical Volumes (2 of 2)

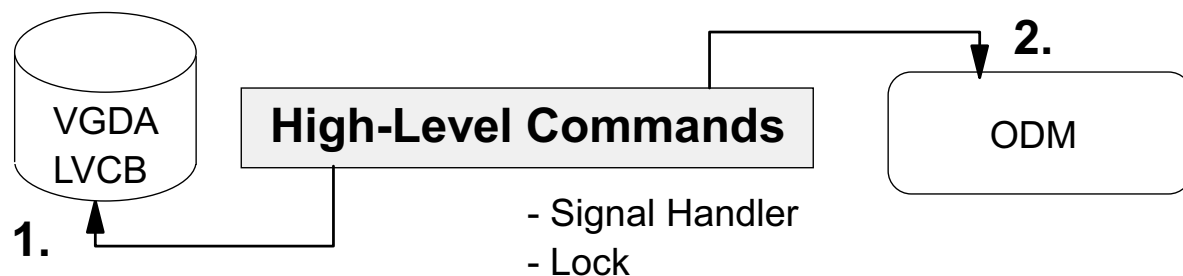
AU1612.0

### Notes:

All logical volumes have an object in **CuDvDr** that is used to create the special file entries in **/dev**.

The ODM class **CuDep** (customized dependencies) stores dependency information for software devices, for example, the logical volume **hd2** is contained in the **rootvg** volume group.

## ODM-Related LVM Problems



### What can cause problems ?

- kill -9, shutdown, system crash
- Improper use of low-level commands
- Hardware changes without or with wrong software actions

© Copyright IBM Corporation 2004

Figure 5-26. ODM-Related LVM Problems

AU1612.0

### Notes:

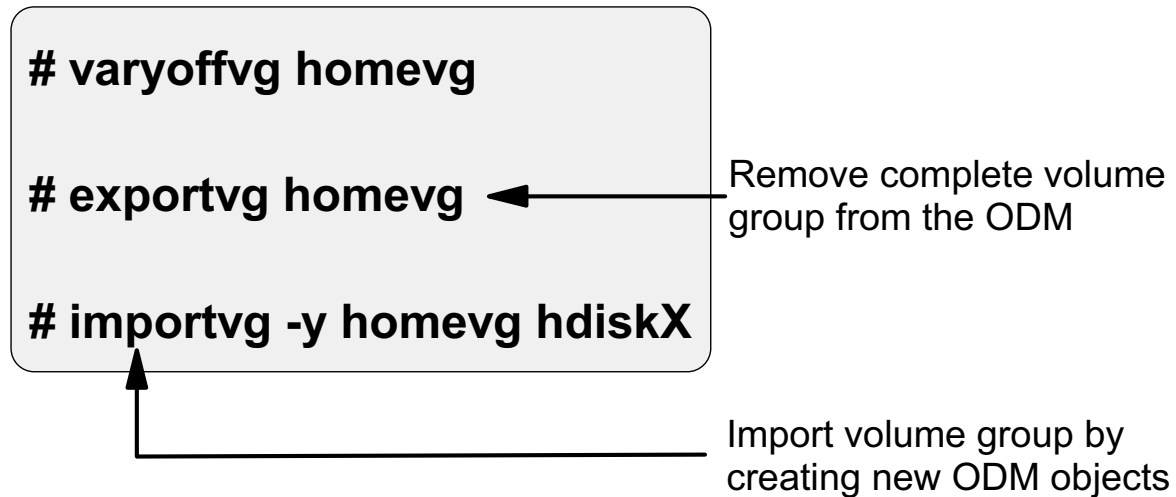
As already mentioned, most of the time administrators use high-level commands to create or update volume groups or logical volumes. These commands use signal handlers to set up a proper cleanup in case of an interruption. Additionally, LVM commands create a locking mechanism to block other commands while a change is in progress.

These signal handlers do not work with a **kill -9**, a system shutdown, or a system crash. You might end up in a situation where the VGDA has been updated, but the change has not been stored in the ODM.

The same situation might come up by the improper use of low-level commands or hardware changes that are not followed by correct administrator actions.

## Fixing ODM Problems (1 of 2)

If the ODM problem is **not in the rootvg**, for example in volume group **homevg**, do the following:



© Copyright IBM Corporation 2004

Figure 5-27. Fixing ODM Problems (1 of 2)

AU1612.0

### Notes:

If you detect ODM problems you must identify whether the volume group is the **rootvg** or not.

Because the **rootvg** cannot be varied off, this procedure applies only to non-rootvg volume groups.

1. In the first step, you vary off the volume group, which requires that all file systems must be unmounted first. To vary off a volume group, use the **varyoffvg** command.
2. In the next step, you export the volume group by using the **exportvg** command. This command removes the complete volume group from the ODM. The VGDA and LVCB are not touched by **exportvg**.
3. In the last step, you import the volume group by using the **importvg** command. Specify the volume group name with option **-y**, otherwise AIX creates a new volume group name.



You need to specify only one intact physical volume of the volume group that you import. The **importvg** command reads the VGDA and LVCB on that disk and creates completely new ODM objects.

We will return to the export and import functions later in this course.

## Fixing ODM Problems (2 of 2)

If the ODM problem is in the **rootvg**, use **rvgrecover**:

```
PV=hdisk0
VG=rootvg
cp /etc/objrepos/CuAt /etc/objrepos/CuAt.$$
cp /etc/objrepos/CuDep /etc/objrepos/CuDep.$$
cp /etc/objrepos/CuDv /etc/objrepos/CuDv.$$
cp /etc/objrepos/CuDvDr /etc/objrepos/CuDvDr.$$
lqueryvg -Lp $PV | awk '{print $2}' | while read LVname;
do
    odmdelete -q "name=$LVname" -o CuAt
    odmdelete -q "name=$LVname" -o CuDv
    odmdelete -q "value3=$LVname" -o CuDvDr
done
odmdelete -q "name=$VG" -o CuAt
odmdelete -q "parent=$VG" -o CuDv
odmdelete -q "name=$VG" -o CuDv
odmdelete -q "name=$VG" -o CuDep
odmdelete -q "dependency=$VG" -o CuDep
odmdelete -q "value1=10" -o CuDvDr
odmdelete -q "value3=$VG" -o CuDvDr
importvg -y $VG $PV    # ignore lvaryoffvg errors
varyonvg $VG
```

- Export rootvg by odmdeletes
- Import rootvg by importvg

© Copyright IBM Corporation 2004

Figure 5-28. Fixing ODM Problems (2 of 2)

AU1612.0

### Notes:

If you detect ODM problems in **rootvg** use the shell script **rvgrecover**. This procedure is described in the *AIX 4.3 Problem Solving Guide and Reference*. Create this script in **/bin** and mark it executable.

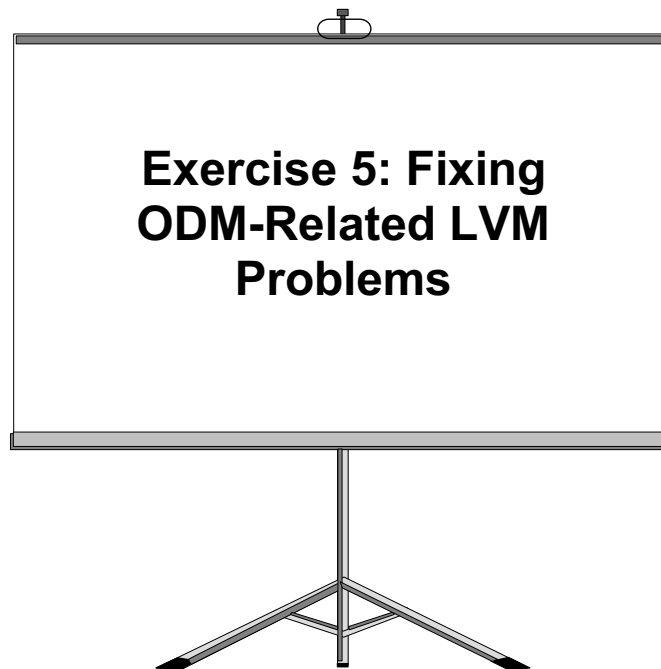
The script **rvgrecover** removes all ODM entries that belong to your **rootvg** by using **odmdelete**. That's the same way **exportvg** works.

After deleting all ODM objects from **rootvg** it imports the **rootvg** by reading the VGDA and LVCB from the boot disk. This results in completely new ODM objects that describe your **rootvg**.

---

## Next Step

---



© Copyright IBM Corporation 2004

Figure 5-29. Next Step

AU1612.0

### **Notes:**

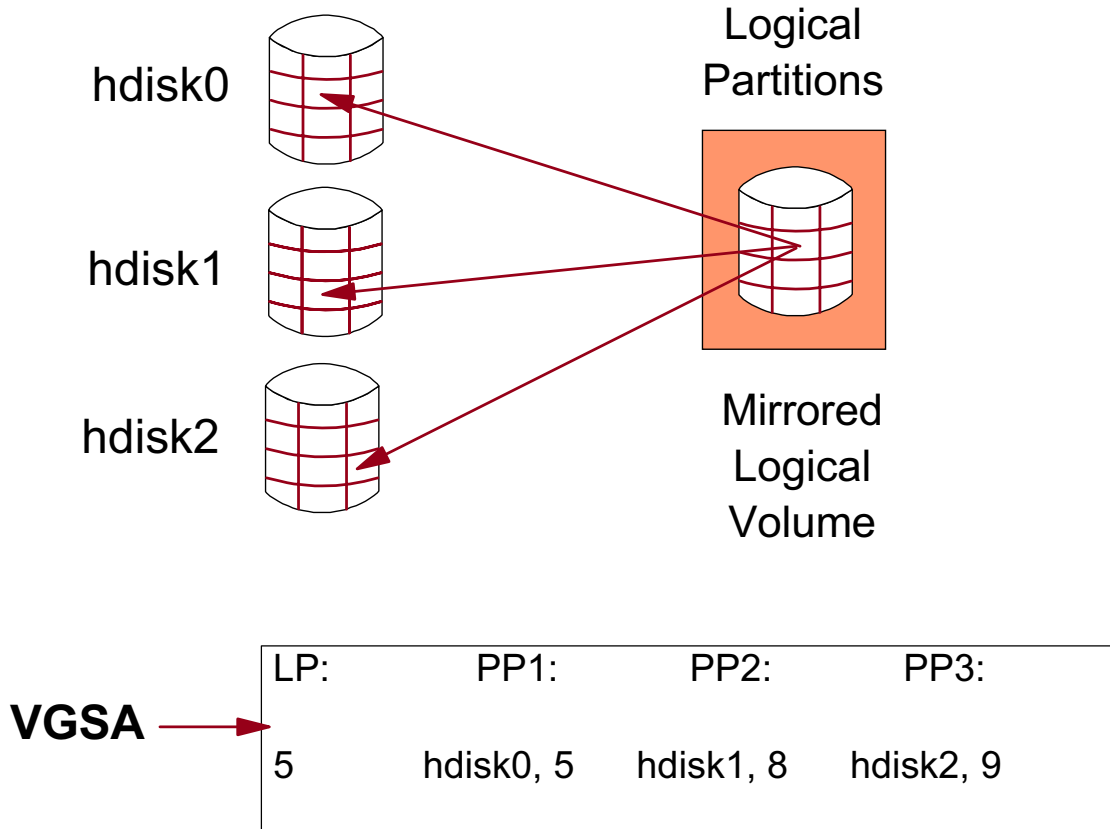
At the end of this exercise, you should be able to:

- Analyze an LVM-related ODM problem
- Fix an LVM-related ODM problem associated with the rootvg



## 5.3 Mirroring and Quorum

# Mirroring



© Copyright IBM Corporation 2004

Figure 5-30. Mirroring

AU1612.0

## Notes:

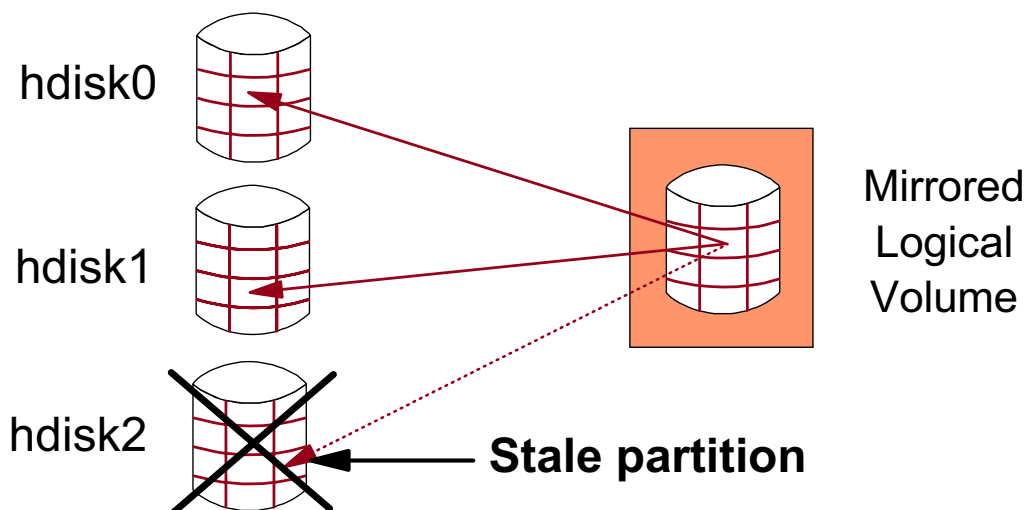
This page shows a mirrored logical volume, where each logical partition is mirrored to three physical partitions. More than three copies are not possible.

If one of the disks fails, there are at least two copies of the data available. That means mirroring is used to increase the availability of a system or a logical volume.

The information about the mirrored partitions is stored in the **VGSA (Volume Group Status Area)**, which is contained on each disk. In the example, we see logical partition 5 points to physical partition 5 on hdisk0, physical partition 8 on hdisk1 and physical partition 9 on hdisk2.

In AIX 4.1/4.2 the maximum number of mirrored partitions on a disk was 1016. AIX 4.3 and subsequent releases allow more than 1016 mirrored partitions on a disk. This maximum depends on the number of disks that can reside in the volume group.

## Stale Partitions



After repair of hdisk2:

- **varyonvg VGName** (calls `syncvg -v VGName`)
- Only stale partitions are updated

© Copyright IBM Corporation 2004

Figure 5-31. Stale Partitions

AU1612.0

### Notes:

If a disk failure occurs, for example **hdisk2** fails, which contains a mirrored logical volume, the data on the failed disk becomes **stale**.

The state information is kept per physical partition. A physical volume is shown as stale (**lsvg VGName**), as long as it has one stale partition.

If the disk has been repaired (for example after a power failure), you should issue the **varyonvg** command which starts the **syncvg** command to synchronize the stale partitions. The **syncvg** command is started as a background job that updates all stale partitions from the volume group.

Always use the **varyonvg** command to update stale partitions. After a power failure, a disk forgets its reservation. The **syncvg** command cannot reestablish the reservation, whereas **varyonvg** does this before calling **syncvg**. The term *reservation* means that a disk is reserved for one system. The disk driver puts the disk in a state where you can work with the disk (at the same time the control LED of the disk turns on).

**varyonvg** works if the volume group is already varied on or if the volume group is the **rootvg**.



## Creating Mirrored LVs (smit mklv)

Add a Logical Volume	
Type or select values in entry fields. Press Enter AFTER making all desired changes.	
[TOP]	[Entry Fields]
Logical volume NAME	[lv01]
VOLUME GROUP name	rootvg
Number of LOGICAL PARTITIONS	[50]
PHYSICAL VOLUME names	[hdisk2 hdisk4]
Logical Volume TYPE	[]
POSITION on physical volume	edge
RANGE of physical volumes	minimum
MAXIMUM NUMBER of PHYSICAL VOLUMES to use for allocation	[]
Number of COPIES of each logical partition	[2]
Mirror Write Consistency?	active
Allocate each logical partition copy on a SEPARATE physical volume?	yes
...	
SCHEDULING POLICY for reading/writing logical partition copies	parallel

© Copyright IBM Corporation 2004

Figure 5-32. Creating Mirrored LVs (smit mklv)

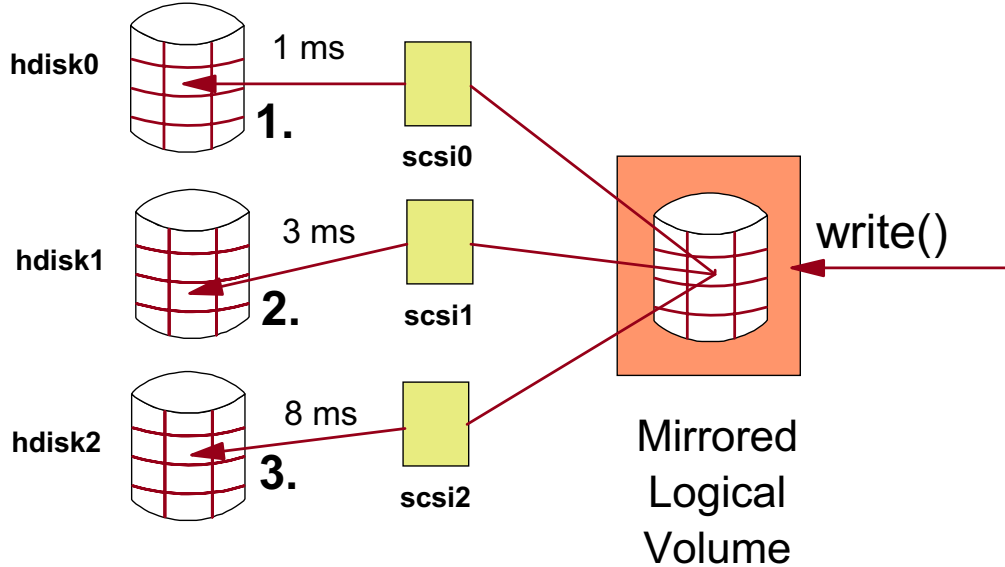
AU1612.0

### Notes:

A very easy way to create a mirrored logical volume is to use the smit fastpath **mklv**.

- Specify the logical volume name, for example **lv01**.
- Specify the number of logical partitions, for example 50.
- Specify the disks where the physical partitions reside. If you want mirroring on separate adapters, choose disk names that reside on different adapters.
- Specify the number of copies, for example two for a single mirror or three for a double mirror.
- Do not change the default entry for **Allocate each logical partition copy on a SEPARATE physical volume**, which is **yes**. Otherwise you would mirror on the same disk, which makes no sense. If you leave the default entry of **yes** and no separate disk is available, **mklvcopy** will fail.
- The terms **Mirror Write Consistency** and **Scheduling Policy** are explained on the next page.

# Scheduling Policies: Sequential



- Second physical write operation is not started unless the first has completed successfully
- In case of a total disk failure there is always a "good copy"
- Increases availability, but decreases performance
- In this example the write operation takes 12 ms

© Copyright IBM Corporation 2004

Figure 5-33. Scheduling Policies: Sequential

AU1612.0

## Notes:

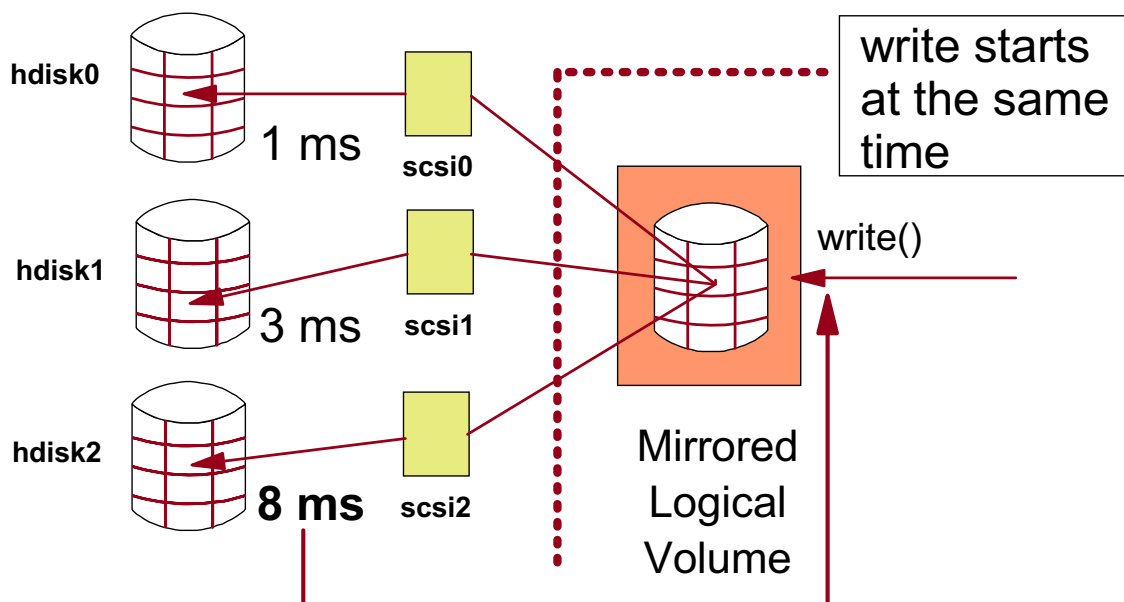
The sequential scheduling performs writes to multiple copies in order. The multiple physical partitions representing the mirrored copies of a single logical partition are designated primary, secondary, and tertiary.

In sequential scheduling, the physical partitions are written to in sequence; the system waits for the write operation for one physical partition to complete before starting the write operation for the next one.

The write()-operation of the application must wait until all three partitions are written to the disk. This decreases the performance but increases availability. In case of a total disk failure (for example, due to a power loss), there will always be a good copy.

For read operations on mirrored logical volumes with a sequential scheduling policy, only the primary copy is read. If that read operation is unsuccessful, the next copy is read. During the read-retry operation on the next copy, the failed primary copy is corrected by the LVM with a hardware relocation. Thus, the bad block that prevented the first read from completing is patched for future access.

## Scheduling Policies: Parallel



- Write operations for physical partitions starts at the same time: When the longest write (8 ms) finishes, the write operation is complete
- Improves performance (especially READ-Performance)

© Copyright IBM Corporation 2004

Figure 5-34. Scheduling Policies: Parallel

AU1612.0

### Notes:

The parallel scheduling policy starts the write operation to all copies at the same time. When the write operation that takes the longest to complete finishes (for example, the one that takes 8 milliseconds), the write() from the application completes.

Specifying mirrored logical volumes with a parallel scheduling policy may increase overall performance due to a common read/write ratio of 3:1 or 4:1. With sequential policy, the primary copy is always read; with parallel policy, the copy that's best reachable is used. On each read, the system checks whether the primary is busy. If it is not busy, the read is initiated on the primary. If the primary is busy, the system checks the secondary. If it is not busy, the read is initiated on the secondary. If the secondary is busy, the read is initiated on the copy with the least number of outstanding I/Os.

The parallel/sequential policy always initiates reads from the primary copy, but initiates writes concurrently.

The parallel/round-robin policy alternates reads between the copies. This results in equal utilization for reads even when there is more than one I/O outstanding at a time. Writes are performed concurrently.

A parallel policy offers the best performance if you mirror on separate adapters.

## Mirror Write Consistency (MWC)

### Problem:

- Parallel scheduling policy and ...
- ... system crashes **before the write to all mirrors** have been completed
- Mirrors of the logical volume are in an **inconsistent** state

### Solution: Mirror Write Consistency

- Allows identifying the correct physical partition after reboot
- Separate area of each disk (outer edge)
- Place logical volumes with mirror write consistency on the outer edger

© Copyright IBM Corporation 2004

Figure 5-35. Mirror Write Consistency (MWC)

AU1612.0

### Notes:

When working with parallel scheduling policy, the LVM starts the write operation for the physical partition at the same time. If a system crashes (for example, due to a power failure) **before the write to all mirrors** has been completed, the mirrors of the logical volume are in an inconsistent state.

To avoid this situation, always use **mirror write consistency** when working with parallel scheduling policy.

When the volume group is varied back online for use, this information is used to make logical partitions consistent again.

Active mirror write consistency is implemented as a cache on the disk and behaves similarly to the JFS and JFS2 log devices. The physical write operation proceeds when the MWC cache has been updated. The disk cache resides in the outer edge area. Therefore, always try to place a logical volume that uses active MWC in the same area as the MWC. This improves disk access times.

AIX 5L introduces the new **passive** option to the mirror write consistency (MWC) algorithm for mirrored logical volumes. This option only applies to big volume groups. Big volume groups allow up to 512 logical volumes and 128 physical volumes per volume group. Without the big volume group format up to 256 logical volumes and 32 physical volumes can exist within a volume group.

Passive MWC reduces the problem of having to update the MWC log on the disk. This method logs that the logical volume has been opened but does not log writes. If the system crashes, then the LVM starts a forced synchronization of the entire logical volume when the system restarts.

The following syntax is used with either the `mk1v` or `ch1v` command to set MWC options:

```
mk1v -w y|a|p|n
```

```
ch1v -w y|a|p|n
```

Here is a description of the MWC options:

Option	Description
y or a	Logical partitions that might be inconsistent if the system or the volume group is not shut down properly are identified. When the volume group is varied back online, this information is used to make logical partitions consistent.
p	The volume group logs that the logical volume has been opened. After a crash when the volume group is varied on, an automatic forced synchronization of the logical volume is started. Consistency is maintained while the forced synchronization is in progress by using a copy of the read recovery policy that propagates the blocks being read to the other mirrors in the logical volume.
n	<p>The mirrors of a mirrored logical volume can be left in an inconsistent state in the event of a system or volume group crash. There is no automatic protection of mirror consistency. Writes outstanding at the time of the crash can leave mirrors with inconsistent data the next time the volume group is varied on. After a crash, any mirrored logical volume that has MWC turned OFF should perform a forced synchronization before the data within the logical volume is used. For example,</p> <pre>syncvg -f -l LVname</pre> <p>An exception to forced synchronization is logical volumes whose content is only valid while the logical volume is open, such as paging spaces.</p>

## Adding Mirrors to Existing LVs (mklvcopy)

Add Copies to a Logical Volume	
Type or select values in entry fields. Press Enter AFTER making all desired changes.	
	[Entry Fields]
Logical volume NAME	[hd2]
NEW TOTAL number of logical partition copies	<b>2</b>
PHYSICAL VOLUME names	<b>[hdisk1]</b>
POSITION on physical volume	<b>edge</b>
RANGE of physical volumes	minimum
MAXIMUM NUMBER of PHYSICAL VOLUMES to use for allocation	[32]
Allocate each logical partition copy on a SEPARATE physical volume?	<b>yes</b>
File containing ALLOCATION MAP	[]
SYNCHRONIZE the data in the new logical partition copies?	<b>no</b>

© Copyright IBM Corporation 2004

Figure 5-36. Adding Mirrors to Existing LVs (mklvcopy)

AU1612.0

### Notes:

Using the **mklvcopy** command or the smit fastpath **smit mklvcopy** you can add mirrors to existing logical volumes. You need to specify the new total number of logical partition copies and the disks where the physical partitions reside. If you work with active MWC, use **edge** as the position policy to increase performance.

If there are many LVs to synchronize it's better not to synchronize the new copies immediately after the creation (that's the default).

Here are some examples for the **mklvcopy** command:

1. Add a copy for logical volume **lv01** on disk **hdisk7**:  
**# mklvcopy lv01 2 hdisk7**
2. **Add a copy for logical volume lv02** on disk **hdisk4**. The copies should reside in the outer edge area. The synchronization will be done immediately:  
**# mklvcopy -a e -k lv02 2 hdisk4**

**To remove copies from a logical volume use rmlvcopy** or the smit fastpath **smit rmlvcopy**.

# Mirroring rootvg



1. extendvg
2. chvg -Qn
3. mirrorvg -s
4. syncvg -v

5. bosboot -a
6. bootlist
7. shutdown -Fr
8. bootinfo -b

- Make a copy of all rootvg LVs via **mirrorvg** and place copies on the second disk
- Execute **bosboot** and change your **bootlist**

© Copyright IBM Corporation 2004

Figure 5-37. Mirroring rootvg

AU1612.0

## Notes:

What is the reason to mirror the rootvg?

If your rootvg is on one disk, you get a **single point of failure**; that means, if this disk fails, your machine is not available any longer.

If you mirror rootvg to a second (or third) disk, and one disk fails, there will be another disk that contains the mirrored rootvg. You increase the availability of your system.

The following steps show how to mirror the rootvg.

- Add the new disk to the volume group (for example, **hdisk1**):  

```
# extendvg [ -f ] rootvg hdisk1
```
- If you use one mirror disk, be sure that a quorum is not required for varyon:  

```
# chvg -Qn rootvg
```
- **Add the mirrors for all rootvg logical volumes:**



```
# mklvcopy hd1 2 hdisk1
# mklvcopy hd2 2 hdisk1
# mklvcopy hd3 2 hdisk1
# mklvcopy hd4 2 hdisk1
# mklvcopy hd5 2 hdisk1
# mklvcopy hd6 2 hdisk1
# mklvcopy hd8 2 hdisk1
# mklvcopy hd9var 2 hdisk1
# mklvcopy hd10opt 2 hdisk1
OR better
# mirrorvg -s rootvg
```

If you have other logical volumes in your rootvg, be sure to create copies for them as well.

An alternative to running multiple **mklvcopy** commands is to use **mirrorvg**. This command was added in version 4.2 to simplify mirroring VGs. The **mirrorvg** command by default will disable quorum and mirror the existing LVs in the specified VG. To mirror rootvg, use the command:

```
mirrorvg -s rootvg
```

- Now synchronize the new copies you created:

```
# syncvg -v rootvg
```

- As we want to be able to boot from different disks, we need to do a bosboot:

```
# bosboot -a
```

As **hd5** is mirrored, there is no need to do it for each disk.

- Update the **boot list**. In case of a disk failure we must be able to boot from different disks.

```
# bootlist -m normal hdisk1 hdisk0
```

```
# bootlist -m service hdisk1 hdisk0
```

- Reboot the system because we disabled the quorum to take effect

```
# shutdown -Fr
```

- Check that the system boots from the first boot disk.

```
# bootinfo -b
```

# Mirroring Volume Groups (mirrorvg)

## Mirror a Volume Group

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

VOLUME GROUP name	[Entry Fields] rootvg
Mirror sync mode	[Foreground]
PHYSICAL VOLUME names	[hdisk1]
Number of COPIES of each logical partition	2
Keep Quorum Checking On?	no
Create Exact LV Mapping?	no

For rootvg, you need to execute:

- **bosboot**
- **bootlist -m normal ...**

© Copyright IBM Corporation 2004

Figure 5-38. Mirroring Volume Groups (mirrorvg)

AU1612.0

### Notes:

Another way to mirror a volume group is to use the **mirrorvg** command or the smit fastpath **smit mirrorvg**.

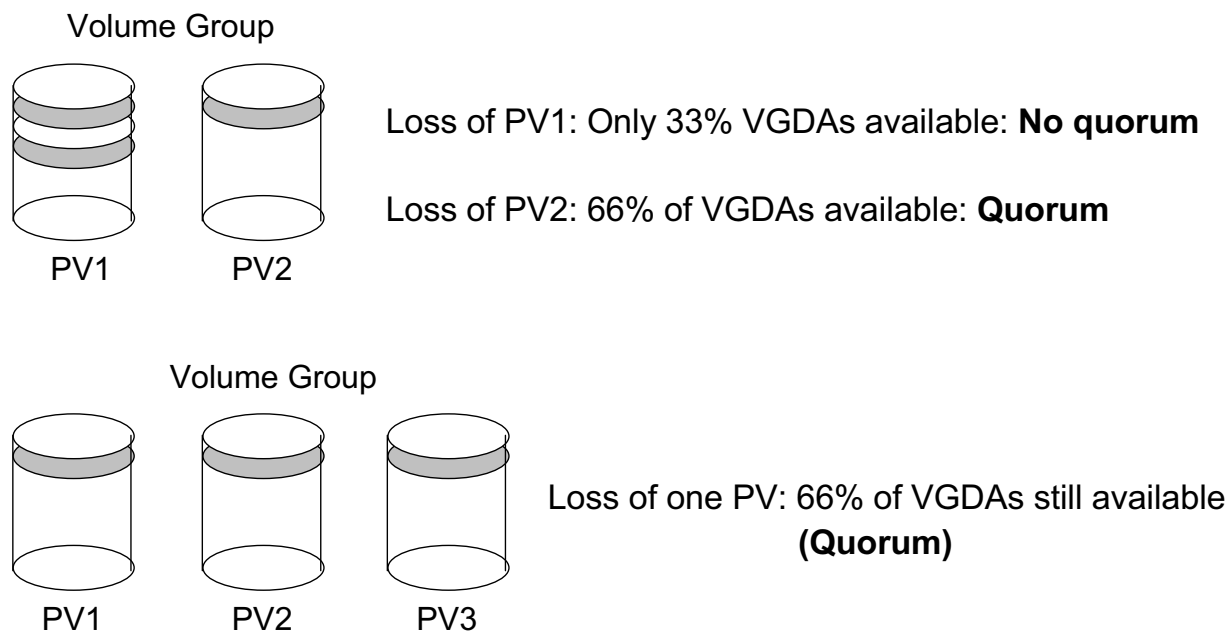
**Note:** If you mirror the rootvg with the **mirrorvg** command you need to execute a **bosboot** afterwards. Additionally, you need to change your **boot list**.

The **mirrorvg** command was introduced with AIX 4.2.1.

The opposite of the **mirrorvg** command is **unmirrorvg** which removes mirrored copies for an entire volume group.

As you see the quorum checking is disabled by default. Let's review what the term quorum means.

## VGDA Count



© Copyright IBM Corporation 2004

Figure 5-39. VGDA Count

AU1612.0

### Notes:

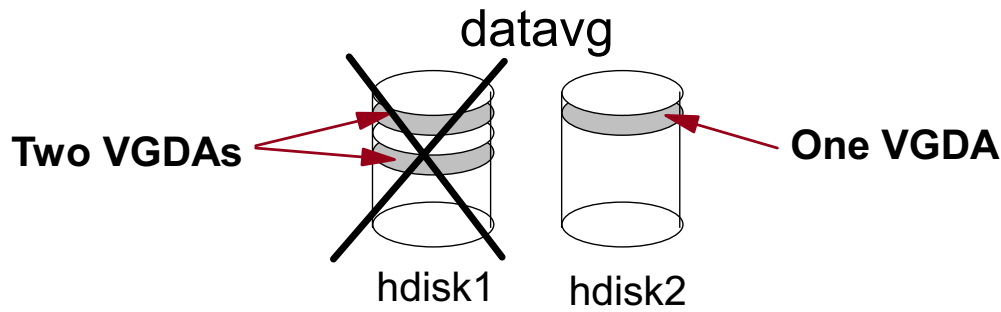
Each disk that is contained in a volume group contains at least one VGDA. The LVM always reserves space for two VGDA slots on each disk.

If a volume group consists of two disks, one disk contains two VGDA slots, the other one contains only one. If the disk with the two VGDA slots fails, we have only 33 percent of VGDA slots available, that means we have less than 50 percent of VGDA slots. In this case the quorum, which means that more than 50 percent of VGDA slots must be available, is not fulfilled.

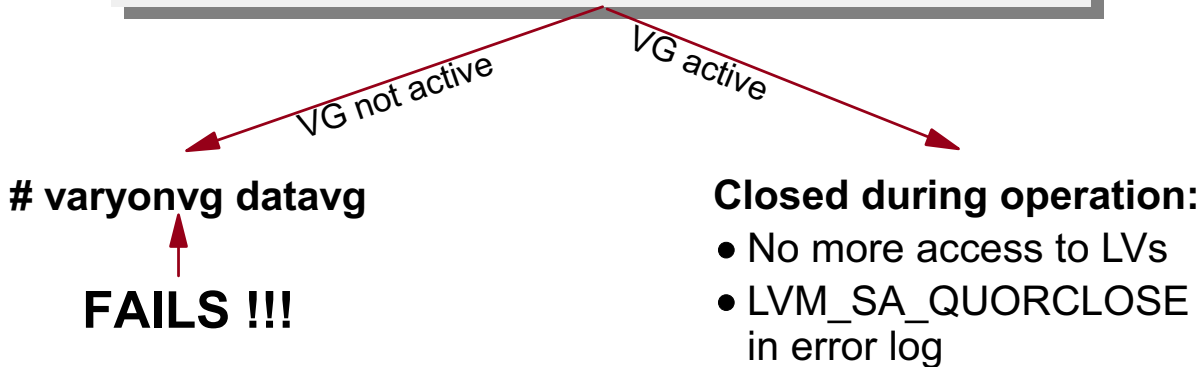
If a volume group consists of more than two disks, each disk contains one VGDA slot. If one disk fails, we still have 66 percent of VGDA slots available and the quorum is fulfilled.

What happens if a quorum is not available?

# Quorum = yes



**If hdisk1 fails, datavg has no quorum**



© Copyright IBM Corporation 2004

Figure 5-40. Quorum

AU1612.0

## Notes:

What happens if a quorum is not available in a volume group? Consider the following example.

In a two-disk volume group **datavg**, the disk **hdisk1** is not available due to a hardware defect. **hdisk1** is the disk that contains the two VGDA's; that means the volume group does not have a quorum of VGDA's. If the volume group is not varied on and the administrator tries to vary on **datavg**, the **varyonvg** command will fail.

If the volume group is already varied on when losing the quorum, the LVM will deactivate the volume group. There is no more access to any logical volume that is part of this volume group. At this point the system sometimes shows strange behavior. This situation is posted to the error log, which shows an error entry **LVM\_SA\_QUORCLOSE**. After losing the quorum, the volume group may still be listed as active (**lsvg -o**), however, all application data access and LVM functions requiring data access to the volume group will fail. The volume group is dropped from the active list as soon as the last logical volume is closed. You can still use **fuser -k /dev/LVname** and **umount /dev/LVname**, but no data is actually written to the disk.

## Nonquorum Volume Groups

With single mirroring, always disable the quorum:

- `chvg -Qn datavg`
- `varyoffvg datavg`
- `varyonvg datavg`

Additional considerations for `rootvg`:

- `chvg -Qn rootvg`
- `bosboot -ad /dev/hdiskX`
- `reboot`

- Turning off the quorum does not allow a normal `varyonvg` without a quorum
- It prevents closing the volume group when losing the quorum

© Copyright IBM Corporation 2004

Figure 5-41. Nonquorum Volume Groups

AU1612.0

### Notes:

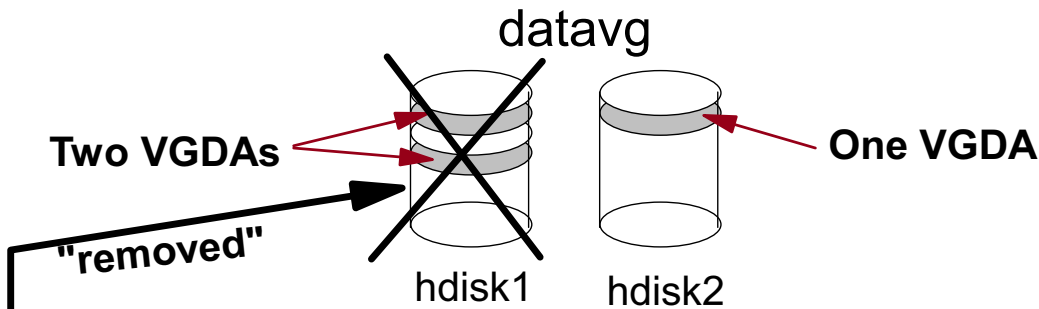
When a nonquorum volume group loses its quorum it will not be deactivated, it will be active until it loses all of its physical volumes.

When working with single mirroring, always disable the quorum using the command **chvg -Qn**. For data volume groups you must vary off and vary on the volume group to make the change work.

When turning off the quorum for **rootvg**, you must do a **bosboot** (or a **savebase**), to reflect the change in the ODM in the boot logical volume. Afterwards reboot the machine.

It's important that you know that turning off the quorum does not allow a **varyonvg** without a quorum. It just prevents the closing of an active volume group when losing its quorum.

## Forced Varyon (**varyonvg -f**)



```
# varyonvg datavg FAILS !!! (even when quorum disabled)
```

Check the reason for the failure (cable, adapter, power), before doing ...

```
# varyonvg -f datavg
```

Failure accessing hdisk1. Set PV STATE to removed.

Volume group datavg is varied on.

© Copyright IBM Corporation 2004

Figure 5-42. Forced Varyon (varyonvg -f)

AU1612.0

### Notes:

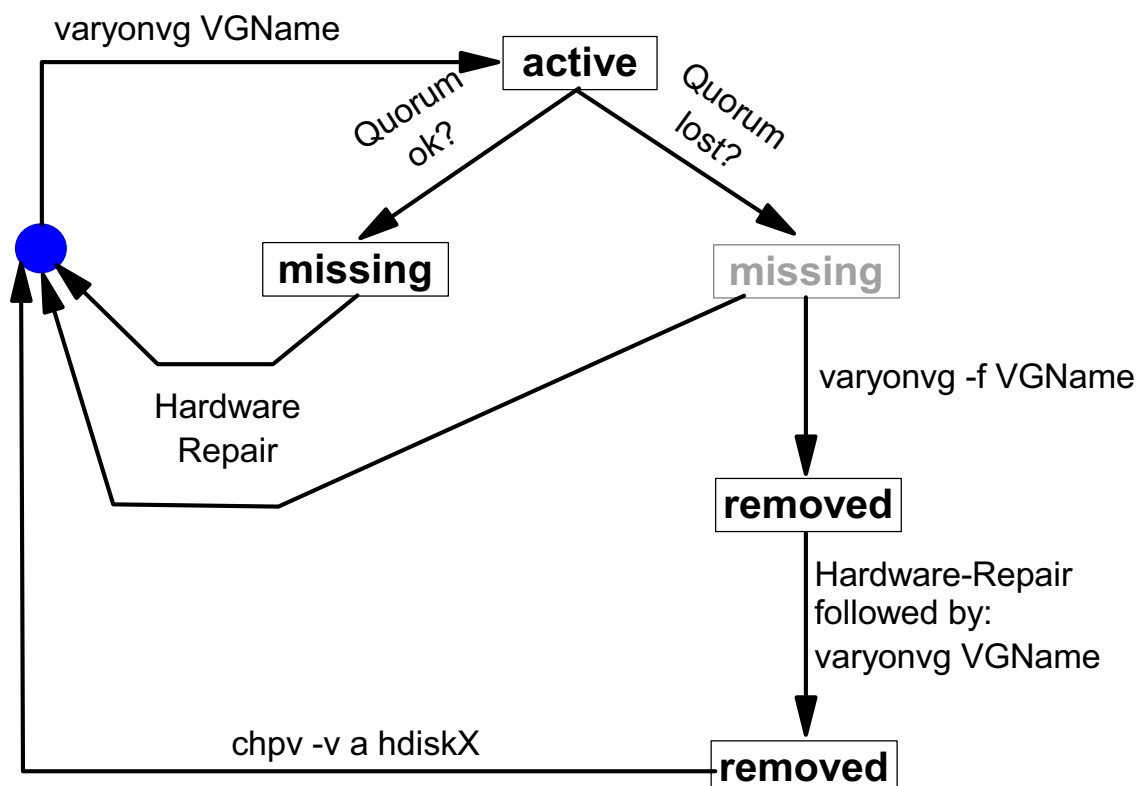
If the quorum of VGDA's is not available during vary on, the **varyonvg** command fails, even when quorum is disabled.

Before doing a forced vary on (**varyonvg -f**) always check the reason of the failure. If the physical volume appears to be permanently damaged use a forced **varyonvg**.

All physical volumes that are missing during this forced vary on will be changed to physical volume state **removed**. This means that all the VGDA and VGSA copies will be removed from these physical volumes. Once this is done, these physical volumes will no longer take part in quorum checking, nor will they be allowed to become active within the volume group until you return them to the volume group.

In our example, the active disk **hdisk2** becomes the disk with the two VGDA's. This does not change, even if the failed disk can be brought back.

## Physical Volume States and Quorum = yes



© Copyright IBM Corporation 2004

Figure 5-43. Physical Volume States

AU1612.0

### Notes:

This page introduces **physical volume states** (not device states!) Physical volume states can be displayed with **lsvg -p VGName**.

What physical volume states must you know about?

- If a disk can be accessed during a **varyonvg** it gets a PV state of **active**.
- If a disk can not be accessed during a **varyonvg**, but quorum is available, the failing disk gets a PV state **missing**.

If the disk can be repaired, for example, due to a power failure, you just have to issue a **varyonvg VGName** to bring the disk into the **active** state again. Any stale partitions will be synchronized.

- If a disk cannot be accessed during a **varyonvg** and the quorum of disks is not available, you can issue a **varyonvg -f VGName**, a forced vary on of the volume group.

The failing disk gets a PV state of **removed** and it will not be used for quorum checks anymore.

If you are able to repair the disk (for example after a power failure), executing a **varyonvg** alone does not bring the disk back into the **active** state. It maintains the **removed** state.

At this stage you have to announce the fact that the failure is over by using the following command:

**# chpv -va hdiskX**

**This defines the disk hdiskX as active.**

Note that you have to do a **varyonvg VGName** afterwards to synchronize any stale partitions.

The opposite of **chpv -va** is **chpv -vr** which brings the disk into the **removed** state. This works only when all logical volumes have been closed on the disk that will be defined as removed. Additionally, **chpv -vr** does not work when the quorum will be lost in the volume group after removing the disk.



## Summary Quorum

	Quorum ON	Quorum OFF
<b>rootvg (active)</b>	> 50 %	>=1
<b>datavg (active)</b>	> 50 %	>=1
<b>rootvg (varyon)</b>	>=1	>=1
<b>datavg (varyon)</b>	> 50 %	100 % or varyonvg -f or MISSINGPV_VARYON=TRUE

© Copyright IBM Corporation 2004

Figure 5-44. Summary Quorum

AU1612.0

### Notes:

With **Quorum turned ON** you always need > 50% of the VGDA's available (except of rootvg varyon).

If **Quorum is turned OFF** you have to make a difference between an already active volume group and between performing a varyon.

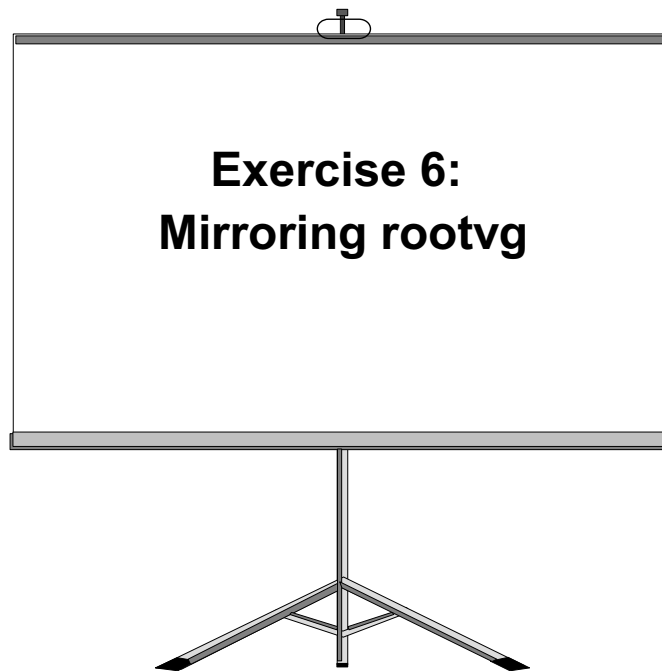
An active Volume Group will be kept open as long as there is at least one VGDA available.

Set **MISSINGPV\_VARYON=true** in /etc/environment if volume group needs to be varied on with missing disks at the boot time.

When using **varyonvg -f** or using **MISSINGPV\_VARYON=true** you take full responsibility for the volume group integrity.

## Next Step

---



© Copyright IBM Corporation 2004

Figure 5-45. Next Step...

AU1612.0

### **Notes:**

At the end of the exercise, you should be able to:

- Mirror the rootvg
- Describe physical volume states
- Unmirror the rootvg

## Checkpoint

---

Answer True or False to the following statements:

1. All LVM information is stored in the ODM.
2. You detect that a physical volume hdisk1 that is contained in your rootvg is missing in the ODM. This problem can be fixed by exporting and importing the rootvg.
3. The LVM supports RAID-5 without separate hardware.

© Copyright IBM Corporation 2004

---

Figure 5-46. Checkpoint

AU1612.0

### **Notes:**

## Unit Summary

---

- The LVM information is held in a number of different places on the disk, including the ODM and the VGDA
- ODM related problems can be solved by:
  - exportvg/importvg (non rootvg VGs)
  - rvgrecovery (rootvg)
- Mirroring improves the availability of a system or a logical volume
- Striping improves the performance of a logical volume
- Quorum means that more than 50% of VGDA's must be available

© Copyright IBM Corporation 2004

Figure 5-47. Unit Summary

AU1612.0

### **Notes:**

# Unit 6. Disk Management Procedures

## What This Unit Is About

This unit describes different disk management procedures:

- Disk replacement procedures
- Procedures to solve problems caused by an incorrect disk replacement
- Export and import of volume groups

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Replace a disk under different circumstances
- Recover from a total volume group failure
- Rectify problems caused by incorrect actions that have been taken to change disks
- Export and import volume groups

## How You Will Check Your Progress

Accountability:

- Lab exercises
- Checkpoint questions

## References

Online *Commands Reference*  
GG24-4484-00 *AIX Storage Management*

## Unit Objectives

---

After completing this unit, students should be able to:

- Replace a disk under different circumstances
- Recover from a total volume group failure
- Rectify problems caused by incorrect actions that have been taken to change disks
- Export and import volume groups

© Copyright IBM Corporation 2004

Figure 6-1. Unit Objectives

AU1612.0

### **Notes:**

This unit presents many disk management procedures that are very important for any AIX system administrator.

## 6.1 Disk Replacement Techniques

# Disk Replacement: Starting Point

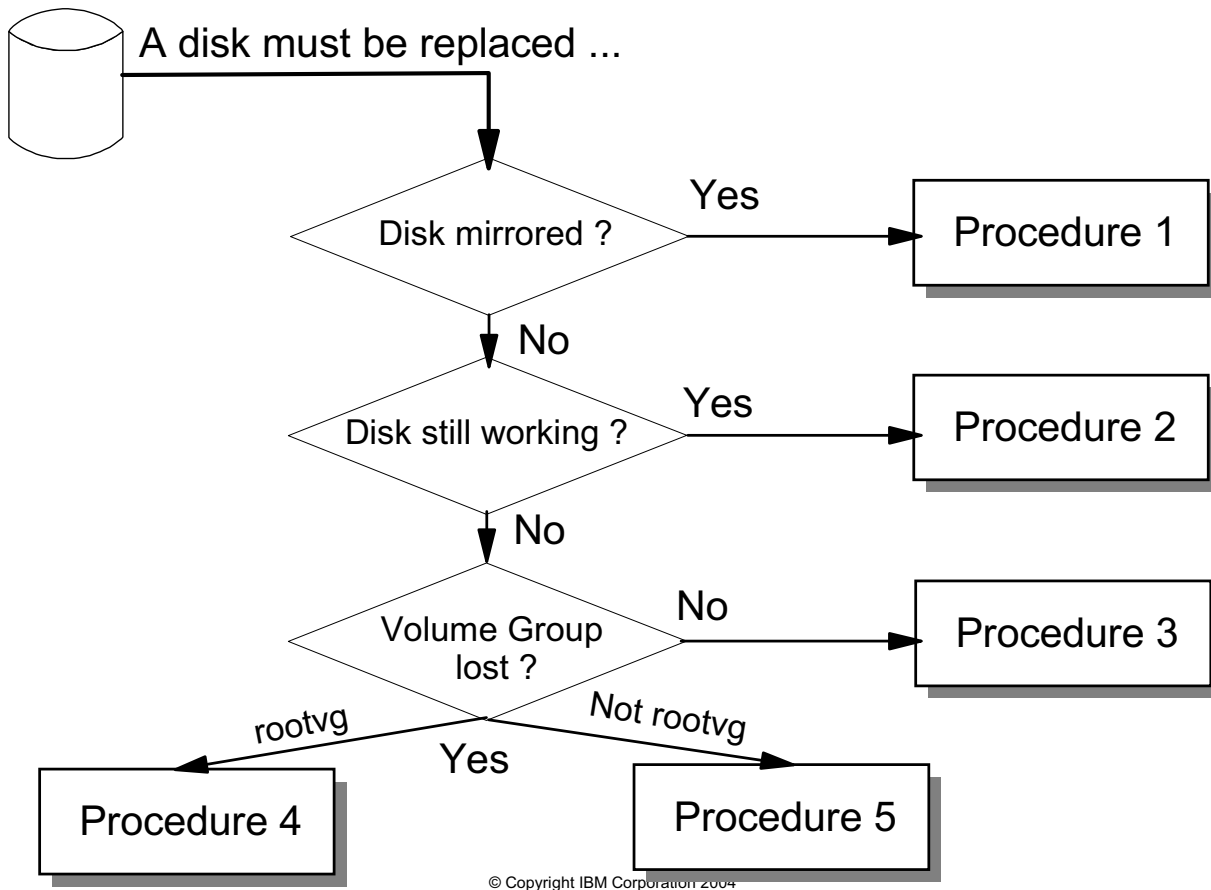


Figure 6-2. Disk Replacement: Starting Point

AU1612.0

## Notes:

Many reasons might require the replacement of a disk, for example:

- Disk too small
- Disk too slow
- Disk produces many **DISK\_ERR4** log entries

Before starting the disk replacement, always follow the flowchart that is shown on this page. This will help you whenever you have to replace a disk.

1. If the disk that must be replaced is completely mirrored onto another disk, follow **procedure 1**.
2. If a disk is not mirrored, but still works, follow **procedure 2**.
3. If you are absolutely sure that a disk failed and you are not able to repair the disk, do the following:  
If the volume group can be varied on (normal or forced), use **procedure 3**.



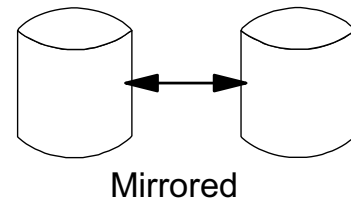
If the volume group is totally lost after the disk failure, that means the volume group could not be varied on (either normal or forced), follow **procedure 4** if the volume group is the **rootvg**.

If the volume group that is lost is **not** the **rootvg** follow **procedure 5**.

Let's start with **procedure 1**.

## Procedure 1: Disk Mirrored

1. Remove all copies from disk:  
# unmirrorvg *vg\_name* *hdiskX*
2. Remove disk from volume group:  
# reducevg *vg\_name* *hdiskX*
3. Remove disk from ODM:  
# rmdev -l *hdiskX* -d
4. Connect new disk to system:  
# reboot (if not hot-swappable)
5. Add new disk to volume group:  
# extendvg *vg\_name* *hdiskY*
6. Create new copies:  
# mirrorvg *vg\_name* *hdiskY*  
# varyonvg *vg\_name*



© Copyright IBM Corporation 2004

Figure 6-3. Procedure 1: Disk Mirrored

AU1612.0

### Notes:

Use **procedure 1** when the disk that must be replaced is mirrored.

This procedure requires that the disk state of the failed disk be either **missing** or **removed**. Refer to Physical Volume States in Unit 5: Disk Management Theory for more information on disk states. Use **lspv *hdiskX*** to check the state of your physical volume. If the disk is still in the **active** state you cannot remove any copies or logical volumes from the failing disk. In this case one way to bring the disk into a **removed** or **missing** state is to run the **reducevg -d** command or to do a varyoffvg and a varyonvg on the volume group by rebooting the system.

Remember to disable the quorum check if you have only two disks in your volume group.

The goal of each disk replacement is to remove all logical volumes from a disk.

1. Start removing all logical volume copies from the disk. Use either the smit fastpath **smit unmirrorvg** or the **unmirrorvg** command as shown. This must be done for each logical volume that is mirrored on the disk.

---

If you have additional unmirrored logical volumes on the disk you have to either move them to another disk (**migratepv**), or remove them if the disk cannot be accessed (**rmlv**). As mentioned the latter will only work if the disk state is either **missing** or **removed**.

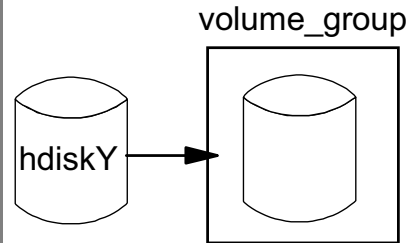
2. If the disk is completely empty, remove the disk from the volume group. Use smit fastpath **smit reducevg** or the **reducevg** command.
3. After the disk has been removed from the volume group, you can remove it from the ODM. Use the **rmdev** command as shown.

If the disk must be removed from the system, shut down the machine and then remove it.

4. Connect the new disk to the system and reboot your system. The **cfgmgr** will configure the new disk. If using hot-swappable disks, a reboot is not necessary.
5. Add the new disk to the volume group. Use either the smit fastpath **smit extendvg** or the **extendvg** command.
6. Finally create the new copies for each logical volume on the new disk. Use either the smit fastpath **smit mirrorvg** or the **mirrorvg** command. Synchronize the volume group (or each logical volume) afterwards, using the **varyonvg** command.

## Procedure 2: Disk Still Working

1. Connect new disk to system
2. Add new disk to volume group:  
# `extendvg vg_name hdiskY`
3. Migrate old disk to new disk: (\*)  
# `migratepv hdiskX hdiskY`
4. Remove old disk from volume group:  
# `reducevg vg_name hdiskX`
5. Remove old disk from ODM:  
# `rmdev -l hdiskX -d`



(\*) : Is the disk in rootvg?  
See next foil for further considerations!

© Copyright IBM Corporation 2004

Figure 6-4. Procedure 2: Disk Still Working

AU1612.0

### Notes:

**Procedure 2** applies to a disk replacement where the disk is **unmirrored** but could be accessed.

The goal is the same as always. Before we can replace a disk we must remove everything from the disk.

1. Shut down your system if you need to physically attach a new disk to the system. Boot the system so that **cfgmgr** will configure the new disk.
2. Add the new disk to the volume group. Use either the smit fastpath **smit extendvg** or the **extendvg** command.
3. Before executing the next step it is necessary to distinguish between the **rootvg** and a **non-rootvg** volume group.

If the disk that is replaced is in **rootvg** execute the steps that are shown on page *Procedure 2: Special Steps for rootvg*.

If the disk that is replaced is **not** in the **rootvg**, use the **migratepv** command:

**# migratepv hdisk\_old hdisk\_new**

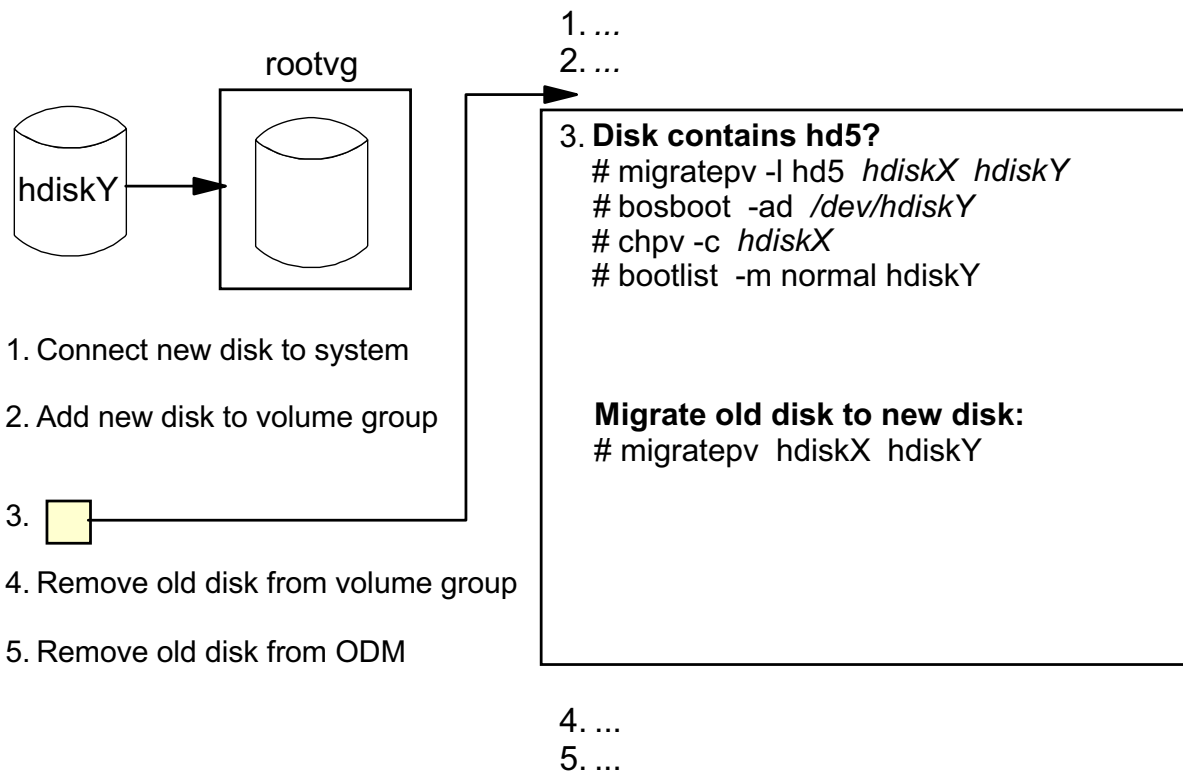
**This command moves all logical volumes from one disk to another. You can do this during normal system activity. The command migratepv** requires that the disks are in the same volume group.

4. If the old disk has been completely migrated, remove it from the volume group. Use either the smit fastpath **smit reducevg** or the **reducevg** command.
5. If you need to remove the disk from the system, remove it from the ODM using the **rmdev** command as shown. Finally remove the physical disk from the system.

**Note:**

If the disk that must be replaced is in **rootvg**, follow the instructions on the next page.

## Procedure 2: Special Steps for rootvg



© Copyright IBM Corporation 2004

Figure 6-5. Procedure 2: Special Steps for rootvg

AU1612.0

### Notes:

**Procedure 2** requires some additional steps if the disk that must be replaced is in **rootvg**.

1. Connect the new disk to the system as described in procedure 2.
2. Add the new disk to the volume group. Use **smit extendvg** or the **extendvg** command.
3. This step requires special considerations for **rootvg**:
  - Check whether your disk contains the **boot logical volume** (default is **/dev/hd5**).

Use command **lspv -l** to check the logical volumes on the disk that must be replaced.

If the disk contains the **boot logical volume**, migrate the logical volume to the new disk and update the boot logical volume on the new disk. To avoid a potential boot from the old disk, clear the old boot record, by using the **chpv -c** command. Then change your boot list:

```
# migratepv -l hd5 hdiskX hdiskY
# bosboot -ad /dev/hdiskY
```

```
# chpv -c hdiskX  
# bootlist -m normal hdiskY
```

If the disk contains the **primary dump device**, you must deactivate the dump before migrating the corresponding logical volume:

```
# sysdumpdev -p /dev/sysdumpnull
```

- **Migrate the complete old disk to the new one:**

```
# migratepv hdiskX hdiskY
```

If the **primary dump device** has been deactivated, you have to activate it again:

```
# sysdumpdev -p /dev/hdX (Default is /dev/hd6 in AIX 4)
```

4. After the disk has been migrated, remove it from the root volume group.

```
# reducevg rootvg hdiskX
```

5. If the disk must be removed from the system, remove it from the ODM (use the **rmdev** command), shut down your AIX, and remove the disk from the system afterwards.

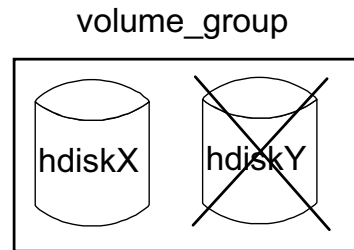
```
# rmdev -l hdiskX -d
```

---

## Procedure 3: Total Disk Failure

---

1. Identify all LVs and file systems on failing disk:  
`# lspv -l hdiskY`
2. Unmount all file systems on failing disk:  
`# umount /dev/lv_xx`
3. Remove all file systems and LVs from failing disk:  
`# smit rmfs                                       # rmlv lv_xx`
4. Remove disk from volume group:  
`# reducevg vg_name hdiskY`
5. Remove disk from system:  
`# rmdev -l hdiskY -d`
6. Add new disk to volume group:  
`# extendvg vg_name hdiskZ`
7. Re-create all LVs and file systems on new disk:  
`# mklv -y lv_xx                               # smit crfs`
8. Restore file systems from backup:  
`# restore -rvqf /dev/rmt0`



```
# lspv hdiskY
...
PV STATE: removed

# lspv hdiskY
...
PV STATE: missing
```

© Copyright IBM Corporation 2004

Figure 6-6. Procedure 3: Total Disk Failure

AU1612.0

### Notes:

**Procedure 3** applies to a disk replacement where a disk **could not be accessed** but the volume group is intact. The failing disk is either in a state (not device state) of **missing** (normal varyonvg worked) or **removed** (forced varyonvg was necessary to bring the volume group online).

If the failing disk is in an **active** state (this is **not** a device state), this procedure will not work. In this case one way to bring the disk into a **removed** or **missing** state is to run the **reducevg -d** command or to do a varyoffvg and a varyonvg on the volume group by rebooting the system. The reboot is necessary because you cannot vary off a volume group with open logical volumes. Because the failing disk is **active** there is no way to unmount file systems.

If the failing disk is in a **missing** or **removed** state, start the procedure:

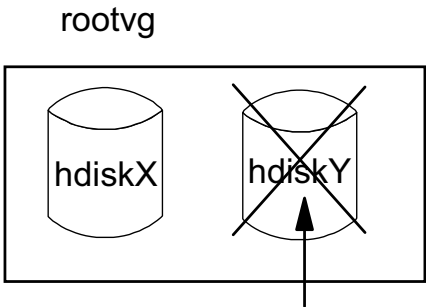
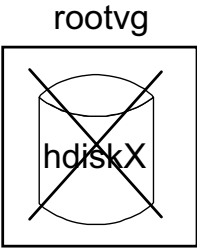
1. Identify all logical volumes and file systems on the failing disk. Use commands like **lspv**, **lslv** or **lsfs** to provide this information. These commands will work on a failing disk.



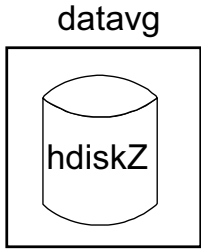
2. If you have **mounted** file systems on a logical volume on the failing disk, you must unmount them. Use the **umount** command.
3. Remove all file systems from the failing disk, using **smit rmfs** or the **rmfs** command. If you remove a file system, the corresponding logical volume and stanza in **/etc/filesystems** is removed as well.
4. Remove the remaining logical volumes (those not associated with a file system) from the failing disk using **smit rmlv** or the **rmlv** command.
5. Remove the disk from the volume group, using the smit fastpath **smit reducevg** or the **reducevg** command.
6. Remove the disk from the ODM (**rmdev**) and from the system.
7. Add the new disk to the system and extend your volume group. Use **smit extendvg** or the **extendvg** command.
8. Re-create all logical volumes and file systems that have been removed due to the disk failure. Use **smit mklv**, **smit crfs** or the commands directly.
9. Due to the total disk failure, you lost all data on the disk. This data has to be restored, either by the **restore** command or any other tool you use to restore data (for example, TSM).

# Procedure 4: Total rootvg Failure

1. Replace bad disk.
2. Boot in maintenance mode.
3. Restore from a **mksysb** tape.
4. Import each volume group into the new ODM (importvg) if needed.



contains OS logical volumes



© Copyright IBM Corporation 2004

Figure 6-7. Procedure 4: Total rootvg Failure

AU1612.0

## Notes:

**Procedure 4** applies to a total **rootvg** failure.

This situation might come up when your **rootvg** consists of one disk that fails. Or your **rootvg** is installed on two disks and the disk fails that contains **operating system** logical volumes (for example, **/dev/hd4**).

1. Replace the bad disk and boot your system in **maintenance mode**.
2. Restore your system from a **mksysb** tape.

Remember that if any rootvg file systems were not mounted when the mksysb was made, those file systems are not included on the backup image. You will need to create and restore those as a separate step.

If your **mksysb** tape does not contain user volume group definitions (for example, you created a volume group after saving your **rootvg**), you have to import the user volume group after restoring the **mksysb**. For example:

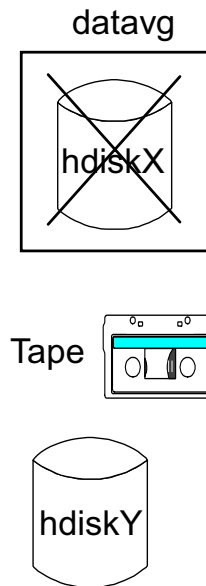
```
# importvg -y datavg hdisk9
```

Only one disk from the volume group (in our example **hdisk9**), needs to be selected.

Export and import of volume groups is discussed in more detail in the next topic.

## Procedure 5: Total non-rootvg Failure

1. Export the volume group from the system:  
# exportvg vg\_name
  2. Check /etc/filesystems.
  3. Remove bad disk from ODM and the system:  
# rmdev -l hdiskX -d
  4. Connect new disk
  5. If volume group backup available (savevg):  
# restvg -f /dev/rmt0 hdiskY
  6. If **no** volume group backup available: Recreate ...
    - volume group (mkvg)
    - logical volumes and filesystems (mklv, crfs).
- Restore data from a backup:  
# restore -rqvf /dev/rmt0



© Copyright IBM Corporation 2004

Figure 6-8. Procedure 5: Total non-rootvg Failure

AU1612.0

### Notes:

**Procedure 5** applies to a total failure of a non-rootvg volume group. This situation might come up if your volume group consists of only one disk that fails. Before starting this procedure make sure this is not just a temporary disk failure (for example, a power failure).

1. To fix this problem, export the volume group from the system. Use the command **exportvg** as shown. During the export of the volume group all ODM objects that are related to the volume group will be deleted.
2. Check your **/etc/filesystems**. There should be no references to logical volumes or file systems from the exported volume group.
3. Remove the bad disk from the ODM (Use **rmdev** as shown). Shut down your system and remove the physical disk from the system.
4. Connect the new drive and boot the system. The **cfgmgr** will configure the new disk.
5. If you have a **volume group backup** available (created by the **savevg** command), you can restore the complete volume group with the **restvg** command (or the smit fastpath **smit restvg**). All logical volumes and file systems are recovered.

If you have more than one disk that should be used during **restvg** you must specify these disks:

```
# restvg -f /dev/rmt0 hdiskY hdiskZ
```

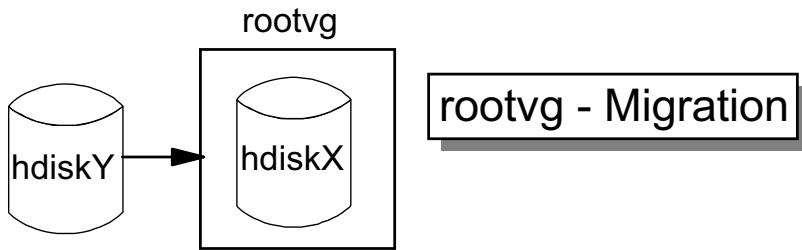
We will talk more about **savevg** and **restvg** in a future chapter.

6. If you have **no** volume group backup available, you have to re-create everything that was part of the volume group.

Re-create the volume group (**mkvg** or **smit mkvg**), all logical volumes (**mklv** or **smit mklv**) and all file systems (**crfs** or **smit crfs**).

Finally, restore the lost data from backups, for example with the **restore** command or any other tool you use to restore data in your environment.

## Frequent Disk Replacement Errors (1 of 4)



### Boot problems after migration:

- Firmware LED codes cycle

### Fix:

- Check bootlist (SMS Menu)
- Check bootlist (bootlist)
- Re-create boot logical volume (bosboot)

© Copyright IBM Corporation 2004

Figure 6-9. Frequent Disk Replacement Errors (1 of 4)

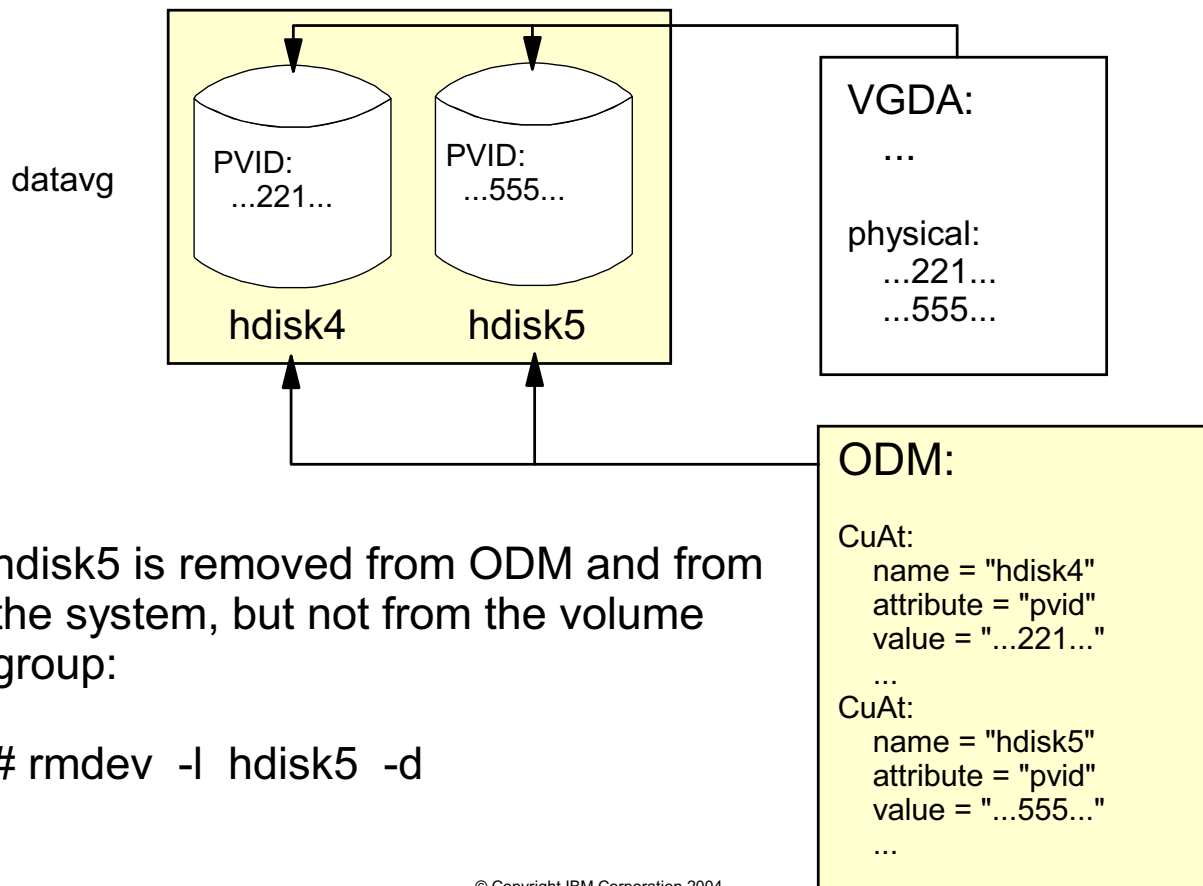
AU1612.0

### Notes:

A common problem seen after a migration of the **rootvg** is that the machine will not boot. On a microchannel system you get alternating LED codes **223-229**, on a PCI system the LED codes cycle. This loop indicates that the firmware is not able to find a bootstrap code to boot from.

This problem is usually easy to fix. Boot in **SMS Menu (F1)** and check your bootlist (use Multi-boot menu) or boot in **maintenance mode** and check your boot list (use the **bootlist** command). If the boot list is correct, update the **boot logical volume** (use the **bosboot** command).

## Frequent Disk Replacement Errors (2 of 4)



hdisk5 is removed from ODM and from the system, but not from the volume group:

```
# rmdev -l hdisk5 -d
```

© Copyright IBM Corporation 2004

Figure 6-10. Frequent Disk Replacement Errors (2 of 4)

AU1612.0

### Notes:

**Note:** Throughout this discussion the physical volume ID is abbreviated in the visuals for simplicity. The physical volume id is actually 32 characters.

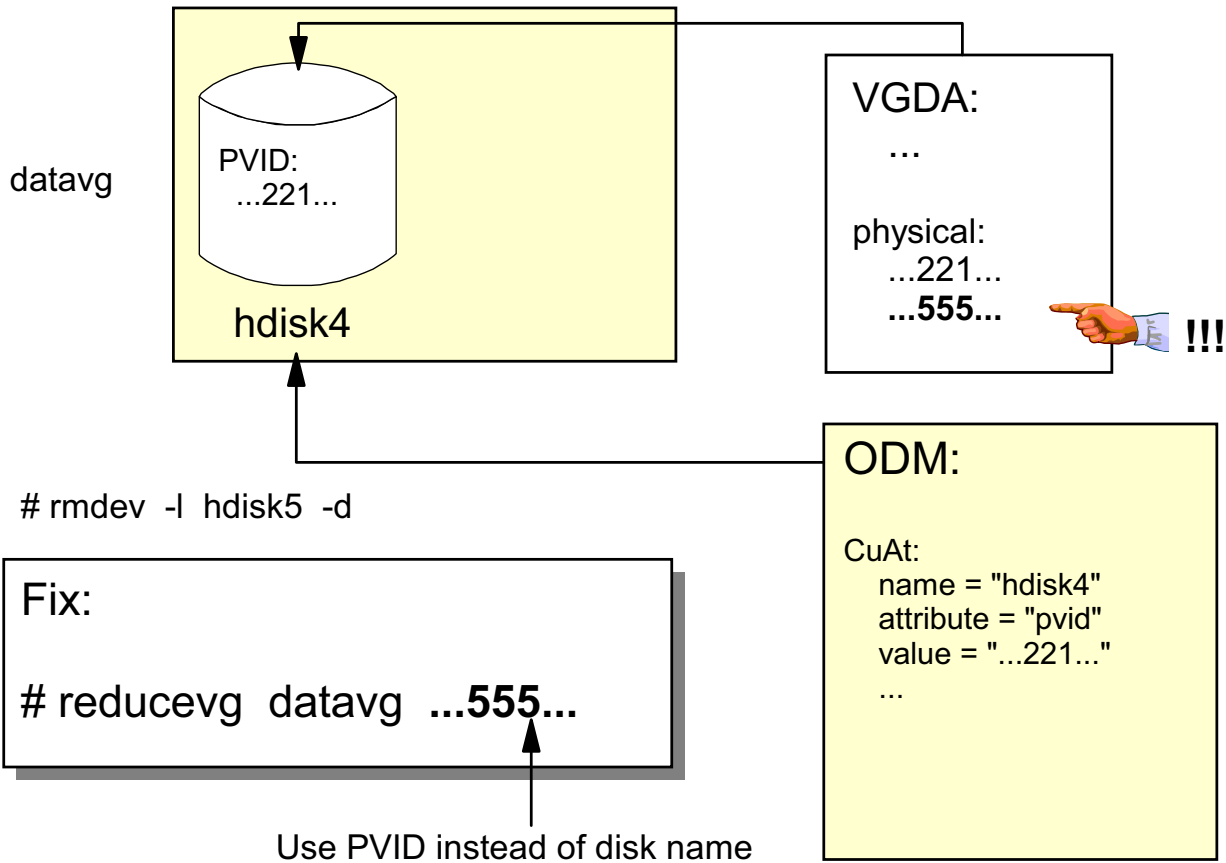
Another frequent error comes up when administrators remove a disk from the ODM (by executing **rmdev**) and physically remove the disk from the system, but do not remove entries from the volume group descriptor area.

Before discussing the fix for this problem, remember that the **VGDA** stores information about all physical volumes of the volume group. Each disk has at least one **VGDA**.

Disk information is also stored in the ODM, for example, the physical volume identifiers are stored in the ODM class **CuAt**.

What happens if a disk is removed from the ODM but not from the volume group?

# Frequent Disk Replacement Errors (3 of 4)



© Copyright IBM Corporation 2004

Figure 6-11. Frequent Disk Replacement Errors (3 of 4)

AU1612.0

## Notes:

After removing the disk from the ODM you still have a reference in the **VGDA** to the removed disk. In early AIX versions the fix for this problem was difficult. You had to add ODM objects that described the attributes of the removed disk.

Fix this problem by executing the **reducevg** command. Instead of passing the disk name you pass the **physical volume ID** of the removed disk.

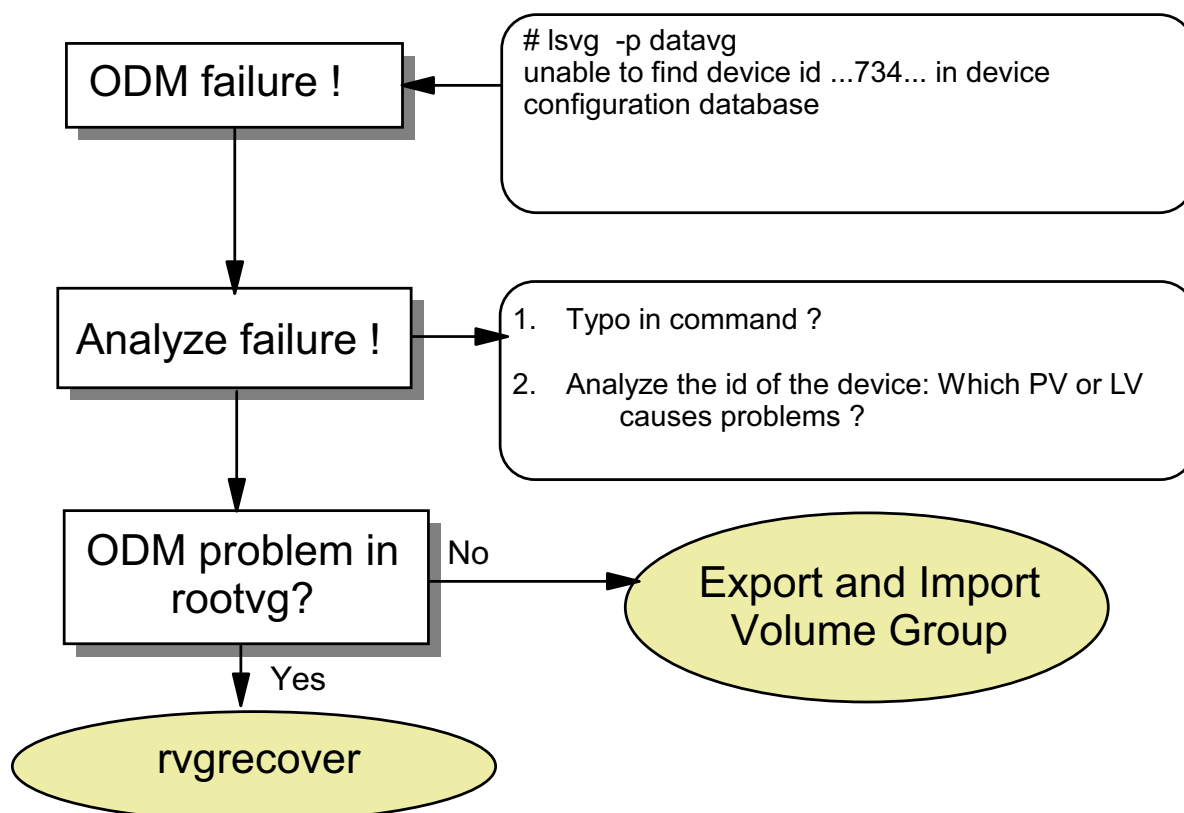
Execute the **lspv** command to identify the **missing disk**. Write down the **physical volume ID** of the missing disk and compare this id with the contents of the **VGDA**. Use the following command to query the **VGDA** on a disk:

**# lqueryvg -p hdisk4 -At (Use any disk from the volume group)**

**If you are sure that you found the missing pvid, pass this pvid to the reducevg command.**



## Frequent Disk Replacement Errors (4 of 4)



© Copyright IBM Corporation 2004

Figure 6-12. Frequent Disk Replacement Errors (4 of 4)

AU1612.0

### Notes:

After an incorrect disk replacement you might detect ODM failures. A typical error message is shown:

**unable to find device id 00837734 in device configuration database**

**In this case a device could not be found in the ODM. Before starting any fixes check the command you typed in. Maybe it just contains a typo.**

**Analyze the failure. Find out what device corresponds to the ID that is shown in the error message.**

**If you are not sure what caused the problem, remember the two ways you learned already to fix an ODM problem.**

- **If the ODM problem is related to the rootvg**, execute the **rvgrecover** procedure. If the ODM problem is **not** related to the **rootvg**, export the volume group and import it again.

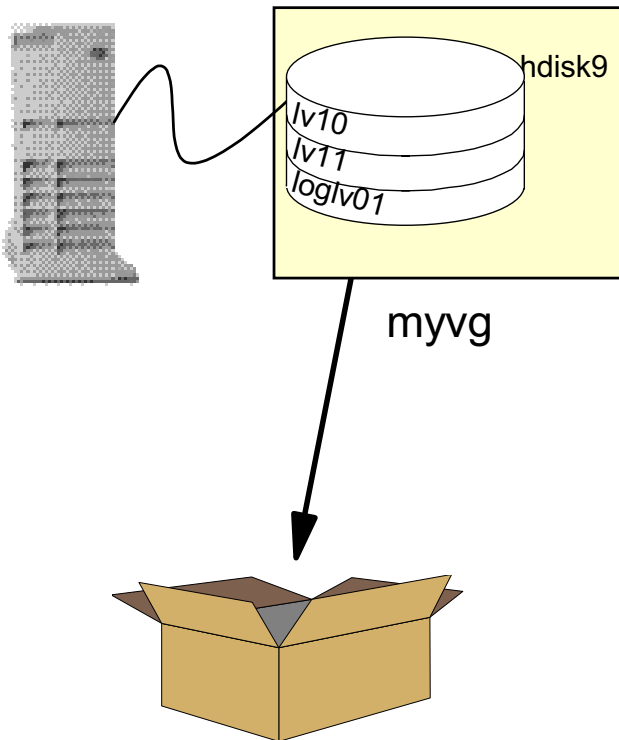
Export and import will be explained in more detail in the next topic.



## 6.2 Export and Import

# Exporting a Volume Group

moon



To export a volume group:

1. Unmount all filesystems:  
# umount /dev/lv10  
# umount /dev/lv11
2. Vary off the volume group:  
# varyoffvg myvg
3. Export volume group:  
# exportvg myvg

The complete volume group is removed from the ODM.

© Copyright IBM Corporation 2004

Figure 6-13. Exporting a Volume Group

AU1612.0

## Notes:

As you learned already, **exportvg** and **importvg** can be used to fix ODM problems. Additionally, these commands provide a way to transfer data between different AIX systems. This page provides an example of how to export a volume group:

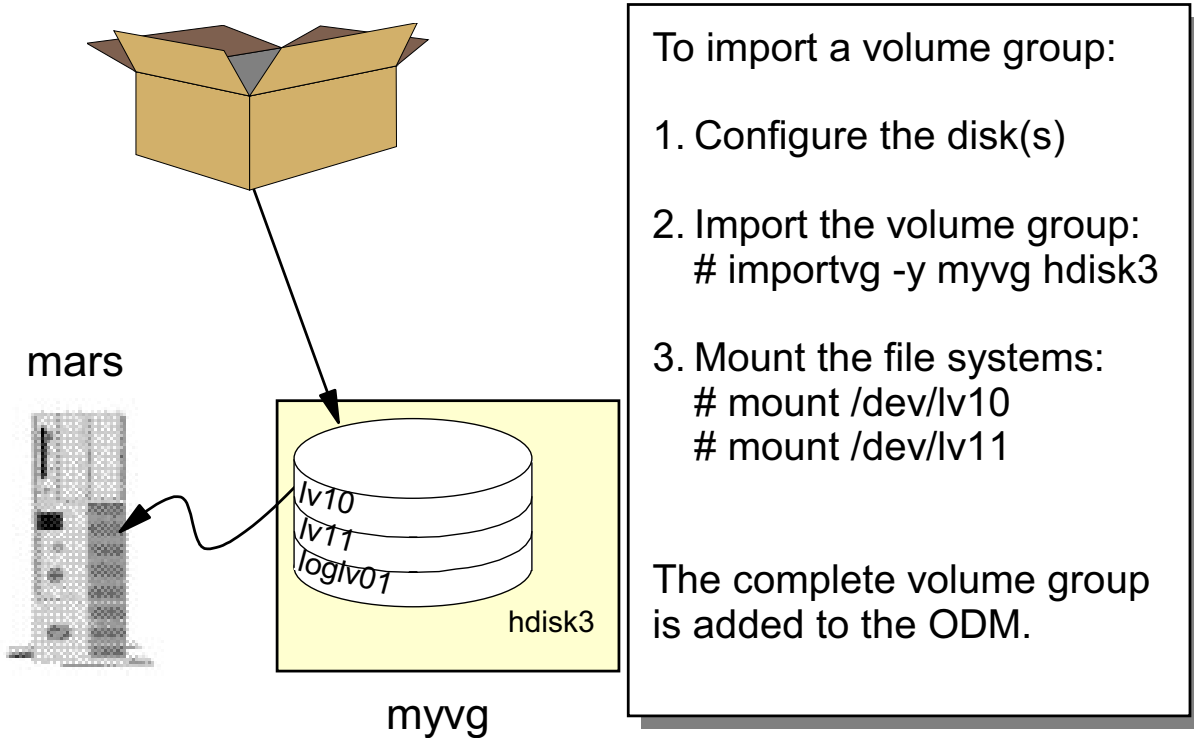
On a system named **moon** a disk **hdisk9** is connected. This disk belongs to a volume group **myvg**. This volume group needs to be transferred to another system. Execute the following steps to export this volume group:

1. **Unmount** all file systems from the volume group. As you see we have two logical volumes **lv10** and **lv11** in **myvg**. Another logical volume **loglv01** exists in the volume group **myvg**. This logical volume is the JFS log device for the file systems in **myvg**, which is closed when all file systems are unmounted.
2. When all logical volumes are closed, we vary off the volume group. Execute the **varyoffvg** command as shown.

3. Finally export the volume group, using the **exportvg** command. After this point the complete volume group (including all file systems and logical volumes) is removed from the ODM.

After exporting the volume group you can transfer the disk to another system.

# Importing a Volume Group



© Copyright IBM Corporation 2004

Figure 6-14. Importing a Volume Group

AU1612.0

## Notes:

To import a volume group into a system, for example into a system named **mars**, execute the following steps.

1. Connect all disks (in our example we have only one disk) and reboot the system so that **cfgmgr** will configure the added disks.
2. Notice that you only have to specify one disk (using either **hdisk#** or **PVID**) during.....). If you do not specify the option **-y** the command will generate a new volume group name.

Notice that you only have to specify **one** disk during the **importvg**. Because all disks contain the same **VGDA** information, the system can determine this information by querying any **VGDA** from any disk in the VG.

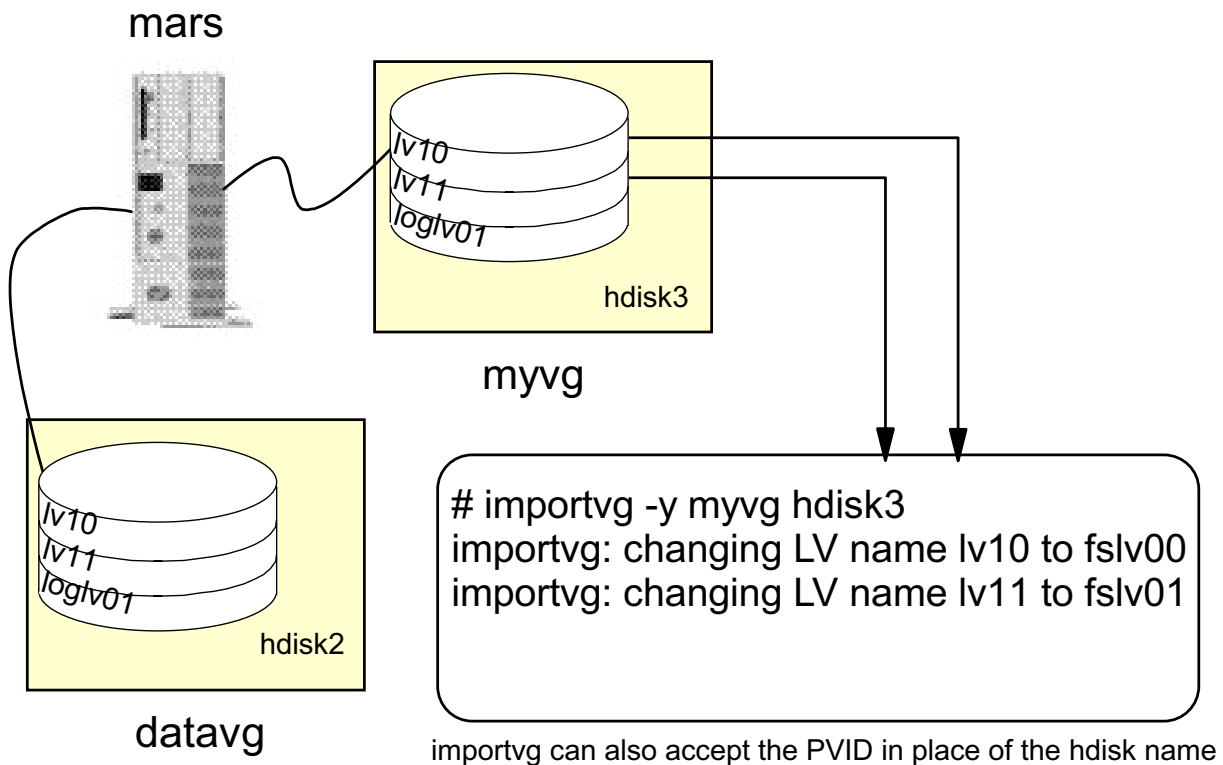
The command **importvg** generates completely new ODM entries.

3. In AIX 4.3 and subsequent releases of the operating system the volume group is automatically varied on. If you are using another AIX version, you have to check whether the volume group is varied on after the **importvg**.

If the volume group is **not automatically** varied on, execute the **varyonvg** command to vary on the volume group.

4. Finally mount the file systems.

# importvg and Existing Logical Volumes



© Copyright IBM Corporation 2004

Figure 6-15. importvg and Existing Logical Volumes

AU1612.0

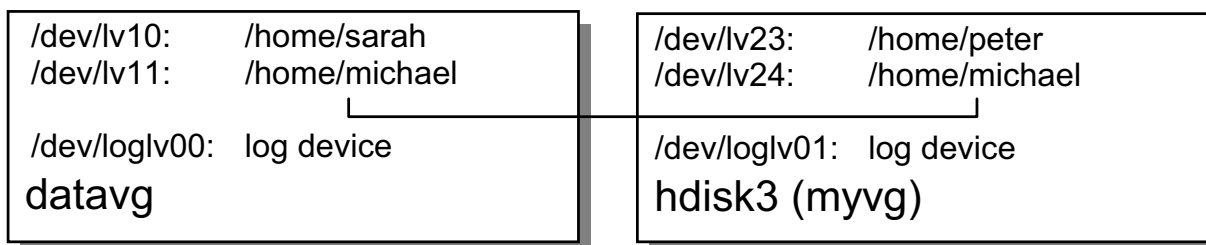
## Notes:

If you are importing a volume group with logical volumes that already exist on the system, the **importvg** command renames the **logical volumes** from the volume group that is imported.

The logical volumes **/dev/lv10** and **/dev/lv11** exist in both volume groups. During the **importvg** command the logical volumes from **myvg** are renamed to **/dev/fslv00** and **/dev/fslv01**.



## importvg and Existing Filesystems (1 of 2)



```
# importvg -y myvg hdisk3
```

Warning: mount point /home/michael already exists in /etc/filesystems

```
# umount /home/michael
```

```
# mount -o log=/dev/loglv01 /dev/lv24 /home/michael
```

© Copyright IBM Corporation 2004

Figure 6-16. importvg and Existing Filesystems (1 of 2)

AU1612.0

### Notes:

If a file system (for example **/home/michael**) already exists on a system, you run into problems when you **mount** the file system that was imported.

This page explains the one thing you can do:

- Unmount the file system that exists on the system (**/home/michael** from **datavg**).
- Mount the imported file system. Note that you have to specify the **log device** (-o log=/dev/lvlog01), the **logical volume name** (/dev/lv24) and the **mount point** (/home/michael). If the filesystem type is jfs2 you have to specify this as well ( **-V jfs2** ). You can get all this informations by running the command **getlvcb lv24 -At**

Another possibility is to add a new stanza to the **/etc/filesystems** file. This is covered on the next page.

## importvg and Existing Filesystems (2 of 2)

```
# vi /etc/filesystems
```

```
/home/michael:
dev    = /dev/lv11
vfs    = jfs
log    = /dev/loglv00
mount  = false
options = rw
account= false
```

```
/home/michael_moon:
dev    = /dev/lv24
vfs    = jfs
log    = /dev/loglv01
mount  = false
options = rw
account= false
```

```
/dev/lv10:  /home/sarah
/dev/lv11:  /home/michael

/dev/loglv00: log device
datavg
```

```
/dev/lv23:  /home/peter
/dev/lv24:  /home/michael

/dev/loglv01: log device
hdisk3 (myvg)
```

```
# mount /home/michael
# mount /home/michael_moon
```

Mount point must exist !

© Copyright IBM Corporation 2004

Figure 6-17. importvg and Existing Filesystems (2 of 2)

AU1612.0

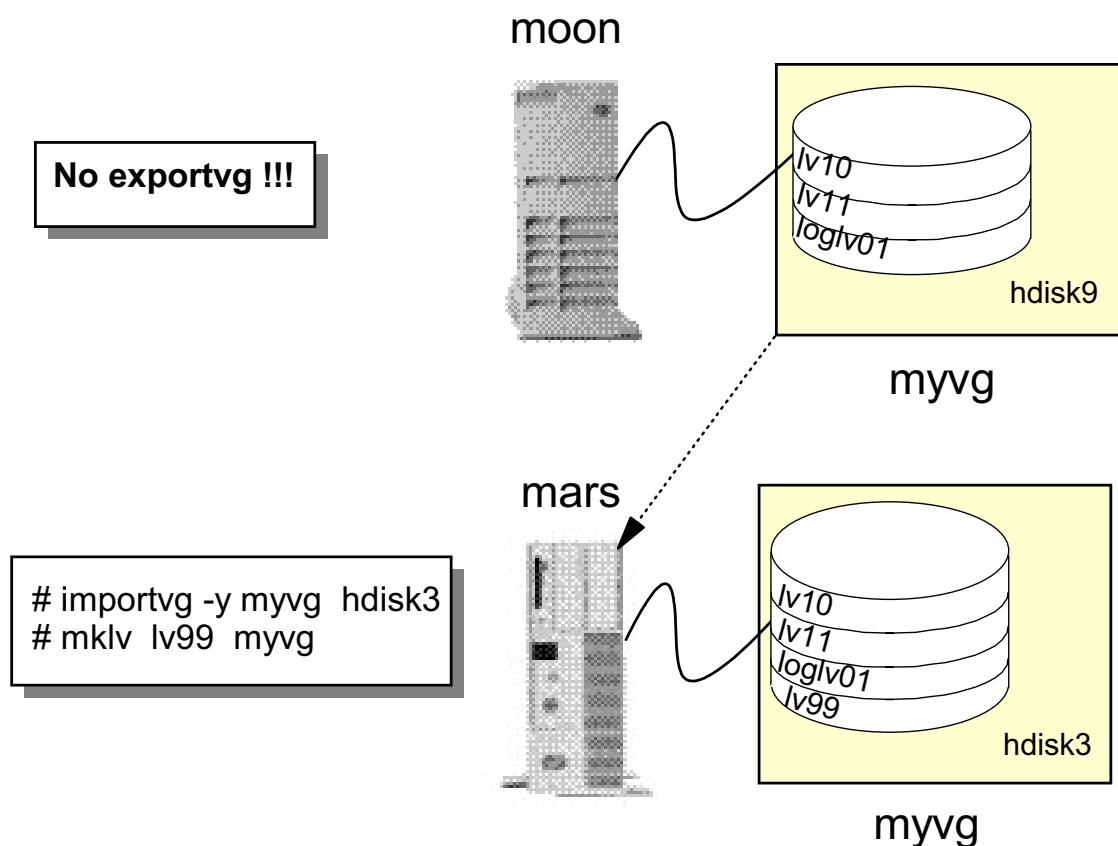
### Notes:

If you need both file systems (the imported and the one that already exists) mounted at the same time, you need to create a new stanza in **/etc/filesystems**. In our example we create a second stanza for our imported logical volume, **/home/michael\_moon**:

- **dev** specifies the logical volume, in our example **/dev/lv24**.
- **vfs** specifies the file system type, in our example a **journaled file system**.
- **log** specifies the **JFS log device** for the file system.
- **mount** specifies whether this file system should be mounted by default. The value **false** specifies no default mounting during boot. The value **true** indicates that a file system should be mounted during the boot process.
- **options** specifies that this file system should be mounted with read and write access.
- **account** specifies whether the file system should be processed by the accounting system. A value of false indicates no accounting.

Before mounting the file system **/home/michael\_moon**, the corresponding mount point must be created.

## importvg -L (1 of 2)



© Copyright IBM Corporation 2004

Figure 6-18. importvg -L (1 of 2)

AU1612.0

### Notes:

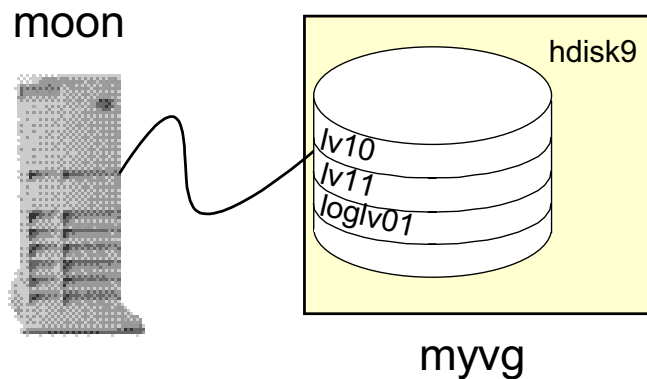
The command **importvg** has a very interesting option, **-L**, which stands for *learn about possible changes*. What does this mean?

Let's discuss an example:

- On system **moon** a volume group **myvg** exists, which contains three logical volumes: **lv10**, **lv11**, **loglv01**.
- The volume group resides on one disk **hdisk9**, which is now moved to another system, **mars**. Note that we do not export **myvg** on system **moon**!
- The volume group **myvg** is now imported on system **mars**, by executing the **importvg** command. Additionally, a new logical volume, **lv99** is created in **myvg**.
- The disk that contains the volume group **myvg**, plus the newly created logical volume **lv99** is now moved back to the system **moon**.

Because we did not export the volume group **myvg** on **moon**, we cannot import the volume group again. Now, how can we fix this problem? This is shown on the next visual.

## importvg -L (2 of 2)



*"Learn about possible changes!"*

```
# importvg -L myvg hdisk9
# varyonvg myvg

==> importvg -L fails, if a name clash is detected
```

© Copyright IBM Corporation 2004

Figure 6-19. importvg -L (2 of 2)

AU1612.0

### Notes:

To import an existing volume group, the command **importvg** offers the option **-L**.

In our example, the following command must be executed to import the volume group **myvg**:

```
# importvg -L myvg hdisk9
```

After executing this command, the new logical volume **lv99** will be recognized by the system.

The volume group must not be active. Additionally the volume group is not automatically varied on, which is a difference to a normal **importvg**.

The command **importvg -L** fails, if a logical volume name clash is detected.

## Next Step

---



© Copyright IBM Corporation 2004

Figure 6-20. Next Step

AU1612.0

### **Notes:**

At the end of the exercise, you should be able to:

- Export a volume group
- Import a volume group

# Checkpoint

---

1. Although everything seems to be working fine, you detect error log entries for disk **hdisk0** in your **rootvg**. The disk is not mirrored to another disk. You decide to replace this disk. Which procedure would you use to migrate this disk?

---

---

2. You detect an unrecoverable disk failure in volume group **datavg**. This volume group consists of two disks that are completely mirrored. Because of the disk failure you are not able to vary on **datavg**. How do you recover from this situation?

---

---

3. After disk replacement you recognize that a disk has been removed from the system but not from the volume group. How do you fix this problem?

---

---

© Copyright IBM Corporation 2004

Figure 6-21. Checkpoint

AU1612.0

## Notes:

---

## Unit Summary

---

- Different procedures are available that can be used to fix disk problems under any circumstance:
  - Procedure 1: Mirrored Disk
  - Procedure 2: Disk still working (rootvg specials)
  - Procedure 3: Total disk failure
  - Procedure 4: Total rootvg failure
  - Procedure 5: Total non-rootvg failure
- exportvg and importvg can be used to easily transfer volume groups between systems

© Copyright IBM Corporation 2004

Figure 6-22. Unit Summary

AU1612.0

### **Notes:**





# Unit 7. Saving and Restoring Volume Groups and Online JFS/JFS2 Backups

## What This Unit Is About

This unit describes how to back up and restore different kinds of volume groups. Additionally, alternate disk installation techniques are introduced.

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Back up and restore the root volume group
- Back up and restore user volume groups
- List different ways of alternate disk installation
- Split an LV mirror to perform an online JFS or JFS2 backup

## How You Will Check Your Progress

Accountability:

- Checkpoint questions
- Activities
- Lab exercise

## Unit Objectives

---

After completing this unit, students should be able to:

- Create, verify, and restore **mksysb** images
- Set up **cloning** using **mksysb** images
- **Shrink** file systems and logical volumes
- Provide **alternate disk installation** techniques
- **Backup** and **restore** non-rootvg volume groups
- Perform an online JFS or JFS2 backup

© Copyright IBM Corporation 2004

Figure 7-1. Unit Objectives

AU1612.0

### **Notes:**

## 7.1 Saving and Restoring the rootvg

# Creating a System Backup: mksysb

```
# smit mksysb
```

## Back Up the System

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

[Entry Fields]

WARNING: Execution of the mksysb command will result in the loss of all material previously stored on the selected output medium. This command backs up only rootvg volume group.

* Backup DEVICE or FILE	[ ]	+/
Create MAP files?	no	+
EXCLUDE files?	no	+
List files as they are backed up?	no	+
Generate new /image.data file?	yes	+
EXPAND /tmp if needed?	no	+
Disable software packing of backup?	no	+
Number of BLOCKS to write in a single output (Leave blank to use a system default)	[ ]	#

© Copyright IBM Corporation 2004

Figure 7-2. Creating a System Backup: mksysb

AU1612.0

## Notes:

The **mksysb** command is used to back up the **rootvg** volume group. It is considered a system backup. You can use this backup to reinstall a system to its original state after it has been corrupted. If you create the backup on tape, the tape is bootable and includes the programs needed to boot into maintenance mode. In maintenance mode, you can access the rootvg and its files.

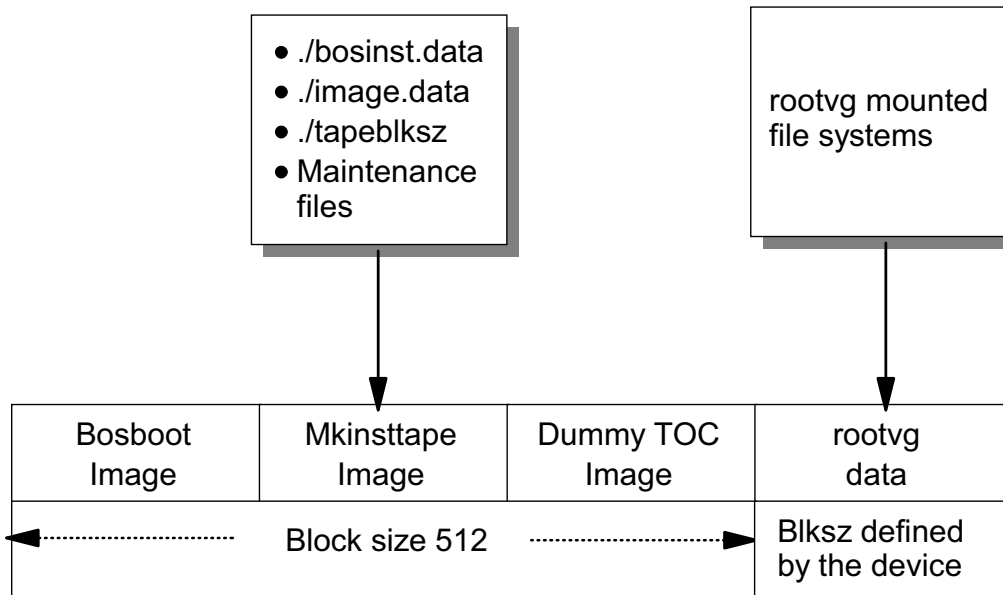
When creating the **mksysb** image, the **/tmp** file system must have at least 8.8 MB free space.

After creating the **mksysb** image, note how many **volume groups** the system has, what disks they are located on, and the **location** of each disk. Hdisk#'s are not retained when restoring the **mksysb** image.

Creating a **mksysb** to a file will create a non-bootable, single-image backup and restore archive containing ONLY **rootvg** jfs and jfs2 mounted file systems.

In AIX Version 5.2, mksysb can be used with the -V option to verify the backup. It verifies the file header of each file on the backup tape and reports any read errors as they occur.

# mksysb Tape Images



© Copyright IBM Corporation 2004

Figure 7-3. mksysb Tape Images

AU1612.0

## Notes:

There will be four images on the **mksysb** tape, and the fourth image will contain only **rootvg** jfs and jfs2 mounted file systems. The following is a description of **mksysb**'s four images.

1. Image #1: The **bosboot** image contains a copy of the system's kernel and specific device drivers, allowing the user to boot from this tape.
2. Image #2: The **mkinsttape** image contains files to be loaded into the RAM file system when booting in maintenance. Example files in this image are **bosinst.data**, **image.data** or **tapeblksz**, which contains the blocksize for the fourth image.
3. Image #3: The dummy image contains a single file containing the words "dummy toc". This image is used to make the **mksysb** tape contain the same number of images as a BOS install tape.
4. Image #4: The **rootvg** image contains data from the **rootvg** volume group (mounted jfs and jfs2 file systems only).

The blocksize for the first three images is set to **512 bytes**. The blocksize for the **rootvg** image is determined by the tape device.

If you are not sure what blocksize is used for the **rootvg** image, restore the file **tapeblksz** from the second image:

```
# chdev -l rmt0 -a block_size=512
# tctl -f /dev/rmt0 rewind
# restore -s2 -xqvf /dev/rmt0.1 ./tapeblksz
# cat tapeblksz
1024
```

In this example the blocksize used in the fourth image is **1024**.

---

## CD or DVD mksysb

---

- Personal system backup
  - Will only boot and install the system where it was created
- Generic backup
  - Will boot and install any platform (rspc, rs6k, chrp)
- Non-bootable VG backup
  - Contains only a VG image (rootvg and non-rootvg)
  - Can install AIX after boot from product CD-ROM (rootvg)
  - Can be source for alt\_disk\_install
  - Can be restored using restvg (non-rootvg)

© Copyright IBM Corporation 2004

Figure 7-4. CD or DVD mksysb

AU1612.0

### **Notes:**

CD (CD-R, CD-RW), DVD (DVD-R, DVD-RAM) are devices supported as mksysb media on AIX 5L.

The three types of CDs (or DVDs) that can be created are listed above.

## Required Hardware and Software for Backup CDs and DVDs

Software	Hardware
GNU & Free Software Foundation, Inc. cdrrecord Version 1.8a5 mkisofs Version 1.5	Yamaha CRW4416S - CD=RW Yamaha CRW8424S - CD-RW Ricoh MP6201SE 6XR-2X - CD-R Panasonic CW-7502-B - CD-R
Jodian System and Software, Inc. CDWrite Version 1.3 mkcdimg Version 2.0	Yamaha CRW4416S - CD=RW Ricoh MP6201SE 6XR-2X - CD-R Panasonic CW-7502-B - CD-R
Youngminds, Inc. MakeDisk Version 1.3-Beta2	Young Minds CD Studio - CD-R
Youngminds, Inc.	Young Minds Turbo Studio - DVD-R
GNU Software	Matsushita LF-D291 - DVD-RAM IBM DVD-RAM

© Copyright IBM Corporation 2003

Figure 7-5. Required Hardware and Software for Backup CDs and DVDs

AU1612.0

### Notes:

Because IBM does not sell or support the software to create CDs, they must be obtained from independent vendors.

The listed drives have been tested by IBM.

The listed software is used in conjunction with the **mkcd** command.



---

## The mkcd Command

---

- **mksysb** and **savevg** images are written to CD-Rs and DVDs using **mkcd**
- Supports ISO09660 and UDF formats
- Requires third party code to create the Rock Ridge file system and write the backup image

© Copyright IBM Corporation 2004

Figure 7-6. The mkcd Command

AU1612.0

### Notes:

This code must be linked to `/usr/sbin/mkrr_fs` (for creating the Rock Ridge format image) and `/usr/sbin/burn_cd` (for writing to the CD-R or DVD-RAM device). For example, if you are using Jodian software, you will need to create the following links:

```
ln -s /usr/samples/oem_cdwriters/mkrr_fs_gnu /usr/sbin/mkrr_fs
ln -s /usr/samples/oem_cdwriters/burn_cd_gnu_dvdram /usr/sbin/burn_cd
```

The process for creating a mksysb CD using the mkcd command is:

1. If file systems or directories are not specified, they will be created by mkcd and removed at the end of the command (unless the `-R` or `-S` flags are used). mkcd will create following file systems:
  - `/mkcd/mksysb_image`  
Contains a mksysb image. Enough space must be free to hold the mksysb.

- /mkcd/cd\_fs  
Contains CD file systems structures. At least 645 MB of free space is required (up to 8.8 GB for DVD).
- /mkcd/cd\_image  
Contains final the CD image before writing to CD-R. At least 645 MB of free space is required (up to 8.8 GB for DVD).

The /mkcd/cd\_fs and /mkcd/cd\_image may be required to have 8.8 GB of free space each, depending how big the mksysb is.

**Note:** The /mkcd/cd\_images (with an 's') may need to be even larger than 8.8 GB or 645 MB if the -R or -S flags were specified (if it is multi-volume), because there must be sufficient space to hold each volume.

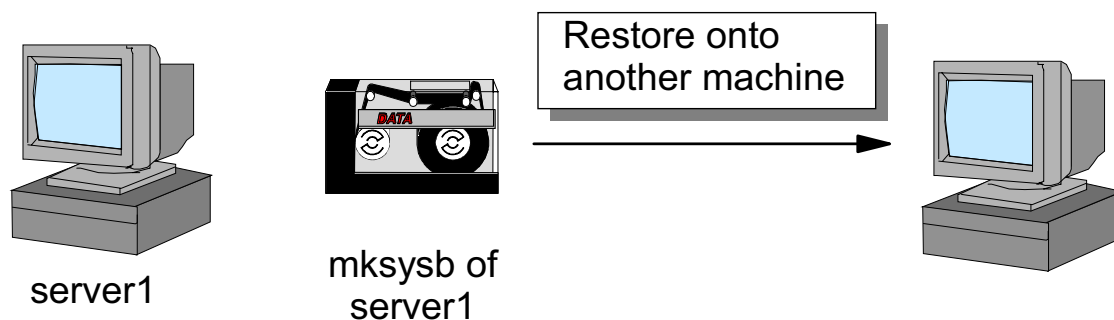
User provided file systems or directories can be NFS mounted.

The file systems provided by the user will be checked for adequate space and an error will be given if there is not enough space. Write access will also be checked.

2. If a mksysb image is not provided, mkcd calls mksysb, and stores the image in the directory specified with the -M flag or in /mkcd/mksysb\_image.
3. The mkcd command creates the directory structure and copies files based on the cdfs.required.list and the cdfs.optional.list files.
4. Device images are copied to ./install/ppc or ./installp if the -G flag is used or the -I flag is given (with a list of images to copy).
5. The mksysb image is copied to the file system. It determines the current size of the CD file system at this point, so it knows how much space is available for the mksysb. If the mksysb image is larger than the remaining space, multiple CDs are required. It uses dd to copy the specified number of bytes of the image to the CD file system. It then updates the volume ID in a file. A variable is set from a function that determines how many CDs are required to hold the entire mksysb image.
6. The mkcd command then calls the mkrr\_fs command to create a RockRidge file system and places the image in the specified directory.
7. The mkcd command then calls the burn\_cd command to create the CD.

If multiple CDs are required, the user is instructed to remove the CD and put the next one in and the process continues until the entire mksysb image is put on the CDs. Only the first CD supports system boot.

## Verifying a System Backup After mksysb Completion (1 of 2)



- The only method to verify that a system backup will correctly restore with no problems is to actually restore the mksysb onto another machine.
- This should be done to test your company's DISASTER RECOVERY PLAN.

© Copyright IBM Corporation 2003

Figure 7-7. Verifying a System Backup After mksysb Completion (1 of 2)

AU1612.0

### Notes:

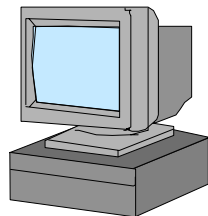
After creating the **mksysb** tape, you must verify that the image will correctly restore with no problems.

The ONLY method to verify this is to restore the **mksysb** onto another machine.

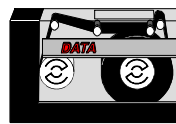
This must be part of a company's **disaster recovery plan**. A disaster is a situation where you have to reinstall a system from scratch. The first step will be to reinstall the operating system, that means to restore the **mksysb** image.

How can you verify the **mksysb** tape if you do not have a second machine available?

## Verifying a System Backup (2 of 2)



server1

mksysb of  
server1

- Data Verification:

```
# tctl -f /dev/rmt0 rewind
# restore -s4 -Tqvf /dev/rmt0.1 > /tmp/mksysb.log
```

- Boot Verification:

Boot from the tape without restoring any data.  
WARNING: Check the PROMPT field in bosinst.data!

© Copyright IBM Corporation 2004

Figure 7-8. Verifying a System Backup After mksysb Completion (2 of 2)

AU1612.0

### Notes:

If you cannot test the installability of your image, execute the following tasks:

1. Do a **data verification**. Test that you can access the **rootvg** image without any errors. The option **-T** in the **restore** command indicates that a **table of contents** should be created.
2. Do a **boot verification**. Shut down a system and boot from the **mksysb** tape. Do not restore any data from the **mksysb** tape.

Having the **PROMPT** field in the **bosinst.data** file set to **no**, causes the system to begin the **mksysb** restore automatically using preset values with no user invention.

If you want to check the state of the **PROMPT** field, restore the **bosinst.data** file from the image:

```
# chdev -l rmt0 -a block_size=512
# tctl -f /dev/rmt0 rewind
# restore -s2 -xqvf /dev/rmt0 ./bosinst.data
```

If the state is **no** it can be changed to **yes** during the boot process. After answering the prompt to select a console during the startup process, a **rotating character** will be seen in the lower left of the screen. As soon as this character appears, type **000** and press Enter. This will set the prompt variable to **yes**.

## mksysb Control File: bosinst.data

```

control_flow:
    CONSOLE =
    INSTALL_METHOD = overwrite
    PROMPT = yes
    EXISTING_SYSTEM_OVERWRITE = yes
    INSTALL_X_IF_ADAPTER = yes
    RUN_STARTUP = yes
    RM_INST_ROOTS = no
    ERROR_EXIT =
    CUSTOMIZATION_FILE =
    TCB = no
    INSTALL_TYPE =
    BUNDLES =
    SWITCH_TO_PRODUCT_TAPE =
    RECOVER_DEVICES = yes
    BOSINST_DEBUG = no

target_disk_data:
    LOCATION =
    SIZE_MB =
    HDISKNAME =

locale:
    BOSINST_LANG =
    CULTURAL_CONVENTION =
    MESSAGES =
    KEYBOARD = . . .

```

© Copyright IBM Corporation 2004

Figure 7-9. mksysb Control File: bosinst.data

AU1612.0

### Notes:

The **bosinst.data** file controls the restore process on the target system. It allows the administrator to specify requirements at the target system and how the user interacts with the target system.

The system backup utilities copy the **/bosinst.data** as the first file in the **rootvg** image on the **mksysb** tape. If this file is **not** in the root directory, the **/usr/lpp/bosinst/bosinst.template** is copied to **/bosinst.data**.

Normally there is no need to change the stanzas from **bosinst.data**. One exception is to enable an **unattended** installation:

To enable an unattended installation process of the **mksysb** tape, edit the **bosinst.data** as follows:

- Specify the console on the **CONSOLE** line, for example **CONSOLE=/dev/tty0** or **CONSOLE=/dev/lft0**.
- Set **PROMPT=no**, to disable installation menus.

Three lines were added to the control\_flow stanza in AIX 4.2: to Other lines in the control\_flow stanza include:

- The option **SWITCH\_TO\_PRODUCT\_TAPE** must be set to **yes** if you are **cloning** a system from a **product tape**. Cloning is introduced later in this unit.
- The option **RECOVER\_DEVICES** allows the choice to recover the **CuAt** (customized attributes) ODM class, which contains attributes like network addresses, static routes, tty settings and more. If the **mksysb** tape is used to clone systems, this stanza could be set to **no**. In this case, the **CuAt** will not be restored on the target system. If you are restoring the **mksysb** on the same system, do not change the default value, which is **yes**.
- The option **BOSINST\_DEBUG** specifies whether to show debug information during the installation process. The value **yes** will send **set -x** debug output to the screen during the installation. Possible values are **no** (default) and **yes**.

You can overwrite the default value of **no** debug information during the installation process. If the rotating character appears on the lower left screen during the installation, type in **911**. This number indicates to the installation routines to turn on debug information.

If you do not want to use the **mksysb's bosinst.data** during the installation, you can create one that can be read from a floppy. Execute the following steps:

1. Create a file named **signature** in the following way:

```
# echo "data" > signature
```

2. Edit your **bosinst.data** file and change the appropriate stanzas

3. Create a floppy diskette with the following command:

```
# ls ./bosinst.data ./signature | backup -iqv
```

Before restoring the **mksysb** insert this diskette into the floppy drive.

## Restoring a mksysb (1 of 2)

Boot from AIX bootable media

Welcome to Base Operating System  
Installation and Maintenance

Type the number of your choice and press Enter. Choice is indicated by >>.

- 1 Start Install Now With Default Settings
- 2 Change/Show Installation Settings and Install
- >> 3 Start Maintenance Mode for System Recovery

Maintenance

Type the number of your choice and press Enter.

- 1 Access A Root Volume Group
- 2 Copy a System Dump to Removable Media
- 3 Access Advanced Maintenance Functions
- >> 4 Install from a System Backup

Choose Tape Drive

Type the number of the tape drive containing the system backup to be installed and press Enter.

- |      | Tape Drive        | Path Name |
|------|-------------------|-----------|
| >> 1 | tape/scsi/4mm/2GB | /dev/rmt0 |

© Copyright IBM Corporation 2004

Figure 7-10. Restoring a mksysb (1 of 2)

AU1612.0

### Notes:

Restoring a **mksysb** is very easy. Follow these steps:

- Boot the system (as you learned in this course) from an AIX CD, an AIX product tape or the **mksysb** tape.
- From the **Installation and Maintenance** menu, select option **3**.
- From the **Maintenance** menu, select option **4**.

Choose the drive that contains the **mksysb** image.

The AIX 5.3 mksysb screen will also include an “Erase Disks” option on the “Maintenance” menu. This will take the user to the **Select Disk(s) That You Want to Erase** menu.

Continuation from there will take the user to the **Erasure Options for Disks** menu. This menu allows you to select erasure pattern ( write ‘0’, write ‘ff’, and so forth) then run the erase utility and exit.



## Restoring a mksysb (2 of 2)

```

Welcome to Base Operating System
Installation and Maintenance

Type the number of your choice and press Enter. Choice is indicated by >>.
  1 Start Install Now With Default Settings
>> 2 Change/Show Installation Settings and Install
  3 Start Maintenance Mode for System Recovery

```

```

System Backup Installation and Settings

Type the number of your choice and press Enter.

  1 Disk(s) where you want to install          hdisk0
  2 Use Maps                                  No
  3 Shrink File systems                       No
  0 Install with the settings listed above

```

© Copyright IBM Corporation 2004

Figure 7-11. Restoring a mksysb (2 of 2)

AU1612.0

### Notes:

- After selecting the tape drive (and a language, which is not shown on the visuals) you will return to the **Installation and Maintenance** menu. Now select option **2**.
- From the **System Backup Installation and Settings** menu, select **1** and select the disks where you want to install.

Be sure to select all physical volumes required for the volume group. This is especially important if mirroring has been set up.

Two other options can be enabled in this menu:

1. The option **Use Maps** indicates that map files must be used. These map files allow an exact placement of the physical partitions from **rootvg** on the disks, as specified in the **mksysb** image. The default value is no.
2. The option **Shrink Filesystems** allows you to install the file systems using the minimum required space. The default value is no. If yes, all file systems are shrunk. So remember after the restore, evaluate the current file system sizes. You might need to increase their sizes. You will learn later, how to shrink selected file systems and logical volumes.

- At the end, select option **0** (Install with the settings above). Your **mksysb** image will be restored.
- After the restore is complete, the system reboots.

The total restore time varies from system to system. A good rule of thumb is twice the amount of time it took to create the **mksysb**.

The AIX 5.3 option 2 screen will look like the following:

### System Backup Installation and Settings

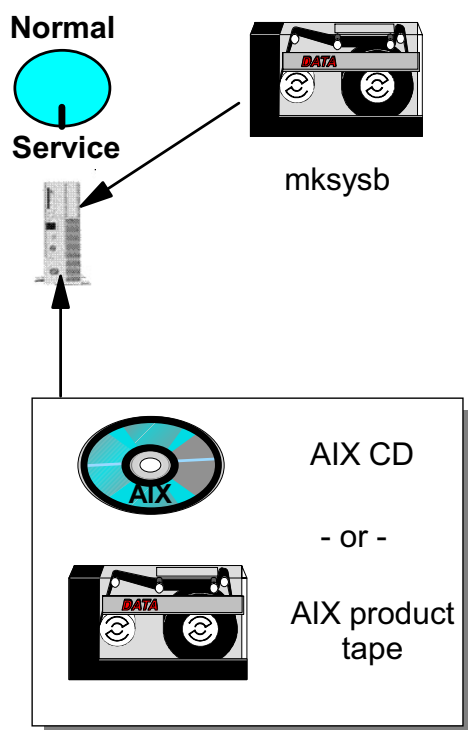
Either type 0 and press Enter to install with the current settings, or type the number of the setting you want to change and press Enter.

Setting:	Current Choice(s):
----------	--------------------

- |   |        |
|---|--------|
| 1 Disk(s) where you want to install ..... | hdisk0 |
| Use Maps.....                             | No     |
| 2 Shrink File Systems.....                | No     |
| 3 Import User Volume Groups.....          | No     |
| 4 Recover Devices.....                    | No     |

>>> 0 Install with the settings listed above.

# Cloning Systems Using mksysb Tapes



1. Insert the mksysb tape and the AIX CD (same AIX level!)
2. Boot from the AIX CD (\*)
3. "Install from a System Backup":  
Missing device support is installed from the AIX CD

(\*): If no AIX CD available, use an AIX product tape, but check bosinst.data:

```
bosinst.data:
SWITCH_TO_PRODUCT_TAPE=yes
```

© Copyright IBM Corporation 2004

Figure 7-12. Cloning Systems Using mksysb Tapes

AU1612.0

## Notes:

Beginning in AIX 5.2, all devices and kernel support are installed by default during the base operating system (BOS) installation process. If the "Enable System Backups to install any system" selection in the Install Software menu is set to yes, you can create a mksysb image that boots and installs supported systems. Verify that your system is installed with all devices and kernel support by typing the following command:

```
# grep ALL_DEVICES_KERNELS /bosinst.data
```

Output similar to the following displays:

```
ALL_DEVICES_KERNELS = yes
```

If all device and kernel support was not installed, you will need to boot from the appropriate product media for your system at the same maintenance level of BOS as the installed source system on which the mksysb tape was created.

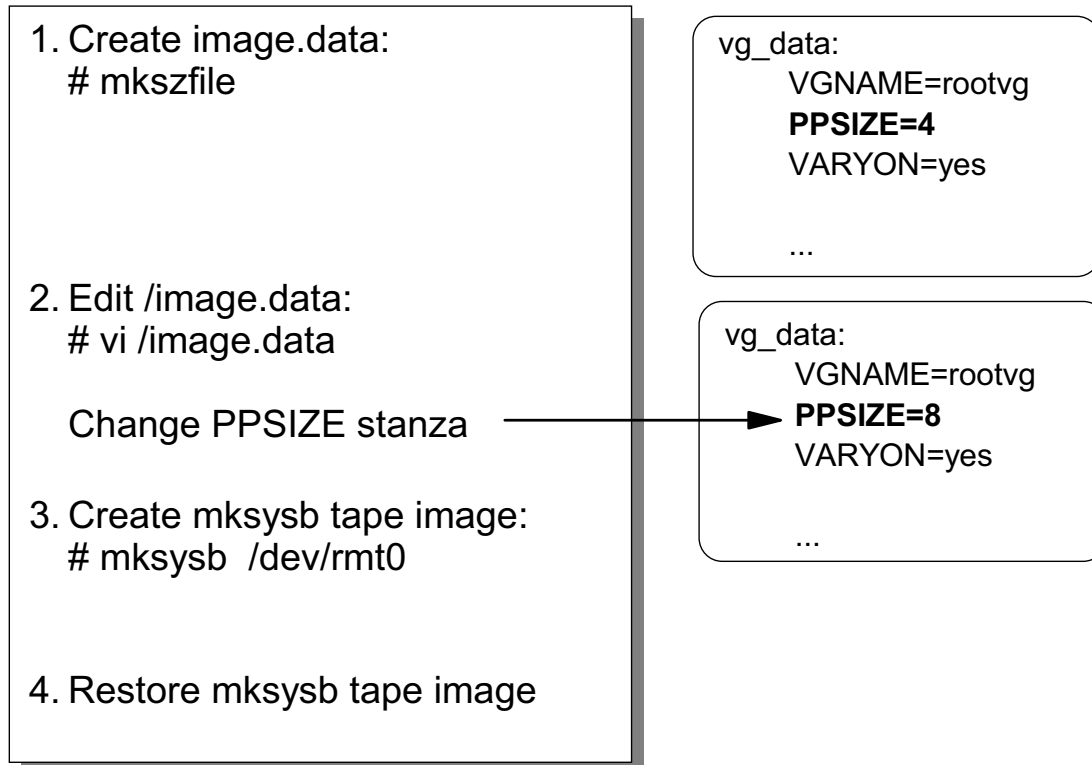
In this scenario, you will do the following:

1. Insert the **mksysb** tape and the AIX CD into the target system. Note that both **must** have the same AIX level. If you have, for example, an AIX 5.2.0 **mksysb** image, you must use the AIX 5.2.0 CD.
2. Boot your system from the CD, **not** from the **mksysb** image.
3. Start the maintenance mode and install the system from the system backup (the menus have been shown on the last two pages).

After the mksysb installation completes, the installation program automatically installs additional devices and the kernel (uniprocessor or multiprocessor) on your system, using the original product media you booted from.

If you work with an AIX product tape, you need to set the stanza **SWITCH\_TO\_PRODUCT\_TAPE** in **bosinst.data** to **yes**. Anyway it is preferable to use the AIX CD. If the installation tape is used, the installation tape and the **mksysb** tape may need to be switched back and forth a few times during the restoration.

## Changing the Partition Size in rootvg



© Copyright IBM Corporation 2004

Figure 7-13. Changing the Partition Size in rootvg

AU1612.0

### Notes:

What can you do if you have to increase the **physical partition size** in your **rootvg**? Remember: if your **rootvg** has a physical partition size of **4 MB**, the maximum disk space is **4 GB** (4 MB \* 1016 partitions). In this case you cannot use a **8 GB** disk (you can, but you waste 50 percent of the disk space).

To solve this situation, execute the following steps:

1. Execute the command **mkszfile**:

```
# mkszfile
```

This command creates a file **image.data** in the root directory.

2. Edit the file **/image.data**. Locate the stanza **vg\_data** and change the attribute **PPSIZE** to the desired value, for example to **8 MB**.
3. Create a new **mkysyb** image with the following command:

```
# mkysyb /dev/rmt0 (or whatever your tape device is)
```

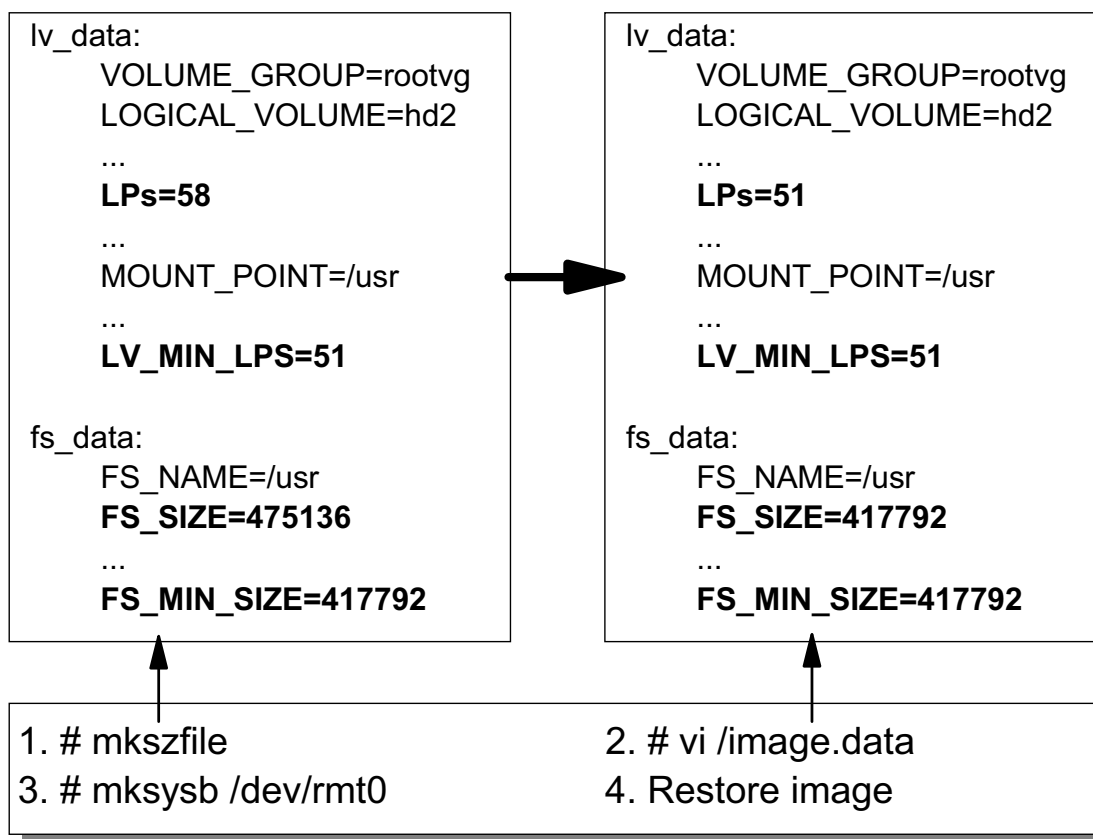
If you use smit to create the **mksysb** image, be sure to answer “no” to “Generate new /image.data file?” Reason:

Smit will use mksysb -i otherwise which will create a new image.data file overwriting your modifications.

When the **mksysb** image is complete, verify the image, as learned in this unit, before restoring it.

4. Restore the **mksysb** image on the system. Your **rootvg** will be allocated with the changed partition size.

## Reducing a File System in rootvg



© Copyright IBM Corporation 2004

Figure 7-14. Reducing a File System in rootvg

AU1612.0

### Notes:

Another very nice thing you can do with **mksysb** images is to reduce the file system size of **one** file system. Remember that you can shrink **all** file systems when restoring the **mksysb**. The advantage of this technique is that you shrink only one selected file system.

In the following example, we change the **/usr** file system:

1. Execute the **mkszfile** command to create a file **/image.data**:

```
# mkszfile
```

2. Change the file **/image.data** in the following way:

- You can either increase or decrease the number of logical partitions needed to contain the file system data.

In the example we decrease the number of logical partitions (LPs=58 to LPs=51) to the minimum required size (LV\_MIN\_LPS=51). **Note:** If you enter a value that is less than the minimum size, the reinstallation process will fail.

- After reducing the number of logical partitions, you must change the file system size. In our example we change the file system size to the minimum required size (FS\_SIZE=475136 to FS\_SIZE=417792), indicated by FS\_MIN\_SIZE. Note that FS\_SIZE and FS\_MIN\_SIZE are in 512-byte blocks.
3. After changing **/image.data**, create a new **mksysb** tape image. Verify the image as you learned earlier in this unit.
  4. Finally restore the image.



---

# Let's Review: Working with mksysb Images

---



© Copyright IBM Corporation 2004

Figure 7-15. Let's Review: Working with mksysb Images

AU1612.0

## **Notes:**

### **Please answer the following questions:**

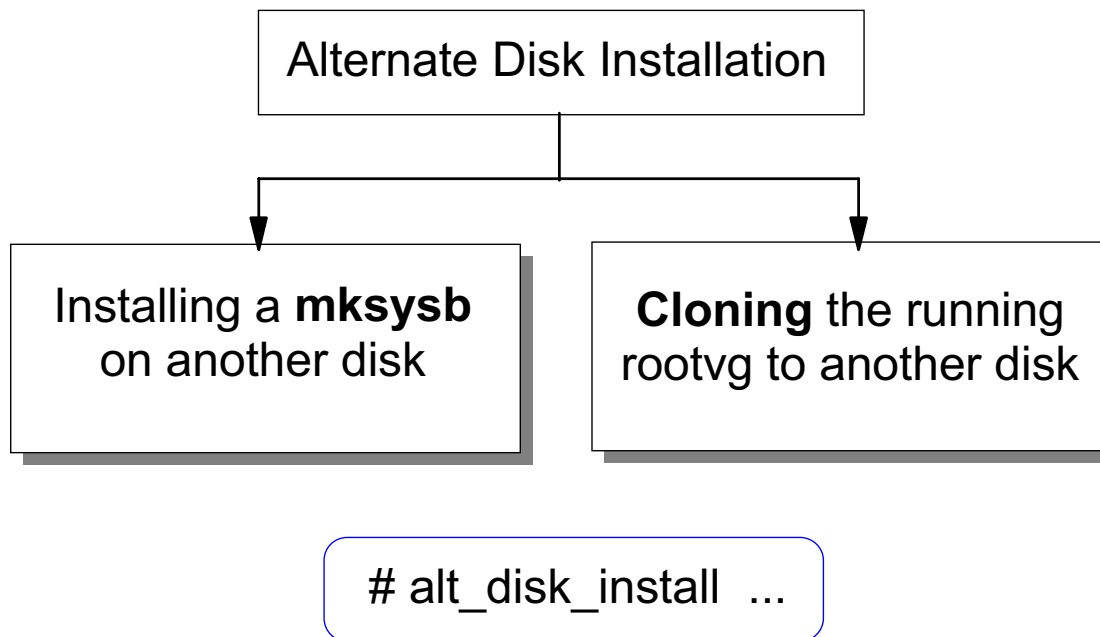
- \_\_\_ 1. True or False: A **mksysb** image contains a backup of all volume groups.  
\_\_\_\_\_
- \_\_\_ 2. How can you determine the blocksize of the fourth image in a **mksysb** tape image?  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_
- \_\_\_ 3. Describe the meaning of the attribute **RECOVER\_DEVICES** from **bosinst.data**.  
\_\_\_\_\_  
\_\_\_\_\_
- \_\_\_ 4. True or False: Cloning AIX systems is only possible if the source and target system use the **same** hardware architecture.

\_\_\_ 5. What happens if you execute the command **mkszfile**?

---

## 7.2 Alternate Disk Installation

# Alternate Disk Installation



© Copyright IBM Corporation 2004

Figure 7-16. Alternate Disk Installation

AU1612.0

## Notes:

Alternate disk installation, available in AIX 4.3 and subsequent versions of the operating system, allows installing the system while it is still up and running, allowing installation or upgrade time to be decreased considerably. It also allows large facilities to manage an upgrade because systems can be installed over a longer period of time while the systems are running at the same version. The switchover to the new version can then happen at the same time.

Alternate disk installation can be used in one of two ways:

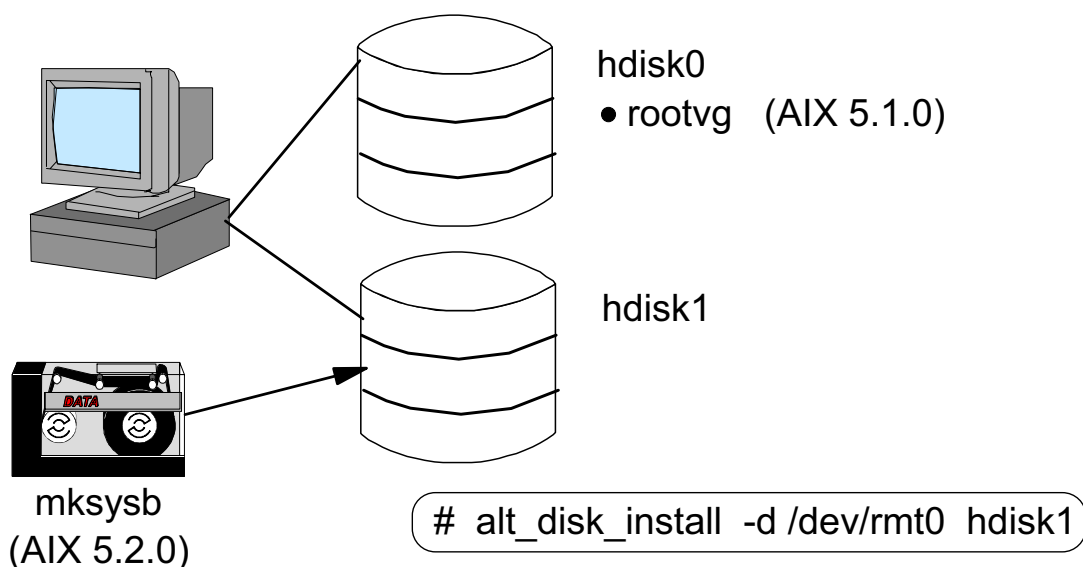
1. Installing a **mksysb** image on another disk.
2. Cloning the current running **rootvg** to an alternate disk.

The command that is used for alternate disk installation is **alt\_disk\_install**. This command runs on AIX 4.1.4 and higher systems.

Both techniques are introduced on the following pages.

The fileset **bos.alt\_disk\_install** must be installed on the system.

## Alternate mksysb Disk Installation (1 of 2)



- Installs a 5.2.0 **mksysb** on **hdisk1** ("second rootvg")
- Bootlist will be set to alternate disk (**hdisk1**)
- Changing the bootlist allows to boot different AIX levels (**hdisk0** boots AIX 5.1.0, **hdisk1** boots AIX 5.2.0)

© Copyright IBM Corporation 2004

Figure 7-17. Alternate mksysb Disk Installation (1 of 2)

AU1612.0

### Notes:

Alternate **mksysb** installation involves installing a **mksysb** image that has already been created from another system onto an alternate disk of the target system.

In the example, an AIX 5.2.0 **mksysb** tape image is installed on an alternate disk, **hdisk1** by executing the following command:

```
# alt_disk_install -d /dev/rmt0 hdisk1
```

The system contains now two **rootvgs** on different disks. In the example, one **rootvg** has an AIX level 5.1.0 (**hdisk0**), one has an AIX level 5.2.0 (**hdisk1**).

The **alt\_disk\_install** command changes the boot list by default. During the next reboot, the system will boot from the new **rootvg**. If you do not want to change the boot list, use the option **-B** from **alt\_disk\_install**.

By changing the boot list you determine, which AIX level you want to boot.

Alternate **mksysb** disk installation requires a **mksysb** image created on a system running AIX 4.3 or subsequent versions of the operating system.

The AIX 5L Version 5.3 has implemented a number of changes to make the alt\_disk\_install operations easier to use, document, and maintain. The following functional changes have been implemented:

alt\_disk\_install has been partitioned into separate modules with separate syntax based on operation and functionality.

A library of common functions that can be accessed by the modules has been implemented.

Error checking and robustness of existing alt\_disk\_install operations has been improved.

Documentation has been improved by creating a separate man page for each module (currently there is one extremely large man page).

The following three new commands have been added:

alt\_disk\_copy will create copies of rootvg on an alternate set of disks.

alt\_disk\_mksysb will install an existing mkysb on an alternate set of disks.

alt\_rootvg\_op will perform Wake, Sleep, and Customize operations.

Also, a new library, alt\_lib, has been added that serves as a common library shared by all alt\_disk\_install commands. The alt\_disk\_install module will continue to ship as a wrapper to the new modules. However, it will not support any new functions, flags or features.

The following table displays how the existing operation flags for alt\_disk\_install will map to the new modules. The alt\_disk\_install command will now call the new modules after printing an attention notice that it is obsolete. All other flags will apply as currently defined.

alt_disk_install Command arguments	New commands
-C <args> <disks>	alt_disk_copy <args> -d <disks>
-d <mkysb> <args> <disks>	alt_disk_mkysb -m <mkysb> <args> -d <disks>
-W <args> <disk>	alt_rootvg_op -W <args> -d <disk>
-S <args>	alt_rootvg_op -S <args>
-P2 <args> <disks>	alt_rootvg_op -C <args> -d <disks>
-X <args>	alt_rootvg_op -X <args>
-v <args> <disk>	alt_rootvg_op -v <args> -d <disk>
-q <args> <disk>	alt_rootvg_op -q <args> -d <disk>

## Alternate mksysb Disk Installation (2 of 2)

```
# smit alt_mksysb
```

Install mksysb on an Alternate Disk

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

	[Entry Fields]	
* Target Disk(s) to install	[hdisk1]	+
* Device or image name	[/dev/rmt0]	+
Phase to execute	all	+
image.data file	[]	/
Customization script	[]	/
Set bootlist to boot from this disk on next reboot?	yes	+
Reboot when complete?	no	+
Verbose output?	no	+
Debug output?	no	+
resolv.conf file	[]	/

© Copyright IBM Corporation 2004

Figure 7-18. Alternate mksysb Disk Installation (2 of 2)

AU1612.0

### Notes:

To execute alternate **mksysb** disk installation, you can either work with the command **alt\_disk\_install** or the smit fastpath **smit alt\_mksysb**.

The installation on the alternate disk is broken into three phases:

1. **Phase 1** creates the **altinst\_rootvg** volume group, the **alt\_logical** volumes, the **/alt\_inst** file systems and restores the **mksysb** data.
2. **Phase 2** runs any specified **customization script** and copies a **resolv.conf** file if specified.
3. **Phase 3** umounts the **/alt\_inst** file systems, renames the file systems and logical volumes and varies off the **altinst\_rootvg**. It sets the **boot list** and reboots if specified.

You can run each phase separately. You must use phase 3 to get a volume group that is a usable rootvg.

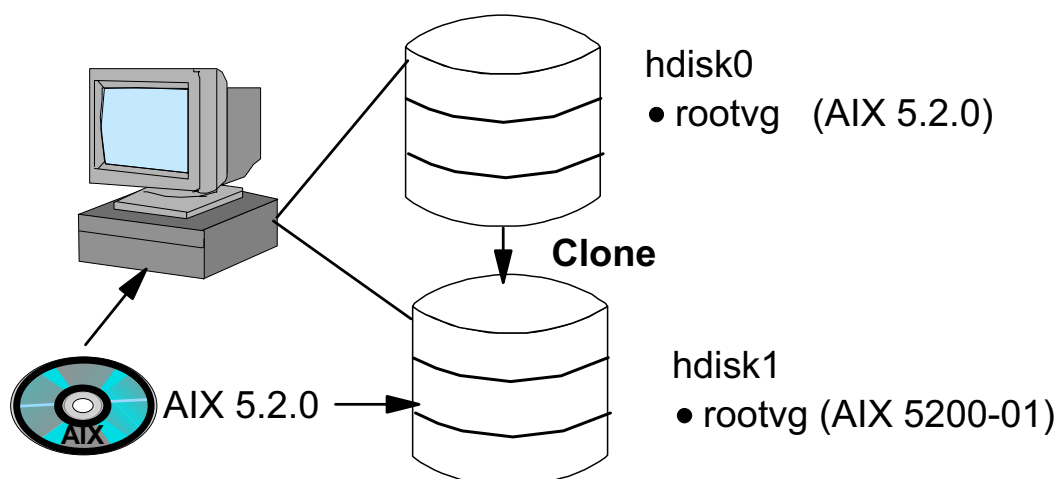
**Important:**

The **mksysb** image used for the installation must be created on a system that has either the same hardware configuration as the target system, or must have all the device and kernel support installed for a different machine type or platform. In this case the following filesets must be contained in the **mksysb**:

- devices.\*
- bos.mp
- bos.up
- bos.64bit (if necessary)



## Alternate Disk rootvg Cloning (1 of 2)



```
# alt_disk_install -C -b update_all -l /dev/cd0 hdisk1
```

- Creates a copy of the current rootvg ("clone") on hdisk1
- Installs a maintenance level on clone (AIX 5200-01)
- Changing the bootlist allows you to boot different AIX levels (hdisk0 boots AIX 5.2.0, hdisk1 boots AIX 5200-01)

© Copyright IBM Corporation 2004

Figure 7-19. Alternate Disk rootvg Cloning (1 of 2)

AU1612.0

### Notes:

Cloning the **rootvg** to an alternate disk can have many advantages. One advantage is having an online backup available, in case of a disaster. Another benefit of **rootvg** cloning is in applying new maintenance levels or updates. A copy of the **rootvg** is made to an alternate disk (in our example **hdisk1**), then a maintenance level is installed on the copy. The system runs uninterrupted during this time. When it is rebooted, the system will boot from the newly updated **rootvg** for testing. If the maintenance level causes problems, the old **rootvg** can be retrieved by simply resetting the **boot list** and rebooting.

In the example we clone the current **rootvg** which resides on **hdisk0** to the alternate disk **hdisk1**. Additionally, a new maintenance level will be applied to the cloned version of AIX.

# Alternate Disk rootvg Cloning (2 of 2)

# smit alt\_clone

Clone the rootvg to an Alternate Disk

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

	[Entry Fields]	
* Target Disk(s) to install	[hdisk1]	+
Phase to execute	all	+
image.data file	<input type="text"/>	/
Exclude list	<input type="text"/>	/
Bundle to install	[update_all]	+
Filesets to install	<input type="text"/>	
...		
Fixes to install	<input type="text"/>	
Directory or Device with images	[/dev/cd0]	
Customization script	<input type="text"/>	/
Set bootlist to boot from this disk on next reboot?	yes	+
Reboot when complete?	no	+
...		

© Copyright IBM Corporation 2004

Figure 7-20. Alternate Disk rootvg Cloning (2 of 2)

AU1612.0

## Notes:

The smit fastpath for alternate disk rootvg cloning is **smit alt\_clone**.

The target disk in the example is **hdisk1**, that means the **rootvg** will be copied to that disk. When you specify a bundle, a fileset or a fix, the installation or the update takes place on the clone, not in the original **rootvg**.

By default the **boot list** will be set to the new disk.

Changing the boot list allows you to boot from the original **rootvg** or the cloned **rootvg**.

## Removing an Alternate Disk Installation

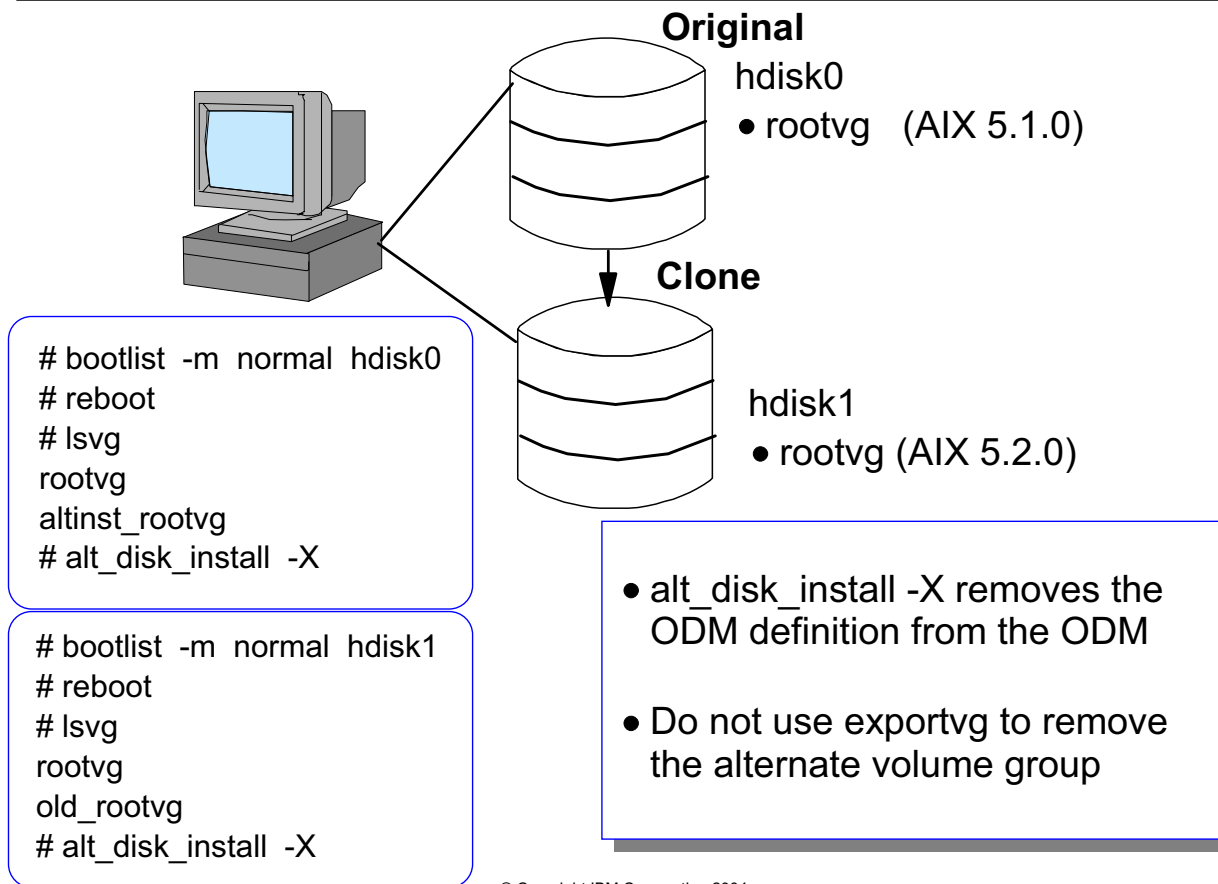


Figure 7-21. Removing an Alternate Disk Installation

AU1612.0

### Notes:

If you have created an alternate **rootvg** with **alt\_disk\_install**, but no longer wish to use it, boot your system from the original disk (in our example, **hdisk0**).

When executing **lsvg** to list the volume groups in the system, the alternate **rootvg** is shown with the name **altinst\_rootvg**.

If you want to remove the alternate **rootvg**, do not use the **exportvg** command. Simply run the following command:

```
# alt_disk_install -X
```

**This command removes the altinst\_rootvg definition from the ODM database.**

If **exportvg** is run by accident, you must re-create the **/etc/filesystems** file before rebooting the system. The system will not boot without a correct **/etc/filesystems**.

If you have created an alternate **rootvg** with **alt\_disk\_install**, and no longer wish to use the original disk, boot your system from the cloned disk (in our example, **hdisk1**).

When executing **lsvg** to list the volume groups in the system, the alternate **rootvg** is shown with the name **old\_rootvg**.

If you want to remove the original **rootvg**, do not use the **exportvg** command. Simply run the following command:

```
# alt_disk_install -X
```

**This command removes the old\_rootvg definition from the ODM database.**

If **exportvg** is run by accident, you must re-create the **/etc/filesystems** file before rebooting the system. The system will not boot without a correct **/etc/filesystems**.

---

# Let's Review: Alternate Disk Installation

---



© Copyright IBM Corporation 2004

Figure 7-22. Let's Review: Alternate Disk Installation

AU1612.0

## **Notes:**

Answer the following review questions:

1. Name the two ways alternate disk installation can be used.

---

---

2. At what version of AIX can an alternate mkysb disk installation occur?

---

3. What are the advantages of alternate disk rootvg cloning?

---

---

---

4. How do you remove an alternate rootvg?

---

5. Why not use **exportvg**?

---

---

## 7.3 Saving and Restoring non-rootvg Volume Groups

# Saving a non-rootvg Volume Group

```
# smit savevg
```

## Back Up a Volume Group to Tape/File

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

[Entry Fields]

WARNING: Execution of the savevg command will result in the loss of all material previously stored on the selected output medium.

* Backup DEVICE or FILE	[/dev/rmt0]	+/
* VOLUME GROUP to back up	[datavg]	+
List files as they are backed up?	no	+
Generate new vg.data file?	yes	+
Create MAP files?	no	+
EXCLUDE files?	no	+
EXPAND /tmp if needed?	no	+
Disable software packing of backup?	no	+
Number of BLOCKS to write in a single output (Leave blank to use a system default)	[ ]	#

© Copyright IBM Corporation 2004

Figure 7-23. Saving a non-rootvg Volume Group

AU1612.0

## Notes:

The **savevg** command allows backups of non-rootvg volume groups. This backup contains the complete definition for all logical volumes and file systems and the corresponding data. In case of a disaster where you have to restore the complete volume group, this backup offers the fastest way to recover the volume group.

When executing the **savevg** command, the volume group must be varied-on and all file systems must be mounted.

In the example we save the volume group **datavg** to the tape device **/dev/rmt0**. The command that **smit** executes is the following:

```
# savevg -i -f/dev/rmt0 datavg
```

The option **-i** indicates the **mkvgdata** command is executed before saving the data. This command behaves like **mkszfile**. It creates a file **vgname.data** (in our example the name is **datavg.data**) that contains information about the volume group. This file is located in **/tmp/vgdata/vgname**, for example, **/tmp/vgdata/datavg**.



## savevg/restvg Control File: vgname.data

```
# mkvgdata datavg
# vi /tmp/vgdata/datavg/datavg.data
```

```
vg_data:
  VGNAME=datavg
  PPSIZE=8
  VARYON=yes

lv_data:

  LPS=128

  LV_MIN_LPS=128

fs_data:

  ...
```

```
# savevg -f /dev/rmt0 datavg
```

© Copyright IBM Corporation 2004

Figure 7-24. savevg/restvg Control File: vgname.data

AU1612.0

### Notes:

If you want to change characteristics in a user volume group, execute the following steps:

1. Execute the command **mkvgdata**. This command generates a file **/tmp/vgdata/vgname/vgname.data**. In our example the filename is **/tmp/vgdata/datavg/datavg.data**.
2. Edit this file and change the corresponding characteristic. In the example we change the **number of logical partitions** in a logical volume.
3. Finally save the volume group. If you use **smit**, set “Generate new vg.data file?” to “NO” or **smit** will overwrite your changes.

To make the changes active, this volume group backup must be restored. Here is one way how you handle this:

1. Unmount all file systems.
2. Varyoff the volume group.
3. Export the volume group, using **exportvg**.
4. Restore the volume group, using the **restvg** command.

The **restvg** command is explained on the next page.

# Restoring a non-rootvg Volume Group

```
# smit restvg
```

## Remake a Volume Group

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

	[Entry Fields]	
* Restore DEVICE or FILE	[/dev/rmt0]	/+
SHRINK the file systems?	no	+
Recreate logical volumes and filesystems only?	no	+
PHYSICAL VOLUME names	[ ]	+
(Leave blank to use the PHYSICAL VOLUMES listed in the vname.data file in the backup image)		
Use existing MAP files?	yes	+
Physical partition SIZE in megabytes	[ ]	+#
(Leave blank to have the SIZE determined based on disk size)		
Number of BLOCKS to read in a single input	[ ]	#
(Leave blank to use a system default)		
Alternate vg.data file	[ ]	/
(Leave blank to use vg.data stored in backup image)		

© Copyright IBM Corporation 2004

Figure 7-25. Restoring a non-rootvg Volume Group

AU1612.0

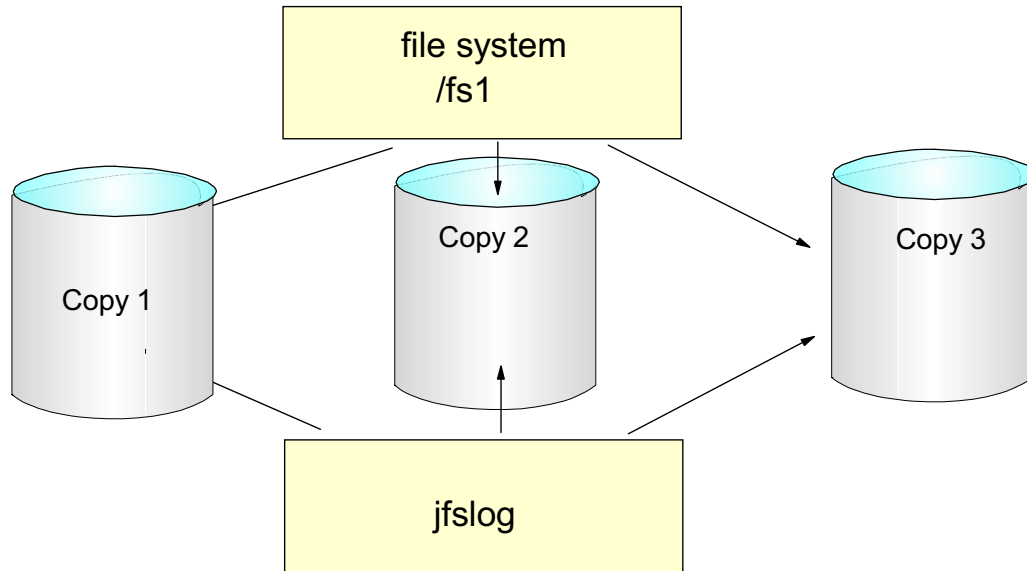
## Notes:

The **restvg** command restores the user volume group and all its containers and files, as specified in **/tmp/vgdata/vgname/vgname.data**. In our example we restore the volume group from the tape device.

Note that you can specify a partition size for the volume group. If not specified, **restvg** uses the best value for the partition size, dependent upon the largest disk being restored to. If this is not the same as the size specified in the **vgname.data** file, the number of partitions in each logical volume will be appropriately altered with respect to the new partition size.

## 7.4 Online JFS and JFS2 Backup; JFS2 Snapshot; VG Snapshot

# Online JFS Backup



```
# lsvg -l newvg
newvg:
LV NAME      TYPE    LPs  PPs  PVs  LV STATE  MOUNT
POINT
loglv00     jfslog   1    3    3    open/syncd  N/A
lv03        jfs      1    3    3    open/syncd  /fs1
```

© Copyright IBM Corporation 2004

Figure 7-26. Online jfs and jfs2 Backup

AU1612.0

## Notes:

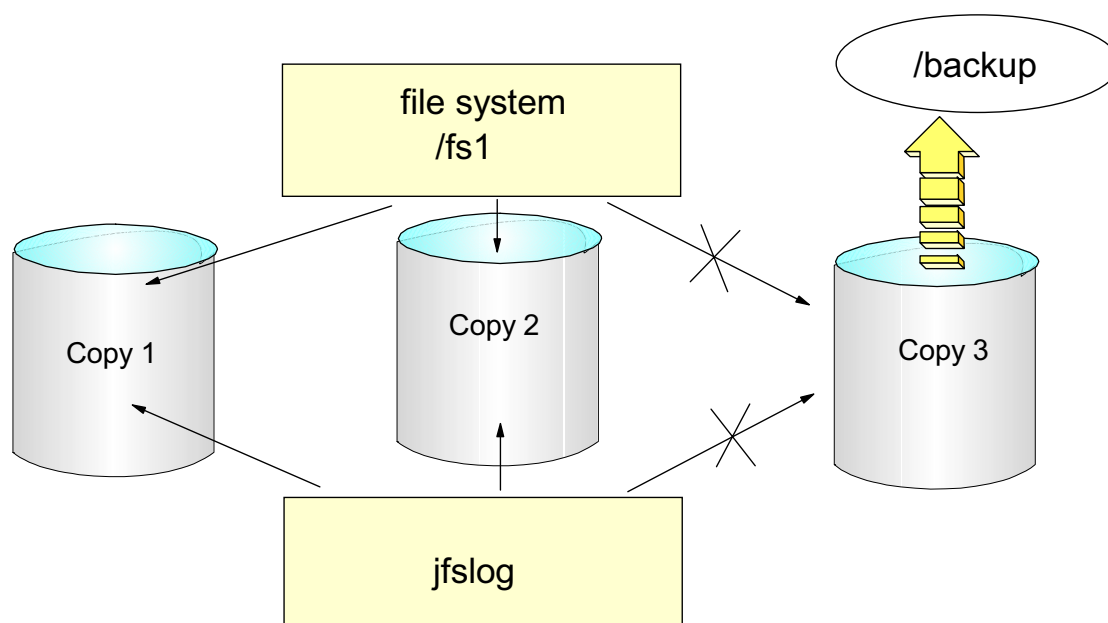
By splitting a mirror, you can perform a backup of the mirror that is not changing while the other mirrors remain online.

To do this, it is best to have three copies of your data. You will need to stop one of the copies but the other two will continue to provide redundancy for the online portion of the logical volume.

You are also required to have the log mirrored.

The picture above shows the output from **lsvg -l** indicating that the logical volume and the log are both mirrored.

## Splitting the Mirror



```
# chfs -a splitcopy = /backup -a copy=3 /fs1
```

© Copyright IBM Corporation 2004

Figure 7-27. Splitting the Mirror

AU1612.0

### Notes:

The command **chfs** is used to split the mirror to form a “snapshot” of the file system. This creates a read-only file system called **/backup** that can be accessed to perform a backup.

```
# lsvg -l newvg
```

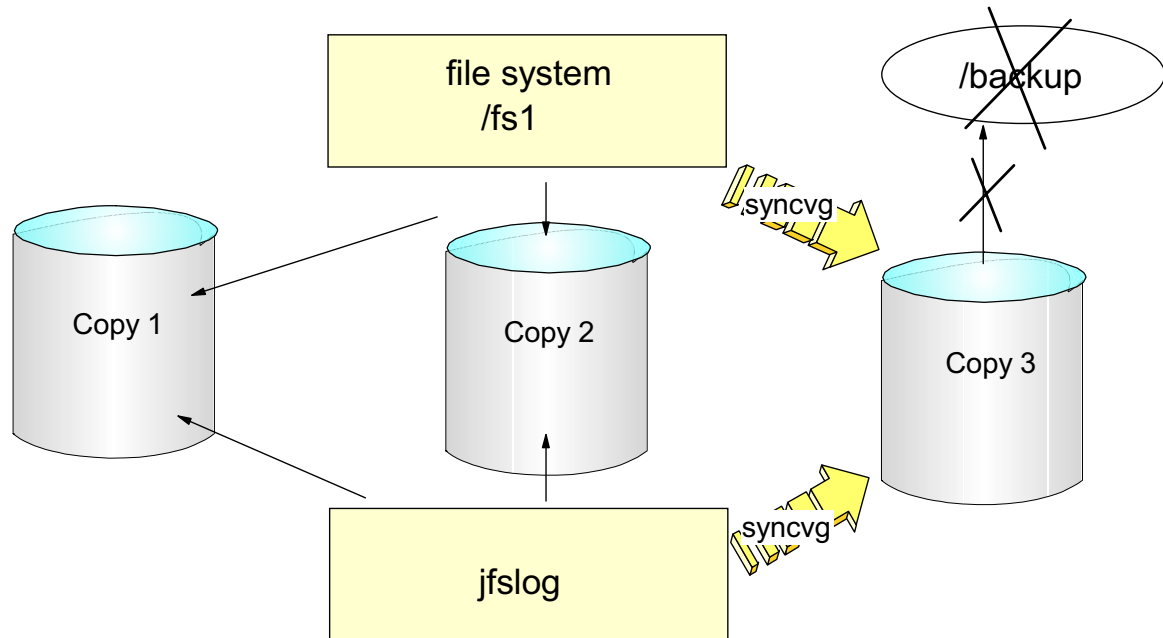
```
newvg:
```

LV NAME	TYPE	LPs	PPs	PVs	LV STATE	MOUNT
loglv00	jfslog	1	3	3	open/syncd	N/A
lv03	jfs	1	3	3	open/stale	/fs1
lv03copy00	jfs	0	0	0	open/syncd	/backup

The **/fs1** file system still contains 3 PPs but the mirror is now stale. The “stale” copy is now accessible by the newly created read-only file system **/backup**. That file system is contained on a newly created logical volume **lv03copy00**. This LV is not sync’ed or stale and it does not indicate any LP’s since the LP’s really belong to **lv03**.

You can look at the content and interact with the **/backup** file system just like any other read-only file system.

## Reintegrate a Mirror Backup Copy



```
# unmount /backup
# rmfs /backup
```

© Copyright IBM Corporation 2004

Figure 7-28. Reintegrate a Mirror Backup Copy

AU1612.0

### Notes:

To reintegrate the “snapshot” into the file system, unmount the **/backup** file system and remove it.

The third copy will automatically re-sync and come online.

---

## JFS2 Snapshot Image

---

- For a JFS2 file system, the point-in-time image is called a snapshot.
- A snapshot image of a JFS2 file system can be used to:
  - create a backup of the filesystem at the given point in time the snapshot was created
  - provide the capability to access files or directories as they were at the time of the snapshot
  - backup removable media
- The snapshot stays stable even if the file system that the snapshot was taken from continues to change.

© Copyright IBM Corporation 2004

Figure 7-29. JFS2 Snapshot Image

AU1612.0

### **Notes:**

Beginning with AIX 5.2, you can make a snapshot of a mounted JFS2 that establishes a consistent block-level image of the file system at a point in time.

The snapshot image remains stable even as the file system that was used to create the snapshot, called the snappedFS, continues to change.

The snapshot retains the same security permissions as the snappedFS had when the snapshot was made.

## Creation of a JFS2 Snapshot

---

- JFS2 snapshots can be created on the command line, through SMIT or the Web-based System Manager
- Some of the new commands included in Version 5.2 that support the JFS2 snapshot function are:
  - Snapshot - create, delete, and query a snapshot
  - Backsnap - create and backup a snapshot
  - fsdb - examine and modify snapshot superblock and snapshot map

© Copyright IBM Corporation 2004

Figure 7-30. Creation of a JFS2 Snapshot

AU1612.0

### **Notes:**

To create a snapshot of the /home/abc/test file system and back it up (by name) to the tape device /dev/rmt0, use the following command:

```
backsnap -m /tmp/snapshot -s size=16M -i f/dev/rmt0 /home/abc/test
```

This command creates a logical volume of 16 MB for the snapshot of the JFS2 file system (/home/abc/test). The snapshot is mounted on /tmp/snapshot and then a backup by name of the snapshot is made to the tape device. After the backup completes, the snapshot remains mounted. Use the -R flag with the backsnap command if you want the snapshot removed when the backup completes.



---

## Using a JFS2 Snapshot

---

- When a file becomes corrupted, you can replace it if you have an accurate copy in an online JFS2 snapshot.
- Use the following procedure to recover one or more files from a JFS2 snapshot image:
  - Mount the snapshot. For example:
    - `mount -v jfs2 -o snapshot /dev/mysnaplv /home/aaa/mysnap`
  - Change to the directory that contains the snapshot. For example:
    - `cd /home/aaa/mysnap`
  - Copy the accurate file to overwrite the corrupted one. For example:
    - `cp myfile /home/aaa/myfs` (copies only the file named myfile)
- The following example copies all files at once:
  - `cp -R /home/aaa/mysnap /home/aaa/myfs`

© Copyright IBM Corporation 2004

Figure 7-31. Using a JFS2 Snapshot

AU1612.0

### **Notes:**

This shows the procedure for using a JFS2 snapshot to recover a corrupted enhanced file system.

## Snapshot Support for Mirrored VGs

---

- Split a mirrored copy of a fully mirrored VG into a snapshot VG
- All LVs must be mirrored on disks that contains only those mirrors
- New LVs and mount points are created in the snapshot VG
- Both VGs keep track of changes in PPs
  - Writes to PP in original VG causes corresponding PP in snapshot VG to be marked stale
  - Writes to PP in snapshot VG causes that PP to be marked stale
- When the VGs are rejoined the stale PPs are resynchronized
- The user will see the same data in the rejoined VG as was in the original VG before the rejoin.

© Copyright IBM Corporation 2004

Figure 7-32. Snapshot Support for Mirrored VGs

AU1612.0

### **Notes:**

Snapshot support for a mirrored volume group is provided to split a mirrored copy of a fully mirrored volume group into a snapshot volume group.

When the VG is split the original VG will stop using the disks that are now part of the snapshot volume group.

Both volume groups will keep track of changes in physical partitions within the VG so that when the snapshot volume group is rejoined with the original VG, consistent data is maintained across the rejoined mirror copies.

## Snapshot VG Commands

```
splitvg [ -y SnapVGname ] [-c copy] [-f] [-i] Vgname
-y specifies the name of the snapped VG
-c specifies which mirror to use (1, 2 or 3)
-f forces the split even if there are stale partitions
-i creates an independent VG which cannot be rejoined into the original
```

- Example: File system /data is in the VG datavg. These commands split the VG, creates a backup of the /data file system and then rejoins the snapshot VG with the original.
  1. `splitvg -y snapvg datavg`
    - The VG datavg is split and the VG snapvg is created. The mount point /fs/data is created.
  2. `backup -f /dev/rmt0 /fs/data`
    - An i-node based backup of the unmounted file system /fs/data is created on tape.
  3. `joinvg datavg`
    - snapvg is rejoined with the original VG and synced in the background.

© Copyright IBM Corporation 2004

Figure 7-33. Snapshot VG Commands

AU1612.0

### Notes:

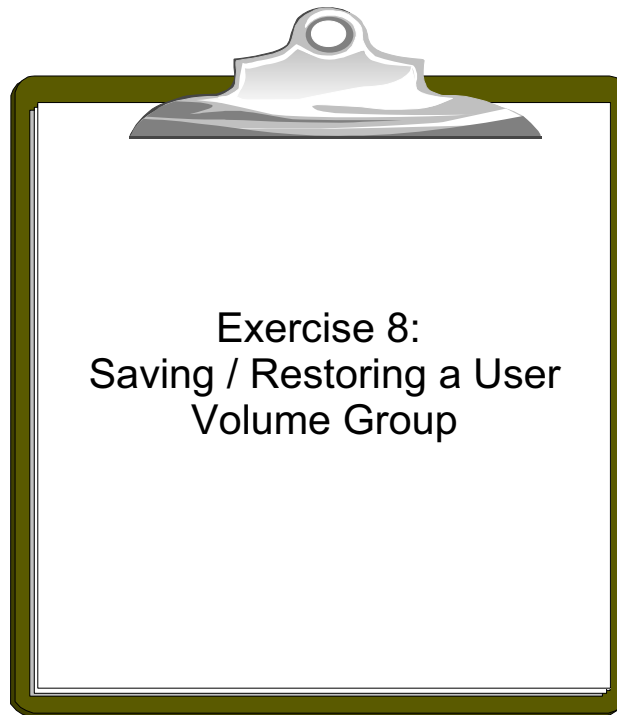
The `splitvg` command will fail if any of the disks to be split are not active within the original volume group.

In the event of a system crash or loss of quorum while running this command, the `joinvg` command must be run to rejoin the disks back to the original volume group.

You must have root authority to run this command.

## Next Step

---



© Copyright IBM Corporation 2004

Figure 7-34. Next Step

AU1612.0

### **Notes:**

After the exercise, you should be able to:

- Use the **savevg** command to back up a user volume group
- Use the **restvg** command to restore a user volume group
- Change volume group characteristics

---

## Checkpoint

---

1. **T/F:** After restoring a **mksysb** image all passwords are restored as well.  

---
2. The **mkszfile** will create a file named:
  - a. /bosinst.data
  - b. /image.data
  - c. /vgname.data

---
3. Which two alternate disk installation techniques are available?  

---
4. What are the commands to backup and restore a non-rootvg volume group?  

---
5. If you want to shrink one file system in a volume group **myvg**, which file must be changed before backing up the user volume group?  

---
6. How many mirror copies should you have before performing an online JFS or JFS2 backup?  

---

© Copyright IBM Corporation 2004

Figure 7-35. Checkpoint

AU1612.0

### Notes:

## Unit Summary

---

- Backing up rootvg is performed with the **mksysb** command. A **mksysb** image should always be verified before using it.
- **mksysb** control files are **bosinst.data** and **image.data**
- Two alternate disk installation techniques are available:
  - Installing a **mksysb** onto an **alternate** disk
  - **Cloning** the current **rootvg** onto an **alternate** disk
  - Changing the **bootlist** allows booting different AIX levels
- Backing up a non-rootvg volume group is performed with the **savevg** command.
- Restoring a non-rootvg volume group is done using the **restvg** command.
- Online JFS and JFS2 backups can be done using **chfs**.

© Copyright IBM Corporation 2004

Figure 7-36. Unit Summary

AU1612.0

### **Notes:**

# Unit 8. Error Log and syslogd

## What This Unit Is About

This unit is an overview of the error logging facility available in AIX and shows how to work with the syslogd daemon.

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Analyze error log entries
- Identify and maintain the error log components
- Provide different **error notification** methods
- Log system messages using the **syslogd** daemon

## How You Will Check Your Progress

Accountability:

- Activities
- Lab exercise
- Checkpoint questions

## References

- |               |   |
|---------------|---|
| <i>Online</i> | <i>General Programming Concepts: Writing and Debugging Programs Chapter 4. Error Notification</i> |
| <i>Online</i> | <i>Commands Reference</i>   |

## Unit Objectives

---

After completing this unit, students should be able to:

- Analyze **error log entries**
- Identify and maintain the **error log components**
- Provide different **error notification** methods
- Log system messages using the **syslogd** daemon

© Copyright IBM Corporation 2004

---

Figure 8-1. Unit Objectives

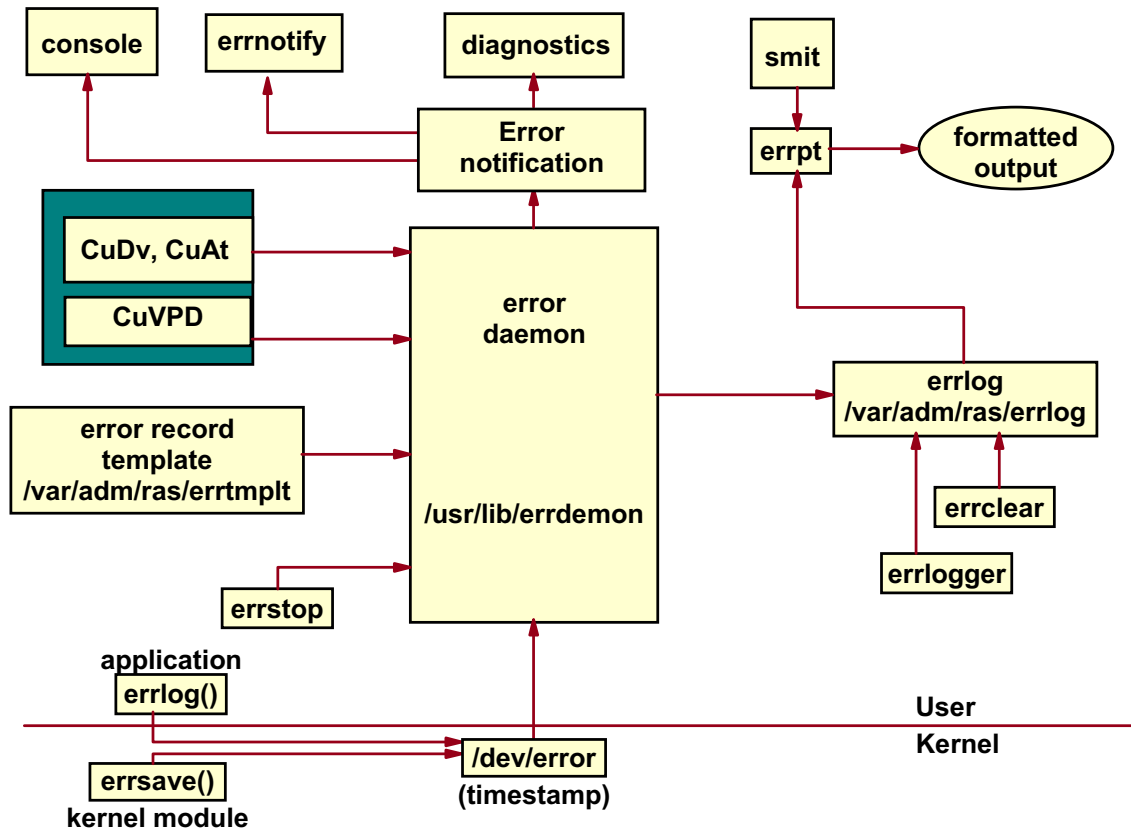
AU1612.0

### **Notes:**



## 8.1 Working With Error Log

# Error Logging Components



© Copyright IBM Corporation 2004

Figure 8-2. Error Logging Components

AU1612.0

## Notes:

The error logging process begins when an operating system module detects an error. The error detecting segment of code then sends error information to either the **errsave()** kernel service or the **errlog()** application subroutine, where the information is in turn written to the **/dev/error** special file. This process then adds a timestamp to the collected data. The **errdemon** daemon constantly checks the **/dev/error** file for new entries, and when new data is written, the daemon conducts a series of operations.

Before an entry is written to the error log, the **errdemon** daemon compares the label sent by the kernel or the application code to the contents of the Error Record Template Repository. If the label matches an item in the repository, the daemon collects additional data from other parts of the system.

To create an entry in the error log, the **errdemon** daemon retrieves the appropriate template from the repository, the resource name of the unit that caused the error, and the detail data. Also, if the error signifies a hardware-related problem and hardware vital product data (VPD) exists, the daemon retrieves the VPD from the ODM. When you access the error log, either through SMIT or with the **errpt** command, the error log is formatted

according to the error template in the error template repository and presented in either a summary or detailed report. Most entries in the error log are attributable to hardware and software problems, but informational messages can also be logged, for example, by the system administrator.

The **errlogger** command allows the system administrator to record messages of up to 1024 bytes in the error log. Whenever you perform a maintenance activity, such as clearing entries from the error log, replacing hardware, or applying a software fix, it is a good idea to record this activity in the system error log.

For example:

**# errlogger system hard disk '(hdisk0)' replaced.**

This message will be listed as part of the error log.

Error log hardening

Under very rare circumstances, like powering off the system exactly while the errdemon is writing into the error log, the error log may get corrupted. In AIX 5L V5.3, there are minor modifications done to the errdemon to improve its robustness and to the recovery of the error log file at its start.

When the errdemon starts, it checks for error log consistency. First, it makes a backup copy of the existing error log file to /tmp/errlog.save and then it corrects the error log file, while preserving consistent error log entries. The difference from the previous versions of AIX is that the errdemon is used to reset the log file if it was corrupted, instead of repairing it.

# Generating an Error Report via smit

# smit errpt

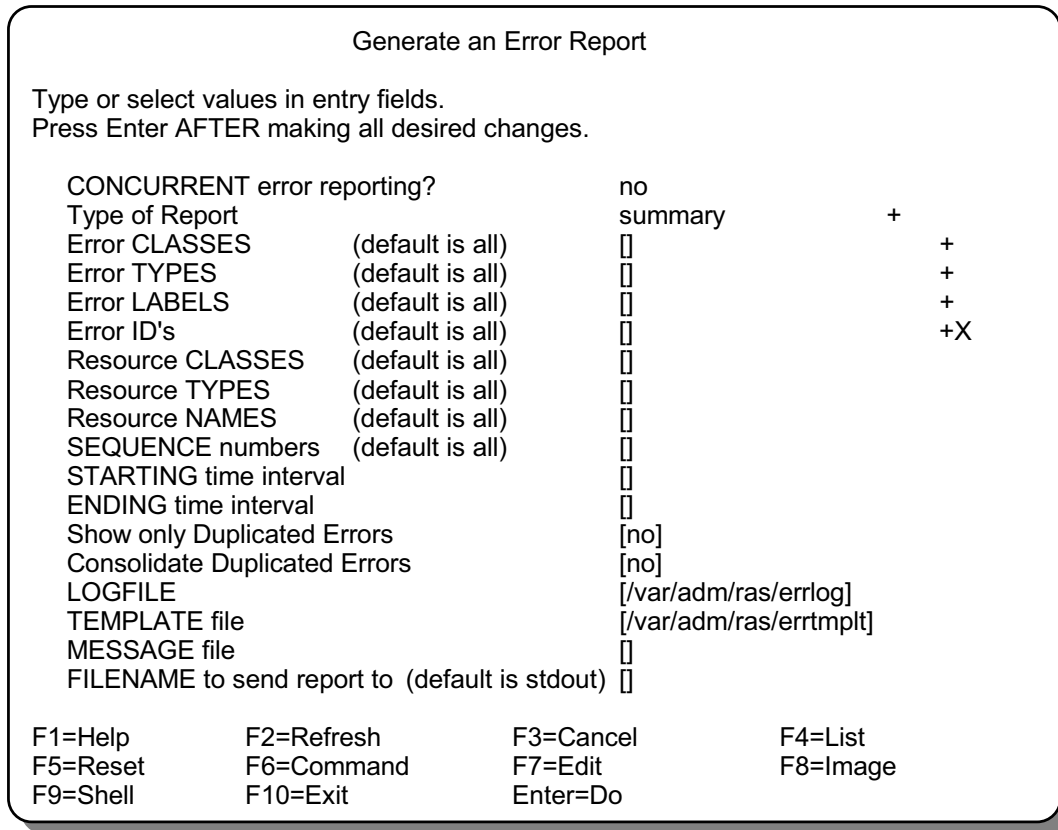


Figure 8-3. Generating an Error Report via smit

AU1612.0

## Notes:

Any user can use this screen. The fields can be specified as:

**CONCURRENT error reporting** Yes means you want errors displayed or printed as the errors are entered into the error log - a sort of **tail -f**

**Type of Report** Summary, intermediate and detailed reports are available. Detailed reports give comprehensive information. Intermediate reports display most of the error information. Summary reports contain concise descriptions of errors.

**Error CLASSES** Values are H (hardware), S (software) and O (operator messages created with **errlogger**). You can specify more than one error class.

**Resource CLASSES** Means device class for hardware errors (for example, disk).

**Error TYPES**

<b>PEND</b>	The loss of availability of a device or component is imminent
<b>PERF</b>	The performance of the device or component has degraded to below an acceptable level
<b>TEMP</b>	Recovered from condition after several attempts
<b>PERM</b>	Unable to recover from error condition. Error types with this value are usually most severe error and imply that you have a hardware or software defect. Error types other than PERM usually do not indicate a defect, but they are recorded so that they can be analyzed by the diagnostic programs
<b>UNKN</b>	Severity of the error cannot be determined.
<b>INFO</b>	The error type is used to record informational entries
<b>Resource TYPES</b>	Device type for hardware (for example 355 MB)
<b>Resource NAMES</b>	Common device name (for example hdisk0)
<b>ID</b>	Is the error identifier
<b>STARTING and ENDING dates</b>	Format mmddhhmmyy can be used to select only errors from the log that are time stamped between the two values.
<b>Show only Duplicated Errors</b>	Yes will report only those errors that are exact duplicates of previous errors generated during the interval of time specified. The default time interval is 100 milliseconds. This value can be changed with the <code>errdemon -t</code> command. The default for the Show only Duplicated Errors option is no.
<b>Consolidate Duplicated Errors</b>	Yes will report only the number of duplicate errors and timestamps of the first and last occurrence of that error. The default for the Consolidate Duplicated Errors option is no.

# The errpt Command

- Summary report:  
# errpt
- Summary report of all hardware errors:  
# errpt -d H
- Intermediate report:  
# errpt -A
- Detailed report:  
# errpt -a
- Detailed report of all software errors:  
# errpt -a -d S
- Concurrent error logging ("Real-time" error logging):  
# errpt -c > /dev/console

© Copyright IBM Corporation 2004

Figure 8-4. The errpt Command

AU1612.0

## Notes:

The **errpt** command generates a report of logged errors. Three different layouts are produced dependent on the options that are used:

- A **summary** report, which gives an overview (default).
- An **intermediate** report, which only displays the values for the LABEL, Date/Time, Type, Resource Name, Description and Detailed Data fields. Use the option **-A** to specify an intermediate report.
- A **detailed** report, which shows a detailed description of all the error entries. Use the option **-a** to specify a detailed report.

The **errpt** command queries the error log file **/var/adm/ras/errlog** to produce the error report.

If you want to display the error entries concurrently, that is, at the time they are logged, you must execute **errpt -c**. In the example, we direct the output to the system console.

Duplicate errors can be consolidated using **errpt -D**. When used with the **-a** option, **errpt -D** reports only the number of duplicate errors and the timestamp for the first and last occurrence of the identical error.

The **errpt** command has many options. Refer to your AIX commands reference for a complete description.

# A Summary Report (errpt)

```
# errpt
IDENTIFIER  TIMESTAMPT  C  RESOURCE_NAME  DESCRIPTION
94537C2E    0430033899 P  H  tok0           WIRE FAULT
35BFC499    0429090399 P  H  hdisk1         DISK OPERATION ERROR
...
1581762B    0428202699 T  H  hdisk0         DISK OPERATION ERROR
...
E85C5C4C    0428043199 P  S  LFTDD         SOFTWARE PROGRAMM ERROR
2BFA76F6    0427091499 T  S  SYSPROC       SYSTEM SHUTDOWN BY USER
B188909A    0427090899 U  S  LVDD          PHYSICAL PARTITION MARKED STALE
...
9DBCDFDEE   0427090699 T  O  errdemon      ERROR LOGGING TURNED ON
...
2BFA76F6    0426112799 T  S  SYSPROC       SYSTEM SHUTDOWN BY USER
```

- Error Type:**
- P: Permanent, Performance or Pending
  - T: Temporary
  - I: Informational
  - U: Unknown

- Error Class:**
- H: Hardware
  - S: Software
  - O: Operator
  - U: Undetermined

© Copyright IBM Corporation 2004

Figure 8-5. A Summary Report (errpt)

AU1612.0

## Notes:

The **errpt** command creates by default a **summary** report which gives an overview about the different error entries. One line per error is fine to get a feel for what is there, but you need more details to understand problems.

The example shows different hardware and software errors that occurred. To get more information about these errors you must create a **detailed** report.



## A Detailed Error Report (errpt -a)

```

LABEL:                TAPE_ERR4
IDENTIFIER:           5537AC5F

Date/Time:            Thu 27 Feb 13:41:51
Sequence Number:     40
Machine Id:           000031994100
Node Id:              dw6
Class:                H
Type:                 PERM
Resource Name:        rmt0
Resource Class:       tape
Resource Type:        8mm
Location:             00-00-0S-3,0
VPD:
    Manufacturer      EXABYTE
    Machine Type and Model EXB-8200
    Part Number        21F8842
    Device Specific (Z0) 0180010133000000
    Device Specific (Z1) 2680

Description
TAPE DRIVE FAILURE

Probable Causes
ADAPTER
TAPE DRIVE

Failure Causes
ADAPTER
TAPE DRIVE

Recommended Actions
PERFORM PROBLEM DETERMINATION PROCEDURES

Detail Data
SENSE DATA
0603 0000 1700 0000 0000 0000 0000 0000 0200 0800 0000 0000 0000 0000
0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000

```

© Copyright IBM Corporation 2004

Figure 8-6. A Detailed Report (errpt -a)

AU1612.0

### Notes:

The detailed error reports are generated by issuing the **errpt -a** command. The first half of the information is obtained from the ODM (CuDv, CuAt, CuVPD) and is very useful because it shows clearly which part causes the error entry. The next few fields explain probable reasons for the problem, and actions that you can take to correct the problem.

The last field, **SENSE DATA**, is a detailed report about which part of the device is failing. For example, with disks it could tell you which sector on the disk is failing. This information can be used by IBM support to analyze the problem.

Here again is a list of error classes and error types:

1. An error class value of **H** and an error type value of **PERM** indicate that the system encountered a problem with a piece of hardware and could not recover from it.
2. An error class value of **H** and an error type value of **PEND** indicate that a piece of hardware may become unavailable soon due to the numerous errors detected by the system.

3. An error class value of **S** and an error type of **PERM** indicate that the system encountered a problem with software and could not recover from it.
4. An error class value of **S** and an error type of **TEMP** indicate that the system encountered a problem with software. After several attempts, the system was able to recover from the problem.
5. An error class value of **O** indicate that an informational message has been logged.
6. An error class value of **U** indicate that an error could not be determined.

Starting in AIX 5.1, there is a link between the error log and diagnostics. Error reports will include the diagnostic analysis for errors that have been analyzed. Diagnostics, and the diagnostic tool **diag**, will be covered in a later unit.

## Types of Disk Errors

DISK_ERR1	P	Failure of physical volume media Action: Replace device as soon as possible
DISK_ERR2, DISK_ERR3	P	Device does not respond Action: Check power supply
DISK_ERR4	T	Error caused by a bad block or event of a recovered error
SCSI_ERR* (SCSI_ERR10)	P	SCSI Communication Problem Action: Check cable, SCSI addresses, terminator

P = Permanent hardware error

T = Temporary hardware error

Rule of thumb: Replace disk, if it produces more than one DISK\_ERR4 per week

© Copyright IBM Corporation 2004

Figure 8-7. Types Of Disk Errors

AU1612.0

### Notes:

This page explains the most common disk errors you should know about:

1. **DISK\_ERR1** is caused from wear and tear of the disk. Remove the disk as soon as possible from the system and replace it with a new one. Follow the procedures that you've learned earlier in this course.
2. **DISK\_ERR2**, **DISK\_ERR3** error entries are mostly caused by a loss of electrical power.
3. **DISK\_ERR4** is the most interesting one, and the one that you should watch out for, as this indicates bad blocks on the disk. Do not panic if you get a few entries in the log of this type of an error. What you should be aware of is the number of **DISK\_ERR4** errors and their frequency. The more you get, the closer you are getting to a disk failure. You want to prevent this before it happens, so monitor the error log closely.
4. Sometimes **SCSI** errors are logged, mostly with the ID **SCSI\_ERR10**. They indicate that the SCSI controller is not able to communicate with an attached device. In this case, check the cable (and the cable length), the SCSI addresses and the terminator.

A very infrequent error is **DISK\_ERR5**. It is the catch-all (that is, the problem does not match any of the above DISK\_ERR symptoms). You need to investigate further by running the **diagnostic** programs which can detect and produce more information on the problem.

## LVM Error Log Entries

LVM_BBEPOOL, LVM_BBERELMAX, LVM_HWFAIL	S,P	No more bad block relocation. Action: Replace disk as soon as possible
LVM_SA_STALEPP	S,P	Stale physical partition. Action: check disk, synchronize data (syncvg)
LVM_SA_QUORCLOSE	H,P	Quorum lost, volume group closing. Action: Check disk, consider working without quorum

H = Hardware  
S = Software

P = Permanent  
T = Temp

© Copyright IBM Corporation 2004

Figure 8-8. LVM Error Log Entries

AU1612.0

### Notes:

This list shows some very important LVM error codes you should know. All of these errors are permanent errors that cannot be recovered. Very often these errors are accompanied by hardware errors as shown on the previous page.

Errors, like these shown in the list, require your immediate intervention.

# Maintaining the Error Log

# smit errdemon

Change / Show Characteristics of the Error Log

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

LOGFILE	[/var/adm/ras/errlog]			
* Maximum LOGSIZE	[1048576]		#	
Memory Buffer Size	[8192]		#	
...				

# smit errclear

Clean the Error Log

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

Remove entries older than this number of days	[30]			
Error CLASSES	[ ]		+	
Error TYPES	[ ]		+	
...				
Resource CLASSES	[ ]		+	
...				

==> Use the errlogger command as reminder <==

© Copyright IBM Corporation 2004

Figure 8-9. Maintaining the Error Log

AU1612.0

## Notes:

- To change error log attributes like the **error log filename**, the **internal memory buffer size** and the **error log file size** use the **smit** fastpath **smit errdemon**. The error log file is implemented as a **ring**. When the file reaches its limit, the oldest entry is removed to allow adding a new one. The command that **smit** executes is the **errdemon** command. See your AIX command reference for a listing of the different options.
- To clean up error log entries, use the **smit** fastpath **smit errclear**. For example, after removing a bad disk that caused error logs entries, you should remove the corresponding error log entries of the bad disk. The **errclear** command is part of the fileset **bos.sysmgmt.serv\_aid**.

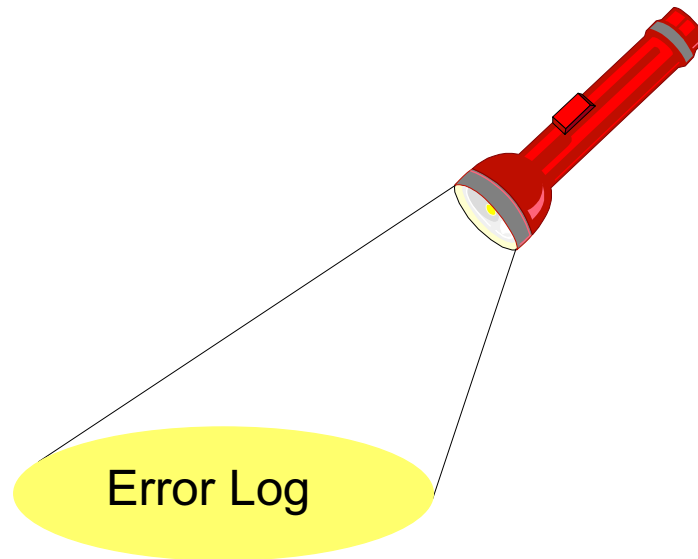
Software and hardware errors are removed by **errclear** using **crontab**. Software and operator errors are purged after 30 days, hardware errors are purged after 90 days.

Follow the reminder from the bottom of the visual. Whenever an important system event takes place, for example the replacement of a disk, log this event using the **errlogger** command.

---

## Activity: Working with the Error Log

---



© Copyright IBM Corporation 2004

Figure 8-10. Activity: Working with the Error Log

AU1612.0

### **Notes:**

This activity allows you to work with the AIX error logging facility.

After the activity, you should be able to:

- Determine what errors are logged on your machine.
- Generate different error reports.
- Start concurrent error notification.

### **Instructions:**

\_\_\_ 1. Generate a **summary report** of your system's error log. Write down the command that you (or smit) used:

\_\_\_\_\_

\_\_\_ 2. Generate a **detailed report** of your system's error log. Write down the command that you (or smit) used:

\_\_\_\_\_

- \_\_\_ 3. Using **smit**, generate the following reports:
- A **summary report** of all errors that occurred during the past 24 hours. Write down the command that **smit** executes:

\_\_\_\_\_

- A **detailed report** of all hardware errors. Write down the command that **smit** executes:

\_\_\_\_\_

- \_\_\_ 4. This instruction requires that a graphical desktop, for example CDE is active. Start two windows. In one window startup **concurrent** error logging, using the **errpt** command. Write down the command that you used:

\_\_\_\_\_

In the other window, execute the **errlogger** command to generate an error entry. Write down the command you used:

\_\_\_\_\_

Is the complete error text shown in the error report?

\_\_\_\_\_

Stop **concurrent** error logging.

- \_\_\_ 5. Write down the characteristics of your error log:

**LOGFILE:**

**Maximum LOGSIZE:**

**Memory BUFFER SIZE:**

**What command have you used to show these characteristics?**

\_\_\_\_\_

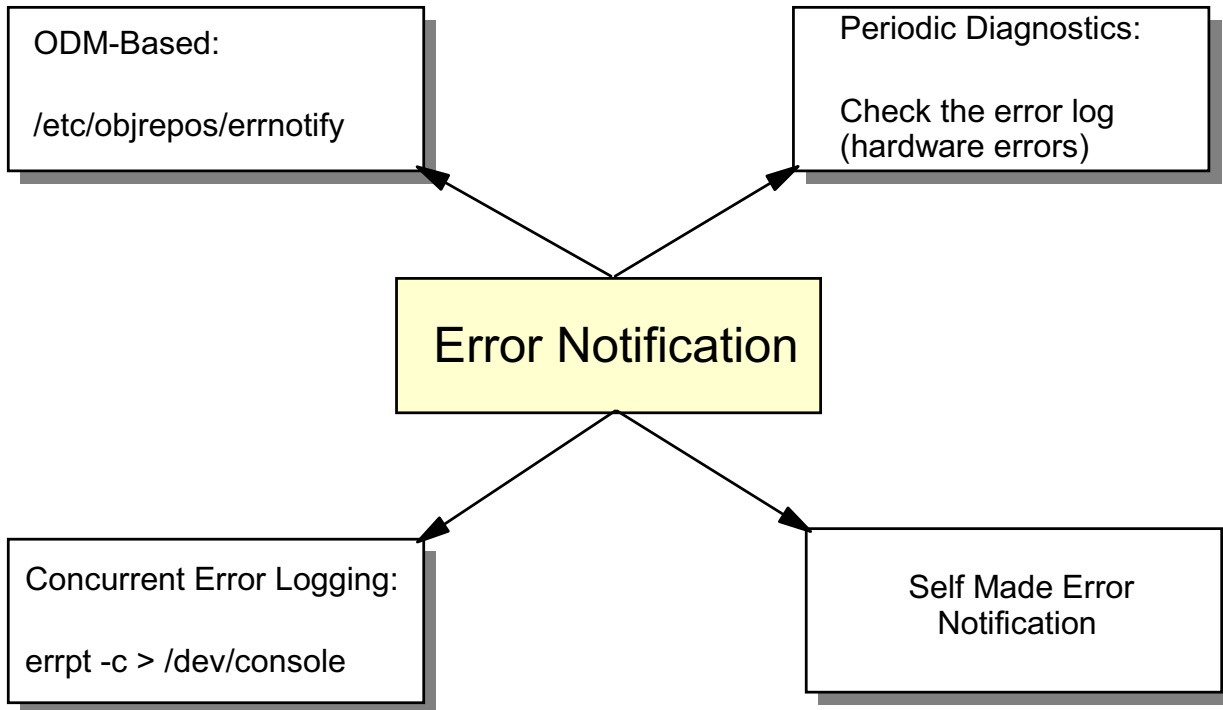
- \_\_\_ 6. **Clean up all error entries that have the error class operator.** Write down the command, you (or smit) used:

\_\_\_\_\_



## 8.2 Error Notification and syslogd

# Error Notification Methods



© Copyright IBM Corporation 2004

Figure 8-11. Error Notification Methods

AU1612.0

## Notes:

The term **error notification** means that the system informs you whenever an error is posted to the error log.

There are different ways to implement **error notification**.

1. **Concurrent Error Logging:** That's the easiest way to implement error notification. By starting **errpt -c** each error is reported when it occurs. By redirecting the output to the console, an operator is informed about each new error entry.
2. **Self-made Error Notification:** Another easy way to implement error notification is to write a shell procedure that regularly checks the error log. This is shown on the next visual.
3. **Periodic Diagnostics:** The **diagnostics** package (**diag command**) contains a periodic diagnostic procedure (**diagela**). Whenever a **hardware error** is posted to the log, all members of the **system group** get a mail message. Additionally a message is sent to the system console. **diagela** has two disadvantages:
  - Since it executes many times a day, the program might slow down your system.

- Only hardware errors are analyzed.
4. **ODM-based error notification:** The **errdemon** program uses an ODM class **errnotify** for error notification. How to work with **errnotify** is introduced later in this topic.

# Self-made Error Notification

```
#!/usr/bin/ksh

errpt > /tmp/errlog.1

while true
do
    sleep 60                # Let's sleep one minute

    errpt > /tmp/errlog.2

    # Compare both files.
    # If no difference, let's sleep again
    cmp -s /tmp/errlog.1 /tmp/errlog.2 && continue

    # Files are different: Let's inform the operator:
    print "Operator: Check error log " > /dev/console

    errpt > /tmp/errlog.1

done
```

© Copyright IBM Corporation 2004

Figure 8-12. Self-made Error Notification

AU1612.0

## Notes:

By using the **errpt** command it's very easy to implement a self-made error notification.

Let's analyze the procedure shown above:

- The first **errpt** command generates a file **/tmp/errlog.1**.
- The construct **while true** implements an infinite loop that never terminates.
- In the loop, the first action is to **sleep** one minute.
- The second **errpt** command generates a second file **/tmp/errlog.2**.
- Both files are compared using the command **cmp -s** (silent compare, that means no output will be reported). If the files are not different, we jump back to the beginning of the loop (continue), and the process will sleep again.
- If there is a difference, a new error entry has been posted to the error log. In this case, we inform the operator that a new entry is in the error log. Instead of **print** you could use the **mail** command to inform another person.

This is a very easy but effective way of implementing error notification.

## ODM-based Error Notification: errnotify

```
errnotify:
  en_pid = 0
  en_name = "sample"
  en_persistenceflg = 1
  en_label = ""
  en_crcid = 0
  en_class = "H"
  en_type = "PERM"
  en_alertflg = ""
  en_resource = ""
  en_rtype = ""
  en_rclass = "disk"
  en_method = "errpt -a -l $1 | mail -s DiskError root"
```

© Copyright IBM Corporation 2004

Figure 8-13. ODM-based Error Notification: errnotify

AU1612.0

### Notes:

The Error Notification object class specifies the conditions and actions to be taken when errors are recorded in the system error log. The user specifies these conditions and actions in an Error Notification object.

Each time an error is logged, the **error notification** daemon determines if the error log entry matches the selection criteria of any of the Error Notification objects. If matches exist, the daemon runs the programmed action, also called a notify method, for each matched object.

The Error Notification object class is located in the **/etc/objrepos/errnotify** file. Error Notification objects are added to the object class by using ODM commands.

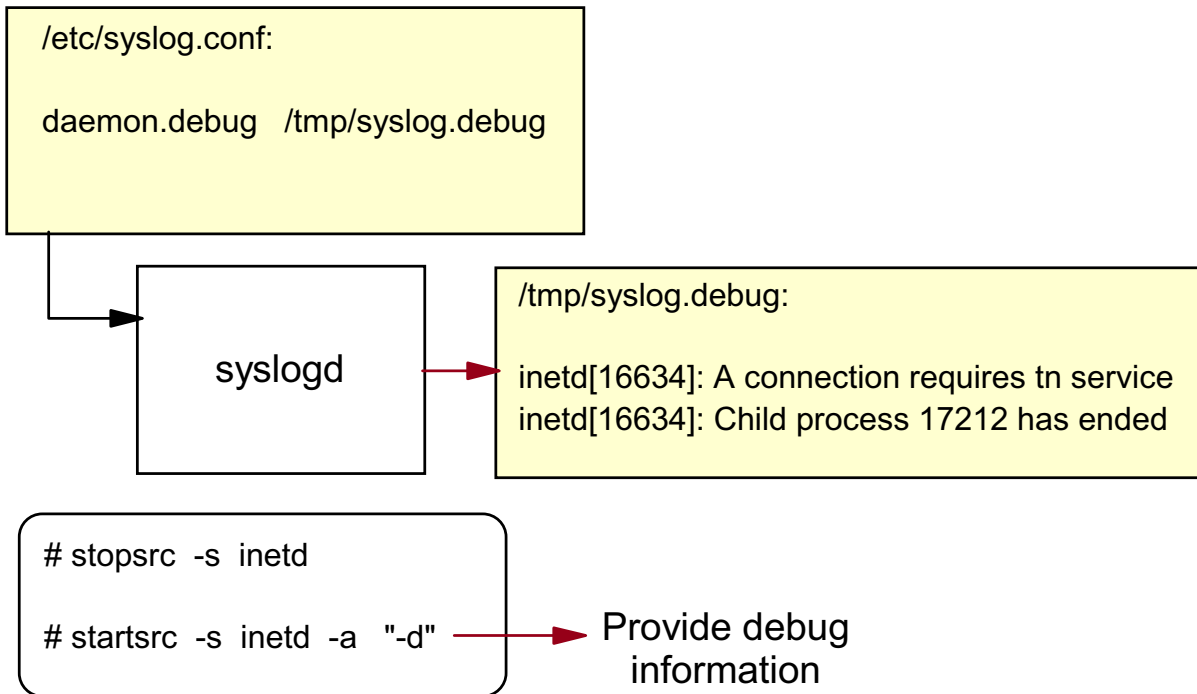
The example shows an object that creates a **mail** message to **root** whenever a **disk** error is posted to the log. Here is a list of all **descriptors**:

**en\_alertflg** Identifies whether the error is alertable. This descriptor is provided for use by alert agents with network management applications. The values are **TRUE** (alertable) or **FALSE** (not alertable).

<b>en_class</b>	Identifies the class of error log entries to match. Valid values are <b>H</b> (hardware errors), <b>S</b> (software errors), <b>O</b> (operator messages) and <b>U</b> (undetermined).
<b>en_crcid</b>	Specifies the error identifier associated with a particular error.
<b>en_label</b>	Specifies the label associated with a particular error identifier as defined in the output of <b>errpt -t</b> (show templates).
<b>en_method</b>	<p>Specifies a user-programmable action, such as a shell script or a command string, to be run when an error matching the selection criteria of this Error Notification object is logged. The error notification daemon uses the <b>sh -c</b> command to execute the notify method.</p> <p>The following keywords are passed to the method as arguments:</p> <ul style="list-style-type: none"><li><b>\$1</b> Sequence number from the error log entry</li><li><b>\$2</b> Error ID from the error log entry</li><li><b>\$3</b> Class from the error log entry</li><li><b>\$4</b> Type from the error log entry</li><li><b>\$5</b> Alert flags from the error log entry</li><li><b>\$6</b> Resource name from the error log entry</li><li><b>\$7</b> Resource type from the error log entry</li><li><b>\$8</b> Resource class from the error log entry</li><li><b>\$9</b> Error label from the error log entry</li></ul>
<b>en_name</b>	Uniquely identifies the object.
<b>en_persistenceflg</b>	Designates whether the Error Notification object should be removed when the system is restarted. <b>0</b> means removed at boot time, <b>1</b> means persists through boot.
<b>en_pid</b>	Specifies a process ID for use in identifying the Error Notification object. Objects that have a PID specified should have the <b>en_persistenceflg</b> descriptor set to <b>0</b> .
<b>en_rclass</b>	Identifies the class of the failing resource. For hardware errors, the resource class is the device class (see PdDv). Not used for software errors.
<b>en_resource</b>	Identifies the name of the failing resource. For hardware errors, the resource name is the device name. Not used for software errors.
<b>en_rtype</b>	Identifies the type of the failing resource. For hardware errors, the resource type is the device type (see PdDv). Not used for software errors.

<b>en_symptom</b>	Enables notification of an error accompanied by a symptom string when set to <b>TRUE</b> .
<b>en_type</b>	Identifies the severity of error log entries to match. Valid values are: <b>INFO</b> : Informational <b>PEND</b> : Impending loss of availability <b>PERM</b> : Permanent <b>PERF</b> : Unacceptable performance degradation <b>TEMP</b> : Temporary <b>UNKN</b> : Unknown <b>TRUE</b> : Matches alertable errors <b>FALSE</b> : Matches non-alertable errors <b>0</b> : Removes the Error Notification object at system restart <b>non-zero</b> : Retains the Error Notification object at system restart

# syslogd Daemon



© Copyright IBM Corporation 2004

Figure 8-14. syslogd Daemon

AU1612.0

## Notes:

The **syslogd** daemon logs system messages from different software components (Kernel, daemon processes, system applications).

When started, the **syslogd** reads a configuration file `/etc/syslog.conf`. Whenever you change this configuration file you need to refresh the **syslogd** subsystem:

```
# refresh -s syslogd
```

The visual shows a configuration that is often used when a daemon process causes a problem. The line:

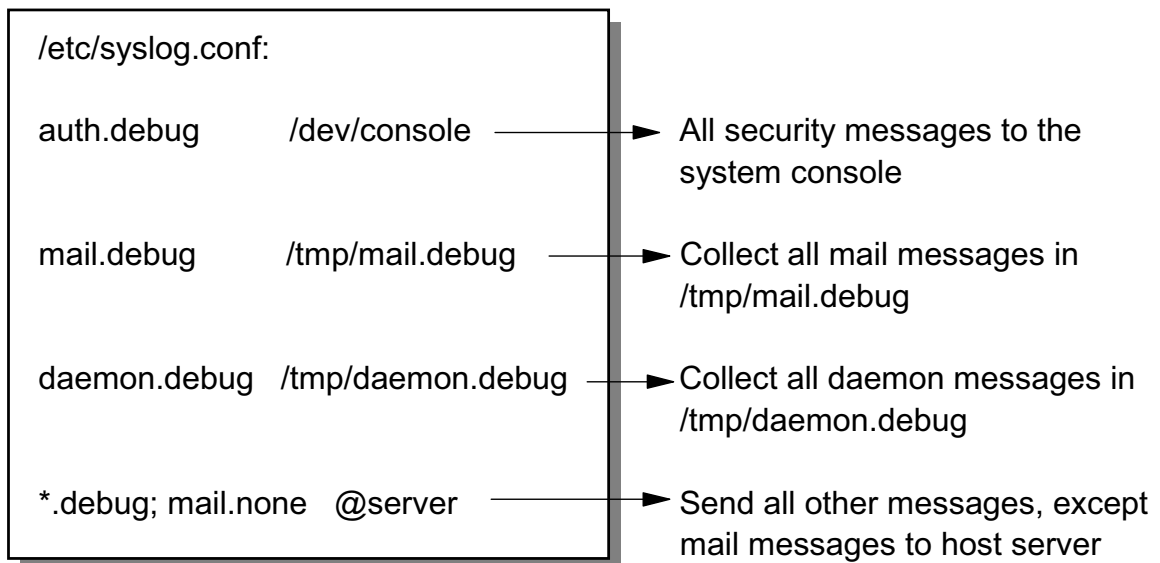
```
daemon.debug /tmp/syslog.debug
```

**indicates that facility daemon** should be controlled. All messages with the priority level **debug** and higher, should be written to the file `/tmp/syslog.debug`. Note that this file **must** exist.

The daemon process that causes problems (in our example the **inetd**) is started with option `-d` to provide debug information. This debug information is collected by the **syslogd** daemon, which writes the information to the log file `/tmp/syslog.debug`.



# syslogd Configuration Examples



After changing `/etc/syslog.conf`:

- `refresh -s syslogd`

© Copyright IBM Corporation 2004

Figure 8-15. syslogd Configuration Examples

AU1612.0

## Notes:

The visual shows some configuration examples in `/etc/syslog.conf`:

- **auth.debug /dev/console** specifies that all security messages are directed to the system console.
- **mail.debug /tmp/mail.debug** specifies that all mail messages are collected in file `/tmp/mail.debug`.
- **daemon.debug /tmp/daemon.debug** specifies that all messages produced from daemon processes are collected in file `/tmp/daemon.debug`.
- **\*.debug; mail.none @server** specifies that all other messages, except messages from the mail subsystem, are sent to the **syslogd** daemon on host **server**.

As you see, the general format in `/etc/syslog.conf` is:

**selector action**

The **selector field** names a **facility** and a **priority level**. Separate facility names with a comma (,). Separate the facility and priority level portions of the selector field with a period (.). Separate multiple entries in the same selector field with a semicolon (;). To select all facilities use an asterisk (\*).

The action field identifies a destination (file, host or user) to receive the messages. If routed to a remote host, the remote system will handle the message as indicated in its own configuration file. To display messages on a user's terminal, the destination field must contain the name of a valid, logged-in system user. If you specify an asterisk (\*) in the action field, a message is sent to all logged-in users.

### Facilities

Use the following system facility names in the selector field:

<b>kern</b>	Kernel
<b>user</b>	User level
<b>mail</b>	Mail subsystem
<b>daemon</b>	System daemons
<b>auth</b>	Security or authorization
<b>syslog</b>	<b>syslogd</b> messages
<b>lpr</b>	Line-printer subsystem
<b>news</b>	News subsystem
<b>uucp</b>	uucp subsystem
<b>*</b>	All facilities

### Priority Levels

Use the following levels in the selector field. Messages of the specified level and all levels above it are sent as directed.

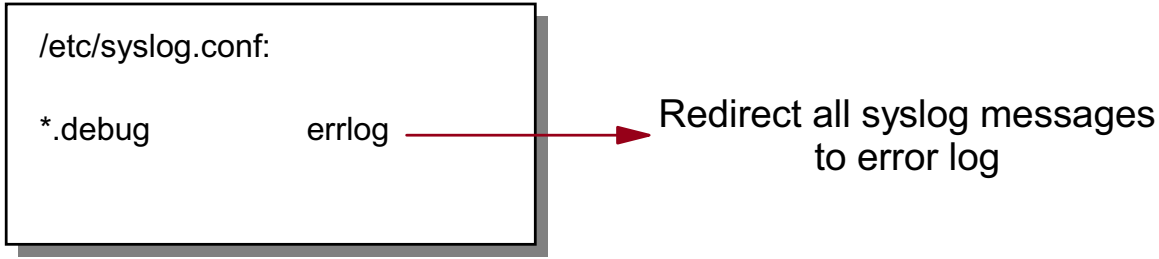
<b>emerg</b>	Specifies emergency messages. These messages are not distributed to all users.
<b>alert</b>	Specifies important messages such as serious hardware errors. These messages are distributed to all users.
<b>crit</b>	Specifies critical messages, not classified as errors, such as improper login attempts. These messages are sent to the system console.
<b>err</b>	Specifies messages that represent error conditions.
<b>warning</b>	Specifies messages for abnormal, but recoverable conditions.
<b>notice</b>	Specifies important informational messages.
<b>info</b>	Specifies information messages that are useful to analyze the system.

**debug** Specifies debugging messages. If you are interested in all messages of a certain facility, use this level.

**none** Excludes the selected facility.

Whenever changing **/etc/syslog.conf**, you must refresh the **syslogd** subsystem.

# Redirecting syslog Messages to Error Log



```
# errpt
```

```
IDENTIFIER  TIMESTAMPT  C  RESOURCE_NAME  DESCRIPTION
...
C6ACA566    0505071399  U  S  syslog          MESSAGE REDIRECTED FROM SYSLOG
...
```

© Copyright IBM Corporation 2004

Figure 8-16. Redirecting syslog Messages to Error Log

AU1612.0

## Notes:

Some applications use **syslogd** for logging errors and events. Some administrators find it desirable to list all errors in one report.

The visual shows how to redirect messages from **syslogd** to the error log.

By setting the action field to **errlog**, all messages are redirected to the AIX error log.

## Directing Error Log Messages to syslogd

```
errnotify:
  en_name = "syslog1"
  en_persistenceflg = 1
  en_method = "logger Error Log: `errpt -l $1 | grep -v TIMESTAMP`"
```

```
errnotify:
  en_name = "syslog1"
  en_persistenceflg = 1
  en_method = "logger Error Log: $(errpt -l $1 | grep -v TIMESTAMP)"
```

Direct the last error entry (-l \$1) to the syslogd.  
Do not show the error log header (grep -v) or (tail -1).

```
errnotify:
  en_name = "syslog1"
  en_persistenceflg = 1
  en_method = "errpt -l $1 | tail -1 | logger -t errpt -p
  daemon.notice"
```

© Copyright IBM Corporation 2004

Figure 8-17. Directing Error Log Messages to syslogd

AU1612.0

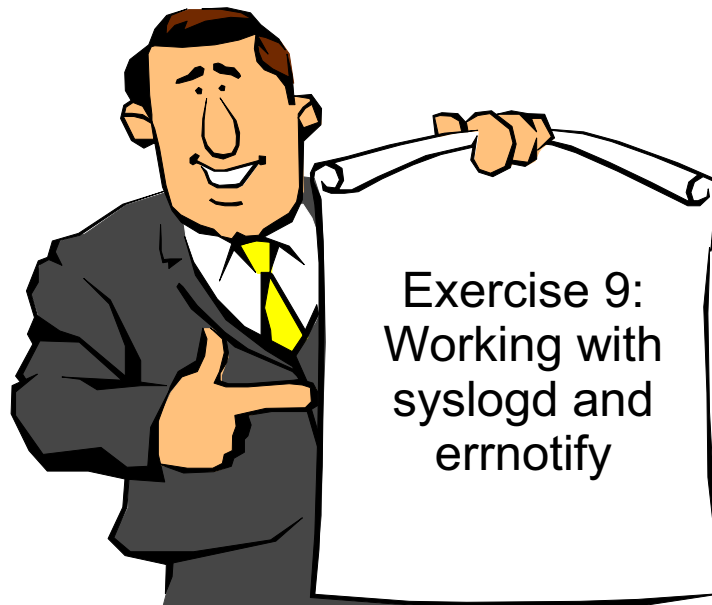
### Notes:

You can log error log events in the **syslog** by using the **logger** command with the **errnotify** ODM class. Whenever an entry is posted to the error log, this last entry will be passed to the **logger** command.

Note that you must use **backquotes** to do a command substitution before calling the **logger** command.

## Next Step

---



© Copyright IBM Corporation 2004

Figure 8-18. Next Step

AU1612.0

### **Notes:**

At the end of the lab, you should be able to:

- Configure the **syslogd** daemon
- Redirect **syslogd** messages to the Error Log
- Implement error notification with **errnotify**

---

## Checkpoint

---

1. Which command generates error reports?

---

---

2. Which type of disk error indicates bad blocks?

---

3. What do the following commands do?

**errclear** \_\_\_\_\_

**errlogger** \_\_\_\_\_

4. What does the following line in /etc/syslog.conf indicate:

**\*.debug errlog**

---

5. What does the descriptor **en\_method** in **errnotify** indicate?

---

---

---

© Copyright IBM Corporation 2004

Figure 8-19. Checkpoint

AU1612.0

### Notes:

## Unit Summary

---

- Use the **errpt** (**smit errpt**) command to generate error reports
- Different **error notification methods** are available
- Use **smit errdemon** and **smit errclear** to maintain the error log
- Some components use **syslogd** for error logging
- **syslogd** configuration file is **/etc/syslog.conf**
- **syslogd** and Error Log Messages could be redirected

© Copyright IBM Corporation 2004

Figure 8-20. Unit Summary

AU1612.0

### **Notes:**



# Unit 9. Diagnostics

## What This Unit Is About

This unit is an overview of diagnostics available in AIX.

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Use the **diag** command to diagnose hardware
- List the **different** diagnostic program modes
- Use the **System Management Services** on RS/6000 PCI models that do not support **diag**

## How You Will Check Your Progress

Accountability:

- Activity
- Checkpoint questions

## References

Online      *Understanding the Diagnostic Subsystem for AIX*

## Unit Objectives

---

After completing this unit, students should be able to:

- Use the **diag** command to diagnose hardware
- List the different **diagnostic** program **modes**
- Use the **System Management Services** on RS/6000 PCI models that do not support **diag**

© Copyright IBM Corporation 2004

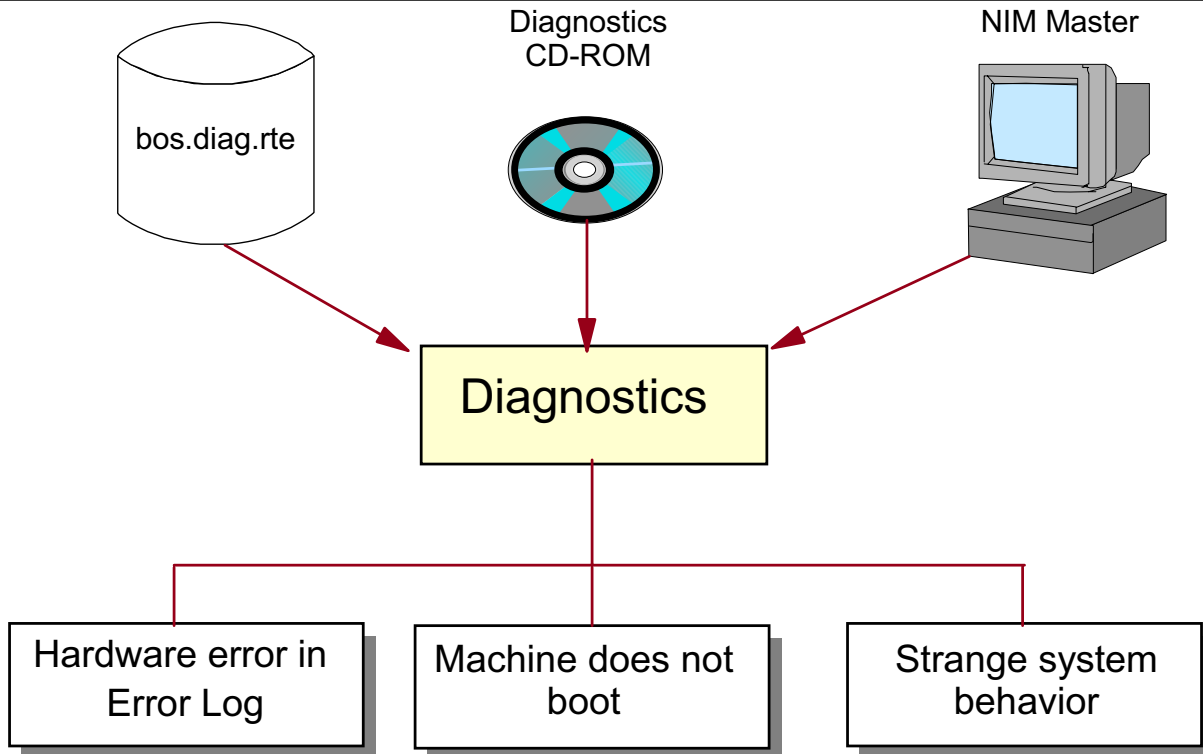
Figure 9-1. Unit Objectives

AU1612.0

### **Notes:**

## 9.1 Diagnostics

# When Do I Need Diagnostics?



© Copyright IBM Corporation 2004

Figure 9-2. When Do I Need Diagnostics?

AU1612.0

## Notes:

The lifetime of hardware is limited. Broken hardware leads to hardware errors in the error log, to systems that will not boot or to very strange system behavior.

The **diagnostic** package helps you to analyze your system and discover hardware that is broken. Additionally the **diagnostic** package provides information to service representatives that allows fast error analysis.

**Diagnostics** are available from different sources.

- A diagnostic package is shipped and installed with your AIX operating system. The fileset name is **bos.diag.rte**.
- **Diagnostic CD-ROMs** are available that allow you to diagnose a system that has no AIX installed. Normally the **diagnostic CD-ROM** is not shipped with the system.
- Diagnostic programs can be loaded from a **NIM master** (NIM=Network Installation Manager). This master holds and maintains different resources, for example a diagnostic package. This package could be loaded via the network to a NIM client, that is used to diagnose the client machine.

## The diag Command

```
# errpt
IDENTIFIER  TIMESTAMP    T   C  RESOURCE_NAME  DESCRIPTION
...
BF93B600   0505071399P   H   tok0                ADAPTER ERROR
...

# diag

A PROBLEM WAS DETECTED ON Thu May 6 09:40:22 1999

The Service Request Number(s)/Probable Cause or Causes:

850-902:   Error log analysis indicates hardware failure

60%       tok0         00-02      Token-Ring Adapter
40%       sysplanar0  00-00      System Planar
```

- **diag** allows testing of a device, if it's not busy
- **diag** allows analyzing the error log

© Copyright IBM Corporation 2004

Figure 9-3. The diag Command

AU1612.0

### Notes:

Whenever you detect a hardware problem, for example, a communication adapter error in the error log, use the **diag** command to diagnose the hardware.

The **diag** command allows testing of a device if the device is not busy. If any AIX process uses a device, the diagnostic programs cannot test it; they must have exclusive use of the device to be tested. Methods used to test devices that are busy are introduced later in this unit.

The **diag** command analyses the error log to fully diagnose a problem if run in the correct mode. It provides information that is very useful for the service representative, for example **SRNs** (Service Request Numbers) or probable causes.

Starting in AIX 5.1, there is a cross link between the AIX error log and diagnostics. When the **errpt** command is used to display an error log entry, diagnostic results related to that entry are also displayed.

## Working with diag (1 of 2)

# diag

FUNCTION SELECTION

801002

Move cursor to selection, then press Enter.

Diagnostic Routines

This selection will test the machine hardware. Wrap plugs and other advanced functions will not be used.

...

DIAGNOSTIC MODE SELECTION

801003

Move cursor to selection, then press Enter.

System Verification

This selection will test the system, but **will not analyze the error log**. Use this option to verify that the machine is functioning correctly after completing a repair or an upgrade.

Problem Determination

This selection tests the system and analyzes the error log if one is available. Use this option when a **problem is suspected** on the machine.

Figure 9-4. Working with diag (1 of 2)

AU1612.0

### Notes:

The **diag** command is menu driven, and offers different ways to test hardware devices or the complete system. Here is one method to test hardware devices with **diag**:

- Start the **diag** command. A welcome screen appears, which is not shown on the visual. After pressing Enter, the **FUNCTION SELECTION** menu is shown.
- Select **Diagnostic Routines**, which allows you to test hardware devices.
- The next menu is **DIAGNOSTIC MODE SELECTION**. Here you have two selections:

**System Verification** tests the hardware without analyzing the error log. This option is used after a repair to test the new component. If a part is replaced due to an error log analysis, the service provider must log a repair action to reset error counters and prevent the problem from being reported again. Running Advanced Diagnostics in System Verification mode will log a repair action.

**Problem Determination** tests hardware components **and** analyzes the error log. When you suspect a problem on a machine, use this selection. Do not use this selection after you have repaired a device, unless you remove the error log entries of the broken device.

## Working with diag (2 of 2)

DIAGNOSTIC SELECTION 801006

From the list below, select any number of resources by moving the cursor to the resource and pressing 'Enter'.  
 To cancel the selection, press 'Enter' again.  
 To list the supported tasks for the resource highlighted, press 'List'.

Once all selections have been made, press 'Commit'.  
 To exit without selecting a resource, press the 'Exit' key.

All Resources  
 This selection will select all the resources currently displayed.

sysplanar0	00-00	System Planar
proc0	00-00	Processor
mem0	00-0A	4MB Memory Simm
...		
hdisk0	00-00-0S-0,0	2.0 GB SCSI Disk Drive
...		
+ tok0	00-02	Token-Ring Adapter
...		

F1=Help      F4=List      F7=Commit      F10=Exit  
 F3=Previous Menu

© Copyright IBM Corporation 2004

Figure 9-5. Working with diag (2 of 2)

AU1612.0

### Notes:

In the next **diag** menu select the hardware devices that you want to test. If you want to test the complete system, select **All Resources**. If you want to test selected devices, press **Enter** to select any device, then press **F7** to commit your actions. In our example, we select the token-ring adapter.

If you press **F4** (List), **diag** presents tasks the selected devices support, for example:

- Run diagnostics
- Display hardware vital product data
- Display resource attributes
- Change hardware vital product data
- Run error log analysis

To start diagnostics, press **F7** (Commit).



## What Happens If a Device Is Busy?

ADDITIONAL RESOURCES ARE REQUIRED FOR TESTING

801011

No trouble was found. **However, the resource was not tested** because the device driver indicated that the resource was in use.

The resource needed is:

- tok0            00-02            Token-Ring Adapter

To test this resource, you can:

Free this resource and continue testing

Shutdown the system and run in maintenance mode

Run diagnostics from the Diagnostics Standalone Package

...

F3=Cancel

F10=Exit

© Copyright IBM Corporation 2004

Figure 9-6. What Happens If a Device Is Busy?

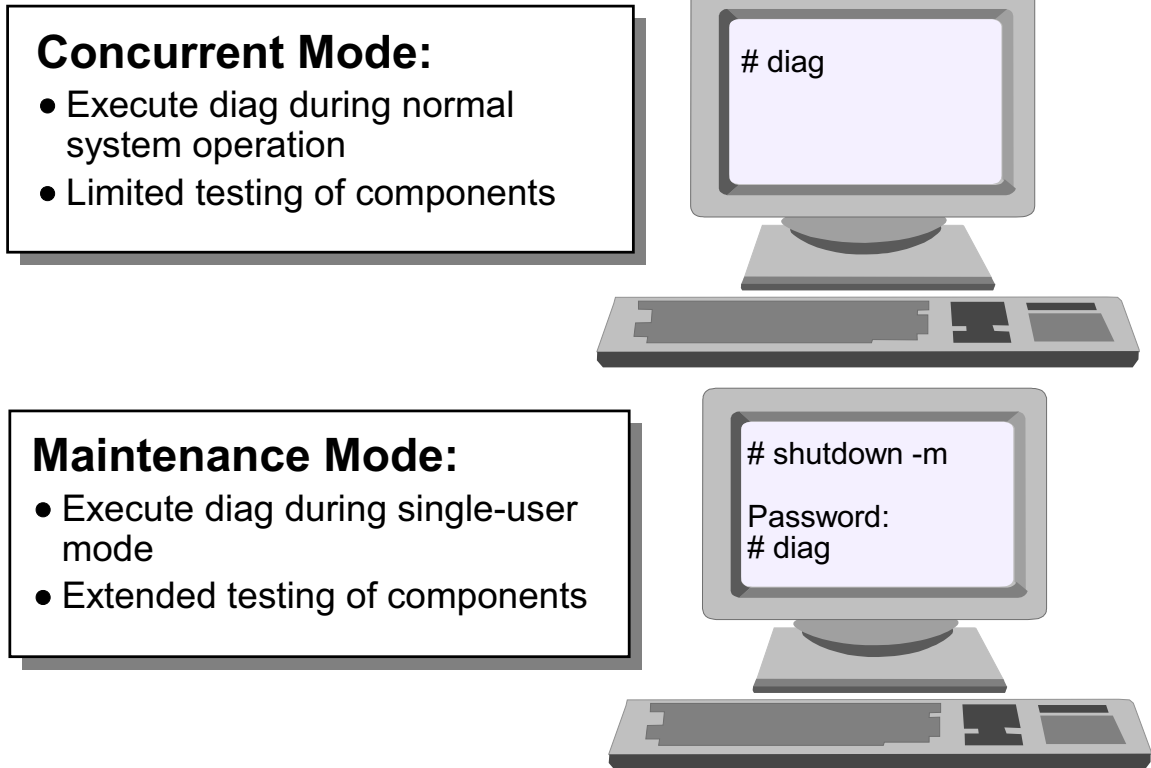
AU1612.0

### Notes:

If a device is busy, meaning the device is in use, the diagnostic programs do not permit testing the device or analyzing the error log.

That's what the visual shows: we selected the token-ring adapter, but the resource was not tested because the device was in use. To test the device we must free the resource. We must use another **diagnostic** mode to test this resource.

## Diagnostic Modes (1 of 2)



© Copyright IBM Corporation 2004

Figure 9-7. Diagnostic Modes (1 of 2)

AU1612.0

### Notes:

Three different diagnostic modes are available: concurrent mode, maintenance (single-user) mode and stand-alone (service) mode (covered on the next foil).

- **Concurrent Mode:**

Concurrent mode means that the diagnostic programs are executed during normal system operation. Certain devices can be tested, for example, a tape device that is currently not in use, but the number of resources that can be tested is very limited. Devices that are in use cannot be tested.

- **Maintenance (Single-User) Mode:**

To expand the list of devices that can be tested, one method is to take the system down to maintenance mode:

```
# shutdown -m
```

Enter the **root** password when prompted, and execute the **diag** command in the shell.

All programs except the operating system itself are stopped. All user volume groups are inactive, which extends the number of devices that can be tested in this mode.

But what do you do if your system does not boot or if you have to test a system without an installed AIX system? In this case you must use the **stand-alone mode**, which is introduced on the next visual.

## Diagnostic Modes (2 of 2)

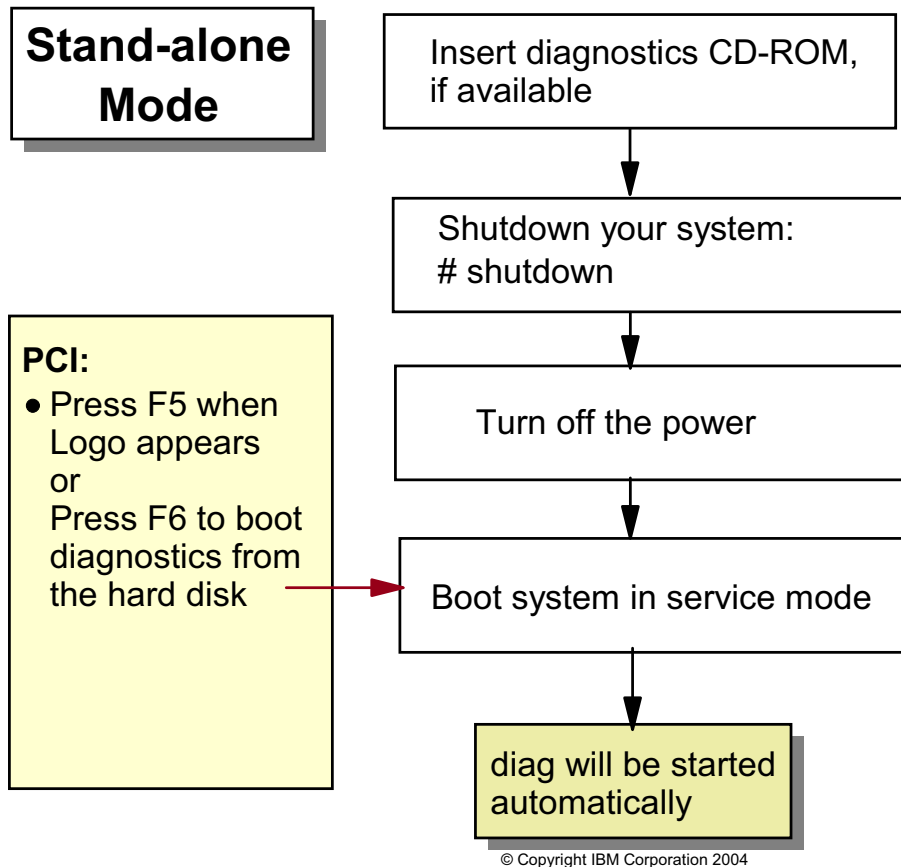


Figure 9-8. Diagnostic Modes (2 of 2)

AU1612.0

### Notes:

The **stand-alone mode** offers the greatest flexibility. You can test systems that do not boot or that have no operating system installed (the latter requires a diagnostic CD-ROM).

Follow these steps to start up diagnostics in **stand-alone mode**:

- If you have a diagnostic CD-ROM (or a diagnostic tape), insert it into the system. If you do not have a diagnostic CD-ROM, you boot diagnostics from the hard disk.
- Shut down the system. When AIX is down, turn off the power.
- Turn on power.
- Press **F5** when an acoustic beep is heard and icons are shown on the display. This simulates booting in service mode (logical key switch).
- The **diag** command will be started automatically, either from the hard disk or the diagnostic CD-ROM.
- At this point you can start your diagnostic routines.

## diag: Using Task Selection

# diag

FUNCTION SELECTION

801002

Move cursor to selection, then press Enter.

...

Task Selection (Diagnostics, Advanced Diagnostics, Service Aids, etc.)

This selection will list the tasks supported by these procedures. Once a task is selected, a resource menu may be presented showing all resources supported by the task.

...

- Run diagnostics
  - Display service hints
  - Display hardware error report
  - Display software product data
  - Display system configuration
  - Display hardware vital product data
  - Display resource attributes
  - Certify media
  - Format media
  - Local area network Analyzer
  - SCSI bus analyzer
  - Download microcode
  - Display or change bootlist
  - Periodic diagnostics
  - Disk maintenance
  - Run error log analysis
- ... and other tasks that are dependent on the devices in the system.

Figure 9-9. diag: Using Task Selection

AU1612.0

### Notes:

The **diag** command offers a wide number of additional tasks that are hardware-related. All these tasks can be found after starting the **diag** main menu and selecting **Task Selection**.

The tasks that are offered are hardware- (or resource) related. For example, if your system has a **service processor**, you will find service processor maintenance tasks, which you don't find on machines without a service processor. Or, on some systems you find tasks to maintain **RAID** and **SSA** storage systems.

# Diagnostic Log

For a summary output:

```
# /usr/lpp/diagnostics/bin/diagrpt -r
```

ID	DATE/TIME	T	RESOURCE_NAME	DESCRIPTION
DC00	Mon Jul 24 18:01:29	I	diag	Diagnostic Session was started
DA00	Mon Jul 24 17:57:16	N	sysplanar0	No Trouble Found
DA00	Mon Jul 24 17:57:12	N	mem0	No Trouble Found
DA00	Mon Jul 24 17:56:49	N	rmt0	No Trouble Found
DC00	Mon Jul 24 17:55:28	I	diag	Diagnostic Session was started

```
# /usr/lpp/diagnostics/bin/diagrpt -a
```

```
IDENTIFIER:      DA00

Date/Time:       Mon Jul 24 17:57:16
Sequence Number: 71
Event type:      No Trouble Found

Resource Name:   sysplanar0
Resource Description: System Planar
Location:        00-00

Diag Session:    13092
Test Mode:       Console,Non-Advanced,Normal IPL,System
                 Verification, System Checkout

Description:     No Trouble Found

-----
IDENTIFIER:      DA00

Date/Time:       Mon Jul 24 17:57:12
Sequence Number: 70
Event type:      No Trouble Found
```

Figure 9-10. Diagnostic Log

AU1612.0

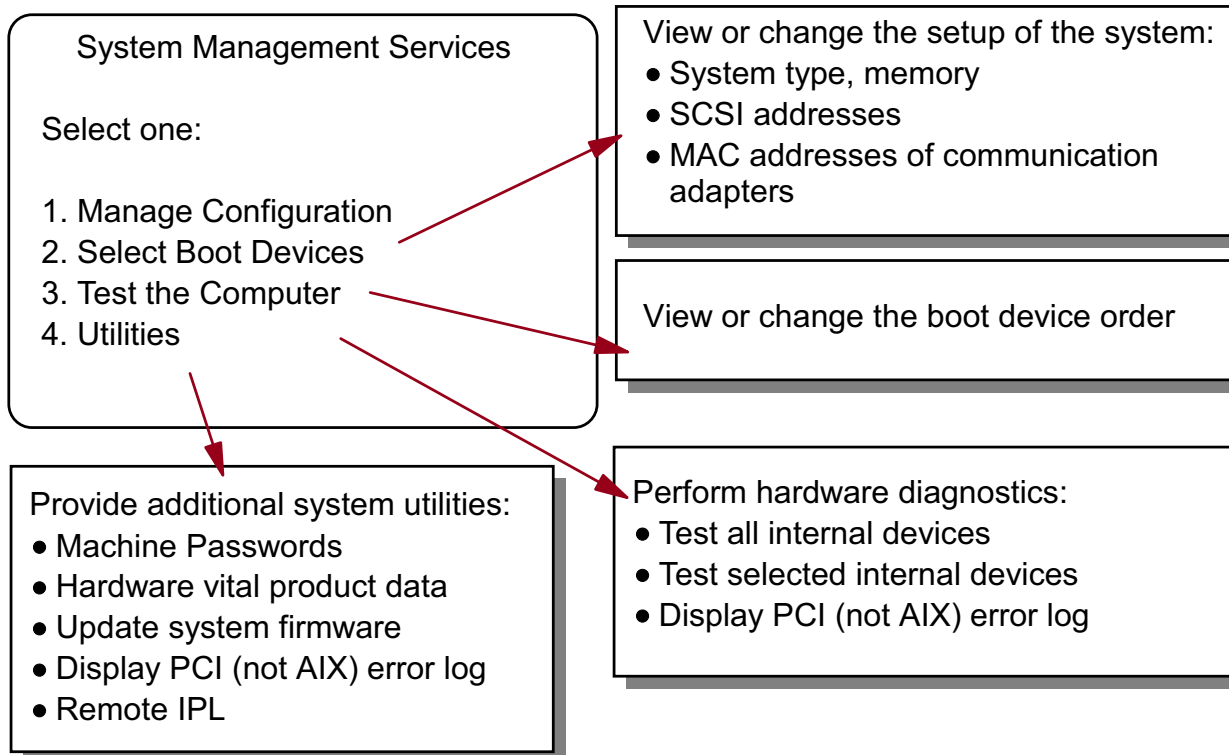
## Notes:

When diagnostics are run, the information is stored into a diagnostics log. The binary file is called `/var/adm/ras/diag_log`. The command `/usr/lpp/diagnostics/bin/diagrpt` is used to read the content of this file.

The **ID** column identifies the event that was logged. In the example above, **DC00** and **DA00** are shown. **DC00** indicated the diagnostics session was started and the **DA00** indicates No Trouble Found (NTF).

The **T** column indicates the type of entry in the log. **I** is for informational messages. **N** is for No Trouble Found. **S** shows the SRN (Service Request Number) for the error that was found. **E** is for an Error Condition.

## PCI: Using SMS for Diagnostics



© Copyright IBM Corporation 2004

Figure 9-11. PCI: Using SMS for Diagnostics

AU1612.0

### Notes:

The AIX **diag** is not supported on older PCI models (40P, 43P without LED). On these systems the **System Management Services** provide a selection, **Test the Computer**. Newer PCI systems that support the **diag** command do not offer this selection.

When you select **Test the Computer**, you can:

- Test all internal devices of the PCI model
- Test selected internal devices (for example memory or keyboard)
- Display the firmware error log

### Note:

Do not confuse the firmware (NVRAM) error log with the AIX error log. The firmware error log contains entries that are logged by the firmware and not from any AIX component. If your PCI system shows hardware errors during boot, always check your firmware error log.

External devices cannot be tested.

Other selections in the SMS are:

- **Manage Configuration:**

Use this selection when you want to view or change the setup of your system. Typical examples are changing a SCSI address or viewing the MAC address from a communication adapter to setup NIM (Network Installation Management).

- **Select Boot Devices:**

Use this selection when you want to view or change the boot order of your system, especially if the **bootlist** command is not supported.

- **Utilities:**

This selection offers a wide number of utilities:

- Manage machine passwords (normal and supervisory password: must be entered when SMS services are started) and start mode
- View or set hardware vital product data
- Update the system firmware, if newer firmware levels are required
- Display the firmware error log
- Set up booting from a remote NIM master (IPL = initial program load)



---

# Activity: Diagnostics

---



© Copyright IBM Corporation 2004

Figure 9-12. Activity: Diagnostics

AU1612.0

## **Notes:**

At the end of the activity, you should be able to:

- Execute hardware diagnostics in different modes

## **Instructions:**

Complete the following steps.

Only one person per machine can execute these commands.

- \_\_\_ 1. Start up diagnostics routines in **concurrent mode** and test a communication adapter of your system. What happens?

---

- \_\_\_ 2. Write down the difference between **System Verification** and **Problem Determination**:

---

---

- \_\_\_ 3. Using **Task Selection** query the vital product data of your **hdisk0**.  
\_\_\_\_\_
- \_\_\_ 4. Using **Task Selection** enable **Periodic Diagnostics** on your system. Who will be notified when a hardware error is posted to the error log?  
\_\_\_\_\_
- \_\_\_ 5. Start up diagnostic routines in **Maintenance Mode**. Write down the steps you executed:  
\_\_\_\_\_  
\_\_\_\_\_
- \_\_\_ 6. Test the communication adapter again in maintenance mode. What happens now?  
\_\_\_\_\_
- \_\_\_ 7. Start up the diagnostic routines in **stand-alone mode**. Write down the steps you executed:  
\_\_\_\_\_  
\_\_\_\_\_
- \_\_\_ 8. Try to **certify** your **hdisk0**. What happens?  
\_\_\_\_\_
- \_\_\_ 9. View the contents of the diagnostics log using both the summary format and detailed format. Did you find any errors?
- \_\_\_ 10. Exit diagnostics and reboot your system in normal mode.

**END OF ACTIVITY**

## Activity Solution:

Here are the solutions for the activity:

- \_\_\_ 1. Start up diagnostic routines in **concurrent mode** and test a communication adapter of your system. What happens?

**Normally, the adapter is used and could not be tested.**

- \_\_\_ 2. Write down the difference between **System Verification** and **Problem Determination**:

**System Verification: Test a resource. Do not analyze the error log**

**Problem Determination: Test a resource, and analyze the error log**

**Problem Determination should not be used after a hardware repair, unless the error log has been cleaned up.**

- \_\_\_ 3. Using **Task Selection** query the vital product data of your **hdisk0**.

**Task Selection**

- **Display Hardware Vital Product Data**

- **Select hdisk0**

- \_\_\_ 4. Using **Task Selection** enable **Periodic Diagnostics** on your system. Who will be notified, when a hardware error is posted to the error log?

**Task Selection**

- **Periodic Diagnostics**

- **Enable Automatic Error Log Analysis**

**All members of group system will be notified (default: root user)**

- \_\_\_ 5. Start up diagnostic routines in **Maintenance Mode**. Write down the steps you executed:

**# shutdown -m**

- **Enter root password**

**# diag**

- \_\_\_ 6. Test the communication adapter again in maintenance mode. What happens now?

**In single-user mode, the communication adapter is not used. Therefore it could be tested.**

- \_\_\_ 7. Start up the diagnostic routines in **stand-alone mode**. Write down the steps you executed:

**# shutdown -F**

- **Power-Off**

- **Power-on**

- **PCI: Press F6 when logo appears**

- **diag is started automatically**

\_\_\_ 8. Try to certify your **hdisk0**. What happens?

**Certification is not possible, because diagnostics have been started from the disk.**

\_\_\_ 9. Exit diagnostics and reboot your system in normal mode.

\_\_\_ 10. View the contents of the diagnostics log using both the summary format and the detailed format. Did you find any error?

**# /usr/lpp/diagnostics/bin/diagrpt -r | more**

**# /usr/lpp/diagnostics/bin/diagrpt -a | more**

**END OF ACTIVITY**

---

## Checkpoint

---

1. **T/F:** The **diag** command is supported on all RS/6000 models.  

---
2. What diagnostic modes are available on a RS/6000?  

---
3. How can you diagnose a communication adapter that is used during normal system operation?  

---

© Copyright IBM Corporation 2004

Figure 9-13. Checkpoint

AU1612.0

### **Notes:**

## Unit Summary

---

- Diagnostics are supported from **hard disk, diagnostic CD-ROM** and over the **network (NIM)**.
- There are three diagnostic modes: **concurrent, maintenance** and **stand-alone**.
- The **diag** command allows **testing and maintaining** the hardware (Task Selection).

© Copyright IBM Corporation 2004

Figure 9-14. Unit Summary

AU1612.0

### **Notes:**

# Unit 10. The AIX System Dump Facility

## What This Unit Is About

This unit outlines how to maintain the AIX system dump facility.

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Explain the meaning of a system dump
- Determine and change the primary and secondary dump devices
- Create a system dump under different conditions
- Execute the **snap** command
- Use the **kdb** command to check the system dump

## How You Will Check Your Progress

Accountability:

- Checkpoint questions
- Lab exercise

## References

Online      *Commands Reference*

## Unit Objectives

---

After completing this unit, students should be able to:

- Explain the meaning of a **system dump**
- Determine and change the **primary** and **secondary dump devices**
- **Create** a system dump
- Execute the **snap** command
- Use the **kdb** command to check a system dump

© Copyright IBM Corporation 2004

Figure 10-1. Unit Objectives

AU1612.0

### **Notes:**

If an AIX kernel - the major component of your operating system - crashes, a dump is created. This dump can be used to analyze the cause of the system crash.

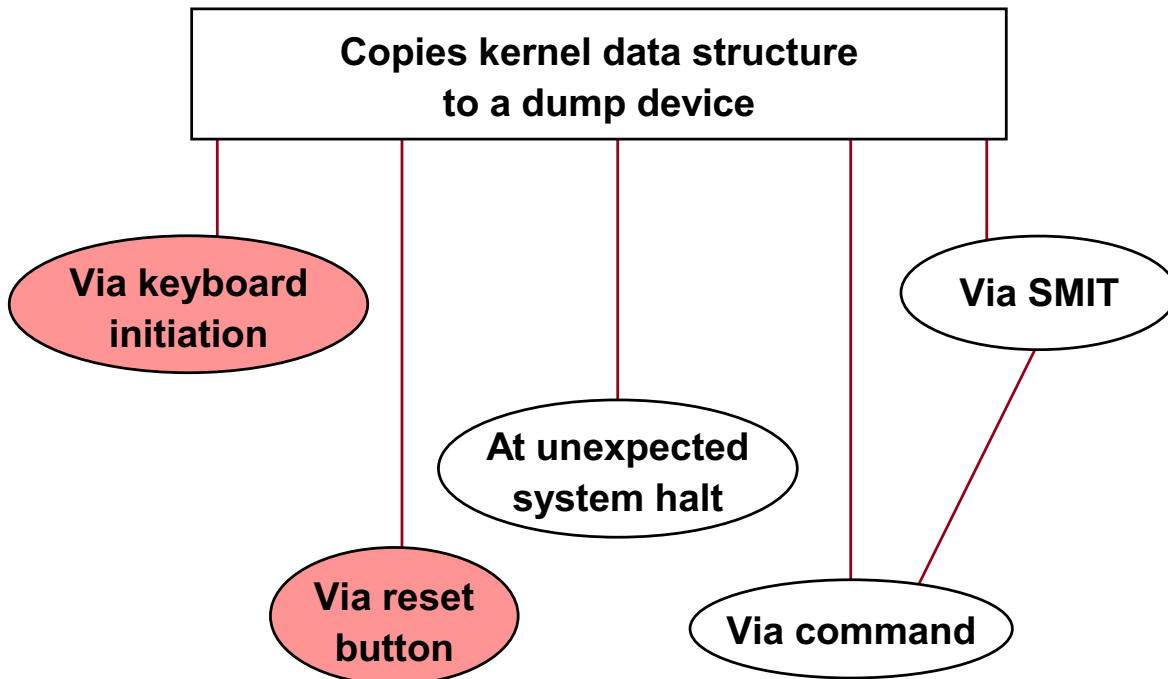
As administrator you have to know what a dump is, how the AIX dump facility is maintained, and how a dump can be started.

Before sending a dump to IBM, use the **snap** command to package the dump.



## 10.1 Working with System Dumps

# How a System Dump Is Invoked



 **By default, with the System Key in Service**

© Copyright IBM Corporation 2004

Figure 10-2. How a System Dump Is Invoked

AU1612.0

## Notes:

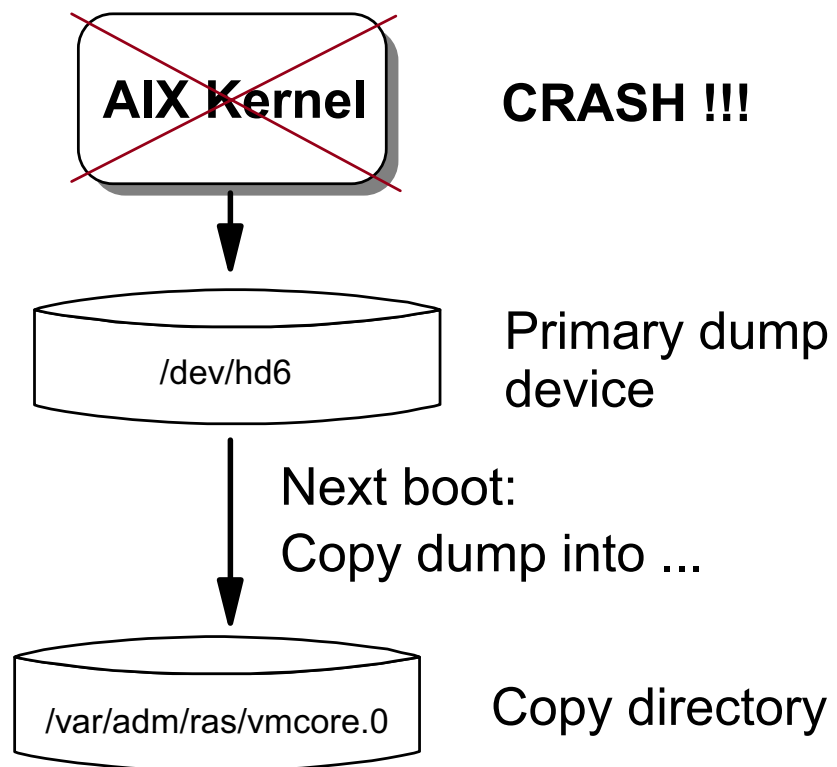
1. A set of special keys on the console keyboard (if it is an lft) can invoke a system dump on a classical RS/6000, when the front panel keylock has been set to service mode.
2. A dump can also be invoked when the reset button is pressed with the front panel keylock set to service mode.
3. If a kernel panic occurs, a dump will be invoked automatically.
4. The superuser can issue a command directly, or through **smit**, to invoke a system dump.

Usually, for persistent problems, the raw dump data is placed on a portable media, such as tape, and sent to a higher level of AIX support for analysis.

The raw dump data can be formatted into readable output via the **kdb** command.

The default setup of the system can be altered with the **sysdumpdev** command. Using this you can configure system dumps to occur regardless of System Key position - which is handy for PCI-bus systems, as they don't have a Switch Key.

## When a Dump Occurs



© Copyright IBM Corporation 2004

Figure 10-3. When a Dump Occurs

AU1612.0

### Notes:

If the AIX kernel crashes (system-initiated or user-initiated) kernel data is written to the primary dump device, which is by default **/dev/hd6**, the paging device. After a kernel crash AIX must be rebooted.

During the next boot, the dump is copied (remember: rc.boot 2) into a dump directory, the default is **/var/adm/ras**. The dump file name is **vmcore.x**, where x indicates the number of the dump (for example 0 indicates the first dump).

# The sysdumpdev Command

```
# sysdumpdev -l ← List dump values
primary                /dev/hd6
secondary              /dev/sysdumpnull
copy directory         /var/adm/ras
forced copy flag       TRUE
always allow dump     FALSE
dump compression      ON

# sysdumpdev -p /dev/sysdumpnull ← Deactivate primary dump device (temporary)

# sysdumpdev -P -s /dev/rmt0 ← Change secondary dump device (Permanent)

# sysdumpdev -L ← Display information about last dump
Device name:           /dev/hd6
Major device number:   10
Minor device number:   2
Size:                  9507840 bytes
Date/Time:             Tue Jun 5 20:41:56 PDT 2001
Dump status:           0
```

© Copyright IBM Corporation 2004

Figure 10-4. The sysdumpdev Command

AU1612.0

## Notes:

Use the **sysdumpdev** command or SMIT to query or change the primary and secondary dump devices. AIX Version 4 and later maintains two system dump devices:

- Primary - usually used when you wish to save the dump data.
- Secondary - can be used to discard dump data (that is, **/dev/sysdumpnull**).

Make sure you know your system and know what your primary and secondary dump devices are set to. Your dump device can be a portable medium, such as a tape drive. AIX Version 4 and later uses **/dev/hd6** (paging) as the default dump device **unless** the system was **migrated** from AIX Version 3, in which case it will continue to use the AIX Version 3's dump device **/dev/hd7**

Flags for the **sysdumpdev** command:

```
-l                list
-e                estimate the size of a dump
-p                primary
```

-C	turns on compression
-c	turns off compression
-s	secondary
-P	make change permanent
-d directory	specifies the directory the dump is copied to at system boot. If the copy fails at boot time, the <b>-d</b> flag ignores the system dump (force copy flag = FALSE)
-D directory	specifies the directory the dump is copied to at system boot. If the copy fails at boot time, using the <b>-D</b> flag allows you to copy the dump to external media (force copy flag = TRUE)
-K	reset button will force a dump with the key in the normal position, or on a machine without a key switch. This option is linked to the “always allow dump” setting.
-z	writes to standard output the string containing the size of the dump in bytes and the name of the dump device, if a new dump is present

Status values, as reported by **sysdumpdev -L**, correspond to dump LED codes (listed in full later) as follows:

<b>0 = 0c0</b>	dump completed
<b>-1 = 0c8</b>	no primary dump device
<b>-2 = 0c4</b>	partial dump
<b>-3 = 0c5</b>	dump failed to start

**Note:** If status is -3, size usually shows as 0, even if some data was written.

System dumps are usually recorded in the error log with the “DUMP\_STATS” label. Here the “Detail Data” section will contain the information that is normally given by the **sysdumpdev -L** command: the major device number, minor device number, size of the dump in bytes, time at which the dump occurred, dump type, that is, primary or secondary, and the dump status code.

### AIX 5.3 Enhancement

AIX 5.3 adds the ability to send the system dump to DVD media. The DVD device could be used as a primary or secondary dump device. In order to get this functionality the target DVD device should be DVD-RAM or writable DVD. Remember to have inserted an empty writable DVD in the drive when using the **sysdumpdev** command, or when you require the dump to be copied to the DVD at boot time after a crash. If the DVD media is not present the commands will give error messages or will not recognize the device as suitable for system dump copy.

During the creation of the system dump, additional information is displayed on the TTY about the progress of the system dump.

```
# sysdumpstart -p
Preparing for AIX System Dump . . .
Dump Started .. Please wait for completion message
AIX Dump .. 23330816 bytes written - time elapsed is 47 secs
Dump Complete .. type=4, status=0x0, dump size:23356416 bytes
Rebooting . . .
```

At this time, the kernel debugger and the 32-bit kernel needs to be enabled to see this function and we tested the functionality only on the S1 port. However, this limitation may change in the future.

Following a system crash there exist scenarios where a system dump may crash or fail without one byte of data written out to the dump device, for example power off or disk errors. For cases where a failed dump does not include the dump minimal table, it is very useful to save some trace back information in the NVRAM. From Version 5.3 the dump procedure is enhanced to use the NVRAM to store minimal dump information. In case that the dump fails, we can use the `sysdumpdev -vL` command (-v is the new verbose flag) to check the reason for the failure.

## Dedicated Dump Device (1 of 2)

- Servers with real memory > 4 GB, will have a dedicated dump device created at installation time

System Memory Size	Dump Device Size
4 GB to, but not including, 12 GB	1 GB
12, but not including, 24 GB	2 GB
24, but not including, 48 GB	3 GB
48 GB and up	4 GB

© Copyright IBM Corporation 2004

Figure 10-5. Dedicated Dump Device (1 of 2)

AU1612.0

### **Notes:**

This dedicated dump device is automatically created and requires no user intervention. The default name of the dump device is lg\_dumplv.

## Dedicated Dump Device (2 of 2)

---

```
/bosinst.data
```

```
·  
·  
·
```

```
large_dump:
```

```
DUMPDEVICE = /dev/lg_dumplv  
SIZE_GB = 1
```

© Copyright IBM Corporation 2004

Figure 10-6. Dedicated Dump Device (2 of 2)

AU1612.0

### **Notes:**

This stanza has been added to the bosinst.data file.

The dedicated dump device size is determined by the amount of memory at system install time.

The dump device name and size can be changed by using the businst.date file on a diskette of boot time.



## The sysdumpdev Command

```
# sysdumpdev -e ← Estimate dump size  
0453-041 estimated dump size in bytes: 52428800
```

```
# sysdumpdev -C ← Turn on dump compression
```

```
# sysdumpdev -e  
0453-041 estimated dump size in bytes: 10485760
```

Use this information to size the /var file system

© Copyright IBM Corporation 2004

Figure 10-7. The sysdumpdev Command

AU1612.0

### Notes:

You should size the /var file system so there is enough space to hold the dump information should your machine ever crash.

The **sysdumpdev -e** command will provide an estimate of the amount of space needed. The size of the dump device is directly related to the amount of RAM on your machine. The more RAM on the machine, the more space that will be needed on the disk. Machines with 16 GB of RAM may need 2 GB of dump space.

In 4.3.2, a option was added to compress the dump data before it is written. To turn on dump compression run **sysdumpdev -C**. This will reduce the amount of space needed by approximately half. To turn off compression use **sysdumpdev -c**.

## dumpcheck Utility

---

- The **dumpcheck** utility will do the following when enabled:
  - Estimate the dump or compressed dump size using **sysdumpdev -e**
  - Find the dump logical volumes and copy directory using **sysdumpdev -l**
  - Estimate the primary and secondary dump device sizes
  - Estimate the copy directory free space
  - Report any errors in the error log file

© Copyright IBM Corporation 2004

Figure 10-8. dumpcheck Utility

AU1612.0

### **Notes:**

A new utility in AIX 5L is the **/usr/lib/ras/dumpcheck** utility. It is used to check the disk resources used by the system dump facility. The command logs an error if either the largest dump device is too small to receive the dump or there is insufficient space in the copy directory when the dump device is a paging space.

If the dump device is a paging space, **dumpcheck** will verify if the free space in the copy directory is large enough to copy the dump.

If the dump device is a logical volume, **dumpcheck** will verify it is large enough to contain a dump.

If the dump device is a tape, **dumpcheck** will exit without message.

Any time a problem is found, **dumpcheck** will log an entry in the error log. If the **-p** flag is present, it will display a message to stdout and mail the information to the root user.

In order to be effective, the **dumpcheck** utility must be enabled. Verify that **dumpcheck** has been enabled by using the following command:

```
# crontab -l | grep dumpcheck
0 15 * * * /usr/lib/ras/dumpcheck >/dev/null 2>&1
```

By default it is set to run at 3 p.m. each afternoon.

Enable the dumpcheck utility by using the **-t** flag. This will create an entry in the root crontab if none exists. In this example the **dumpcheck** utility is set to run at 2 p.m.:

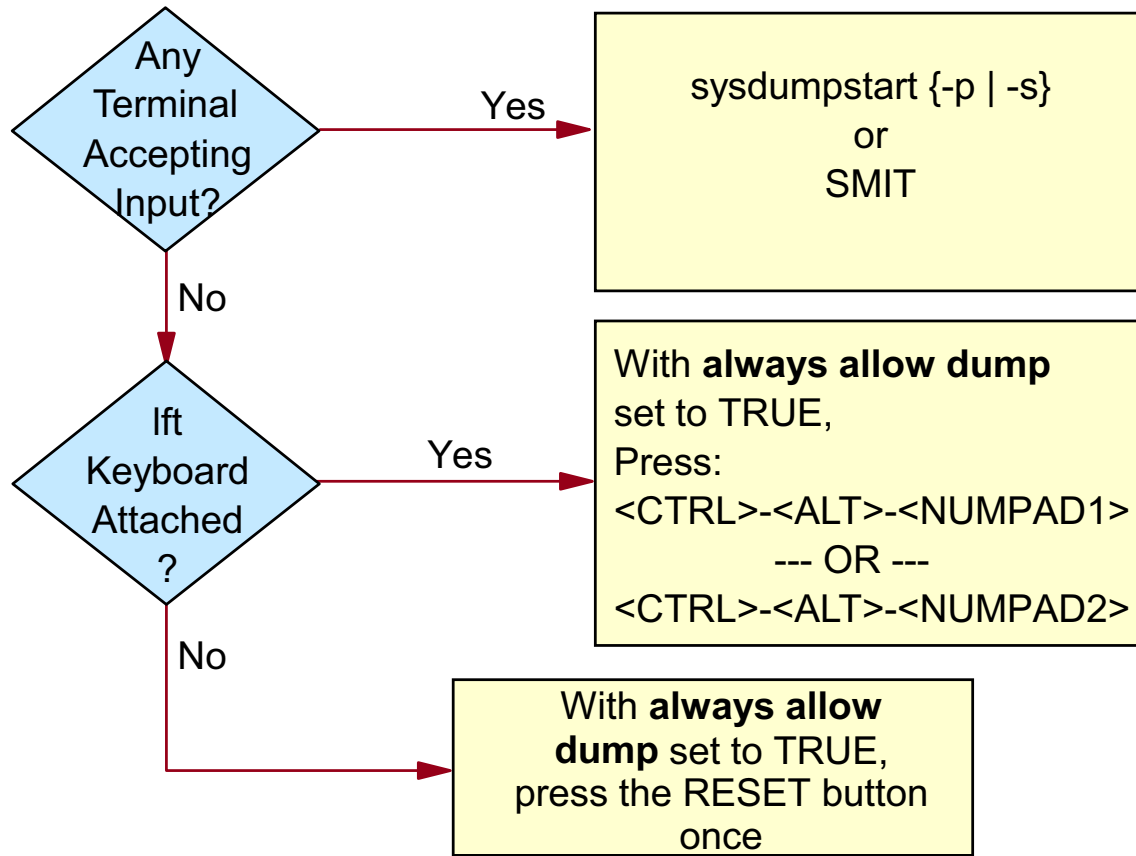
```
# /usr/lib/ras/dumpcheck -t "0 14 * * *"
```

For best results, set **dumpcheck** to run when the system is heavily loaded. This will identify the maximum size the dump will take. The default time is set for 3 p.m.

If you use the **-p** flag in the crontab entry, root will be sent a mail with the standard output of the dumpcheck command:

```
#/usr/lib/ras/dumpcheck -p
```

## Methods of Starting a Dump



© Copyright IBM Corporation 2004

Figure 10-9. Methods of Starting a Dump

AU1612.0

### Notes:

There are three ways for a user to invoke a system dump. Which method is used depends on the condition of the system.

If there is a kernel panic, the system will automatically dump the contents of real memory to the primary dump device.

If the system has halted, but the keyboard will still accept input, a dump can be forced by pressing the **<ctrl-alt-NUMPAD1>** key sequence. This is only possible with an lft keyboard. An ASCII keyboard does not have an "alt" key.

If the keyboard is no longer accepting input, a dump can be created by turning the key to the service position and pressing the reset button. (Pressing the reset button twice will cause the system to reboot.)

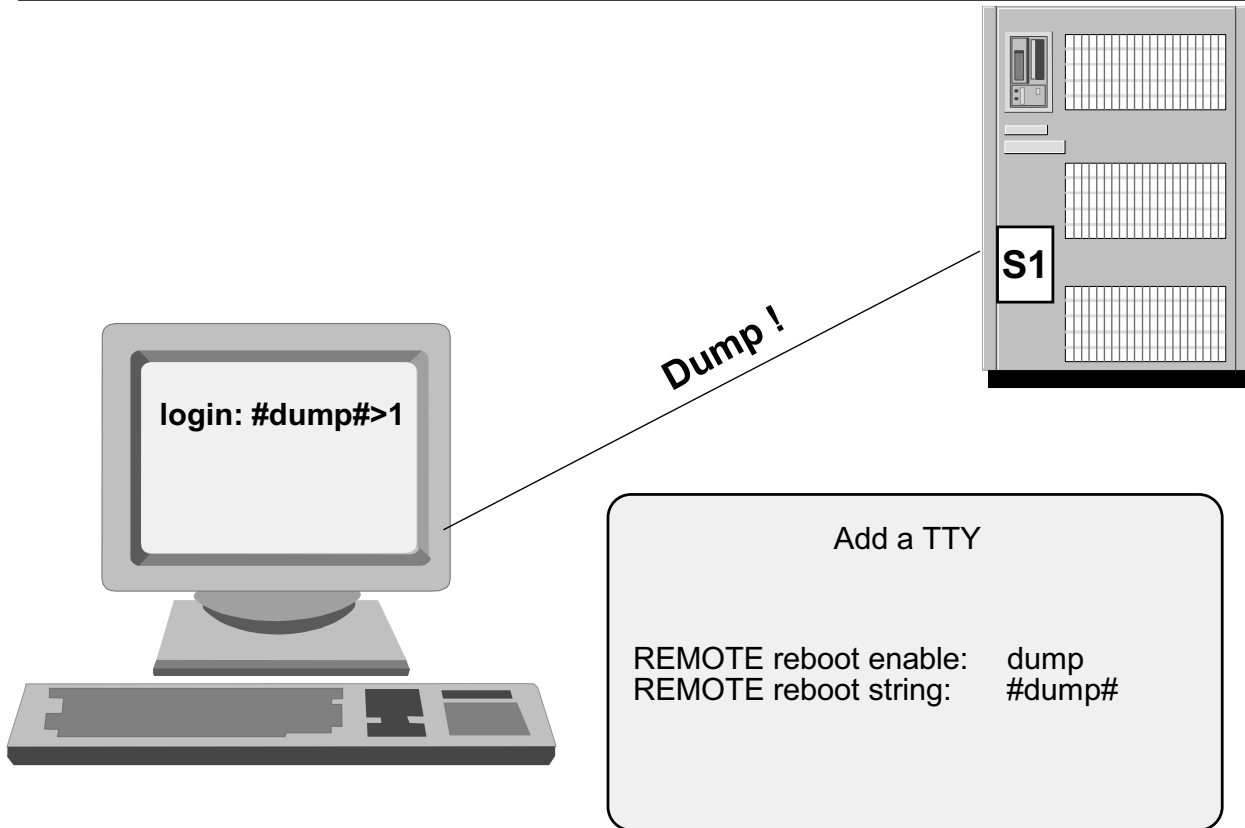
The third method for a user to invoke a dump is to run the **sysdumpstart** command or invoke it through SMIT (fastpath **dump**).

To invoke the dump using the keyboard or the reset button, the "Always allow dump" option must be set to TRUE. This can be done using **sysdumpdev -K**.

Bear in mind that if your system is still operational, a dump taken at this time will not assist in problem determination. A relevant dump is one taken at the time of the system halt.

Now, what can you do if you have no lft terminal available and your machine is a PCI model? This is covered on the next page.

# Start a Dump from a TTY



© Copyright IBM Corporation 2004

Figure 10-10. Start a Dump from a TTY

AU1612.0

## Notes:

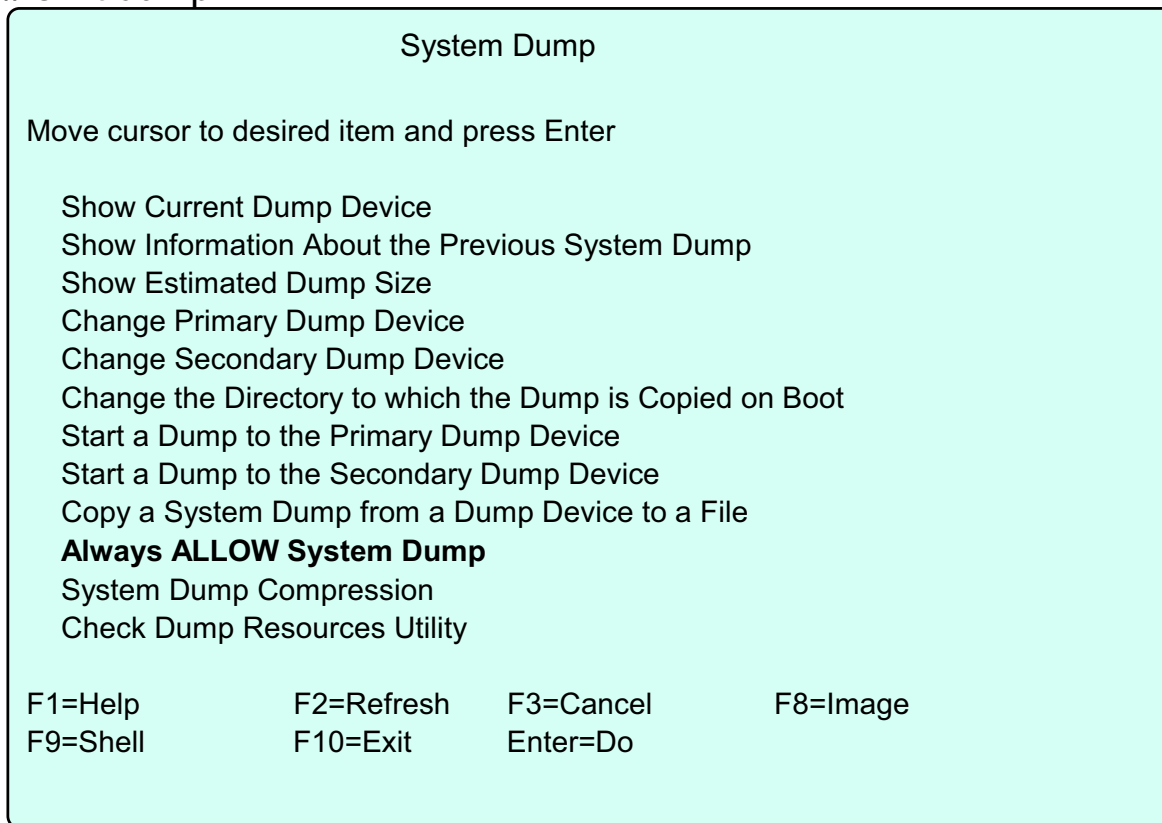
Another possibility allows starting a dump from a terminal. This might be very important if your system does not have an lft terminal attached.

To enable a terminal for starting a dump, you must set **REMOTE reboot enable** to a value of **dump**, when adding or changing a tty. Then specify a self-defined string, for example, **#dump#** to start the dump from a terminal.

This string must be entered at the login line on the terminal, and the string must be followed by a **1** key. Any character other than '1' aborts the dump process.

## Generating Dumps with smit

# smit dump



© Copyright IBM Corporation 2004

Figure 10-11. Generating Dumps with smit

AU1612.0

### Notes:

You can use the SMIT dump interface to work with the dump facility. The menu items that show or change the dump information use the **sysdumpdev** command.

A very important item is **Always Allow System Dump**. If you set this option to yes, the **CTRL-ALT-1** (numpad) and **CTRL-ALT-2** (numpad) key sequence will start a dump even when the key switch is in **normal** position. The reset button also starts a dump when this item is set to yes.

## Dump-related LED Codes

<b>0c0</b>	<b>Dump completed successfully</b>
0c1	An I/O error occurred during the dump
<b>0c2</b>	<b>Dump started by user</b>
0c4	Dump completed unsuccessfully. Not enough space on dump device. Partial dump available
0c5	Dump failed to start. Unexpected error occurred when attempting to write to dump device - e.g. tape not loaded
0c6	Secondary dump started by user
0c8	Dump disabled. No dump device configured
<b>0c9</b>	<b>System-initiated panic dump started</b>
0cc	Failure writing to primary dump device. Switched over to secondary

© Copyright IBM Corporation 2004

Figure 10-12. Dump-related LED Codes

AU1612.0

### Notes:

If a system dump is initiated via a kernel panic, the LEDs on a RS/6000 will display **0c9** while the dump is in progress, and then either a flashing **888** or a steady **0c0**.

All of the LED codes following the flashing **888** (remember: you must use the reset button) should be recorded and passed to IBM. While rotating through the **888** sequence, you will encounter one of the shown codes. The code you want is **0c0**, indicating that the dump completed successfully.

For user-initiated system dumps to the primary dump device, the LED codes should indicate **0c2** for a short period, followed by **0c0** upon completion.

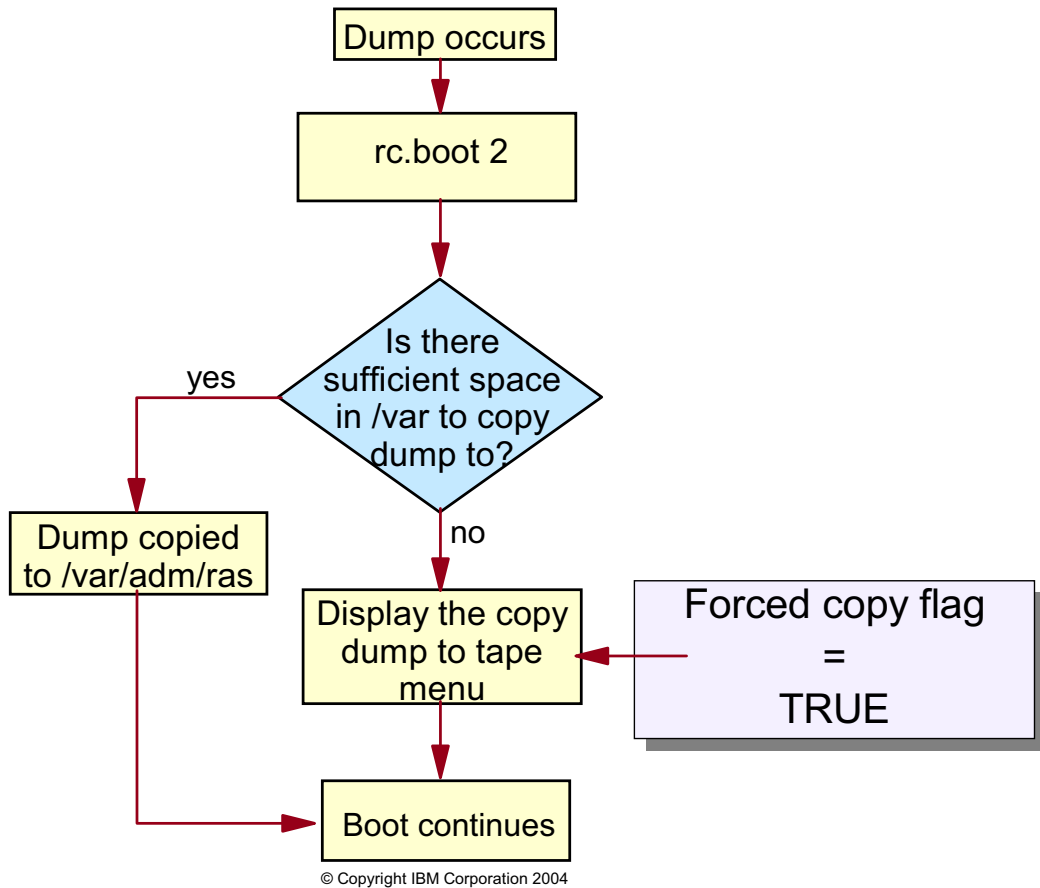
Other common codes include:

- 0c1**                      An I/O error occurred during the dump.
- 0c4**                      Indicates that the dump routine ran out of space on the specified device. It may still be possible to examine and use the data on the dump device, but this tells you that you should increase the size of your dump device.



- 0c5** Check the availability of the medium to which you are writing the dump (for example, whether the tape is in the drive and write enabled).
- 0c6** This is used to indicate a dump request to the secondary device.
- 0c7** A network dump is in progress, and the host is waiting for the server to respond. The value in the three-digit display should alternate between **0c7** and **0c2** or **0c9**. If the value does not change, then the dump did not complete due to an unexpected error.
- 0c8** You have not defined a primary or secondary dump device. The system dump option is not available. Enter the **sysdumpdev** command to configure the dump device.
- 0c9** A dump started by the system did not complete. Wait for one minute for the dump to complete and for the three-digit display value to change. If the three-digit display value changes, find the new value on the list. If the value does not change, then the dump did not complete due to an unexpected error.
- 0cc** This code indicates that the dump could not be written to the primary dump device. Therefore the secondary dump device will be used. This code was introduced with AIX 4.2.1.

# Copying System Dump



© Copyright IBM Corporation 2004

Figure 10-13. Copying System Dump

AU1612.0

## Notes:

For RS/6000s with LED, after a crash, if the LED displays 0c0, then you know that a dump occurred and it completed successfully. At this point you have to reboot your system. If there is enough space to copy the dump from the paging space to the **/var/adm/ras** directory, then it will be copied directly.

If, however, at bootup the system determines that there is not enough space to copy the dump to **/var**, the **/sbin/rc.boot** script (which is executed at bootup) will call the **/lib/boot/srvboot** script. This script in turn calls on the **copydumpmenu** command, which is responsible for displaying the following menu which can be used to copy the dump to removable media:

### Copy a System Dump to Removable Media

The system dump is 583973 bytes and will be copied from /dev/hd6 to media inserted into the device from the list below.

Please make sure that you have sufficient blank, formatted media before proceeding.

Step One:        Insert blank media into the chosen drive.

Step Two:        Type the number for that device and press Enter.

	Device type	Path Name
>>> 1	tape/scsi/8mm	/dev/rmt0
2	Diskette Drive	/dev/fd0
88	Help?	
99	Exit	
>>> Choice	[1]	

# Automatically Reboot After a Crash

```
# smit chgsys
```

Change/Show Characteristics of Operating System

Type or select values in entry fields.  
Press Enter AFTER making all desired changes.

Maximum number of PROCESSES allowed per user	[128]
Maximum number of pages in block I/O BUFFER CACHE	[20]
<b>Automatically REBOOT system after a crash</b>	<b>false</b>
...	
Enable full CORE dump	false
Use pre-430 style CORE dump	false

F1=Help	F2=Refresh	F3=Cancel	F4=List
F5=Reset	F6=Command	F7=Edit	F8=Image
F9=Shell	F10=Exit	Enter=Do	

© Copyright IBM Corporation 2004

Figure 10-14. Automatically Reboot After a Crash

AU1612.0

## Notes:

If you want your system to reboot automatically after a dump, you must change the kernel parameter **autostart** to **true**. This can be easily done by the smit fastpath **smit chgsys**. The corresponding menu item is **Automatically REBOOT system after a crash**. Note that the default value is **true** in V 5.2.

If you do not want to use smit, execute the following command:

```
# chdev -l sys0 -a autorestart=true
```

**If you specify an automatic reboot, verify that the /var file system is large enough to store a system dump.**

## Sending a Dump to IBM

- Copy all system configuration data including a dump onto tape:

**snap -a -o /dev/rmt0**

- The are some AIX 5.3 enhancements
- Label tape with:
  - Problem Management Record (PMR) number
  - Command used to create tape
  - Block size of tape
- Support Center uses **kdb** to examine the dump

© Copyright IBM Corporation 2004

Figure 10-15. Sending a Dump to IBM

AU1612.0

### Notes:

Before sending a dump to the IBM Support Center, use the **snap** command to collect system data. **/usr/sbin/snap -a -o /dev/rmt0** will collect all the necessary data. In AIX 5.2, **pax** is used to write the data to tape. The Support Center will need the information collected by **snap** in addition to the dump and kernel. Do not send just the dump file **vmcore.x** without the corresponding AIX kernel. Without it, the analysis is not possible.

The AIX Systems Support Center will analyze the contents of the dump using the **kdb** command. The **kdb** command uses the kernel that was active on the system at the time of the halt.

The **snap** command was developed by IBM to simplify gathering configuration information. It provides a convenient method of sending **lspp** and **errpt** output to the support centers. It gathers system configuration information and compresses the information to a **pax** file. The file can then be downloaded to disk, or tape.

Some useful flags with the **snap** command are:

- a** Copies all system configuration information to **/tmp/ibmsupt** directory tree
- c** Creates a compressed tar image (snap.tar.Z) of all files in the **/tmp/ibmsupt** directory tree or other named output directory
- f** gather file system information
- g** gather general information
- k** gather kernel information
- D** gather dump and **/unix**
- t** creates tcpip.snap file; gather TCP/IP information

## Snap command AIX 5.3 Enhancements

AIX 5L V5.3 extends its functionality in using external scripts, letting the snap split up the output pax file into smaller pieces or extending the collected data.

### Extending snap to Run External Scripts

The scripts can either be parameters to the snap command, or in case that all is specified, the scripts are expected to be in the **/usr/lib/ras/snapscripts** directory.

### snapsplit Command

The snapsplit command is introduced in AIX 5L V5.3. The command splits the snap.pax.Z file into smaller files or rejoins back the splitted snap files. The command expects to be run from the **/tmp/ibmsupt** directory, where the snap.pax.Z file resides.

The snapsplit command creates the snap.hostname.timestamp.pax.Zxx files of size 4 MB or less. The snapsplit -u command rejoined the files to snap.hostname.timestamp.pax.Z file. You can take the timestamp for the -T flag from the name of the splitted files. The -T or -h flags available that enable you to handle snaps from different systems taken at different times. The -f flag enables you to handle renamed snap files.

### Splitting the snap Output File From the snap Command

The size of the snap output file can be an issue sometimes. There is a new flag, -O megabytes, introduced in AIX 5L V5.3 that enables you to split the snap output file. The snap command calls the snapsplit command. You can use the flag as follows to split the large snap output into smaller 4 MB files.

```
# snap -a -c -O 4
```

## Use kdb to Analyze a Dump

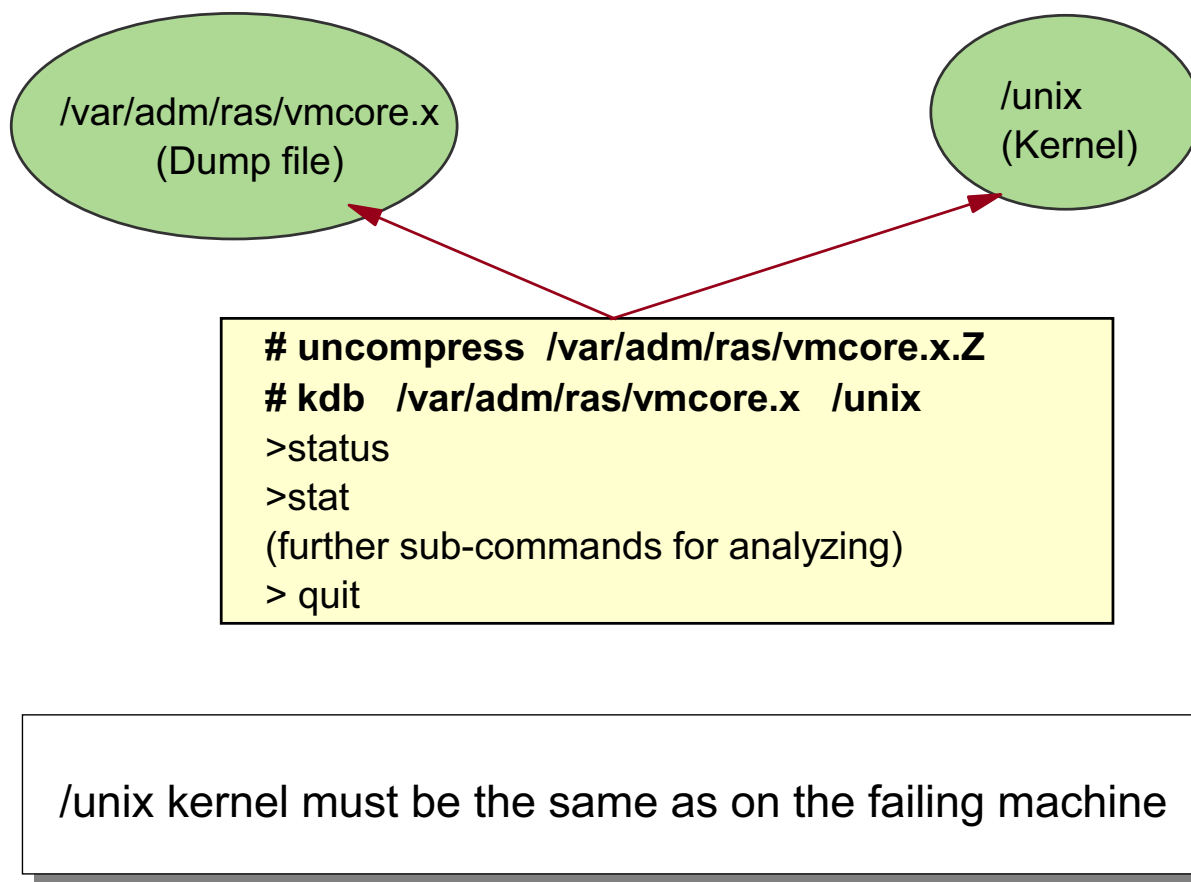


Figure 10-16. Use kdb to Analyze a Dump

AU1612.0

### Notes:

The **kdb** command is an interactive tool for the symbolic visualization of the operating system. Typically, **kdb** is used to examine kernel dumps in a system postmortem state. However, a live running system can also be examined with **kdb**, although due to the dynamic nature of the operating system, the various tables and structures often change while they are being examined, and this precludes extensive analysis.

Prior to AIX 5.1, the **crash** command was used instead of **kdb**.

To examine an active system, you would simply run the **kdb** command without any arguments.

For a dead system, a dump is analyzed using the **kdb** command with file name arguments.

To use **kdb**, the vmcore file must be uncompressed. After a crash it is typically named vmcore.x.Z which indicates it is in a compressed format. Use the **uncompress** command before using **kdb**. To analyze a dump file, enter:

```

# uncompress /var/adm/ras/vmcore.x.Z
# kdb /var/adm/ras/vmcore.x /unix

```

If the copy of **/unix** does not match the dump file, the following output will appear on the screen:

```
WARNING: dumpfile does not appear to match namelist  
>
```

If the dump itself is corrupted in some way, then the following will appear on the screen:

```
...  
dump /var/adm/ras/vmcore.x corrupted
```

Examining a system dump requires an in-depth knowledge of the AIX kernel. However there are two subcommands that might be useful.

The sub-command **status** displays the process that was active at the CPU when the crash occurred. The subcommand **stat** shows the machine status when the dump occurred.

To exit the **kdb** debug program, type **quit** at the > prompt.

The following example stops your running machine and creates a system dump

**Do not execute this in your production environment.**

```
# cat /unix > /dev/mem
```

The LEDs are 888, 102, 300, 0C0. (See Unit 2, Flashing 888)

LED 102 indicates that “a dump has occurred”.

LED 300 stands for crash code “Data Storage Interrupt (DSI)”

LED 0C0 means “Dump completed successfully”

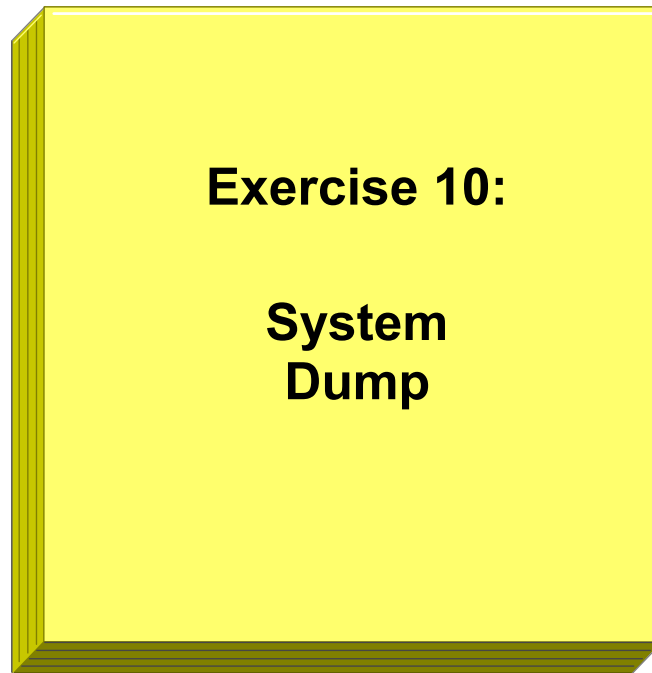
Afterwards you have to power-on your machine and you can analyze your dump.



---

## Next Step

---



© Copyright IBM Corporation 2004

Figure 10-17. Next Step

AU1612.0

### **Notes:**

At the end of the exercise, you should be able to:

- Initiate a dump
- Identify LED codes associated with the dump facility
- Use the **snap** command

# Checkpoint

---

1. What is the default primary dump device? Where do you find the dump file after reboot?

---

---

2. How do you turn on dump compression?

---

3. How do you start a dump from an attached LFT terminal?

---

---

---

4. If the copy directory is too small, will the dump, which is copied during the reboot of the system, be lost?

---

---

5. Which command should you execute before sending a dump to IBM?

---

© Copyright IBM Corporation 2004

Figure 10-18. Checkpoint

AU1612.0

## Notes:

## Unit Summary

---

- When a dump occurs kernel and system data are copied to the primary dump device.
- The system by default has a primary dump device (/dev/hd6) and a secondary device (/dev/sysdumpnull).
- During reboot the dump is copied to the copy directory (/var/adm/ras).
- A system dump should be retrieved from the system using the snap command.
- The support center uses the kdb debugger to examine the dump.

© Copyright IBM Corporation 2004

Figure 10-19. Unit Summary

AU1612.0

### **Notes:**



# Unit 11. Performance and Workload Management

## What This Unit Is About

This unit helps system administrators to identify the cause for performance problems. Workload management techniques will be discussed.

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Provide basic performance concepts
- Provide basic performance analysis
- Manage the workload on a system
- Work with the Performance Diagnostic Tool (PDT)

## How You Will Check Your Progress

Accountability:

- Checkpoint questions
- Exercises

## References

Online      *AIX Performance Tools Guide and Reference*

## Unit Objectives

---

After completing this unit, students should be able to:

- Provide basic performance concepts
- Provide basic performance analysis
- Manage the workload on a system
- Work with the Performance Diagnostic Tool (PDT)

© Copyright IBM Corporation 2004

Figure 11-1. Unit Objectives

AU1612.0

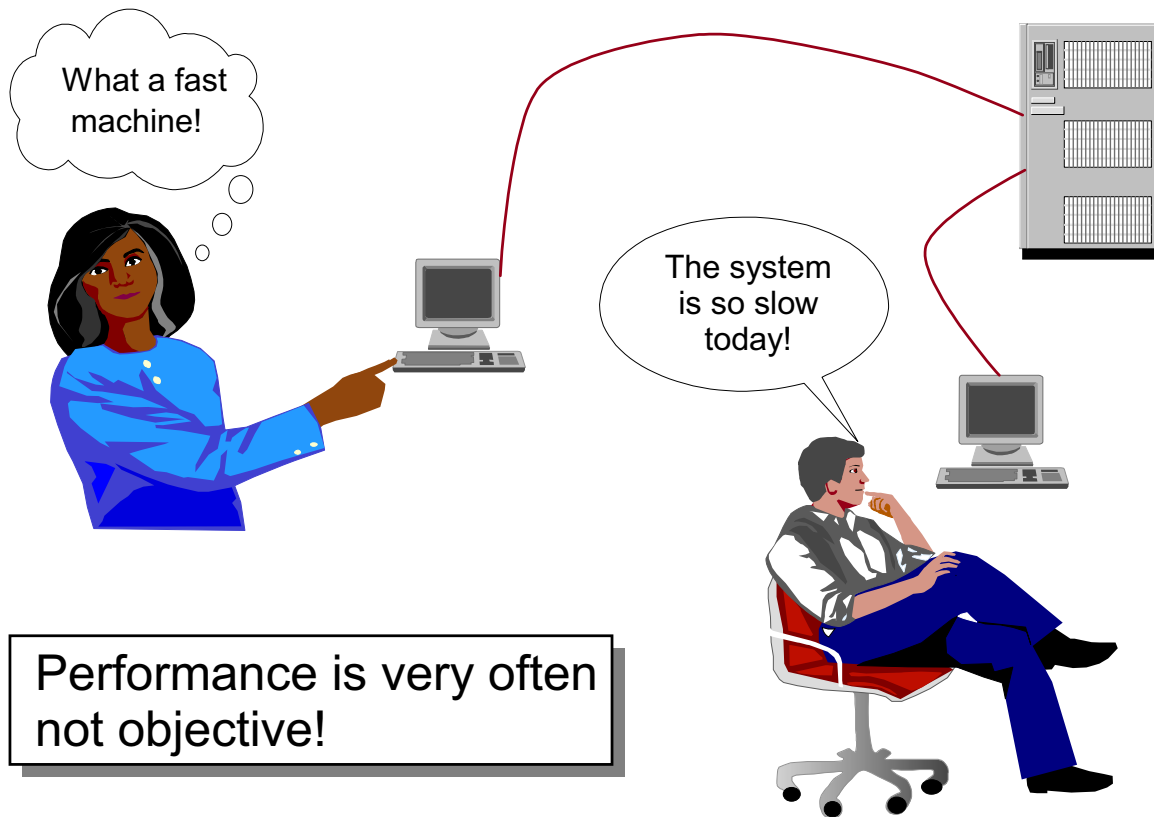
### **Notes:**

This course can only provide an introduction to **performance concepts and tools**. For a more thorough understanding of the subject you should take the AIX Performance Management class.

We will not be covering network monitoring, application development issues, or matters pertaining to SMP and SP machines. Also, this section will not explain the myriad of performance tuning techniques. All of that is addressed by the AIX Performance Management course.

## 11.1 Basic Performance Analysis and Workload Management

# Performance Problems



© Copyright IBM Corporation 2004

Figure 11-2. Performance Problems

AU1612.0

## Notes:

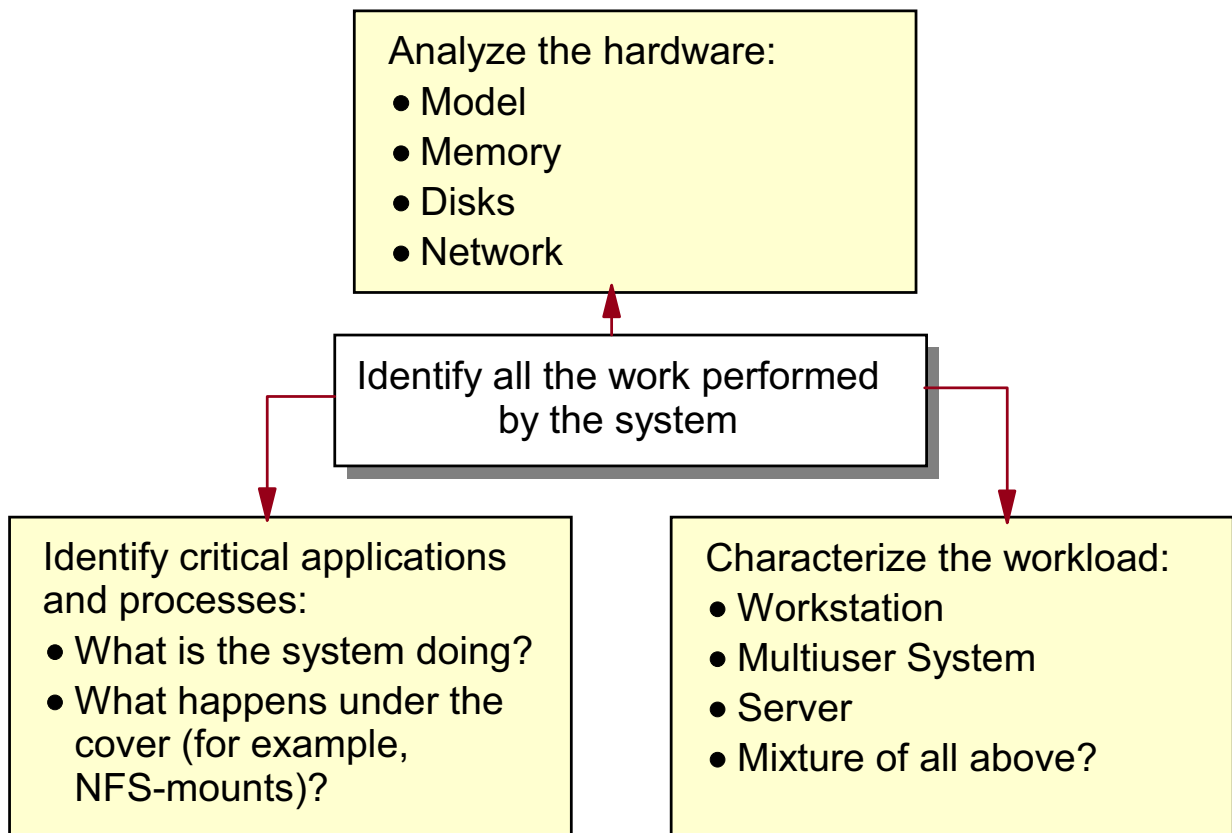
Everyone who uses a computer has an opinion about its performance. Unfortunately, these opinions are often completely different.

Whenever you get performance complaints from users, you must check if this is caused by a system problem or a user (application) problem. If you detect that the system is fast, that means you indicate the problem is user or application-related, check the following:

- What application is running slowly? Has this application always run slowly? Has the source code of this application been changed or a new version installed?
- Check the system's environment. Has something changed? Have files or programs been moved to other directories, disks or systems? Check the file systems to see if they are full.
- Finally, you should check the user's environment. Check the **PATH** variable to determine if it contains any **NFS-mounted** directories. They could cause a very long search time for applications or shared libraries.



# Understand the Workload



© Copyright IBM Corporation 2004

Figure 11-3. Understand the Workload

AU1612.0

## Notes:

If you detect the performance problem is system related, you must analyze the workload of your system. An accurate definition of the system's workload is critical to understanding its performance and performance problems. The workload definition must include not only the type and rate of requests to the system but also the exact software packages and application programs to be executed.

1. **Identify critical applications and processes.** Analyze and document what the system is doing and when the system is executing these tasks. Make sure that you include the work that your system is doing under the cover, for example providing NFS directories to other systems.
2. **Characterize the workload.** Workloads tend to fall naturally into a small number of classes:

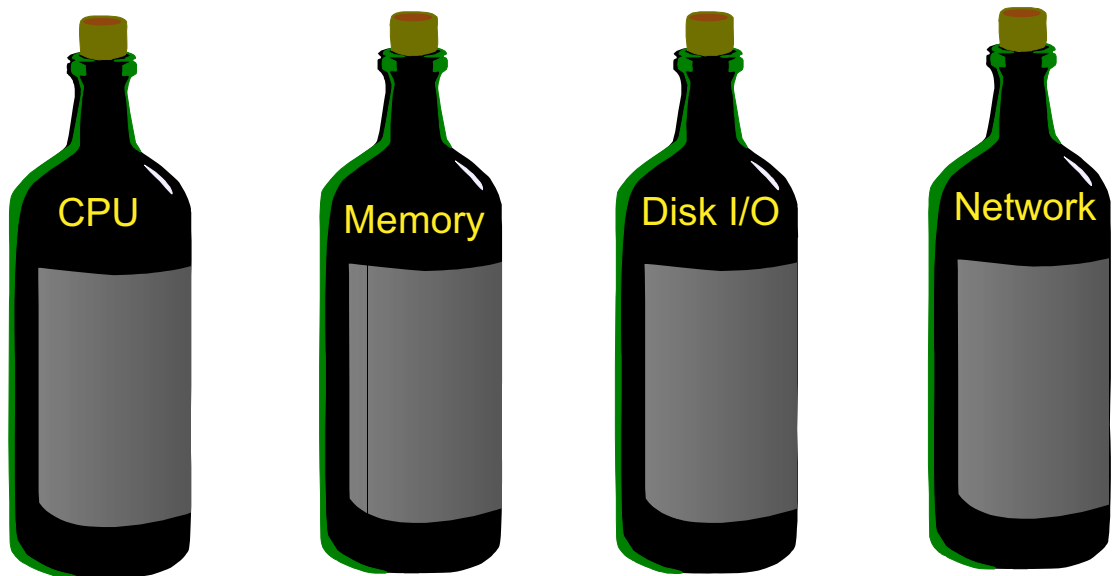
**Workstation** A single user works on a system, submitting work through the keyboard and receiving results on the native display of the system. The highest-priority

	performance objective of such a workload is minimum response time to the user's request.
<b>Multiuser</b>	A number of users submit their work through individual terminals that are connected to one system. The performance objective of such a workload is to maximize system throughput while preserving a specified worst-case response time.
<b>Server</b>	A workload that consists of requests from other systems, for example a file-server workload. The performance objective of such a system is maximum throughput within a given response time.

With multiuser or server workloads, the performance specialist must quantify both the typical and peak request rates.

When you have a clear understanding of the workload requests, analyze and document the physical hardware (what kind of model, how much memory, what kind of disks, what network is used).

# Critical Resources: The Four Bottlenecks



- Number of processes
- Process-Priorities
- Real memory
- Paging
- Memory leaks
- Disk balancing
- Types of disks
- LVM policies
- NFS used to load applications
- Network type
- Network traffic

© Copyright IBM Corporation 2004

Figure 11-4. Critical Resource: The Four Bottlenecks

AU1612.0

## Notes:

The performance of a given workload is determined by the availability and speed of different system resources. These resources that most often affect performance are:

- **CPU** (Central Processing Unit):

Is the CPU able to handle all the processes or is the CPU overloaded? Are there any processes that run with a very high priority that manipulates the system performance in general? Is it possible to run certain processes with a lower priority?

- **Memory:**

Is the real memory sufficient or is there a high paging rate? Are there faulty applications with memory leaks?

- **Disk I/O:**

Is the CPU often waiting for disk I/O? Are the disks in good balance? How good is the disk performance? Can I change LVM policies, to improve the performance (for example, to use striping)?

- **Network:**

How much is NFS used on the system? What kind of networks are used? How much network traffic takes place? Any faulty network cards?

Note that we cannot cover any network-related performance issues in this course. This goes beyond the scope of the class.

Now that we have identified the critical resources, we'll show how to measure the utilization of these resources.

# Identify CPU-Intensive Programs: ps aux

```
# ps aux
USER      PID    %CPU    %MEM    ...    STIME      TIME      COMMAND
root      516    98.2    0.0    ...    13:00:00   1329:38   wait
johnp    7570    1.2    1.0    ...    17:48:32    0:01    -ksh
root     1032    0.8    0.0    ...    15:13:47    78:37    kproc
root      1      0.1    1.0    ...    15:13:50    13:59    /etc/init
```

Percentage of time the process has used the CPU

Percentage of real memory

Total Execution Time

© Copyright IBM Corporation 2004

Figure 11-5. Identify CPU-Intensive Programs: ps aux

AU1612.0

## Notes:

For many performance-related problems a simple check with **ps** may reveal the reason. Execute **ps aux** to identify the CPU and memory usage of your processes. Concentrate on the following two columns:

- **%CPU**: This column indicates the percentage of time the process has used the CPU since the process started. The value is computed by dividing the time the process uses the CPU by the elapsed time of the process. In a multiprocessor environment, the value is further divided by the number of available CPUs.
- **%MEM**: The percentage of real memory used by this process.

By running **ps aux** identify your top applications related to CPU and memory usage.

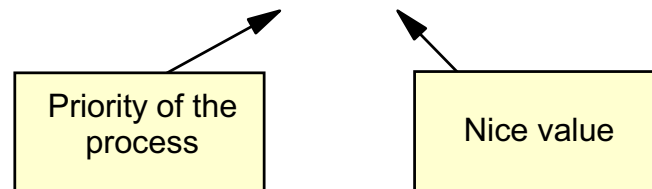
Many administrators use the **ps aux** command to create an alias definition that sorts the output according to the CPU usage:

```
alias top="ps aux | tail +2 | sort -k 1.15,1.19nr"
```

**In the visual a process with PID 516** is shown. That's the **wait** process that is assigned to the CPU, if the system is idle. With AIX, the CPU must always be doing work. If the system is idle, the **wait** process will be executed.

## Identify High-Priority Processes: ps -elf

```
# ps -elf
  F   S  UID  PID  PPID  C  PRI   NI   ...  TIME  CMD
200003  A    0    1    0    0  60   20   ... 13.59  init
240001  A    0 3860    1    0  60   20   ...  6:06  syncd
200001  A   299 7852 7570 24  72   20   ...  0:00  ps
```



- The smaller the PRI value, the higher the priority of the process. The average process runs a priority around 60.
- The NI value is used to adjust the process priority. The higher the nice value is, the lower the priority of the process.

© Copyright IBM Corporation 2004

Figure 11-6. Identify High-Priority Processes: ps -elf

AU1612.0

### Notes:

After identifying CPU and memory-intensive processes, check the priorities of your processes.

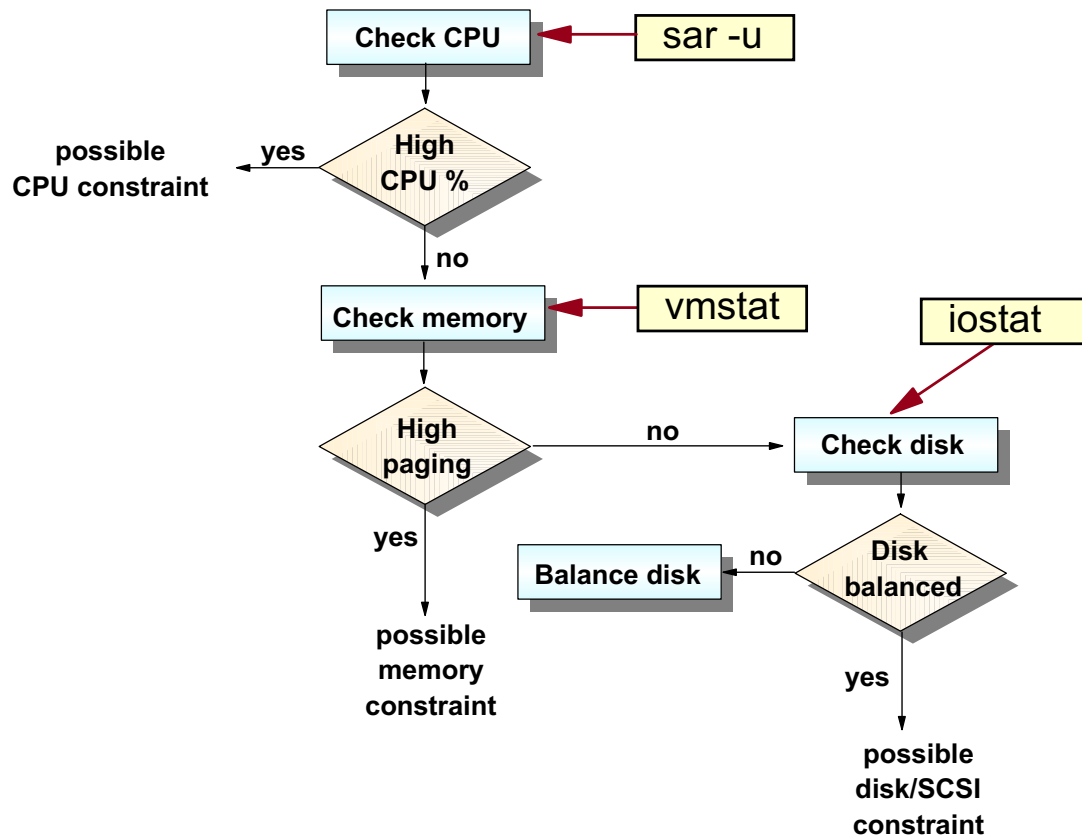
The priority of a process controls when a process will be executed.

AIX distinguishes **fixed** and **non-fixed** priorities. If a process uses a **fixed** priority, this priority will be unchanged throughout the whole lifetime of the process. Default priorities are **non-fixed**, that means after a certain timeslice, the priority will be recalculated. The new priority is determined by the amount of CPU time used and the **nice** value.

The nice value is shown in column **NI**. The default nice value is **20**. The higher the nice value is, the lower the priority of the process. We will learn later how to change the nice value.

The actual priority of the process is shown in the **PRI** column. The smaller this value, the higher the priority. Note that processes generally run with a PRI in the 60s. Keep an eye on processes that use a higher priority than this value.

# Basic Performance Analysis



© Copyright IBM Corporation 2004

Figure 11-7. Basic Performance Analysis

AU1612.0

## Notes:

There is a basic methodology that can make it easier to identify performance problems. The steps are as follows:

Look at the big picture. Is the problem CPU, I/O, or memory related?

- If you have a high CPU utilization, this could mean that there is a CPU bottleneck.
- If it's I/O-related, then is it paging or normal disk I/O?
- If it's paging, then increasing memory might help. You may also want to try to isolate the program and/or user causing the problem.
- If it's disk, then is disk activity balanced?
- If not, perhaps logical volumes should be reorganized to make more efficient use of the subsystem. Tools are available to determine which logical volumes to move.
- If balanced, then there may be too many physical volumes on a bus. More than three or four on a single SCSI bus may create problems. You may need to install another SCSI adapter. Otherwise, more disks may be needed to spread out the data.



## Monitoring CPU Usage: sar -u

```
# sar -u 60 30
AIX www 1 5 000400B24C00 06/06/01
08:24:10 %usr %sys %wio %idle
08:25:10 48 52 0 0
08:26:10 63 37 0 0
08:27:10 59 41 0 0
.
.
Average 57 43 0 0
```

A system is CPU bound, if:  
 $\%usr + \%sys > 80\%$

© Copyright IBM Corporation 2004

Figure 11-8. Monitoring CPU Usage: sar -u

AU1612.0

### Notes:

The **sar** command collects and reports system activity information.

The **sar** parameters on the visual indicate:

- **-u** collect CPU usage data
- **60** interval in seconds
- **30** number of intervals

The columns provide the following information:

- **%usr:**

Reports the percentage of time the CPU spent in execution at the user (or application) level

- **%sys:**

Reports the percentage of time the CPU spent in execution at the system (or kernel) level. This is the time the CPU spent in execution of system functions.

- **%wio:**

Reports the percentage of time the CPU was idle waiting for disk I/O to complete. This does not include waiting for remote disk access.

- **%idle:**

Reports the percentage of time the CPU was idle with no outstanding disk I/O requests.

The CPU usage report from **sar** is a good place to begin narrowing down whether a bottleneck is a CPU problem or an I/O problem. If the **%idle** time is high, it is likely there is no problem in either.

If the sum from **%usr** and **%sys** is always greater than 80%, it indicates that the CPU is approaching its limits. In other words, your system is **CPU bound**.

If you detect that your CPU always has outstanding disk I/Os, you must further investigate in this area. The system could be **I/O bound**.

Those with LPAR based systems should be aware that in AIX 5.3 there can be additional information in the output of all of the performance commands. If the POWER5 LPAR has Shared CPU resource allocated, the **sar** command output could look something like the following:

```
# sar -u 2 10
```

```
AIX console59 3 5 00C0288E4C00 11/19/04
```

```
System configuration: lcpu=2 ent=0.40
```

11:13:03	%usr	%sys	%wio	%idle	physc	%entc
11:13:05	0	1	0	99	0.01	1.4
11:13:07	0	0	0	100	0.00	0.8
11:13:09	0	0	0	100	0.00	0.8
11:13:11	0	0	0	100	0.00	0.8
11:13:13	0	0	0	100	0.00	0.8
11:13:15	0	0	0	100	0.00	0.8
11:13:17	0	0	0	100	0.00	0.8
11:13:19	0	0	0	100	0.00	0.8
11:13:21	0	0	0	100	0.00	0.8
11:13:23	0	0	0	100	0.00	0.8
Average	0	0	0	100	0.00	0.9

in the “System configuration: lcpu=2 ent=0.40” line, the “lcpu” means logical cpus and the “ent” means the LPAR’s entitled capacity.

Notice the “phyc” and “entc” columns. “phyc” reports the number of physical processors consumed. This will be reported only if the partition is running with shared processors or simultaneous multi-threading enabled. “entc” reports the percentage of entitled capacity consumed.

## Simultaneous Multi-Threading (SMT)

---

- Each chip appears as a two-way SMP to software
  - Appear as 2 logical CPUs
  - Performance tools may show number of logical CPUs
- Processor resources optimized for enhanced SMT performance
  - May result in a 25-40% boost and even more.
- Benefits vary - based on workload
- To enable:  
`smtctl [ -m off | on [ -w boot | now]]`

© Copyright IBM Corporation 2004

Figure 11-9. Simultaneous Multi-Threading (SMT)

AU1612.0

### **Notes:**

Modern processors have multiple specialized execution units, each of which is capable of handling a small subset of the instruction set architecture – some will handle integer operations, some floating point, and so on. These execution units are capable of operating in parallel and so several instructions of a program may be executing simultaneously.

However, conventional processors execute instructions from a single instruction stream. Despite microarchitectural advances, execution unit utilization remains low in today's microprocessors. It is not unusual to see average execution unit utilization rates of approximately 25% across a broad spectrum of environments. To increase execution unit utilization, designers use thread-level parallelism, in which the physical processor core executes instructions from more than one instruction stream. To the operating system, the physical processor core appears as if it is a symmetric multiprocessor containing two logical processors.

AIX 5.3 introduces Simultaneous multi-threading (SMT) to handle multiple threads of a POWER processor. If SMT is enabled, the POWER5 uses two separate instruction fetch address registers to store the program counters for the two threads. This implementation

provides the ability to schedule instructions for execution from all threads concurrently. With SMT, the system dynamically adjusts to the environment, allowing instructions to execute from each thread if possible, and allowing instructions from one thread to utilize all the execution units if the other thread encounters a long latency event. The performance benefit of simultaneous multi-threading is workload dependent. Most measurements of commercial workloads have received a 25-40% boost and a few have been even greater. Any workload where the majority of individual software threads highly utilize any resource in the processor or memory will benefit little from simultaneous multi-threading. For example, workloads that are heavily floating-point intensive are likely to gain little from simultaneous multi-threading and are the ones most likely to lose performance.

To enable and disable use `smtctl`.

```
smtctl [ -m off | on [ -w boot | now]]
```

-m off This option will set simultaneous multi-threading mode to disabled.

-m on This option will set simultaneous multi-threading mode to enabled.

-w boot This option makes the simultaneous multi-threading mode change effective on next and subsequent reboots.

-w now This option makes the simultaneous multi-threading mode change immediately but will not persist across reboot.

If neither the -w boot or the -w now options are specified, the mode change is made immediately and will persist across subsequent boots.

# Monitoring Memory Usage: vmstat

Summary report every 5 seconds

```
# vmstat 5
```

kthr		memory			page			...			cpu			
r	b	avm	fre	re	pi	po	fr	sr	cy	...	us	sy	id	wa
0	0	8793	81	0	0	0	1	7	0		1	2	95	2
0	0	9192	66	0	0	16	81	167	0		1	6	77	16
0	0	9693	69	0	0	53	95	216	0		1	4	63	33
0	0	10194	64	0	21	0	0	0	0		20	5	42	33
0	0	4794	5821	0	24	0	0	0	0		5	8	41	46

pi, po: Paging space page ins and outs:  
 • If any paging-space I/O is taking place, the workload is approaching the system's memory limit

wa: I/O wait percentage of CPU  
 • If nonzero, a significant amount of time is being spent waiting on file I/O

© Copyright IBM Corporation 2004

Figure 11-10. Monitoring Memory Usage: vmstat

AU1612.0

## Notes:

The **vmstat** command reports virtual memory statistics. It reports statistics about kernel threads, virtual memory, disks, traps and CPU activity.

In our example, we execute **vmstat 5**, that means every 5 seconds a new report will be written until the command is stopped. Note the first report is always the statistic since system startup.

Because our target in this course is to provide a basic performance understanding, we concentrate on the following columns.

- **pi/po**: These columns indicate the number of 4 KB pages that have been paged in or out.

Simply speaking, paging means that the real memory is not large enough to satisfy all memory requests and uses a secondary storage area on disks. If the systems workload always causes paging, you should consider to increasing real memory. Accessing pages on disk is relatively slow.

- **wa**: This column refers to the %wio column of **sar -u**. It indicates that the CPU has to wait for outstanding disk I/Os to complete. If this value is always non-zero, it might indicate that your system is I/O bound.

# Monitoring Disk I/O: iostat

## # iostat 10 2

```
tty:  tin      tout  avg-cpu:  %user  %sys   %idle  %iowait
      0.0      4.3          0.2   0.6   98.8   0.4
```

```
Disks:  %tm_act      Kbps   tps    Kb_read  Kb_wrtn
hdisk0      0.0      0.2    0.0     7993     4408
hdisk1      0.0      0.0    0.0        0        0
cd0         0.0      0.0    0.0        0        0
```

cumulative activity  
since last reboot

```
tty:  tin      tout  avg-cpu:  %user  %sys   %idle  %iowait
      0.1     110.7          7.0  59.4   0.0  33.7
```

```
Disks:  %tm_act      Kbps      tps    Kb_read  Kb_wrtn
hdisk0   77.9    115.7    28.7     456      8
hdisk1   0.0      0.0     0.0        0      0
cd0      0.0      0.0     0.0        0      0
```

A system is I/O bound, if:  
%iowait > 25%, %tm\_act > 70%

© Copyright IBM Corporation 2004

Figure 11-11. Monitoring Disk I/O: iostat

AU1612.0

## Notes:

The **iostat** command reports statistics for tty devices, disks and CD-ROMs.

**iostat** output:

### tty =

Are the number of characters read from (tin) and sent to (tout) terminals.

### avg-cpu =

Gives the same as **sar -u** and **vmstat** outputs (CPU utilization).

### Disk =

Typically shows the most useful information. This gives I/O statistics for each disk and CD-ROM on the system. **%tm\_act** is the percent of time the device was active over the period. **Kbps** is the amount of data, in kilobytes, transferred (read and written) per second. **tps** is the number of transfers per second. **Kb\_read** and **Kb\_wrtn** are the numbers of kilobytes read and written in the interval.



This information is useful for determining if the disk load is **balanced correctly**. In the above example, for that particular interval, one disk is used nearly 80% of the time where the other is not used at all. If this continues, some disk reorganization should take place.

The %iowait refers to %wio shown when using **sar -u**. If your system always shows waiting for outstanding disk requests, you need to investigate in this particular area.

With **iostat**, like **vmstat**, the first report is since system startup.

# topas

```

Topas Monitor for host:   kca81
Wed Jun  6 14:01:20 2001 Interval:  2

CPU Info → Kernel    0.2  |
              User    0.2  |
              Wait   0.0  |
              Idle   99.5  |#####|

Network      KBPS      I-Pack  O-Pack  KB-In  KB-Out
en0          0.1       0.4     0.4     0.0    0.1
lo0          0.0       0.0     0.0     0.0    0.0

iostat Info → Disk      Busy%      KBPS      TPS  KB-Read  KB-Writ
hdisk0      0.0       0.0     0.0     0.0     0.0
hdisk1      0.0       0.0     0.0     0.0     0.0

Name        PID  CPU%  PgSp  Owner
topas       221512  0.5  0.9  root
syncd       98360  0.0  0.3  root
shdaemon    655468  0.0  32.5  root
dtterm      459818  0.0  1.4  root
dtexec      598126  0.0  0.7  root
dtscreen    330877  0.0  0.6  root
cscope      188008  0.0  0.5  root
gil         57358  0.0  0.1  root
dtfile      409804  0.0  2.3  root
init         1      0.0  0.8  root
dtterm      173427  0.0  1.4  root
ksh         103055  0.0  0.7  root
telnetd     211931  0.0  0.7  root

EVENTS/QUEUES      FILE/TTY
Cswitch            32  Readch           25
Syscall           147  Writech          146
Reads              2  Rawin            0
Writes             2  Ttyout           0
Forks              0  Igets            0
Execs              0  Namei            1
Runqueue          0.0  Dirblk           0
Waitqueue         0.0

PAGING             MEMORY
Faults             0  Real,MB          1023
Steals             0  % Comp           28.0
PgspIn            0  % Noncomp        3.3
PgspOut           0  % Client         3.2
PageIn            0
PageOut           0  PAGING SPACE
Sios              0  Size,MB          512
                  % Used           12.0
                  % Free            87.9

NFS (calls/sec)
ServerV2           0
ClientV2           0  Press:
ServerV3           0  "h" for help
ClientV3           0  "q" to quit
    
```

VMSTAT Info

© Copyright IBM Corporation 2004

Figure 11-12. topas

AU1612.0

## Notes:

In 4.3.3, a new command was added that pulls together pieces of the performance commands and presents them on one screen. This command is **topas**.

**topas** continuously updates the screen to show the current state of the system. In the upper left is the same information that is given with **sar**. The middle of the left side shows the same information as **iostat**. The right lower quadrant show information from the virtual memory manager which can be seen with **vmstat**.

To exit from **topas**, just press “q” for quit. “h” is also available for help.

The **topas** command is only available on the POWER platform.

---

## topas, vmstat, and iostat Enhancements for Micro-Partitioning (AIX 5.3)

---

- Added two new values to the default screen
  - Physc and %Entc
- The vmstat command has two new metrics:
  - pc and ec
- The iostat command has also two new metrics:
  - %physc and %entc

© Copyright IBM Corporation 2003

Figure 11-13. topas, vmstat, and iostat Enhancements for Micro-Partitioning (AIX 5.3)

AU1612.0

### **Notes:**

#### topas Enhancements

If topas runs on a partition with a shared processor partition beneath the CPU utilization, there are two new values displayed:

**Physc:** Number of physical processors granted to the partition (if micro-partitioning)

**%Entc:** Percentage of Entitled Capacity granted to a partition (if micro-partitioning)

The -L flag will switch the output to a logical partition display. You can either use -L when invoking the topas command, or as a toggle when running topas. In this mode, topas displays data similar to mpstat and lparstat commands.

#### vmstat enhancements

The vmstat command has been enhanced to support Micro-Partitioning and can now detect and tolerate dynamic configuration changes.

The `vmstat` command has two new metrics that are displayed. These are Physical Processor Granted and Percentage of Entitlement Granted, which are represented as **pc** and **ec** in the output format. The Physical Processor Granted represents the number of physical processors granted to the partition during an interval. The Percentage of Entitlement Granted is the percentage of Entitled Capacity granted to a partition during an interval. These new metrics will be displayed only when the partition is running as a shared processor partition or with SMT enabled. If the partition is running as a dedicated processor partition and with SMT off, the new metrics will not be displayed.

#### iostat enhancements

Beginning with AIX 5L V5.3, the `iostat` command reports the percentage of physical processors consumed (`%physc`), the percentage of entitled capacity consumed (`%entc`), and the processing capacity entitlement when running in a shared processor partition. These metrics will only be displayed on shared processor partitions.

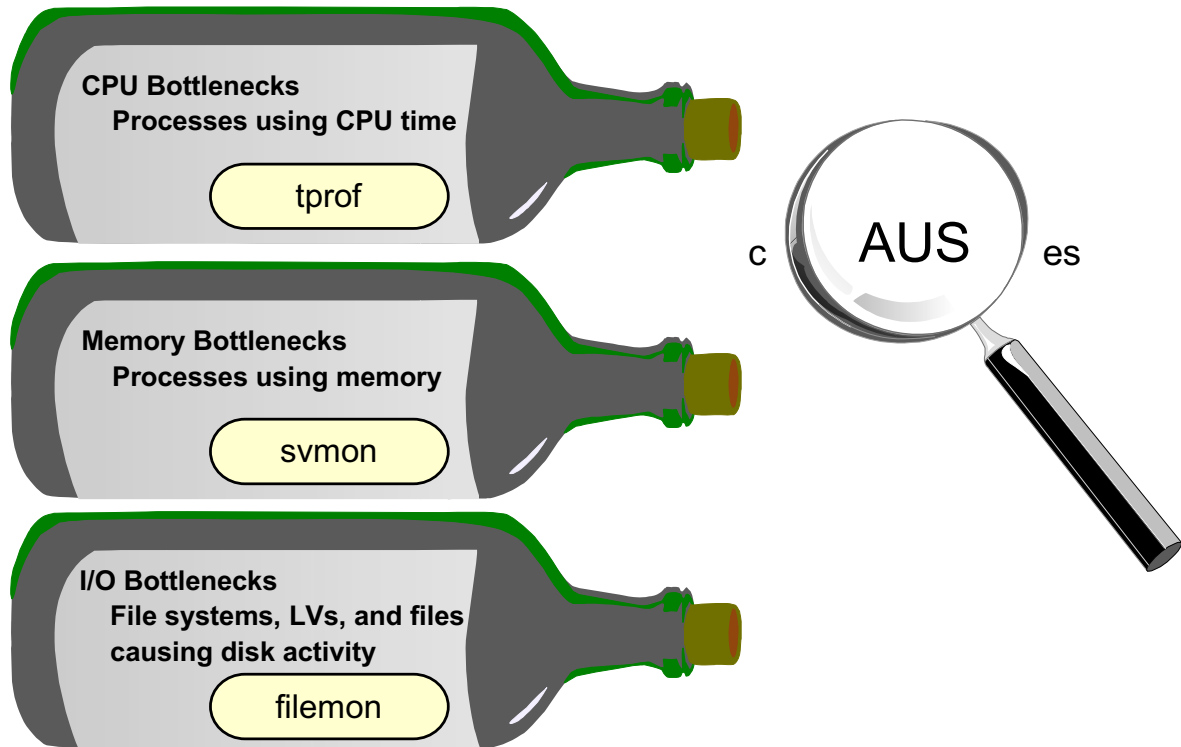
Percentage of entitled capacity consumed is display in the **%entc** column.

Percentage of physical processors consumed is shown in the **%physc** column.

In the system configuration information, you can see the currently assigned processing capacity specified as `ent`.

# AIX Performance Tools

Identify causes of bottlenecks:



© Copyright IBM Corporation 2004

Figure 11-14. AIX Performance Tools

AU1612.0

## Notes:

There are three additional tools that are available in AIX to further determine the cause of the performance bottleneck. **sar**, **vmstat**, and **iostat** are all generic UNIX tools and are good for identifying whether the bottleneck is CPU, memory or disk.

As you try to solve the problem, you need to identify individual applications and processes that put the heaviest workload on the CPU and use the most memory. Also, to solve disk I/O problems, you need to know what file system, logical volumes and file are accessed the most.

This is where **tprof**, **svmon**, and **filemon** are helpful.

The next few graphics are intended as an introduction to these tools. They are extensive in the number of options and the information they can produce. As you learn more about performance and tuning, you should further investigate the capabilities of these tools.

# AIX Tools: tprof

```
# tprof -x sleep 60
# more __prof.all
```

This file is created by tprof

<u>Process</u>	<u>PID</u>	<u>TID</u>	<u>Total</u>	<u>Kernel</u>	<u>User</u>	<u>Shared</u>	<u>Other</u>
wait	516	517	6855	6855	0	0	0
netscape_aix4	23494	40015	201	27	29	145	0
lslpp	17566	43613	11	5	4	2	0

<u>Process</u>	<u>FREQ</u>	<u>Total</u>	<u>Kernel</u>	<u>User</u>	<u>Shared</u>	<u>Other</u>
wait	1	6855	6855	0	0	0
netscape_aix4	5	961	122	139	700	0
ksh	46	77	64	7	6	0

© Copyright IBM Corporation 2004

Figure 11-15. AIX Tools: tprof

AU1612.0

## Notes:

If you have determined that your system is CPU-bound, how do you know what process or processes are using the CPU the most? **tprof** is used to spot those processes.

**tprof** is a trace tool - meaning it monitors the system for a period of time and when it stops, it produces a report. The command **tprof -x sleep 60** analyzes all processes on the system for 60 seconds. It will generate a summary file call **\_\_prof.all** (that is two underscores then prof.all). All files that tprof creates will start with two underscores. By looking at this file, you can see the CPU demand by process in decreasing order.

Our sample output has been reduced to simplify the areas to focus on.

In our sample output, the first section indicates that the process **netscape\_aix4** (pid 23494) used a total of 201 CPU ticks. There are 100 ticks in a second. Therefore, our program used 2.01 seconds of the CPU.

In the second section, you can see there were 5 (FREQ) netscape\_aix4 processes in total running on this system. They took a total of 961 ticks (or 9.61 seconds) of the CPU. This cumulative number can be helpful because one individual process may not be consuming a

significant amount of CPU resources, but together, those similar processes may significantly contribute to the heavy load on the system.

# AIX Tools: svmon

**Global report**

```
# svmon -G
```

	size	inuse	free	pin	virtual
memory	32744	20478	12266	2760	11841
pg space	65536	294			
	work	pers	clnt		
pin	2768	0	0		
in use	13724	6754	0		

Sizes are in # of 4K frames

**Top 3 users of memory**

```
# svmon -Pt 3
```

Pid	Command	Inuse	Pin	Pgsp	Virtual	64-bit	Mthrd
14624	java	6739	1147	425	4288	N	Y
9292	httpd	6307	1154	205	3585	N	Y
3596	X	6035	1147	1069	4252	N	N

\* output has been modified

© Copyright IBM Corporation 2004

Figure 11-16. AIX Tools: svmon

AU1612.0

## Notes:

**svmon** is used to capture and analyze information about virtual memory. This is a very extensive command that can produce a variety of statistics - most of which is beyond our scope for this course.

In both examples, the output has been reduced for simplicity and to show the information that is of interest to this discussion.

In the first example, **svmon -G** provides a global report. You can see the size of memory, how much is in use and the amount that is free. It provides details about how it is being used and it also provide statistics on paging space.

All numbers are reported as the number of frames. A frame is 4 KB in size.

In the second example, **svmon -Pt 3** displays memory usage of the top 3 memory-using processes sorted in decreasing order of memory demand.

**P** - shows processes

**t** - top # to display



# AIX Tools: filemon

# filemon -o fmout ← Starts monitoring disk activity

# trcstop  
# more fmout ← Stops monitoring and creates report

Most Active Logical Volumes					
util	#rblk	#wblk	KB/s	volume	description
0.03	3368	888	26.5	/dev/hd2	/usr
0.02	0	1584	9.9	/dev/hd8	jfslog
0.02	56	928	6.1	/dev/hd4	/
Most Active Physical Volumes					
util	#rblk	#wblk	KB/s	volume	description
0.10	24611	12506	231.4	/dev/hdisk0	N/A
0.02	56	8418	52.8	/dev/hdisk1	N/A

© Copyright IBM Corporation 2004

Figure 11-17. AIX Tools: filemon

AU1612.0

## Notes:

If you have determined your system is I/O bound, you now need to determine how to resolve the problem. You need to identify what is causing your disk activity if you would like to spread the workload among your disks. **filemon** is the tool that can provide that information.

**filemon** is a trace tool. Use the **filemon** command to start the trace. You need to use **trcstop** to stop the trace and generate the report.

In our example, **filemon -o fmout** starts the trace. The **-o** directs the output to the file called fmout. There will be several sections included in this report. The sample output has been reduced to only show two areas: logical volume activity and physical volume activity.

Here is a description of the columns:

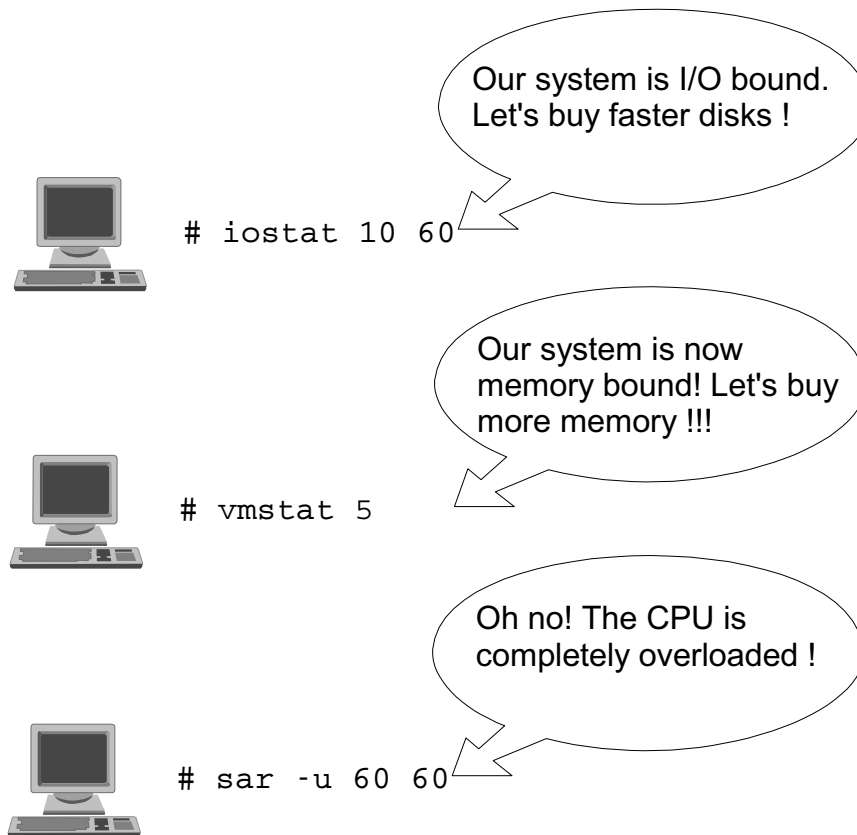
<b>util</b>	utilization over the measured interval (0.03 = 3%)
<b>#rblk</b>	number of 512-byte blocks read
<b>#wblk</b>	number of 512-byte blocks written

<b>KB/s</b>	average data transfer rate
<b>volume</b>	the logical or physical volume name
<b>description</b>	file system name or logical volume type

Since they are ranked by usage, it is very easy to spot the file systems, LV's and disks that are most heavily used.

To break it down even further, you can use **filemon** to see activity of individual files: **filemon -O all -o fmout**

# There Is Always a Next Bottleneck!



© Copyright IBM Corporation 2004

Figure 11-18. There Is Always a Next Bottleneck!

AU1612.0

## Notes:

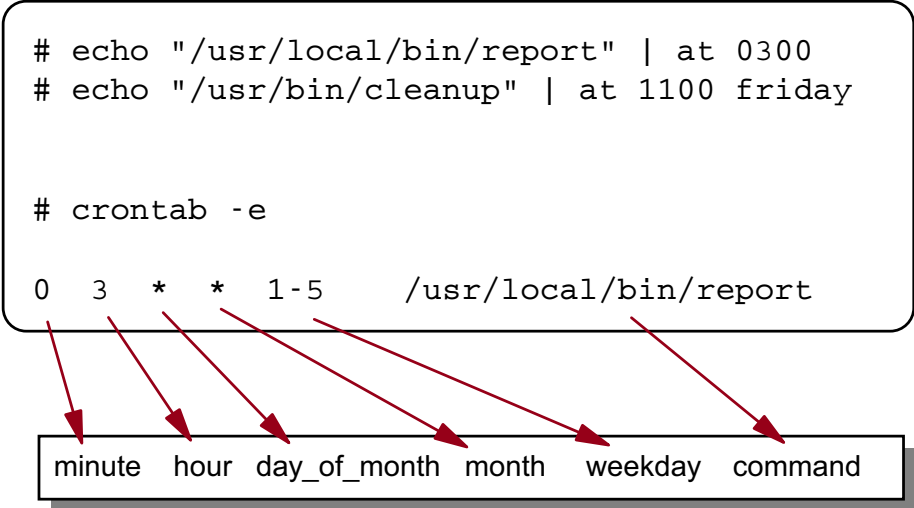
The visual shows a performance truism, “there is always a next bottleneck”. It means that eliminating one bottleneck might lead to another performance bottleneck. For example, eliminating a disk bottleneck might lead to a memory bottleneck. Eliminating the memory bottleneck might lead to a CPU bottleneck.

When you have exhausted all system tuning possibilities and performance is still unsatisfactory, you have one final choice: **Adapt workload-management techniques**

These techniques are provided on the next pages.

# Workload Management Techniques (1 of 3)

Run programs at a specific time



© Copyright IBM Corporation 2004

Figure 11-19. Workload Management Techniques (1 of 3)

AU1612.0

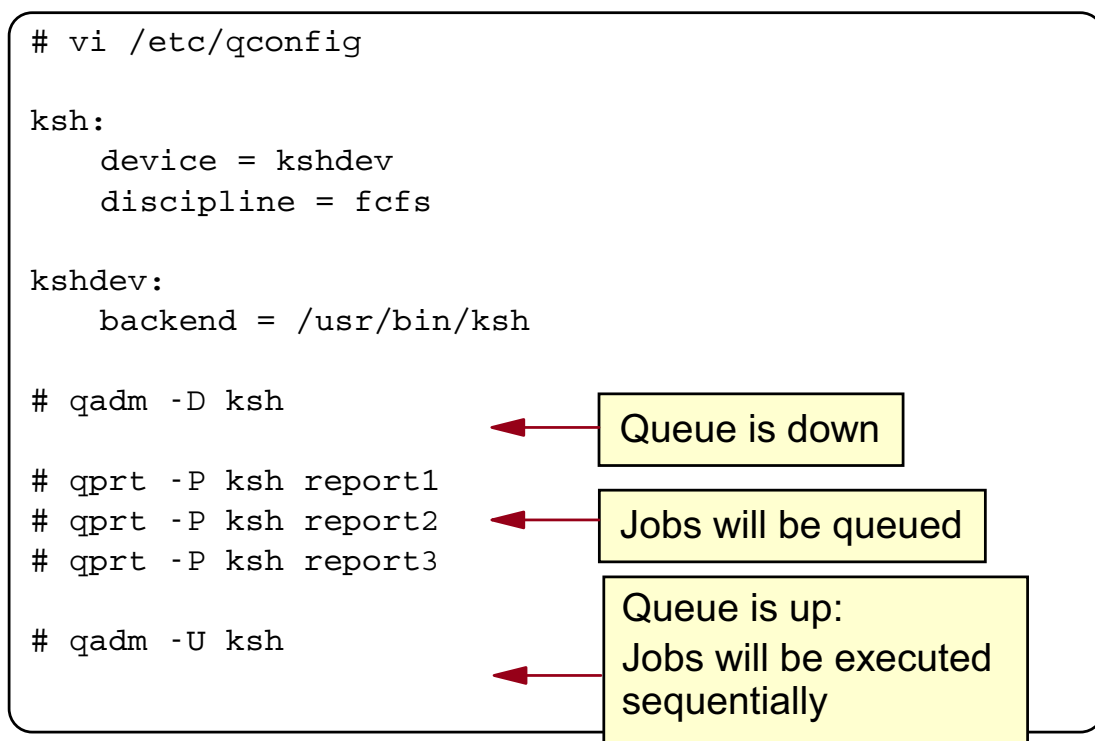
### Notes:

Workload management simply means assessing the components of the workload to determine whether they are all needed as soon as possible. Usually, there is work that can wait for a while. A report that needs to be created for the next morning, could be started at 4 p.m. or at 4 a.m. The difference is that at night the CPU is probably idle.

The cron daemon can be used to spread out the workload by running at different times. To take advantage of the capability, use the **at** command or set up a **crontab** file.

## Workload Management Techniques (2 of 3)

### Sequential execution of programs



© Copyright IBM Corporation 2004

Figure 11-20. Workload Management Techniques (2 of 3)

AU1612.0

### Notes:

Another workload management technique is to put programs or procedures in a **job queue**. In the example we define a **ksh** queue, that uses the **/usr/bin/ksh** as backend (the backend is the program that is called by **qdaemon**).

In the example we bring the queue down:

```
# qadm -D ksh
```

During the day (or when the workload is very high), users put their jobs into this queue:

```
# qprt -P ksh report1
# qprt -P ksh report2
# qprt -P ksh report3
```

During the night (or when the workload is lower), you put the queue up, which leads to a sequential execution of all jobs in the queue:

```
# qadm -U ksh
```

## Workload Management Techniques (3 of 3)

### Run programs at a reduced priority

```
# nice -n 15 backup_all &
# ps -el
  F    S  UID  PID  PPID  C  PRI  NI   ...  TIME  CMD
240001  A    0 3860 2820 30   90   35   ...  0:01  backup_all
```

Very low  
priority

Nice value:  
20+15

```
# renice -n -10 3860
# ps -el
  F    S  UID  PID  PPID  C  PRI  NI   ...  TIME  CMD
240001  A    0 3860 2820 26   78   25   ...  0:02  backup_all
```

© Copyright IBM Corporation 2004

Figure 11-21. Workload Management Techniques (3 of 3)

AU1612.0

### Notes:

Some programs that run during the day can be run with a lower priority. They will take longer to complete, but they will be less in competition with really time-critical processes.

To run a program at a lower priority, use the **nice** command:

```
# nice -n 15 backup_all &
```

This command specifies that the program **backup\_all** runs at a very low priority. The default nice value is 20 (24 for a ksh background process), which is increased here to 35. The nice value can range from 0 to 39, with 39 being the lowest priority.

As **root** user you can use **nice** to start processes with a higher priority. In this case you would use a negative value:

```
# nice -n -15 backup_all &
```

Here the nice value is decreased to 5, which results in a very high priority of the process.

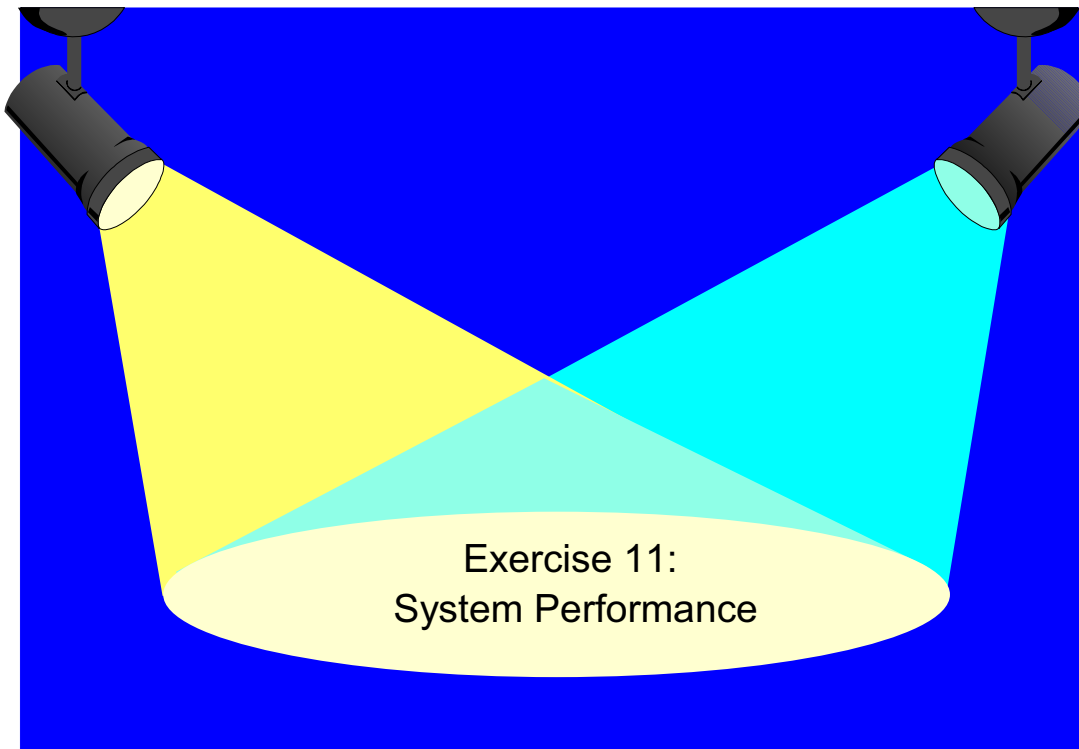
If the process is already running, you can use the **renice** command to reduce or increase the priority:

```
# renice -n -10 3860
```

In the example we decrease the nice value (from 35 to 25), which results in a higher priority. Note that you must specify the process ID when working with **renice**.

## Next Step

---



© Copyright IBM Corporation 2004

Figure 11-22. Next Step

AU1612.0

### **Notes:**

After the exercise you should be able to:

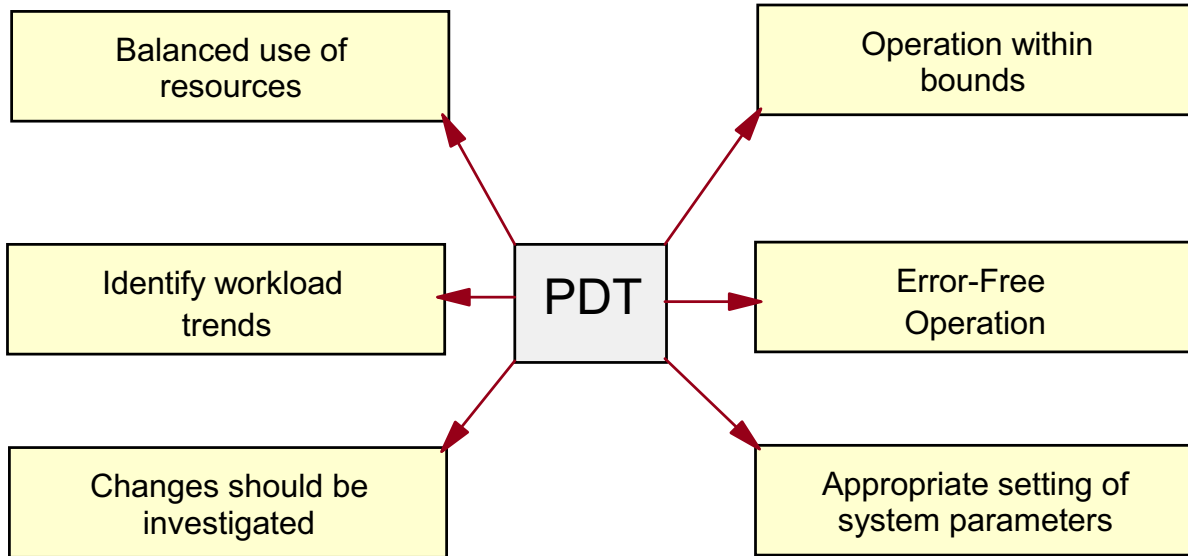
- Use **ps** to identify CPU and memory-intensive programs
- Execute a basic performance analysis
- Implement a korn shell job queue
- Work with **nice** and **renice**



## 11.2 Performance Diagnostic Tool (PDT)

# Performance Diagnostic Tool (PDT)

PDT assesses the current state of a system and tracks changes in workload and performance.



© Copyright IBM Corporation 2004

Figure 11-23. Performance Diagnostic Tool (PDT)

AU1612.0

## Notes:

PDT assesses the current state of a system and tracks changes in workload and performance. It attempts to identify incipient problems and suggest solutions before the problems become critical. PDT is available on all AIX 4 or later systems. It is contained in fileset **bos.perf.diag\_tool**.

PDT attempts to apply some general concepts of well-performing systems to its search for problems. These concepts are:

### 1. **Balanced use of resources:**

In general, if there are several resources of the same type, then a balanced use of those resources produces better performance.

- Comparable numbers of physical volumes on each adapter
- Paging space distributed across multiple physical volumes
- Roughly equal measured load on different physical volumes

**2. Operation within bounds:**

Resources have limits to their use. Trends that would attempt to exceed those limits are reported.

- File system sizes cannot exceed the allocated space
- A disk cannot be utilized more than 100% of the time

**3. Identify workload trends:**

Trends can indicate a change in the nature of the workload as well as increases in the amount of resource used:

- Number of users logged in
- Total number of processes
- CPU-idle percentage

**4. Error-free operation:**

Hardware or software errors often produce performance problems.

- Check the hardware and software error logs
- Report bad VMM pages (pages that have been allocated by applications but have not been freed properly)

**5. Changes should be investigated:**

New workloads or processes that start to consume resources may be the first sign of a problem.

- Appearance of new processes that consume lots of CPU or memory resources

**6. Appropriate setting of system parameters**

There are many parameters in the system, for example the maximum number of processes allowed per user (maxuproc). Are all of them set appropriately?

The PDT data collection and reporting is very easy to implement.

## Enabling PDT

### # /usr/sbin/perf/diag\_tool/pdt\_config

```
-----PDT customization menu-----
1) show current    PDT report recipient and severity level
2) modify/enablePDT reporting
3) disable        PDT reporting
4) modify/enable  PDT collection
5) disable        PDT collection
6) de-install     PDT
7) exit pdt_config

Please enter a number: 4
```

© Copyright IBM Corporation 2004

Figure 11-24. Enabling PDT

AU1612.0

### Notes:

From the PDT menu, option 4 enables the default data collection functions. Actual collection occurs via **cron** jobs run by the **cron** daemon.

The menu is created using the Korn Shell **select** command, and this means the menu options are not reprinted after each selection; however, the program will show the menu again if you press Enter without making a selection.

To alter the recipient of reports use option 2 - the default is the adm user. Reports have severity levels. There are three levels - 1 gives the smallest report, while level 3 will analyze the data in more depth.

Option 6 does not deinstall the program - it simply advises how you might do that.

Analysis by PDT is both static (configuration focused; that is, I/O and paging) and dynamic (over time). Dynamic analysis includes such areas as network, CPU, memory, file size, file system usage, and paging space usage. An additional part of the report evaluates load average, process states, and CPU idle time.

Once PDT is enabled, it maintains data in a historical record for (by default) 35 days. On a daily basis, by default, PDT generates a diagnostic report that is sent to user **adm** and also written to **/var/perf/tmp/PDT\_REPORT**.

## cron Control of PDT Components

```
# cat /var/spool/cron/crontabs/adm
```

```
0 9 * * 1-5 /usr/sbin/perf/diag_tool/Driver_daily
```

Collect system data, each workday at 9:00

```
0 10 * * 1-5 /usr/sbin/perf/diag_tool/Driver_daily2
```

Create a report, each workday at 10:00

```
0 21 * * 6 /usr/sbin/perf/diag_tool/Driver_offweekly
```

Cleanup old data, each saturday evening

© Copyright IBM Corporation 2004

Figure 11-25. cron Control of PDT Components

AU1612.0

### Notes:

The three main components of the PDT system are: collection control, retention control, and reporting control.

When PDT is enabled, by default, it adds entries to the **crontab** file for **adm** to run these functions at certain default times and frequencies. The entries execute a shell script called **Driver\_** in the **/usr/sbin/perf/diag\_tool** directory. This script is passed three different parameters, each representing a collection profile, at three different collection times.

```
# cat /var/spool/cron/crontabs/adm
0 9 * * 1-5 /usr/sbin/perf/diag_tool/Driver_ daily
0 10 * * 1-5 /usr/sbin/perf/diag_tool/Driver_ daily2
0 21 * * 6 /usr/sbin/perf/diag_tool/Driver_ offweekly
```

The **crontab** entries and the **Driver\_** script indicate that daily statistics (**daily**) are collected at 9:00 a.m. and reports (**daily2**) are generated at 10:00 a.m. every work day, and historical data (**offweekly**) is cleaned up every Saturday night at 9:00 p.m.

# PDT Files

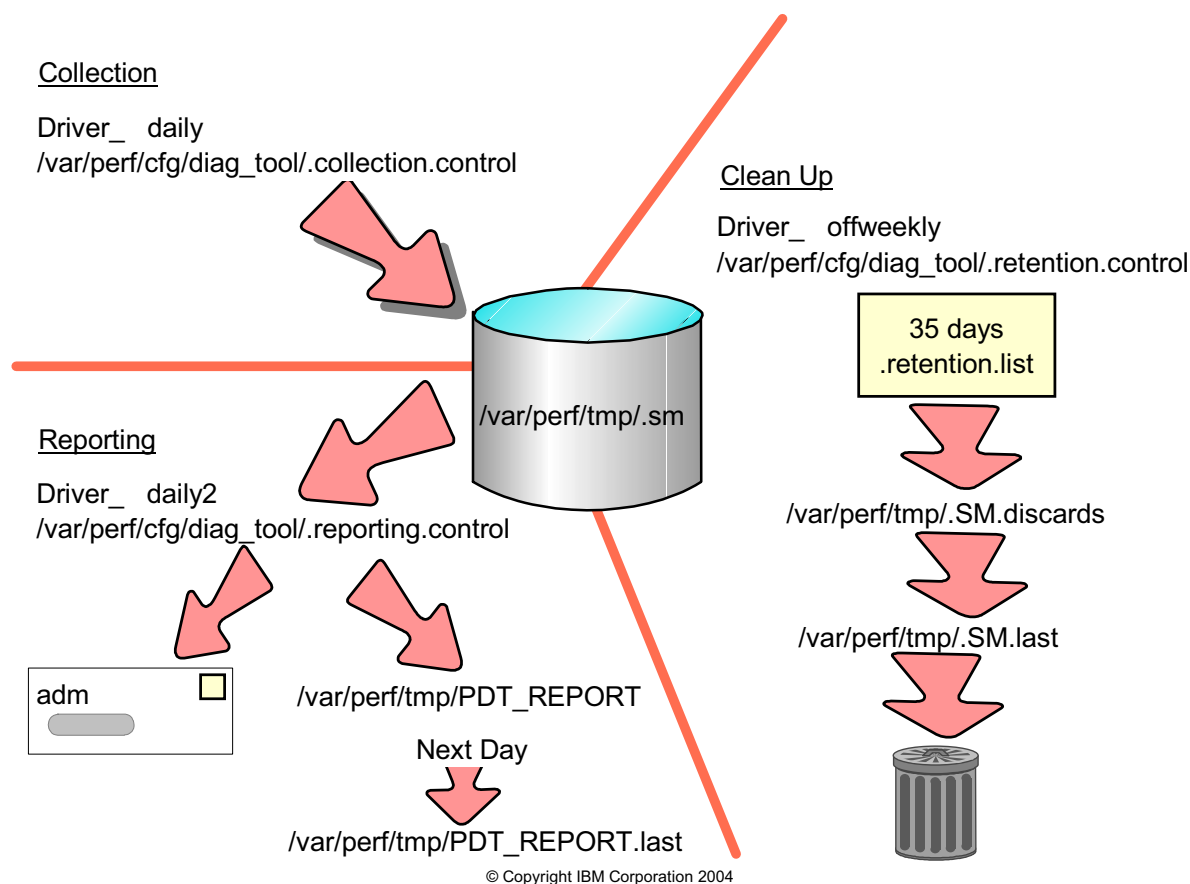


Figure 11-26. PDT Files

AU1612.0

## Notes:

The parameter passed to the **Driver\_** shell script is compared with the contents of the **.control** files found in the **/var/perf/cfg/diag\_tool** directory to find a match. These control files contain the names of scripts to run to collect data and generate reports. When a match is found, the corresponding scripts are run. The scripts that are executed for **daily** are in **.collection.control**, those for **daily2** are in **.reporting.control**, and **offweekly** are in **.retention.control**.

The collection component comprises a set of programs in **/usr/sbin/perf/diag\_tool** that periodically collect and record data on configuration, availability and performance.

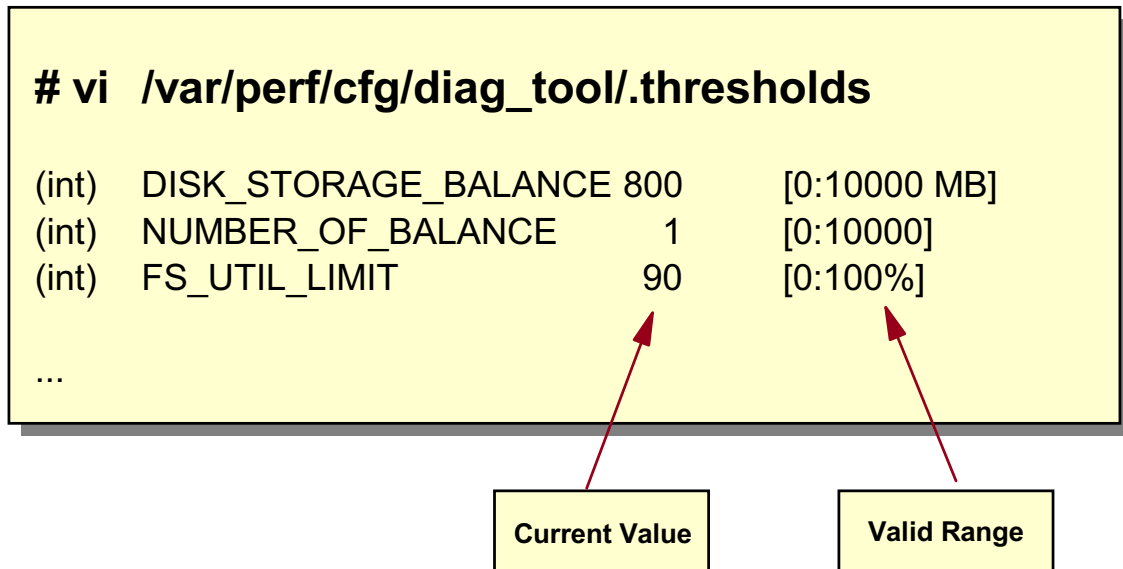
The retention component periodically reviews the collected data and discards data that is out of date. The size of the historical record is controlled by the file **/var/perf/cfg/diag\_tool/.retention.list** - it contains the default number "35" - this number of days may be changed easily. Data that is discarded during the cleanup, is appended to the file **/var/perf/tmp/.SM.discards** - and the cleansed data is kept in **/var/perf/tmp/.SM**, with one last backup held in **/var/perf/tmp/.SM.last**.

Finally, the reporting component periodically produces a diagnostic report from the current set of historical data - on a daily basis PDT generates a diagnostic report and mails the report (by default) to adm and writes it to **/var/perf/tmp/PDT\_REPORT**. The previous day's report is saved to **/var/perf/tmp/PDT\_REPORT.last**.

Any PDT execution errors will be appended to the file **/var/perf/tmp/.stderr**.



# Customizing PDT: Changing Thresholds



© Copyright IBM Corporation 2004

Figure 11-27. Customizing PDT: Changing Thresholds

AU1612.0

## Notes:

The `/var/perf/cfg/diag_tool/.thresholds` file contains the thresholds used in analysis and reporting. The visual shows thresholds that are related to disk balancing (`DISK_STORAGE_BALANCE`, `NUMBER_OF_BALANCE`) and file system utilization. The file may be modified by **root** or **adm**. Here is a complete listing of all thresholds:

### **DISK\_STORAGE\_BALANCE** (MB)

The SCSI controller having the most disk storage space attached to it is identified. The SCSI controller having the smallest disk storage is identified. If the difference (in MB) between these two amounts exceeds `DISK_STORAGE_BALANCE`, then a message will be displayed:

“SCSI Controller **scsiX** has **A.BMB** more storage than **scsiY**”

### **PAGING\_SPACE\_BALANCE**

Not presently used.

### **NUMBER\_OF\_BALANCE**

The SCSI controller having the largest number of disks attached is identified. The SCSI

controller having the least number of disks is identified. If the difference between these two counts exceeds NUMBER\_OF\_BALANCE, then we report:

“SCSI Controller **scsiX** has A more disks than **scsiY**”

The same sort of test is performed on the number of paging areas defined on each physical volume:

“Physical Volume **hdiskX** has A paging areas,  
while Physical Volume **hdiskY** has only B”

### MIN\_UTIL (%)

This threshold is applied to process utilizations. Changes in the top-3 CPU consumers are only reported if the new process had a utilization in excess of MIN\_UTIL:

“First appearance of **PID (process\_name)** on top-3 cpu list”

The same threshold applies to changes in the top-3 memory consumers list:

“First appearance of **PID (process\_name)** on top-3 memory list”

### FS\_UTIL\_LIMIT (%)

Applies to jfs file system utilizations. If a file system is found with percentage use in excess of FS\_UTIL\_LIMIT, then it is identified in the message:

“File system **device\_name (/mount\_point)** is nearly full at **X%**”

The same threshold is applied to paging spaces:

“Paging space **paging\_name** is nearly full at **Y%**”

### MEMORY\_FACTOR

This parameter is employed in the (crude) test to determine if the system has sufficient memory. Conceptually, the objective is to determine if the total amount of memory is adequately backed up by paging space. If real memory size is close to the amount of used paging space, then the system is likely to start paging, and would benefit from the addition of memory. The actual formula is based on experience, and actually compares: MEMORY\_FACTOR \* memory with the mean paging space (+/- 2 standard deviations).

“System has **X** MB memory; may be inadequate.”

The current default is “0.9”; by decreasing this number, the warning will be produced more frequently (and perhaps, unnecessarily). Increasing this number will eliminate the message altogether.

### TREND\_THRESHOLD

This is used in all trending assessments. It is applied after a linear regression is performed on all the available historical data. The slope of the fitted line (assuming a line of significance could be fit, and the regression passes a suite of residuals tests) must exceed (Last Value) \* TREND\_THRESHOLD:

“File system **device\_name (/mount\_point)** is growing,  
now, **X%** full, and growing an avg. of **Y%/day**”

This is purely a heuristic. The objective is to try to ensure that a trend, however strong its statistical significance, actually has some “real world” significance.

So, for example, if we determine that a file system is growing at **A** MB/day, and the last value for the file system size is 100 MB, we require that **A** exceed  $100 \text{ MB} * \text{TREND\_THRESHOLD}$  to be reported as a trend of “real world” significance. The default for TREND\_THRESHOLD is "0.01", so a growth rate of 1 MB per day would be required for reporting. The threshold can be set anywhere between “0.000001” and “100000”.

The assessment applies to trends associated with:

- CPU use by a top-3 process
- Memory use by a top-3 process
- The size of files indicated in the .files file
- FILE SYSTEMS (jfs)
- PAGE SPACES
- Hardware errors
- Software errors
- Workload indicators
- Processes per user
- Ping delay to nodes in the .hosts file
- Packet loss % to nodes in the .hosts file

#### **EVENT\_HORIZON** (Days)

This is used in trending assessments where we report expected time for a given trend to cause a key limit to be reached. For example, in the case of file systems, if we determine that there is a significant (both statistical and “real world”) trend, we estimate the time (at this rate) until the file system is 100% full. If this time is within EVENT\_HORIZON days, we report the estimated full date:

“At this rate, **device\_name** will be full in about **X** days”

This threshold applies to trends associated with:

- FILE SYSTEMS (jfs)
- PAGE SPACES

## Customizing PDT: Specific Monitors

```
# vi /var/perf/cfg/diag_tool/.files
```

```
/var/adm/wtmp
```

```
/var/spool/qdaemon/
```

```
/var/adm/ras/
```

```
/tmp/
```

**Files and directories  
to monitor**

```
# vi /var/perf/cfg/diag_tool/.nodes
```

```
pluto
```

```
neptun
```

```
mars
```

**Machines  
to monitor**

© Copyright IBM Corporation 2004

Figure 11-28. Customizing PDT: Specific Monitors

AU1612.0

### Notes:

By adding files and directories into the file **/var/perf/cfg/diag\_tool/.files** you can monitor the sizes of these files and directories. Here are some examples.

<code>/var/adm/wtmp</code>	is a file used for login recording
<code>/var/spool/qdaemon</code>	is a directory used for print spooler
<code>/var/adm/ras</code>	a directory used for AIX error logging

By adding hostnames to **/var/perf/cfg/diag\_tool/.nodes** you can monitor different systems. By default, no network monitoring takes place, as the **.nodes** file must be created.

# PDT Report Example (Part 1)

## Performance Diagnostic Facility 1.0

Report printed: Wed Jun 6 14:37:07 2001

Host name: master

Range of analysis included measurements from:

Hour 14 on Monday 4th June 2001

to: Hour 9 on Wednesday 6th June

### Alerts

#### I/O CONFIGURATION

- Note: volume hdisk2 has 480 MB available for allocation while volume hdisk1 has 0 MB available

#### PAGING CONFIGURATION

- Physical Volume hdisk1 (type:SCSI) has no paging space defined

#### I/O BALANCE

- Physical volume hdisk0 is significantly busier than others  
volume hdisk0, mean util. = 11.75  
volume hdisk1, mean util. = 0.00

#### NETWORK

- Host sys1 appears to be unreachable

© Copyright IBM Corporation 2004

Figure 11-29. PDT Report Example (Part 1)

AU1612.0

## Notes:

Note that this is a doctored report example. Some sections have been deliberately altered for enhanced dramatic effect; some small parts have been left out for simplicity.

The PDT report consists of several sections. The HEADER section provides information on the time and date of the report, the host name and the time period for which data was analyzed. The content of this section does not differ with changes in the severity level.

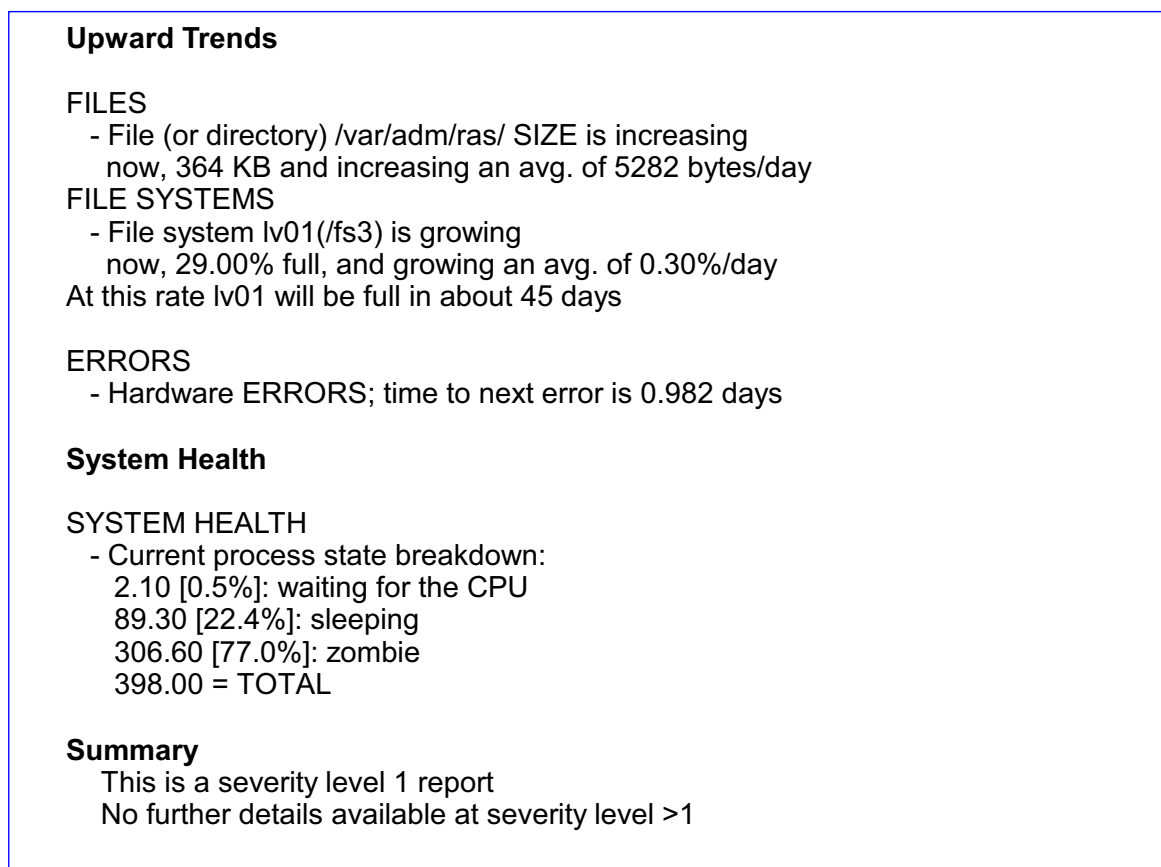
After a HEADER section, the ALERTS section reports on identified violations of concepts and thresholds. If no alerts are found, the section is not included in the report. The ALERTS section focuses on identified violations of applied concepts and thresholds. The following subsystems may have problems and appear in the ALERTS section: file system, I/O configuration, paging configuration, I/O balance, page space, virtual memory, real memory, processes, and network.

For severity 1 levels, ALERTS focus on file systems, physical volumes, paging and memory. If you ask for severity 2 or 3 reporting, it adds information on configuration and processes, as seen here.

Alerts indicate suspicious configuration and load conditions. In this example, it appears that one disk is getting all the I/O activity. Clearly, the I/O load is not distributed to make the best use of the available resources.

The report continues on the next page.

## PDT Report Example (Part 2)



© Copyright IBM Corporation 2004

Figure 11-30. PDT Report Example (Part 2)

AU1612.0

### Notes:

The report then deals with UPWARD TRENDS and DOWNWARD TRENDS. These two sections focus on problem anticipation rather than on the identification of existing problems. The same concepts are applied, but used to project when violations might occur. If no trends are detected, the section does not appear.

PDT employs a statistical technique to determine whether or not there is a trend in a series of measurements. If a trend is detected, the slope of the trend is evaluated for its practical significance. For upward trends, the following items are evaluated: files, file systems, hardware and software errors, paging space, processes, and network. For downward trends the following can be reported: files, file systems, and processes.

The example UPWARD TRENDS section, identifies a possible trend with file system growth on lv01. An estimate is provided for the date at which the file system will be full, based on an assumption of linear growth.

The SYSTEM HEALTH section gives an assessment of the average number of processes in each process state on the system. Additionally, workload indicators are noted for any upward trends.

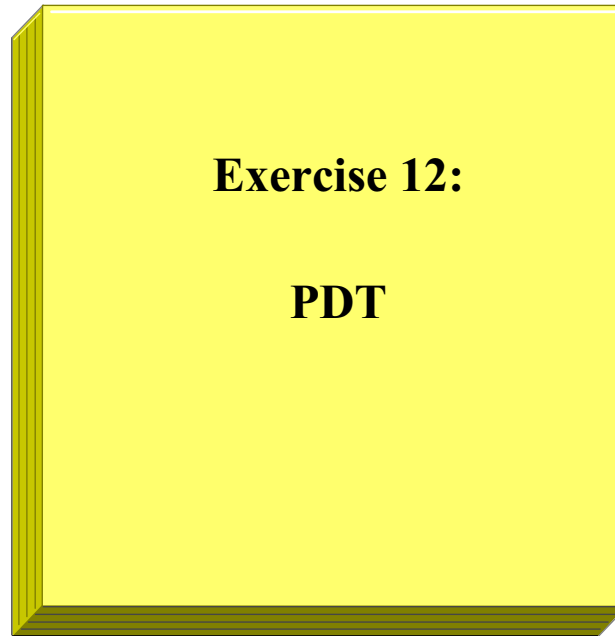
In the Summary section, the severity level of the current report is listed. There is also an indication given as to whether more details are available at higher severity levels. If so, an adhoc report may be generated to get more detail, using the `/usr/sbin/perf/diag_tool/pdt_report` command.



---

## Next Step

---



© Copyright IBM Corporation 2004

Figure 11-31. Next Step

AU1612.0

### **Notes:**

After completing the exercise, you should be able to:

- Use **PDT** for ongoing data capture and analysis of critical system resources

## Checkpoint

---

1. What command can be executed to identify CPU-intensive programs?
2. What command can be executed to start processes with a lower priority?
3. What command can you use to check paging I/O?
4. **T/F:** The higher the PRI value, the higher the priority of a process.

© Copyright IBM Corporation 2004

Figure 11-32. Checkpoint

AU1612.0

### **Notes:**

## Unit Summary

---

- The following commands can be used to identify potential bottlenecks in the system:
  - **ps**
  - **sar**
  - **vmstat**
  - **iostat**
- If you cannot fix a performance problem, **manage your workload** through other means (at, crontab, nice, renice).
- Use **PDT** to assess and control your systems performance.

© Copyright IBM Corporation 2004

Figure 11-33. Unit Summary

AU1612.0

### **Notes:**



# Unit 12. Security

## What This Unit Is About

This unit defines how to configure the auditing subsystem, customize authentication, and work with the Trusted Computing Base (TCB).

## What You Should Be Able to Do

After completing this unit, you should be able to:

- Provide Authentication Procedures
- Specify Extended File Permissions
- Configure the Trusted Computing Base (TCB)

## How You Will Check Your Progress

Accountability:

- Checkpoint questions
- Exercises

## References

- GG24-4433     *Elements of Security: AIX 4.1*
- Online        System Management Guide: Operating System and  
                  Devices: Chapter 2. Security  
                  AIX 5L Version 5.2 Security Guide

# Unit Objectives

---

After completing this unit, students should be able to:

- Provide **Authentication Procedures**
- Specify **Extended File Permissions**
- Configure the **Trusted Computing Base (TCB)**

© Copyright IBM Corporation 2004

---

Figure 12-1. Unit Objectives

AU1612.0

## **Notes:**

## 12.1 Authentication and Access Control Lists (ACLs)

# Protecting Your System

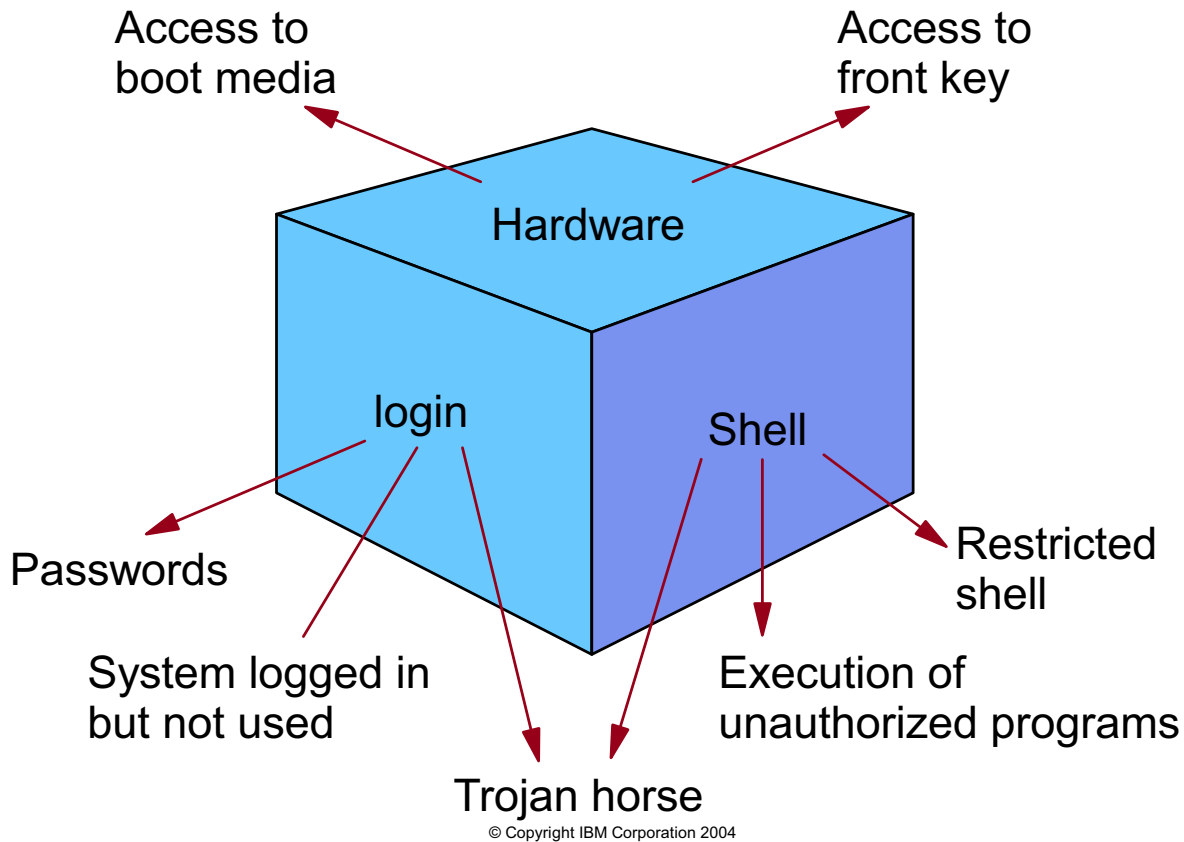


Figure 12-2. Protecting Your System

AU1612.0

## Notes:

If the machine is unattended in an open room, it is at risk from intruders. Anyone who can change the boot list or reboots the machine by pressing F1 or F5 can reboot the machine from an alternate media source (like CD-ROM or tape) and can invade a system. If you are using a PCI-based machine, the boot list is set via the SMS programs or via the bootlist command (run by root) from an AIX prompt. On classical machines, the front panel key selects either the service or normal boot list set by root with the bootlist command (or through service aids).

So what does that mean and how do we protect the systems? The first step is physical security. Without access to the machine, alternate boot media cannot be introduced into the machine. If the intruder has access, they can always shut the machine down by unplugging it. Since bootable media is fairly easy to obtain (especially if your intruder is an administrator), you must protect your machine.

On the classical machines, don't leave the key in the machine. Put the key in the normal or secure positions and take the key out and put it in a secure area. This will either prevent booting (secure) or boot only from a hard disk (normal).



Once logged into a shell, users are able to access, modify or delete any files for which they have permission. If tight control is not kept, they might gain access to unauthorized programs or files which may help them get the access or information they are seeking.

Consider configuring users to use the restricted shells or present them with an application menu instead of a shell prompt. Beware of a user's access to output devices such as printers. They can be used to print confidential material accessible by other users. Watch for "Trojan Horses". This is an executable named and positioned to look like a familiar command. They can perform many tasks without you being aware of it.

Security is the administrator's issue. But, it is also the users' issue. You need to educate your users and hold them accountable when they don't take security seriously. Strongly encourage users to log off when they're finished. Leaving the account logged in and unattended gives anyone access to the machine. It only takes seconds for someone to set up a backdoor. There are several variables that can be used to force a logoff if the session is inactive. In the Korn shell the variable is TMOU and in the Bourne shell it is TIMEOUT.

**Note:** This variable only works at the shell prompt. Remember, variables can be overridden by the user by editing **\$HOME/.profile**.

If a user wishes to lock the terminal but not log out, the **lock** command (or **xlock** command when using Xwindows) can be used. A password is needed to unlock the session.

SUID programs offer users access to the owner's account during the execution of the file. Avoid using them. If the program is poorly written, it could provide inappropriate access to the system. Shell scripts are particularly vulnerable. Fortunately, AIX ignores the SUID bit when used with a shell script. SUID-active files must be machine executable programs, for example, C-programs.

If an account is going to be inactive for a while, lock it. For example, if a user is planning a month long vacation, lock the account. Otherwise a hacker may gain access to the account and no one will notice any problems for the next 30 days. If a user no longer needs access to the system, the account should be locked so that no one can log into it. If the user's data is still required, change the ownership of those files to the new user.

System files, if accessed by an intruder, could be changed to allow the intruder access to the machine after reboot. Monitor the startup scripts which run from **inittab** regularly and ensure that all valid changes are clearly documented.

## How Do You Set Up Your PATH?

---

```
PATH=/usr/bin:/etc:/usr/sbin:/sbin:.
```

- or -

```
PATH=./usr/bin:/etc:/usr/sbin:/sbin
```

???

© Copyright IBM Corporation 2004

Figure 12-3. How Do You Set Up Your PATH?

AU1612.0

### **Notes:**

A common security risk comes up if the **PATH** variable is not set correctly.

At this point, ask yourself which definition do you prefer?

## Trojan Horse: An Easy Example (1 of 3)

```
$ cd /home/hacker
```

```
$ vi ls
```

```
#!/usr/bin/ksh
```

```
cp /usr/bin/ksh /tmp/.hacker
```

```
chown root /tmp/.hacker
```

```
chmod u+s /tmp/.hacker
```

```
rm -f $0
```

```
/usr/bin/ls $*
```

SUID-Bit: Runs under  
root authority



```
$ chmod a+x ls
```

© Copyright IBM Corporation 2004

Figure 12-4. Trojan Horse: An Easy Example (1 of 3)

AU1612.0

### Notes:

What is a **trojan horse**? A trojan horse behaves like an ordinary UNIX command. During the execution of a trojan horse, dangerous actions take place that are intentionally hidden from you. In the example, a user, **hacker**, creates a shell script with the name **ls**. This script really executes an **ls** command, but it does additional things that are not visible during the execution. It copies the shell **/usr/bin/ksh** to a file **/tmp/.hacker**, changes the owner to **root** and sets the **Set-User-Id-Bit**. If the file **/tmp/.hacker** is executed, it runs with **root** authority.

Note that the procedure is destroyed during the execution (`rm -f $0`). The question now is: How can we tell the system administrator to execute the trojan horse?

## Trojan Horse: An Easy Example (2 of 3)

```
$ cd /home/hacker
$ cat > -i
blablabla<CTRL-D>
```

Hello SysAdmin,  
I have a file "-i" and cannot  
remove it. Please help me ...



```
PATH=./usr/bin:/etc:/usr/sbin:/sbin
```

```
# cd /home/hacker
# ls
-i
```

© Copyright IBM Corporation 2004

Figure 12-5. Trojan Horse: An Easy Example (2 of 3)

AU1612.0

### Notes:

The user **hacker** creates a file **-i**. This file is difficult to remove since you cannot run the command **rm -i** without getting a syntax error. The **hacker** sends you a mail requesting your help.

If **root** specifies the **PATH** as shown on the visual, the trojan horse **ls** from user **hacker** will be executed after changing to **/home/hacker**. Note that you do not see the trojan horse itself because it will be destroyed during execution.

## Trojan Horse: An Easy Example (3 of 3)

```
$ cd /tmp
$ .hacker
# passwd root
```

Effective **root** authority



Don't worry, be happy ...

~~PATH=./usr/bin:/etc:/usr/sbin:/sbin~~

When using as root user, **never** specify the working directory in the **PATH** variable.

© Copyright IBM Corporation 2004

Figure 12-6. Trojan Horse: An Easy Example (3 of 3)

AU1612.0

### Notes:

During the execution of the trojan horse the program **/usr/bin/ksh** has been copied to **/tmp/.hacker**. This program has the **SUID-Bit** on.

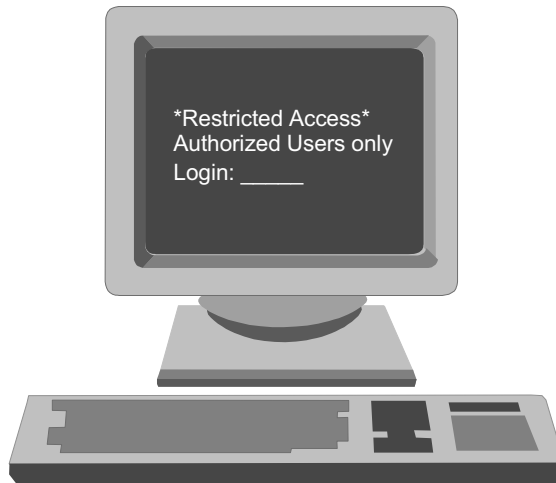
When a normal user executes this program, the user becomes **root** and you might run into big, big problems afterwards.

Never specify the working directory in the **PATH** variable, when working as **root** user.

# login.cfg: login prompts

```
# vi /etc/security/login.cfg
```

```
default:
    sak_enabled = false
    logintimes =
    .
    .
    .
    herald = "\n\n*Restricted Access*\n\nAuthorized Users Only\n\nLogin: "
```



© Copyright IBM Corporation 2004

Figure 12-7. login.cfg: login prompts

AU1612.0

## Notes:

Login prompts present a security issue. Your login prompts should send a clear message that only authorized users should log in and it should not give hackers any additional information about your system. Prompts should not describe your type of system or your company name. This is information that a hacker can use. For example, a login prompt that indicates it is a UNIX machine tells the hacker that there is likely an account call **root**. Now, only a password is needed.

Depending on whether you want to set your ASCII prompt or your graphical login, you will need to alter different files.

For ASCII prompts, edit `/etc/security/login.cfg`. In the **default** stanza, you need to add a line similar to the example:

```
herald = "\n\n*RESTRICTED ACCESS*\n\nAuthorized Users Only\n\nLogin:"
```

The `\n` is a new line and `\r` is a return. These are used to position the text on the screen. Do not use the `<enter>` key inside the quotes. It will not display like you would hope.

For the CDE environment, you need to modify the file **Xresources** in **/etc/dt/config/\$LANG**. If it does not exist, copy **/usr/dt/config/\$LANG/Xresources** to **/etc/dt/config/\$LANG/Xresources**. In this file, locate the lines:

**!! Dtlogin\*greeting.labelString: Welcome to %LocalHost%**

**!! Dtlogin\*greeting.persLabelString: Welcome %s**

Make a copy of both lines before you do any editing. Edit the (copied) lines and remove the comment string “!!”. The information after the colons is what appears on your login screen. **label.String** controls the initial login display when the user is prompted for the login name. **persLabelString** shows when asking for the user’s password. The **%LocalHost** displays the machine name and **%s** displays the user's login name. Modify the message to your liking.

## login.cfg: Restricted Shell

```
# vi /etc/security/login.cfg
```

```
* Other security attributes
```

```
usw:
```

```
shells = /bin/sh,/bin/bsh,/usr/bin/ksh, ...,/usr/bin/Rsh
```

```
# chuser shell=/usr/bin/Rsh michael
```

**michael can't:**

- Change the current directory
- Change the PATH variable
- Use command names containing slashes
- Redirect standard output (>, >>)

© Copyright IBM Corporation 2004

Figure 12-8. login.cfg: Restricted Shell

AU1612.0

### Notes:

All valid login shells are listed in the **usw** stanza in **/etc/security/login.cfg**. If you work on a system where security is a potential problem you can assign a **restricted shell** to users. The effect of these restrictions is to prevent the user from running any command that is not in a directory contained in the **PATH** variable.

To enable a restricted shell on a system, you have to do two things:

1. Add **/usr/bin/Rsh** to the list of shells.
2. Assign the restricted shell to the corresponding users on your system.

If you are going to assign a restricted shell, ensure that the **PATH** variable does not contain directories like **/usr/bin** or **/bin**. Otherwise, the restricted user is able to start other shells (like **ksh**) that are not restricted.

To give a limited set of commands to a user, copy the commands to **/usr/rbin** and add **/usr/rbin** to their **PATH**.



# Customized Authentication

```
# vi /usr/lib/security/methods.cfg
```

```
* Authentication Methods

secondPassword:
  program = /usr/local/bin/getSecondPassword
```

```
# vi /etc/security/user
```

```
michael:
  auth1 = SYSTEM,secondPassword
```

© Copyright IBM Corporation 2004

Figure 12-9. Customized Authentication

AU1612.0

## Notes:

AIX allows you to specify self-written **authentication methods**. These programs are called whenever you log in to your system. To install an additional authentication method, you must do two things:

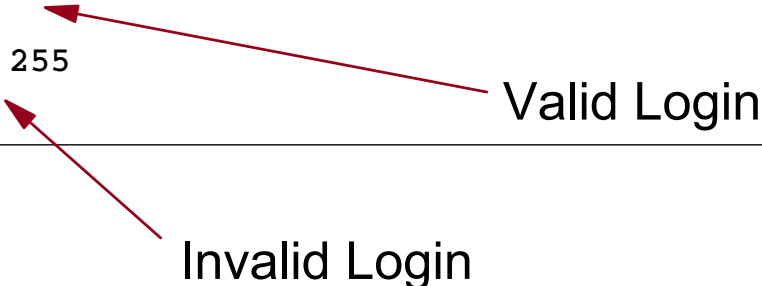
1. Create a stanza for your authentication method in **/usr/lib/security/methods.cfg**. In the example we use the name **secondPassword**. This stanza has only one attribute, **program**. This attribute contains the **full pathname** of the authentication program. Note that this program must be executable.
2. Add the authentication method for the user that should invoke this authentication method during the login-process. To do so, add the **auth1** attribute to the user in **/etc/security/user** as shown on the visual.

The **Common Desktop Environment (CDE)** does not support additional authentication methods.

## Authentication Methods (1 of 2)

```
# vi /usr/local/bin/getSecondPassword
```

```
print "Please enter the second Password: "  
  
stty -echo          # No input visible  
read PASSWORD  
stty echo  
  
if [[ $PASSWORD = "d1f2g3" ]]; then  
    exit 0  
else  
    exit 255  
fi
```



Valid Login

Invalid Login

© Copyright IBM Corporation 2004

Figure 12-10. Authentication Methods (1 of 2)

AU1612.0

### Notes:

The visual shows an **authentication method** that prompts the user for a password. If the correct password (d1f2g3) is entered, the value **0** is returned, indicating a valid log in.

If the password is not correct, a **non-zero value** indicates an invalid login. In this case the user cannot log in.

## Authentication Methods (2 of 2)

```
# vi /usr/local/bin/limitLogins

#!/usr/bin/ksh

# Limit login to one session per user

USER=$1      # User name is first argument

              # How often is the user logged in?
COUNT=$(who | grep "^$USER | wc -l)

              # User already logged in?
if [[ $COUNT -ge 1 ]]; then
    errlogger "$1 tried more than 1 login"
    print "Only one login is allowed"
    exit 128
fi

exit 0      # Return 0 for correct authentication
```

© Copyright IBM Corporation 2004

Figure 12-11. Authentication Methods (2 of 2)

AU1612.0

### Notes:

The visual shows an **authentication method** that **limits the number of login sessions**.

The user name is passed as first argument. For this user the procedure determines via a **command substitution** how often the user is already logged in. If this number is greater or equal to 1, an entry is posted to the error log and the value 128 is returned, indicating an invalid login. Otherwise the value 0 is returned - the login will be successful.

To set this up, add this program name to the authentication methods in **/usr/lib/security/methods.cfg** and set the **auth1** line in the users' stanza in **/etc/security/user**.

# Two-Key Authentication

```
# vi /etc/security/user
```

```
boss:  
  auth1 = SYSTEM;deputy1,SYSTEM;deputy2
```



```
login: boss  
deputy1's Password:  
deputy2's Password: _____
```

© Copyright IBM Corporation 2004

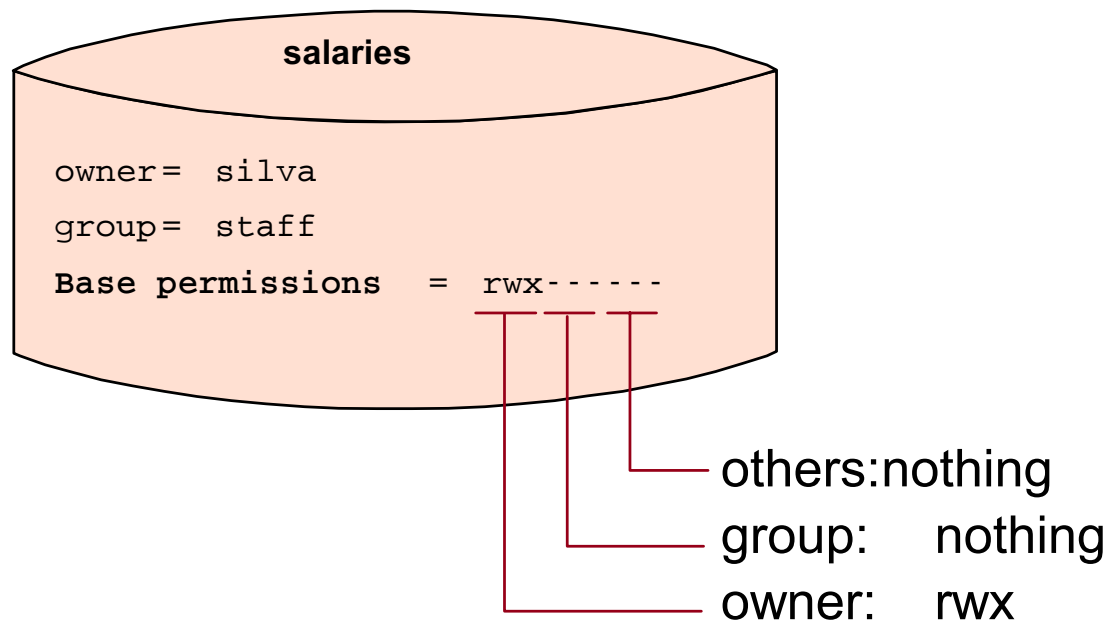
Figure 12-12. Two-Key Authentication

AU1612.0

## Notes:

AIX allows you to create a **two-key** authentication. In the above example, **SYSTEM** is defined as the authentication method twice. **SYSTEM** is supplied with two arguments, **deputy1** and **deputy2**. Therefore, **both** passwords must be entered correctly before the user **boss** may log in.

## Base Permissions



How can **silva** easily give **simon** read access to the file **salaries**?

© Copyright IBM Corporation 2004

Figure 12-13. Base Permissions

AU1612.0

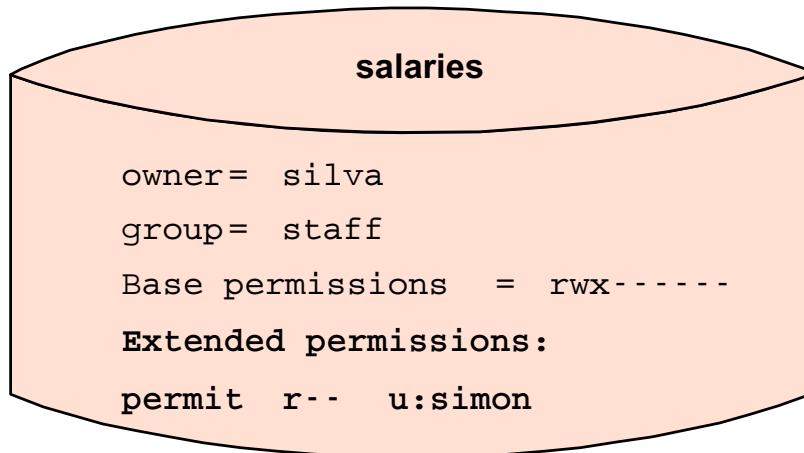
### Notes:

Here is a perplexing problem. If user **silva** owns a file called **salaries**, which contains very sensitive data, how can she easily give user **simon** permission to read the file? Possible solutions:

- **root** could give the file to **simon** (**chown**), but then **silva** won't be able to access it and **simon** can make changes to it.
- **silva** could copy the file for **simon** (**cp**), but then two files would exist, and that causes data integrity problems.
- **silva** could change the group identification for the file (**chgrp**) to a new group and **simon** could have that group added to his list of group membership. But if that were done frequently on the system it would cause a system management nightmare.

The best solution would be if **silva** could add **simon** to a list of those specific users who could read the **salaries** file. This is where **Access Control Lists (ACLs)** come in.

## Extended Permissions: Access Control Lists



```
# acledit salaries
```

EDITOR

```

base permissions
...
extended permissions
enabled
permit r-- u:simon
  
```

© Copyright IBM Corporation 2004

Figure 12-14. Extended Permissions: Access Control Lists

AU1612.0

### Notes:

The **base permissions** control the rights for the **owner**, the **group**, and all others on the system. If you want to specify additional rights, you can use **Access Control Lists** to expand the base permissions.

One way to do this is by executing the **acledit** command, which opens up an editor (specified by the variable **EDITOR**). In the editor session, you must do the following things:

- Enable the extended permissions, by changing the word **disabled** to **enabled**.
- Add additional permissions by using **special keywords**. These keywords are explained on the next visuals. In the example, we **permit** the user **simon** **read** access to file **salaries**.
- Another way to set Extended Permissions is by using the File Manager under CDE.

# ACL Commands

```
# aclget file1 ← Display base/extended permissions
                    ↓ Copy an Access Control List
# aclget status99 | aclput report99

# acledit salaries2 ← To specify extended permissions
```

- `chmod` in the octal format **disables** ACLs
- Only the **backup** command saves ACLs
- **acledit** requires the **EDITOR** variable (full pathname of an AIX editor)

© Copyright IBM Corporation 2004

Figure 12-15. ACL Commands

AU1612.0

## Notes:

Three commands are available to work with **Access Control Lists (ACLs)**:

1. The command **aclget** displays the access control information on standard output.
2. The command **aclput** sets the access control information of a file and is often used in a **pipe** context, to copy the permissions from one file to another as in the above example. Here is another way to copy the ACL from a file:

```
# aclget -o status99.acl status99
# aclput -i status99.acl report99
```

This example works in the same way as the version with the **pipe**. Instead of using a **pipe**, the ACL is written to a file **status99.acl**, that is used by **aclput**.

3. The command **acledit** allows you to edit the access control information of a file. The **EDITOR** variable must be specified with a **complete** pathname, otherwise the command will fail. Note that the entire ACL cannot exceed 4096 bytes.

If you execute a **chmod** in the **octal format**, the ACL will be **disabled**. The extended permissions are still stored, but will not be used. To turn them back on, use **acledit** and

change **disabled** to **enabled**. To prevent this problem, use **chmod** in **symbolic** format if you are working with a file that has extended permissions.

Only the **backup** command saves ACLs. For example, if you use **tar** or **cpio** the ACLs are lost when you restore the corresponding file.

Let's show the **special keywords** that you can use in ACLs.



## ACL Keywords: permit and specify

```
# acledit status99
```

```
attributes:
  base permissions
    owner(fred): rwx
    group(finance): rw-
    others: ---
  extended permissions
  enabled
  permit    - - x    u:michael
  specify   r - -    u:anne,g:account
  specify   r - -    u:nadine
```

- **michael** (member in group finance) gets read, write (base) and execute (extended) permission.
- If **anne** is in group **account**, she gets read permission on file status99.
- **nadine** (member in group finance) gets only read access

Figure 12-16. ACL Keywords: permit and specify

AU1612.0

### Notes:

Extended permissions give the owner of a file the ability to define the access to a file more precisely. **Special keywords** are used to define the access mode:

- The keyword **permit** grants the user or group the specified access to a file. In the example the user **michael** who is a member in group **finance** gets execute privileges. Therefore, **michael** has read, write and execute permission on the file **status99**.
- The keyword **specify** precisely defines the file access for a user or group. In the example the user **anne** gets read permission, but only if she is a member of the group **account**. Putting **u:** and **g:** on the same line requires both conditions to be true for the ACL to apply.
- In the last example, user **nadine** is a member of the finance group which normally has read and write privileges. But, the **specify**, in this case, gives **nadine** only **read** privileges. The base permissions no longer apply to **nadine**.

## ACL Keywords: deny

```
# acledit report99
```

```
attributes:
base permissions
  owner (sarah): rwx
  group (mail): r--
  others: r--
extended permissions
enabled
deny  r--  u:paul g:mail
deny  r--  g:gateway
```

- **deny**: Restricts the user or group from using the specified access to the file
- **deny** overrules **permit** and **specify**

© Copyright IBM Corporation 2004

Figure 12-17. ACL Keywords: deny

AU1612.0

### Notes:

The ACL keyword **deny** restricts the user or group from the specified access to a file.

- In the example, the group **mail** gets a read access to file report99. If the user **paul** is a member of group **mail** then read access is denied for him.
- The rest of the world gets read access to file report99. The exception is group **gateway**; this group has no access rights to the file.

If a user or group is denied a particular access by either a **deny** or **specify** keyword, no other entry can override this access denial.

---

## JFS2 Extended Attributes Version 2 (AIX 5.3)

---

- Extension of normal attributes
- Name and value pairs
- `setea` - to associate name/value pairs
- `getea` - to view

```
#setea -n Author -v DeChalus report1
#getea report1
EAName: Author
EAValue:
DeChalus
```

© Copyright IBM Corporation 2004

Figure 12-18. JFS2 Extended Attributes Version 2 (AIX 5.3)

AU1612.0

### Notes:

Extended attributes are an extension of the normal attributes of a file (such as size and mode). They are (name, value) pairs associated with a file or directory. The name of an attribute is a null-terminated string. The value is arbitrary data of any length. There are two types of extended attribute: extended attribute version 1 (EA<sub>v1</sub>) and extended attribute version 2 (EA<sub>v2</sub>). Starting with AIX 5L Version 5.3, EA<sub>v2</sub> with JFS2 is now available. It should be noted that EA<sub>v2</sub> is required to use an NFS4 ACL (now available with AIX 5.3).

Setting extended attribute version 2

if you created a file named `report1` and want to set attributes to the file such as author, date, revision number, comments and so on (DeChalus as author in this example), you can accomplish this with **setea** to set the value of an extended attribute and **getea** to read the value of an extended attribute shown in the figure.

```
#setea -n Name { -v Value | -d | -f EAFile } FileName ...
```

```
#getea [-n Name] [-e RegExp] [-s] FileName
```

## Next Step

---



© Copyright IBM Corporation 2004

Figure 12-19. Next Step

AU1612.0

### **Notes:**

After the exercise, you should be able to:

- Customize the **login.cfg** file
- Add an additional primary **authentication method** for a user
- Implement **access control lists (ACLs)**

## 12.2 The Trusted Computing Base (TCB)

# The Trusted Computing Base (TCB)

The **TCB** is the part of the system that is responsible for enforcing the **security policies** of the system.

```
# ls -l /etc/passwd
-rw-r--rw- 1 root security ... /etc/passwd

# ls -l /usr/bin/be_happy
-r-sr-xr-x 1 root system ... /usr/bin/be_happy
```

© Copyright IBM Corporation 2004

Figure 12-20. The Trusted Computing Base (TCB)

AU1612.0

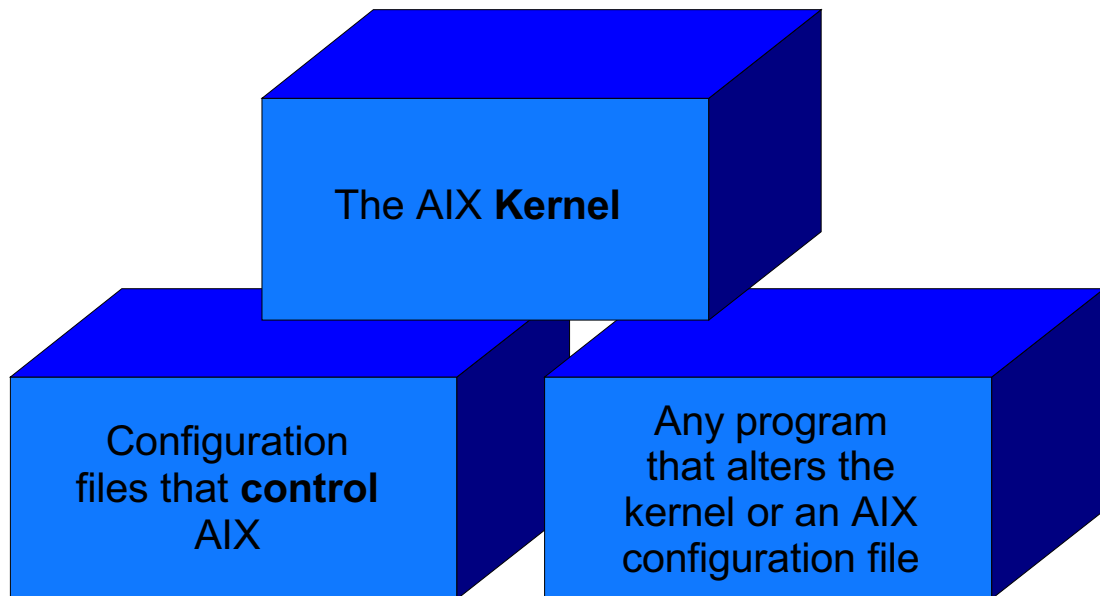
## Notes:

The **Trusted Computing Base** is the part of the system that is responsible for enforcing the information security policies of the system.

The visual shows examples where these security policies have been violated:

- The configuration file **/etc/passwd** allows a write access to all others on the system, which is a big security hole. Somebody has changed the default value of **rw-r--r--** for **/etc/passwd**. If the TCB is enabled on a system, the system administrator will be notified that the file mode for **/etc/passwd** has been changed, when he checks the TCB.
- Somebody has installed a program **/usr/bin/be\_happy**, which is executable for all users. Additionally this program has the **SUID** bit, that means during the execution this program runs with the effective user ID of **root**. If the person who administers the system runs a TCB check, he will be notified that a **SUID**-program has been installed, that is not part of the TCB.

# TCB Components



The TCB can only be enabled at installation time

© Copyright IBM Corporation 2004

Figure 12-21. TCB Components

AU1612.0

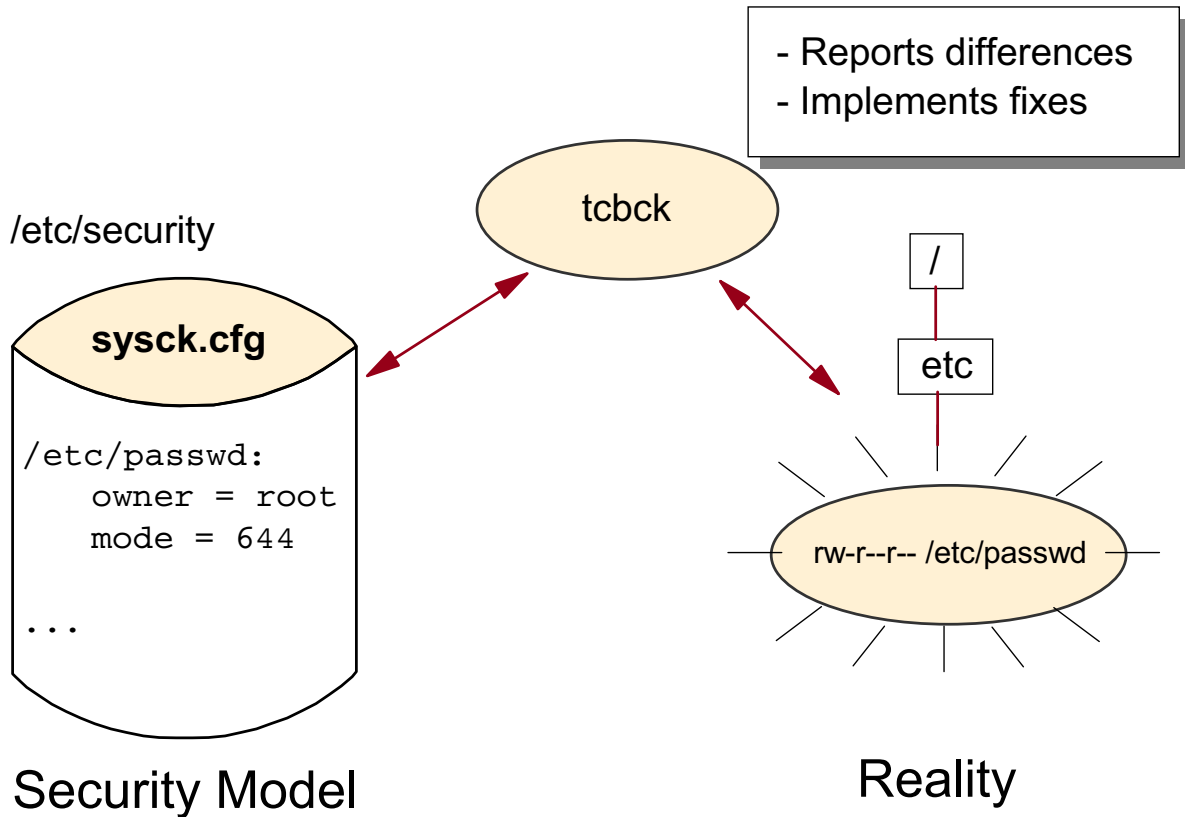
## Notes:

The **Trusted Computing Base (TCB)** consists of:

- The **AIX Kernel** (your operating system)
- All **configuration files** that are used to control AIX (for example: `/etc/passwd`, `/etc/group`)
- Any program that alters the kernel (for example: `mkdev`, `cfgmgr`) or an AIX configuration file (for example: `/usr/bin/passwd`, `/usr/bin/mkuser`)

Many of the TCB functions are optionally enabled at **installation time**. Selecting **yes** for the **Install Trusted Computing Base** option on the *Installation and Settings* menu enables the **TCB**. Selecting **no** disables the **TCB**. The **TCB** can only be enabled at installation time.

# Checking the Trusted Computing Base



© Copyright IBM Corporation 2004

Figure 12-22. Checking the Trusted Computing Base

AU1612.0

## Notes:

To check the security state of your system, the command **tcbck** is used. This command audits the security information by reading the `/etc/security/sysck.cfg`. This file includes a description of all TCB files, configuration files and trusted commands.

If differences between the **security model** as described by `sysck.cfg` and the **reality** occur, the **tcbck** command reports them to standard error. According to the option you use, **tcbck** fixes the differences automatically.

If the **Install Trusted Computing Base** option was not selected during the initial installation, the **tcbck** command will be disabled. The command can be properly enabled only by reinstalling the system.



# The sysck.cfg File

```
# vi /etc/security/sysck.cfg

...

/etc/passwd:
  owner = root
  group = security
  mode = TCB, 644
  type = FILE
  class = apply, inventory, bos.rte.security
  checksum = VOLATILE
  size = VOLATILE

...

# tcbck -t /etc/passwd
```

© Copyright IBM Corporation 2004

Figure 12-23. The sysck.cfg File

AU1612.0

## Notes:

The **tcbck** command reads the **/etc/security/sysck.cfg** file to determine the files to check. Each trusted file on the system should be described by a stanza in the **/etc/security/sysck.cfg** file.

Each file stanza must have the **type** attribute and can have one or more of the following attributes:

<b>acl</b>	Text string representing the <b>access control list</b> for the file. It must be of the same format as the output of the <b>aclget</b> command.
<b>class</b>	Logical name of a <b>group</b> of files. This attribute allows several files with the same class name to be checked by specifying a single argument to the <b>tcbck</b> command.
<b>checksum</b>	Defines the checksum of the file, calculated by the <b>sum -r</b> command.
<b>group</b>	Group ID or name of the file's group.
<b>links</b>	Comma-separated list of path names linked to this file. Defines the absolute paths that have hard links to this object.

<b>mode</b>	Comma-separated list of values. The allowed values are <b>SUID</b> , <b>SGID</b> , <b>SVTX</b> and <b>TCB</b> . The file permissions must be the last value and can be specified either as an octal value or as a 9-character string.
<b>owner</b>	User ID or name of the file owner.
<b>size</b>	Defines the size (in decimal) of the file in bytes. This attribute is only valid for regular files.
<b>program</b>	Comma-separated list of values. The first value is the path name of a <b>checking program</b> . Additional values are passed as arguments to the program when it is executed. The checking program must return 0 to indicate that no errors were found. All errors must be written to standard error. Note that these checker programs run with <b>root</b> authority.
<b>symlinks</b>	Comma-separated list of path names, symbolically linked to this file.
<b>type</b>	The type of the file. One of the following keywords must be used: <b>FILE</b> , <b>DIRECTORY</b> , <b>FIFO</b> , <b>BLK_DEV</b> , <b>CHAR_DEV</b> , <b>MPX_DE</b>

## tcbck: Checking Mode Examples

```
# chmod 777 /etc/passwd
# ls -l /etc/passwd
-rwxrwxrwx 1 root security .../etc/passwd

# tcbck -t /etc/passwd
The file /etc/passwd has the wrong file mode
Change mode for /etc/passwd ?
(yes, no ) yes

# ls -l /etc/passwd
-rw-r--r-- 1 root security .../etc/passwd
```

---

```
# ls -l /tmp/.4711
-rwsr-xr-x 1 root system.../tmp/.4711

# tcbck -t tree
The file /tmp/.4711 is an unregistered set-UID program.
Clear the illegal mode for /tmp/.4711 (yes, no) yes

# ls -l /tmp/.4711
-rwxr-xr-x 1 root system.../tmp/.4711
```

Figure 12-24. tcbck: Checking Mode Examples

AU1612.0

### Notes:

The **tcbck** command audits the security state of a system. The command supplies a **check mode** and an **update mode**. Let's start with the **check mode**:

The visual shows how the check mode of **tcbck** can be used to find any security violations.

- In the first example somebody changed the file mode for **/etc/passwd** to read, write and execute permissions for all users on the system. The command **tcbck -t** specifies checking mode and indicates that errors are to be reported with a prompt asking whether the error should be fixed. In the example we select **yes** and the file mode is restored to its original value as specified in **/etc/security/sysck.cfg**.
- In the second example somebody installed a **SUID** program **/tmp/.4711**. The command **tcbck -t tree** indicates that all files on the system are checked for correct installation. The **tcbck** command discovers any files that are potential threats to system security. It gives you the opportunity to alter the suspected file to remove the offending attribute. The **SUID**-bit is removed after selecting **yes** at the **tcbck** prompt.

## tcbck: Checking Mode Options

	Report:	Fix:
tcbck <b>-n</b> <what>	yes	no
tcbck <b>-p</b> <what>	no	yes
tcbck <b>-t</b> <what>	yes	prompt
tcbck <b>-y</b> <what>	yes	yes

**<what>** can be:

- a *filename* (for example /etc/passwd)
- a *classname*: Logical group of files defined by a **class = name** in sysck.cfg
- **tree**: Check all files in the filesystem tree
- **ALL**: Check all files listed in sysck.cfg

© Copyright IBM Corporation 2004

Figure 12-25. tcbck: Checking Mode Options

AU1612.0

### Notes:

The checking mode of **tcbck** can be enabled by any of the following options:

- n** Indicates that errors are to be reported, but not fixed.
- p** Indicates that errors are to be fixed, but not reported. Be careful with this option.
- t** Indicates that errors are to be reported with a prompt asking whether the error should be fixed.
- y** Indicates that errors are to be fixed and reported. Be careful with this option.

All options that fix automatically should be used with care because the access to system files could be dropped if the **TCB** is not maintained correctly.

The files that must be checked are specified as shown on the visual. After specifying the check mode, you could check:

- One selected file (for example **/etc/passwd**)

- A class of files grouped together by the **class** attribute in **/etc/security/sysck.cfg**
- All files in the file system tree by specifying the word **tree**. In this case, files that are **not** in **/etc/security/sysck.cfg** must *not*:
  - Have the **Trusted Computing Base** attribute set (see **chtcb** for an explanation of this attribute)
  - Be **setuid** or **setgid** to an administrative ID
  - Be linked to a file in the **sysck.cfg** file
  - Be a device special file
- All files listed in **/etc/security/sysck.cfg** by specifying the word **ALL**

## tcbck: Update Mode Examples

```
# tcbck -a /salary/salary.dat class=salary
```

Add salary.dat to  
sysck.cfg
Additional  
class information

```
# tcbck -t salary }
```

Test all files belonging  
to class salary

```
# tcbck -d /etc/cvid }
```

Delete file /etc/cvid from  
sysck.cfg

© Copyright IBM Corporation 2004

Figure 12-26. tcbck: Update Mode Examples

AU1612.0

### Notes:

In the **update mode**, the **tcbck** command adds (-a), deletes (-d) or modifies file definitions in **/etc/security/sysck.cfg**. The visual shows how a file **/salary/salary.dat** is added to **sysck.cfg**. An additional class name **salary** is specified. This class name could be used in the check, to test all files that belong to the class. Here are some more examples where the **update mode** of **tcbck** is used:

1. To add a file **/usr/local/bin/check** with **acl**, **checksum**, **class**, **group** and **owner** attributes to **sysck.cfg**, enter:

```
# tcbck -a /usr/local/bin/check acl checksum class=rocket group owner
```

2. If you remove a file, for example **/etc/cvid**, from the system that is described in **sysck.cfg**, you should also remove the description from this file. To do this, use the option **-d**:

```
# tcbck -d /etc/cvid
```

If you must add **/dev**-files to **sysck.cfg**, you must use the option **-l** (lowercase l). For example to add the newly created **/dev** entries **foo** and **bar**, enter:

```
# tcbck -l /dev/foo /dev/bar
```

## chtcb: Marking Files As Trusted

```
# ls -le /salary/salary.dat
-rw-rw----- root salary ...
salary.dat
```

No "+" indicates not trusted

```
# tcbck -n salary
The file /salary/salary.dat has the wrong
TCB attribute value
```

tcbck indicates a problem!

```
# chtcb on /salary/salary.dat
# ls -le /salary/salary.dat
-rw-rw-----+ root salary ...
salary.dat
```

Now it's trusted !

© Copyright IBM Corporation 2004

Figure 12-27. chtcb: Marking Files As Trusted

AU1612.0

### Notes:

Just adding information about the file to the **sysck.cfg** is not enough. The file must also be marked as **trusted** in the **inode**. To do this, use the **chtcb** command.

In the example, our file **salary.dat** is in the database but is not trusted. If you use the command **ls -le**, a **+** symbol will show in the permissions area, if the file is trusted. When we execute the **tcbck** command to audit the files, it will return an error because our file is not trusted.

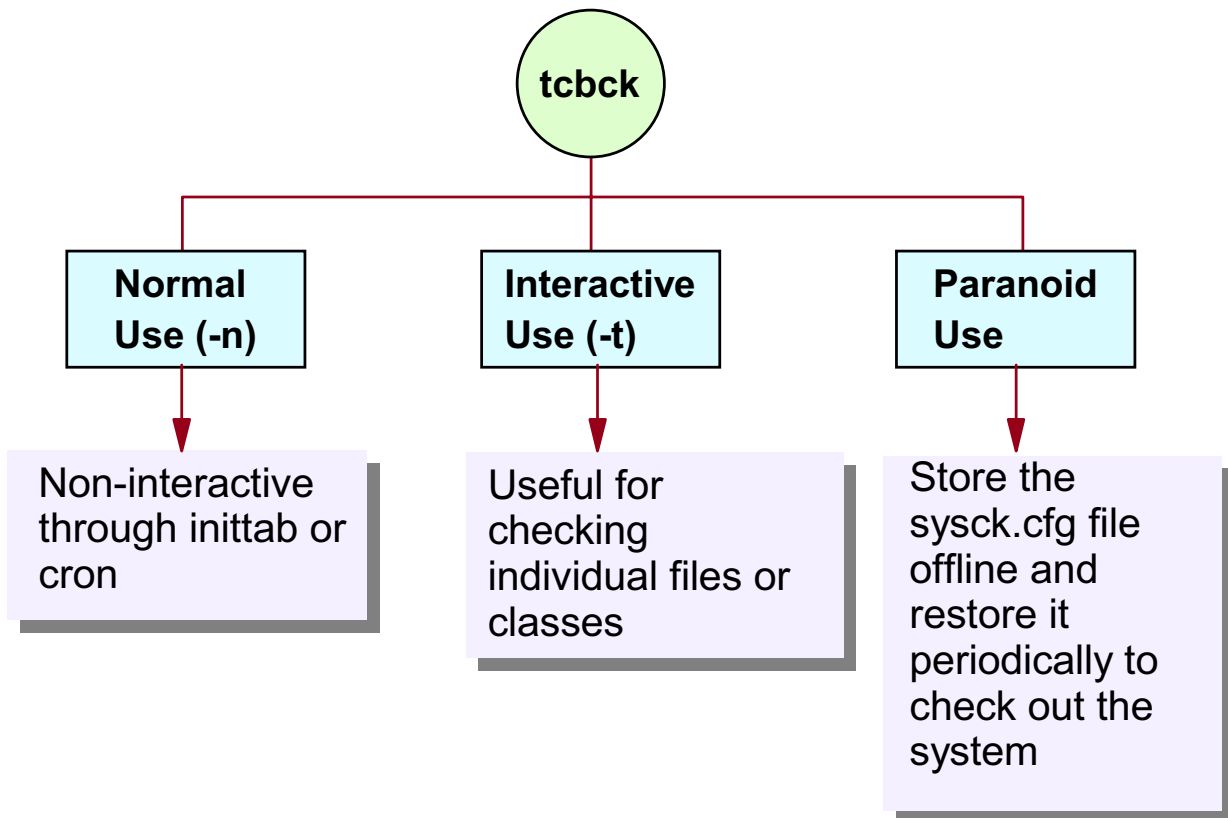
To mark it trusted, run the **chtcb** command with the option of **on**. Now the file is ready.

The **+**-symbol can indicate two things. It can indicate that the file is trusted or that the file contains extended permissions (ACLs). If you are unsure what the **+**-symbol is indicating, you can run **chtcb query** to see if it is a trusted file or **aclget** to see if there are extended permissions.

```
# chtcb query /salary/salary.dat
# aclget /salary/salary.dat
```

We come back to **chtcb** later in this unit.

# tcbck: Effective Usage



© Copyright IBM Corporation 2004

Figure 12-28. tcbck: Effective Usage

AU1612.0

## Notes:

If you decide to use **tcbck**, you should plan and try this very carefully. You need to get some experience with **tcbck**, before you use it in a production environment.

The **tcbck** command can be used in three ways:

- **Normal Use** means that the **tcbck** command is integrated either in an entry in **/etc/inittab** or in **crontab**. In this case, you must redirect standard error to a file that could be analyzed later.
- The **Interactive Use (tcbck -t)** can be used effectively, to check selected files or classes that you've defined.
- **Paranoid Use** means that you store the file **/etc/security/sysck.cfg** offline. The reason for this is if someone successfully hacks into the **root** account, not only can they add programs to the system, but since they have access to everything, they can also update the **sysck.cfg** file. By keeping a copy of **sysck.cfg** offline, you will have a safe copy. Move your offline copy back onto the system and then run the **tcbck** command.



## Trusted Communication Path

The **Trusted Communication Path** allows for secure communication between users and the Trusted Computing Base.

What do you think when you see this screen on a terminal ?



```
AIX Version 4  
(C) Copyrights by IBM and by others 1982, 1996  
login:
```

© Copyright IBM Corporation 2004

Figure 12-29. Trusted Communication Path

AU1612.0

### Notes:

AIX offers an additional feature, the **Trusted Communication Path**, that allows for **secure communication** between users and the **Trusted Computing Base**.

Why do you need this?

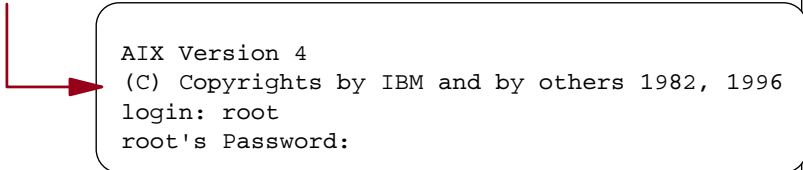
Look on the visual. Imagine you see this prompt on a terminal. What do you think? Surely you think that's a normal login prompt.

Now, look on the next visual.

## Trusted Communication Path: Trojan Horse

```
#!/usr/bin/ksh

print "AIX Version 4"
print "(C) Copyrights by IBM and by others 1982, 1996"
print -n "login: "
read NAME
print -n "$NAME's Password: "
stty -echo
read PASSWORD
stty echo
print $PASSWORD > /tmp/.4711
```



```
AIX Version 4
(C) Copyrights by IBM and by others 1982, 1996
login: root
root's Password:
```

```
$ cat /tmp/.4711
darth22
```

© Copyright IBM Corporation 2004

Figure 12-30. Trusted Communication Path: Trojan Horse

AU1612.0

### Notes:

Look at the shell procedure in the visual. This procedure generates exactly the login prompt that was shown on the last visual. If a system intruder gets the opportunity to start this procedure on a terminal, he can retrieve the password of a user very easily. And if you log in as **root** on this terminal, you are in a very bad position afterwards.

How can you protect yourself against these trojan horses? Request a **trusted communication path** on a terminal, and all trojan horses will be killed.

---

## Trusted Communication Path Elements

---

The **Trusted Communication Path** is based on:

- A **trusted shell** (tsh) that only executes commands that are marked as being trusted
- A **trusted terminal**
- A **reserved key sequence**, called the **secure attention key** (SAK), which allows the user to request a trusted communication path

© Copyright IBM Corporation 2004

Figure 12-31. Trusted Communication Path Elements

AU1612.0

### **Notes:**

The **Trusted Communication Path** is based on:

- A trusted command interpreter (**tsh** command), that only executes commands that are marked as being a member of the **Trusted Computing Base**.
- A terminal that is configured to request a **trusted communication path**.
- A **reserved key sequence**, called the **secure attention key (SAK)**, which allows a user to request a **trusted communication path**.

The **Trusted Communication Path** works only on terminals. In graphical environments (including the **Common Desktop Environment** and commands like **telnet**), the **Trusted Communication Path** is not supported.

# Using the Secure Attention Key (SAK)

## 1. Before logging in at the trusted terminal:

```
AIX Version 4
(C) Copyrights by IBM and by others
1982, 1996
login: <CTRL-x><CTRL-r>
tsh>
```

↑  
Previous login was a trojan horse.

## 2. To establish a **secure environment**:

```
# <CTRL-x><CTRL-r>
tsh>
```

Ensures that no untrusted programs will be run with root authority.

© Copyright IBM Corporation 2004

Figure 12-32. Using the Secure Attention Key (SAK)

AU1612.0

### Notes:

You should use the **Secure Attention Key (SAK)** in two cases:

1. Before you log in on a terminal, press the **SAK**, which is the reserved key sequence **Ctrl-x, Ctrl-r**. If a new login screen scrolls up, you have a **secure path**.

If the **tsh** prompt appears, the initial login was a **trojan horse** that may have been trying to steal your password. Find out who is currently using this terminal with the **who** command, and then log off.

2. When you want to establish a **secure environment**, press the **SAK** sequence, which starts up a **trusted shell**. You may want to use this before you work as **root** user. This ensures that no untrusted programs will be run with **root** user authority.

## Configuring the Secure Attention Key

- Configure a trusted terminal:

```
# vi /etc/security/login.cfg

/dev/tty0:
    sak_enabled = true
```

- Enable a user to use the trusted shell:

```
# vi /etc/security/user

root:
    tpath = on
```

© Copyright IBM Corporation 2004

Figure 12-33. Configuring the Secure Attention Key

AU1612.0

### Notes:

To configure the **SAK**, you should always do two things:

1. Configure your terminals so that pressing the **SAK** sequence creates a **trusted communication path**. This is specified by the **sak\_enabled** attribute in **/etc/security/login.cfg**. If the value of this attribute is **true**, recognition of the **SAK** is enabled.
2. Configure the users that use the **SAK**. This is done by specifying the **tpath** attribute in **/etc/security/user**. Possible values are:

<b>always</b>	The user can only work in the <b>trusted shell</b> . This implies that the user's initial program is <b>/usr/bin/tsh</b> .
<b>notsh</b>	The user cannot invoke the trusted shell on a trusted path. If the user enters the <b>SAK</b> after logging in, the login session ends.
<b>nosak</b>	The <b>SAK</b> is disabled for all processes run by the user. Use this value if the user transfers binary data that might contain the <b>SAK</b> sequence <b>Ctrl-X, Ctrl-R</b> .
<b>on</b>	The user can invoke a trusted shell by entering the <b>SAK</b> on a configured terminal.

## chtcb: Changing the TCB Attribute

```
# chtcb query /usr/bin/ls
/usr/bin/ls is not in the TCB
```

```
tsh>ls *.c
ls: Command must be trusted to run in the tsh
```

```
# chtcb on /usr/bin/ls
```

```
tsh>ls *.c
a.c b.c d.c
```

© Copyright IBM Corporation 2004

Figure 12-34. chtcb: Changing the TCB Attribute

AU1612.0

### Notes:

In a **trusted shell** you can only execute programs that have been marked trusted.

For example, the program **/usr/bin/ls** cannot be executed in a **trusted shell**. It does not have the **TCB** attribute. To enable this attribute, use the keyword **on** as shown in the visual. To disable the **TCB** attribute, use the keyword **off**:

```
# chtcb off /usr/bin/ls
```

If you set the **TCB** attribute for a program, always add the definition for the program to **/etc/security/sysck.cfg** to monitor that the file is not manipulated.

---

## Checkpoint (1 of 2)

---

1. Any programs specified as “auth1” must return a zero in order for the user to log in. True or False?
2. How would you specify that all members of the security group had rwx access to a particular file except for John?
3. In which file must you specify the full path name of the program that is to be used as part of the authentication process when a user logs in?
4. Name the two modes that tcbck supports.

© Copyright IBM Corporation 2004

Figure 12-35. Checkpoint (1 of 2)

AU1612.0

### **Notes:**

## Checkpoint (2 of 2)

---

5. When you execute **<ctrl-x ctrl-r>** at a login prompt and you obtain the **tsh** prompt, what does that indicate?

6. The system administrator must manually mark commands as trusted, which will automatically add the command to the **sysck.cfg** file. True or False?

7. When the **tcbck -p tree** command is executed, all errors are reported and you get a prompt asking if the error should be fixed. True or False?

© Copyright IBM Corporation 2004

Figure 12-36. Checkpoint (2 of 2)

AU1612.0

### **Notes:**



## Unit Summary

---

- The auditing subsystem allows you to capture security-relevant events on a system.
- The authentication process in AIX can be customized by authentication methods.
- Access Control Lists allow a more granular definition of file access modes.
- The Trusted Computing Base is responsible for enforcing the security policies on a system.

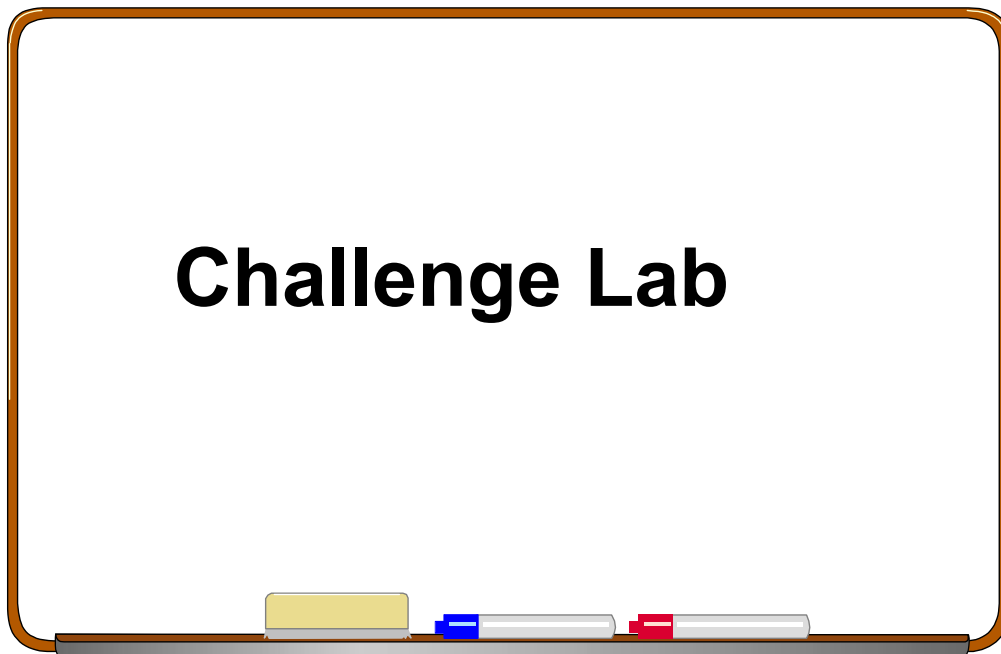
© Copyright IBM Corporation 2004

---

Figure 12-37. Unit Summary

AU1612.0

### **Notes:**



---

Figure 12-38. Challenge LAB

AU1612.0

**Notes:**

This challenge activity presents several “real-world” trouble-shooting problems.

The challenge activity is found in Appendix F. Turn to Appendix F and read the instructions carefully.

---

# Appendix A. Command Summary

## Startup, Logoff, and Shutdown

<Ctrl>d (exit)	log off the system (or the current shell).
shutdown	shuts down the system by disabling all processes. If in single-user mode, may want to use -F option for fast shutdown. -r option will reboot system. Requires user to be root.

## Directories

mkdir	make directory
cd	change directory. Default is \$HOME directory.
rmdir	remove a directory (beware of files starting with “.”)
rm	remove file; -r option removes directory and all files and subdirectories recursively.
pwd	print working directory
ls	list files

- a (all)
- l (long)
- d (directory information)
- r (reverse alphabetic)
- t (time changed)
- C (multi column format)
- R (recursively)
- F (places / after each directory name & \* after each exec file)

## Files - Basic

cat	list files contents (concatenate). Can open a new file with redirection, for example cat > newfile. Use <Ctrl>d to end input.
chmod	change permission mode for files or directories.

- chmod =+- files or directories
- (r,w,x = permissions and u, g, o, a = who)
- can use + or - to grant or revoke specific permissions.
- can also use numerics, 4 = read, 2 = write, 1 = execute.

- can sum them, first is user, next is group, last is other.
- for example “chmod 746 file1” is user = rwx, group = r, other = rw.

chown	change owner of files, for example chown owner file
chgrp	change group of files
cp	copy file
del	delete files with prompting (rm for no prompting)
mv	move and rename file
pg	list files contents by screen (page)
• h (help)	q (quit)
• <cr> (next pg)	f (skip 1 page),
• l (next line)	d (next 1/2 page)
• \$ (last page)	p (previous file),
• n (next file)	. (redisplay current page)
.	Current Directory
.	Parent Directory
	/string (find string forward), ?string (find string backward), - (move backward # pages), +# (move forward # pages)
rm	remove (delete) file(s) (-r option removes directory and all files and subdirectories)
head	print first several lines of a file
tail	print last several lines of a file
wc	report the number of lines (-l), words (-w), characters (-c) in a files. No options gives lines, words, and characters.
su	switch user
id	displays your user ID environment and how it is currently set
tty	displays the device that is currently active. Very useful for Xwindows where there are several pts devices that can be created. It's nice to know which one you have active. <b>who am i</b> will do the same.

## Files - Advanced

awk	programmable text editor / report write
banner	display banner (can redirect to another terminal “nn” with “>/dev/ttynn”)
cal	calendar (cal month year)

- cut cut out specific fields from each line of a file
- diff differences between two files
- find find files anywhere on disks. Specify location by path (will search all subdirectories under specified directory).
- name fl (file names matching fl criteria)
  - user ul (files owned by user ul)
  - size +n (or -n) (files larger (or smaller) than n blocks)
  - mtime +x (-x) (files modified more (less) than x days ago)
  - perm num (files whose access permissions match num)
  - exec (execute a command with results of find command)
  - ok (execute a cmd command interactively with results of find command)
  - o (logical or) print (display results. Usually included)
- find syntax: find path expression action
- for example find / -name "\*.txt" -print
  - or find / -name "\*.txt" -exec li -l {} \;
- (executes li -l where names found are substituted for {})  
; indicates end-of-command to be executed and \ removes usual interpretation as command continuation character)
- grep search for pattern, for example grep pattern file(s). Pattern can include regular expressions.
- c (count lines with matches, but don't list)
  - l (list files with matches, but don't list)
  - n (list line numbers with lines)
  - v (find files without pattern)
- expression metacharacters
- [ ] matches any one character inside.
  - with a - in [ ] will match a range of characters.
  - ^ matches BOL when ^ begins the pattern.
  - \$ matches EOL when \$ ends the pattern.
  - . matches any single character. (same as ? in shell).
  - \* matches 0 or more occurrences of preceding character.
- Note:** ".\*" is the same as "\*" in the shell.
- sed stream (text) editor. Used with editing flat files.
- sort sort and merge files. -r (reverse order); -u (keep only unique lines)

## Editors

ed	line editor
vi	screen editor
INed	LPP editor
emacs	screen editor +

## Shells, Redirection and Pipelining

< (read)	redirect standard input, for example “command < file” reads input for command from file.
> (write)	redirect standard output, for example “command > file” writes output for command to file overwriting contents of file.
>> (append)	redirect standard output, for example “command >> file” appends output for command to the end of file.
2>	redirect standard error (to append standard error to a file, use “command 2>> file”) combined redirection examples: <ul style="list-style-type: none"><li>• command &lt; infile &gt; outfile 2&gt; errfile</li><li>• command &gt;&gt; appendfile 2&gt;&gt; errfile &lt; infile</li></ul>
;	command terminator used to string commands on single line
	pipe information from one command to the next command. For example “ls   cpio -o > /dev/fd0” will pass the results of the ls command to the cpio command.
\	continuation character to continue command on a new line. Will be prompted with > for command continuation.
tee	reads standard input and sends standard output to both standard output and a file. For example “ls   tee ls.save   sort” results in ls output going to ls.save and piped to sort command.

## Metacharacters

*	any number of characters ( 0 or more)
?	any single character
[abc]	[ ] any character from the list
[a-c]	[ ] match any character from the list range
!	not any of the following characters (for example leftbox !abc right box)
;	command terminator used to string commands on a single line

---

&	command preceding and to be run in background mode
#	comment character
\	removes special meaning (no interpretation) of the following character removes special meaning (no interpretation) of character in quotes
"	interprets only \$, backquote, and \ characters between the quotes.
"	used to set variable to results of a command for example now= "date" sets the value of now to current results of the date command.
\$	preceding variable name indicates the value of the variable.

## Physical and Logical Storage

chlv	changes the characteristics of a logical volume.
chpv	changes the state of a physical volume within a volume group.
chvg	changes the characteristics of a volume group.
cplv	makes a copy of a logical volume.
exportvg	exports the definition of a volume group.
importvg	Imports the definition of a volume group
mklvcopy	makes logical partition copies for a logical volume
mkvg	makes a volume group.
reducevg	reduces the size of a volume group and deletes empty groups.
reorgvg	reorganizes the physical partition allocation for a volume group.
rmlv	removes a logical volume
syncvg	synchronizes logical partition copies
copyrawlv	copies the contents of one logical volume to another by directly reading and writing the logical volume devices. The destination logical volume must already exist and must be at least as large as the source.
getlvcb	returns the control block information for the specified logical volume.
getlvname	generates a logical volume name for a new logical volume. This is done using the name provided, or by using the default prefixes as defined in the Predefined ODM object classes.
getlvodm	gets logical volume data from the ODM and writes it to standard output.
getvgname	returns a volume group name. This is done either by using the name supplied by the user, or by using default prefixes as defined in the Predefined ODM.

lvgenmajor	generates a major number for the specified volume group. If a major number already exists for the volume group, that number is returned to standard out.
lvgenminor	generates a minor number for a logical volume or volume group.
lvrelmajor	releases a volume group's major number and removes the device file in the /dev directory.
lvrelminor	releases a logical volumes minor number and removes the /dev entries associated with the minor number.
putlvcb	writes the logical volume control block data into block 0 of the logical volume. The lvcb contains the attributes of the logical volume.
putlvodm	reads data from the command line and writes it to the appropriate ODM data class fields. This includes logical volume attributes, volume group attributes and physical volume attributes.
synclvodm	synchronizes data for the specified volume group or logical volume. The Logical Volume Manager is seen as correct when there are conflicts.
lchangelv	changes the attributes of a logical volume.
lcreatelv	creates an empty logical volume that belongs to the specified volume group.
ldeletelv	deletes a logical volume from its parent volume group.
lextendlv	extends or allocates additional logical partitions to a logical volume.
lquerylv	queries the attributes of a logical volume.
lreducelv	reduces the number of allocated logical partitions in a logical volume.
lresynclv	synchronizes all the mirrored logical partitions in the logical volume.
lchangevpv	changes the attributes of a physical volume.
ldeletepv	deletes a physical volume from its parent volume group.
linstallpv	installs or adds a physical volume to a volume group.
lquerypv	queries the attributes of a physical volume.
lresyncpv	synchronizes all mirrored partitions in a physical volume.
lcreatevg	creates a new physical volume and installs the first physical volume in the volume group.
lqueryvg	queries the attributes of a volume group.
lqueryvgs	queries the ID numbers of all volume groups in the system.
lvaryonvg	varies a volume group online. It can varyon in one of two ways: a) The volume group is varied on but the logical volumes cannot be opened. b) The volume group is varied on and the logical volumes are opened.



---

lvaryoffvg	varies a volume group offline. It is assumed that all Logical Volumes in the volume group must be closed before varyoff can complete.
lresynclp	synchronizes all physical partitions belonging to a logical partition.
lmigratepp	moves a physical partition to a specified physical volume.
chfs	changes file system attributes such as mount point, permissions, and size
compress	reduces the size of the specified file using the adaptive LZ algorithm
crfs	creates a file system within a previously created logical volume
extendlv	extends the size of a logical volume
extendvg	extends a volume group by adding a physical volume
fsck	checks for file system consistency, and allows interactive repair of file systems
fuser	lists the process numbers of local processes that use the files specified
lsattr	lists the attributes of the devices known to the system
lscfg	gives detailed information about the RS/6000 hardware configuration
lsdev	lists the devices known to the system
lsfs	displays characteristics of the specified file system such as mount points, permissions, and file system size
lslv	shows you information about a logical volume
lspv	shows you information about a physical volume in a volume group
lsvg	shows you information about the volume groups in your system
migratepv	used to move physical partitions from one physical volume to another
mkdev	configures a device
mkfs	makes a new file system on the specified device
mklv	creates a logical volume
mkvg	creates a volume group
mount	instructs the operating system to make the specified file system available for use from the specified point
quotaon	starts the disk quota monitor
rmdev	removes a device
rmlv	removes logical volumes from a volume group
rmlvcopy	removes copies from a logical volume
umount	unmounts a file system from its mount point

uncompress	restores files compressed by the compress command to their original size
umount	exactly the same function as the umount command
varyoffvg	deactivates a volume group so that it cannot be accessed
varyonvg	activates a volume group so that it can be accessed

## Variables

=	set a variable (for example d="day" sets the value of d to "day"). Can also set the variable to the results of a command by the ` character; for example now=date sets the value of now to the current result of the date command.
HOME	home directory
PATH	path to be checked
SHELL	shell to be used
TERM	terminal being used
PS1	primary prompt characters, usually \$ or #
PS2	secondary prompt characters, usually >
\$_	return code of the last command executed
set	displays current local variable settings
export	exports variable so that they are inherited by child processes
env	displays inherited variables
echo	echo a message (for example "echo HI" or "echo \$d"). Can turn off carriage returns with \c at the end of the message. Can print a blank line with \n at the end of the message.

## Tapes and Diskettes

dd	reads a file in, converts the data (if required), and copies the file out
fdformat	formats diskettes or read/write optical media disks
flcopy	copies information to and from diskettes
format	AIX command to format a diskette
backup	backs up individual files.

- i reads file names from standard input
- v list files as backed up;

for example "backup -iv -f/dev/rmto file1, file2"

- u backup file system at specified level;
    - for example “backup -level -u filesystem”
    - Can pipe list of files to be backed up into command; for example “find . -print | backup -ivf/dev/rmt0” where you are in directory to be backed up.
- mksysb creates an installable image of the root volume group
- restore restores commands from backup
- x restores files created with “backup -i”
  - v list files as restore
  - T list files stored of tape or diskette
  - r restores file systems created with “backup -level -u”;
    - for example “restore -xv -f/dev/rmt0”
- cpio copies to and from an I/O device. Destroys all data previously on tape or diskette. For input, must be able to place files in the same relative (or absolute) path name as when copied out (can determine path names with -it option). For input, if file exists, compares last modification date and keeps most recent (can override with -u option).
- o (output)
  - i (input),
  - t (table of contents)
  - v (verbose),
  - d (create needed directory for relative path names)
  - u (unconditional to override last modification date)
- for example “cpio -o > /dev/fd0”
- ```
“file1”
“file2”
“<Ctrl-d>”
```
- or “cpio -iv file1 < /dev/fd0”
- tapechk performs simple consistency checking for streaming tape drives
- tcopy copies information from one tape device to another
- tctl sends commands to a streaming tape device
- tar alternative utility to backup and restore files
- pax alternative utility to cpio and tar commands

## Transmitting

|             |                                                                                                                                                                                                                                        |
|-------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| mail        | send and receive mail. With user ID sends mail to userid. Without userid, displays your mail. When processing your mail, at the ? prompt for each mail item, you can:<br>d - delete s - append<br>q - quit enter - skip<br>m - forward |
| mailx       | upgrade of mail                                                                                                                                                                                                                        |
| uucp        | copy file to other UNIX systems (UNIX to UNIX copy)                                                                                                                                                                                    |
| uuto/uupick | send and retrieve files to public directories                                                                                                                                                                                          |
| uux         | execute on remote system (UNIX to UNIX execute)                                                                                                                                                                                        |

## System Administration

|            |                                                                                                    |
|------------|----------------------------------------------------------------------------------------------------|
| df         | display file system usage                                                                          |
| installp   | install program                                                                                    |
| kill (pid) | kill batch process with id or (pid) (find using ps);<br>kill -9 (PID) will absolutely kill process |
| mount      | associate logical volume to a directory; for example "mount device directory"                      |
| ps -ef     | shows process status (ps -ef)                                                                      |
| umount     | disassociate file system from directory                                                            |
| smit       | system management interface tool                                                                   |

## Miscellaneous

|         |                                                                        |
|---------|------------------------------------------------------------------------|
| banner  | displays banner                                                        |
| date    | displays current date and time                                         |
| newgrp  | change active groups                                                   |
| nice    | assigns lower priority to following command (for example "nice ps -f") |
| passwd  | modifies current password                                              |
| sleep n | sleep for n seconds                                                    |
| stty    | show and or set terminal settings                                      |
| touch   | create a zero length file(s)                                           |
| xinit   | initiate X-Windows                                                     |

|          |                                                                  |
|----------|------------------------------------------------------------------|
| wall     | sends message to all logged-in users.                            |
| who      | list users currently logged in (“who am i” identifies this user) |
| man,info | displays manual pages                                            |

## System Files

|                         |                                                                                                                  |
|-------------------------|------------------------------------------------------------------------------------------------------------------|
| /etc/group              | list of groups                                                                                                   |
| /etc/motd               | message of the day, displayed at login.                                                                          |
| /etc/passwd             | list of users and signon information. Password shown as !. Can prevent password checking by editing to remove !. |
| /etc/profile            | system-wide user profile executed at login. Can override variables by resetting in the user’s .profile file.     |
| /etc/security           | directory not accessible to normal users                                                                         |
| /etc/security/environ   | user environment settings                                                                                        |
| /etc/security/group     | group attributes                                                                                                 |
| /etc/security/limits    | user limits                                                                                                      |
| /etc/security/login.cfg | login settings                                                                                                   |
| /etc/security/passwd    | user passwords                                                                                                   |
| /etc/security/user      | user attributes, password restrictions                                                                           |

## Shell Programming Summary

### Variables

|            |                                                                                                                                                                                                                                                         |
|------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| var=string | set variable to equal string. (NO SPACES). Spaces must be enclosed by double quotes. Special characters in string must be enclosed by single quotes to prevent substitution. Piping ( ), redirection (<, >, >>), and “and” symbols are not interpreted. |
| \$var      | gives value of var in a compound                                                                                                                                                                                                                        |
| echo       | displays value of var, for example “echo \$var”                                                                                                                                                                                                         |
| HOME       | = home directory of user                                                                                                                                                                                                                                |
| MAIL       | = mail file name                                                                                                                                                                                                                                        |
| PS1        | = primary prompt characters, usually “\$” or “#”                                                                                                                                                                                                        |
| PS2        | = secondary prompt characters, usually “>”                                                                                                                                                                                                              |
| PATH       | = search path                                                                                                                                                                                                                                           |

|                 |                                                                         |
|-----------------|-------------------------------------------------------------------------|
| TERM            | = terminal type being used                                              |
| export          | exports variables to the environment                                    |
| env             | displays environment variables settings                                 |
| \${var:-string} | gives value of var in a command. If var is null, uses "string" instead. |
| \$1 \$2 \$3...  | positional parameters for variable passed into the shell script         |
| \$*             | used for all arguments passed into shell script                         |
| \$#             | number of arguments passed into shell script                            |
| \$0             | name of shell script                                                    |
| \$\$            | process id (pid)                                                        |
| \$?             | last return code from a command                                         |

## Commands

|              |                                                                                                                                    |
|--------------|------------------------------------------------------------------------------------------------------------------------------------|
| #            | comment designator                                                                                                                 |
| &&           | logical-and. Run command following && only if command preceding && succeeds (return code = 0).                                     |
|              | logical-or. Run command following    only if command preceding    fails (return code < > 0).                                       |
| exit n       | used to pass return code nl from shell script. Passed as variable \$? to parent shell                                              |
| expr         | arithmetic expressions<br>Syntax: "expr expression1 operator expression2"<br>operators: + - \* (multiply) / (divide) % (remainder) |
| for loop     | for n (or: for variable in \$*); for example:<br>do<br>command<br>done                                                             |
| if-then-else | if test expression<br>then command<br>elif test expression<br>then command<br>else<br>then command<br>fi                           |
| read         | read from standard input                                                                                                           |

|            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |
|------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| shift      | shifts arguments 1-9 one position to the left and decrements number of arguments                                                                                                                                                                                                                                                                                                                                                                                                                        |
| test       | used for conditional test, has two formats.<br>if test expression (for example "if test \$# -eq 2")<br>if [expression]<br>(for example "if [# -eq 2]") (spaces req'd)<br>integer operators:<br>-eq (=) -lt (<) -le (=<)<br>-ne (<>) -gt (>) -ge (=>)<br>string operators:<br>= != (not eq.) -z (zero length)<br>file status (for example -opt file1)<br>-f (ordinary file)<br>-r (readable by this process)<br>-w (writable by this process)<br>-x (executable by this process)<br>-s (non-zero length) |
| while loop | while test expression<br>do<br>command<br>done                                                                                                                                                                                                                                                                                                                                                                                                                                                          |

## Miscellaneous

|    |                                                                                                 |
|----|-------------------------------------------------------------------------------------------------|
| sh | execute shell script in the sh shellx (execute step by step - used for debugging shell scripts) |
|----|-------------------------------------------------------------------------------------------------|

## vi Editor

### Entering vi

|               |                                                                                                                                                       |
|---------------|-------------------------------------------------------------------------------------------------------------------------------------------------------|
| vi file       | edits the file named file                                                                                                                             |
| vi file file2 | edit files consecutively (via :n)                                                                                                                     |
| .exrc         | file that contains the vi profile                                                                                                                     |
| wm=nn         | sets wrap margin to nn<br>Can enter a file other than at first line by adding + (last line), +n (line n), or +/pattern (first occurrence of pattern). |
| vi -r         | lists saved files                                                                                                                                     |
| vi -r file    | recover file named file from crash                                                                                                                    |

:n next file in stack  
:set all show all options  
:set nu display line numbers (off when set nonu)  
:set list display control characters in file  
:set wm=n set wrap margin to n  
:set showmode sets display of "INPUT" when in input mode

## Read, Write, Exit

:w write buffer contents  
:w file2 write buffer contents to file2  
:w >> file2 write buffer contents to end of file2  
:q quit editing session  
:q! quit editing session and discard any changes  
:r file2 read file2 contents into buffer following current cursor  
:r! com read results of shell command "com" following current cursor  
:! exit shell command (filter through command)  
:wq or ZZ write and quit edit session

## Units of Measure

h, l character left, character right  
k or <Ctrl>p move cursor to character above cursor  
j or <Ctrl>n move cursor to character below cursor  
w, b word right, word left  
^, \$ beginning, end of current line  
<CR> or + beginning of next line  
- beginning of previous line  
G last line of buffer

## Cursor Movements

Can precede cursor movement commands (including cursor arrow) with number of times to repeat, for example 9--> moves right 9 characters.

0 move to first character in line



---

|         |                                                          |
|---------|----------------------------------------------------------|
| \$      | move to last character in line                           |
| ^       | move to first nonblank character in line                 |
| fx      | move right to character “x”                              |
| Fx      | move left to character “x”                               |
| tx      | move right to character preceding character “x”          |
| Tx      | move left to character preceding character “x”           |
| ;       | find next occurrence of “x” in same direction            |
| ,       | find next occurrence of “x” in opposite direction        |
| w       | tab word (nw = n tab word) (punctuation is a word)       |
| W       | tab word (nw = n tab word) (ignore punctuation)          |
| b       | backtab word (punctuation is a word)                     |
| B       | backtab word (ignore punctuation)                        |
| e       | tab to ending char. of next word (punctuation is a word) |
| E       | tab to ending char. of next word (ignore punctuation)    |
| (       | move to beginning of current sentence                    |
| )       | move to beginning of next sentence                       |
| {       | move to beginning of current paragraph                   |
| }       | move to beginning of next paragraph                      |
| H       | move to first line on screen                             |
| M       | move to middle line on screen                            |
| L       | move to last line on screen                              |
| <Ctrl>f | scroll forward 1 screen (3 lines overlap)                |
| <Ctrl>d | scroll forward 1/2 screen                                |
| <Ctrl>b | scroll backward 1 screen (0 line overlap)                |
| <Ctrl>u | scroll backward 1/2 screen                               |
| G       | go to last line in file                                  |
| nG      | go to line “n”                                           |
| <Ctrl>g | display current line number                              |

## Search and Replace

|          |                               |
|----------|-------------------------------|
| /pattern | search forward for “pattern”  |
| ?pattern | search backward for “pattern” |

|   |                                       |
|---|---------------------------------------|
| n | repeat find in the same direction     |
| N | repeat find in the opposite direction |

## Adding Text

|       |                                                     |
|-------|-----------------------------------------------------|
| a     | add text after the cursor (end with <esc>)          |
| A     | add text at end of current line (end with <esc>)    |
| i     | add text before the cursor (end with <esc>)         |
| I     | add text before first nonblank char in current line |
| o     | add line following current line                     |
| O     | add line before current line                        |
| <esc> | return to command mode                              |

## Deleting Text

|         |                                                            |
|---------|------------------------------------------------------------|
| <Ctrl>w | undo entry of current word                                 |
| @       | kill the insert on this line                               |
| x       | delete current character                                   |
| dw      | delete to end of current word (observe punctuation)        |
| dW      | delete to end of current word (ignore punctuation)         |
| dd      | delete current line                                        |
| d       | erase to end of line (same as d\$)                         |
| d)      | delete current sentence                                    |
| d}      | delete current paragraph                                   |
| dG      | delete current line thru end-of buffer                     |
| d^      | delete to the beginning of line                            |
| u       | undo last change command                                   |
| U       | restore current line to original state before modification |

## Replacing Text

|    |                                                         |
|----|---------------------------------------------------------|
| ra | replace current character with "a"                      |
| R  | replace all characters overtyped until <esc> is entered |
| s  | delete current character and append text until <esc>.   |

---

|         |                                                                                                                                                  |
|---------|--------------------------------------------------------------------------------------------------------------------------------------------------|
| s/s1/s2 | replace s1 with s2 (in the same line only)                                                                                                       |
| S       | delete all characters in the line and append text                                                                                                |
| cc      | replace all characters in the line (same as S)                                                                                                   |
| ncx     | delete “n” text objects of type “x”; w, b = words,) = sentences, } = paragraphs, \$ = end-of-line, ^M = beginning of line) and enter append mode |
| C       | replace all characters from cursor to end-of-line.                                                                                               |

## Moving Text

|      |                                                                                                                                                                                                       |
|------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| p    | paste last text deleted after cursor (xp will transpose 2 characters)                                                                                                                                 |
| P    | paste last text deleted before cursor                                                                                                                                                                 |
| nYx  | yank “n” text objects of type “x” (w, b = words,) = sentences, } = paragraphs, \$ = end-of-line, and no “x” indicates lines. Can then paste them with “p” command. Yank does not delete the original. |
| “ayy | can use named registers for moving, copying, cut/paste with “ayy for register a (use registers a-z). Can then paste them with “ap command.                                                            |

## Miscellaneous

|   |                               |
|---|-------------------------------|
| . | repeat last command           |
| J | join current line w/next line |



# Appendix B. Checkpoint Solutions

## Unit 1

1. What are the four major problem determination steps?

**Correct Answer**

Identify the problem.  
Talk to users.  
Collect system data.  
Resolve the problem.

2. Who should provide information about the problems?

**Correct Answer**

Always talk to the users about the problem to gather as much information as possible.

3. **T/F** If there is a problem with the software, it is necessary to get the next release of the product to resolve the problem.

**Correct Answer**

False. In most cases it is only necessary to apply fixes or upgrade microcode.

4. **T/F** Documentation can be viewed or downloaded from the IBM Web site.

**Correct Answer**

True.

## Unit 2

1. In which ODM class do you find the physical volume IDs of your disks?

**Correct Answer**

CuAt

2. What is the difference between state defined and available?

**Correct Answer**

When a device is defined there is an entry in ODM class CuDv. When a device is available, the device driver has been loaded. The device driver can be accessed by the entries in the /dev directory.

## Unit 3

1. **T/F** During the AIX boot process, the AIX kernel is loaded from the root file system.

**Correct Answer**

False, the AIX kernel is loaded from hd5.

2. Which RS/6000 models do not have a boot list for the service mode?

**Correct Answer**

Some PCI models.

3. How do you boot an AIX machine in maintenance mode?

**Correct Answer**

You need to boot from an AIX CD or mkysb tape.

4. Your machine keeps rebooting and repeating the POST. What could be reasons for this?

**Correct Answer**

Invalid boot list, corrupted boot logical volume, hardware failures of boot device.

## Unit 4

1. From where is rc.boot 3 run?

**Correct Answer**

From rootvg -/etc/inittab file

2. Your system stops booting with LED 557. In which rc.boot phase does the system stop? What can be the reasons for this problem?

**Correct Answer**

rc.boot 2

Corrupted BLV, corrupted JFS log, or rootvg unable to varyon.

3. Which ODM file is used by the cfgmgr during boot to configure the devices in the correct sequence?

**Correct Answer**

Config\_Rules

4. What does the line init:2:initdefault: in /etc/inittab mean?

**Correct Answer**

This line is used by the init process, to determine the initial run level (2=multiuser).

## Unit 5

1. **T/F:** All LVM information is stored in the ODM.

**Correct Answer**

False. There are many other AIX files and disk control blocks (like VGDA and LVCB).

2. **T/F:** You detect that a physical volume `hdisk1` that is contained in your `rootvg` is missing in the ODM. This problem can be fixed by exporting and importing the `rootvg`.

**Correct Answer**

False. Use script `rvgrecover` instead. This script creates complete new `rootvg` ODM entries.

3. **T/F:** The LVM supports RAID-5 without separate hardware.

**Correct Answer**

False. The LVM supports RAID-0 (striping) and RAID-1 (mirroring) without additional hardware.

## Unit 6

1. Although everything seems to be working fine, you detect error log entries for disk **hdisk0** in your **rootvg**. The disk is not mirrored to another disk. You decide to replace this disk. Which procedure would you use to migrate this disk?

**Correct Answer**

Procedure 2: Disk still working.

There are some additional steps necessary for `hd5` and the primary dump device `hd6`.

2. You detect an unrecoverable disk failure in volume group **datavg**. This volume group consists of two disks that are completely mirrored. Because of the disk failure you are not able to vary on **datavg**. How do you recover from this situation?

**Correct Answer**

Forced varyon:`varyonvg -f datavg`

Use Procedure 1 for mirrored disks.

3. After a disk replacement you recognize that a disk has been removed from the system but not from the volume group. How do you fix this problem?

**Correct Answer**

Use PVID instead of disk name:

```
reducevg vg_name PVID
```

## Unit 7

1. **T/F:** After restoring a `mksysb` image all passwords are restored as well.

**Correct Answer**

True

2. The mkszfile will create a file named

- a. /bosinst.data
- b. /image.data
- c. /vgname.data

**Correct Answer**

b

3. Which two alternate disk installation techniques are available?

**Correct Answer**

Installing a mksysb on another disk

Cloning the rootvg to another disk

4. What are the commands to back up and restore a non-rootvg volume group?

**Correct Answer**

savevg

restvg

5. If you want to shrink one file system in a volume group myvg, which file must be changed before backing up the user volume group?

**Correct Answer**

The control file is:

/tmp/vgdata/myvg/myvg.data

6. How many copies should you have before performing an online JFS or JFS2 backup?

**Correct Answer**

3

## Unit 8

1. Which command generates error reports?

**Correct Answer**

errpt

errpt -a

2. Which type of disk error indicates bad blocks?

**Correct Answer**

DISK\_ERR4



3. What do the following commands do?

**errclear**

**errlogger**

**Correct Answer**

Clears entries from the error log.

Used by root to add entries into the error log.

4. What does the following line in /etc/syslog.conf indicate:

**\*.debug errlog**

**Correct Answer**

All syslogd messages are directed to the error log

5. What does the descriptor **en\_method** in **errnotify** indicate?

**Correct Answer**

Specifies a program or a command to be run when an error matching the selection criteria is logged.

## Unit 9

1. **T/F:** The diag command is supported on all RS/6000 models.

**Correct Answer**

False

2. What diagnostic modes are available on a RS/6000?

**Correct Answer**

Maintenance, concurrent and stand-alone modes.

3. How can you diagnose a communication adapter that is used during normal system operation?

**Correct Answer**

Either in maintenance or stand-alone mode.

## Unit 10

1. What is the default primary dump device? Where do you find the dump file after reboot?

**Correct Answer**

Default primary dump device: /dev/hd6

Dump file (default): /var/adm/ras/vmcore.x

2. How do you turn on dump compression?

**Correct Answer**

sysdumpdev -C

3. How do you start a dump from an attached LFT terminal?

**Correct Answer**

You have to specify Always Allow Dump in smit, or you must execute the command sysdumpdev -k, then press <ctrl><alt><num-pad-1>.

4. If the copy directory is too small, will the dump which is copied during the reboot of the system, be lost.

**Correct Answer**

No. A special menu is shown during reboot. From this menu you can copy the dump to portable media.

5. Which command should you execute before sending a dump to IBM?

**Correct Answer**

The snap command.

## Unit 11

1. What command can be executed to identify CPU-intensive programs?

**Correct Answer**

ps aux and tprof

2. What command can be executed to start processes with a lower priority?

**Correct Answer**

The nice command

3. What command can you use to check paging I/O?

**Correct Answer**

vmstat

4. **T/F:** The higher the PRI value, the higher the priority of a process.

**Correct Answer**

False

## Unit 12

1. **T/F:** Any programs specified as “auth1” must return a zero in order for the user to log in.

**Correct Answer**

True

2. How would you specify that all members of the security group had rwx access to a particular file except for John?

**Correct Answer**

Using ACLs  
extended permission  
enabled  
permit rwx g:security  
deny rwx u:john

3. In which file must you specify the full path name of the program that is to be used as part of the authentication process when a user logs in?

**Correct Answer**

/usr/lib/security/methods.cfg

4. Name the two modes that tcbck supports.

**Correct Answer**

Check mode  
Update mode

5. When you execute **<ctrl-x ctrl-r>** at a login prompt and you obtain the **tsh** prompt, what does this indicate?

**Correct Answer**

This indicates that there is someone already logged in running a fake getty program -a Trojan Horse!

6. **T/F:** The system administrator must manually mark commands as trusted, which will automatically add the commands to the **sysck.cfg** file.

**Correct Answer**

False. The system administrator has to also remember to add the commands to the **sysck.cfg** file using the **tcbck -a** command.

7. **T/F:** When the **tcbck -p tree** command is executed, all errors are reported and you get a prompt asking if the error should be fixed.

**Correct Answer**

False. Option -p indicates fixing and no reporting. A very dangerous option!



## Appendix C. RS/6000 Three-Digit Display Values

This appendix is an extract out of the *AIX 4.3 Messages Guide and Reference*.

### 0c0 - 0cc

|     |                                                                                                          |
|-----|----------------------------------------------------------------------------------------------------------|
| 0c0 | A user-requested dump completed successfully.                                                            |
| 0c1 | An I/O error occurred during the dump.                                                                   |
| 0c2 | A user-requested dump is in progress. Wait at least one minute for the dump to complete.                 |
| 0c4 | The dump ran out of space. Partial dump is available.                                                    |
| 0c5 | The dump failed due to an internal failure. A partial dump may exist.                                    |
| 0c7 | Progress indicator. Remote dump is in progress.                                                          |
| 0c8 | The dump device is disabled. No dump device configured.                                                  |
| 0c9 | A system-initiated dump has started. Wait at least one minute for the dump to complete.                  |
| 0cc | (AIX 4.2.1 and later) Error occurred writing to the primary dump device. Switched over to the secondary. |

### 100 - 195

|     |                                                                                                 |
|-----|-------------------------------------------------------------------------------------------------|
| 100 | Progress indicator. BIST completed successfully.                                                |
| 101 | Progress indicator. Initial BIST started following system reset.                                |
| 102 | Progress indicator. BIST started following power-on reset.                                      |
| 103 | BIST could not determine the system model number.                                               |
| 104 | BIST could not find the common on-chip processor bus address.                                   |
| 105 | BIST could not read from the on-chip sequencer EPROM.                                           |
| 106 | BIST detected a module failure.                                                                 |
| 111 | On-chip sequencer stopped. BIST detected a module error.                                        |
| 112 | Checkstop occurred during BIST and checkstop results could not be logged out.                   |
| 113 | The BIST checkstop count equals 3, that means three unsuccessful system restarts. System halts. |
| 120 | Progress indicator. BIST started CRC check on the EPROM.                                        |
| 121 | BIST detected a bad CRC on the on-chip sequencer EPROM.                                         |

- 122 Progress indicator. BIST started CRC check on the EPROM.
- 123 BIST detected a bad CRC on the on-chip sequencer NVRAM.
- 124 Progress indicator. BIST started CRC check on the on-chip sequencer NVRAM.
- 125 BIST detected a bad CRC on the time-of-day NVRAM.
- 126 Progress indicator. BIST started CRC check on the time-of-day NVRAM.
- 127 BIST detected a bad CRC on the EPROM.
- 130 Progress indicator. BIST presence test started.
- 140 BIST was unsuccessful. System halts.
- 142 BIST was unsuccessful. System halts.
- 143 Invalid memory configuration.
- 144 BIST was unsuccessful. System halts.
- 151 Progress indicator. BIST started.
- 152 Progress indicator. BIST started direct-current logic self-test (DCLST) code.
- 153 Progress indicator. BIST started.
- 154 Progress indicator. BIST started array self-test (AST) test code.
- 160 BIST detected a missing Early Power-Off Warning (EPOW) connector.
- 161 The Bump quick I/O tests failed.
- 162 The JTAG tests failed.
- 164 BIST encountered an error while reading low NVRAM.
- 165 BIST encountered an error while writing low NVRAM.
- 166 BIST encountered an error while reading high NVRAM.
- 167 BIST encountered an error while writing high NVRAM.
- 168 BIST encountered an error while reading the serial input/output register.
- 169 BIST encountered an error while writing the serial input/output register.
- 180 Progress indicator. BIST checkstop logout in progress.
- 182 BIST COP bus is not responding.
- 185 Checktop occurred during BIST.
- 186 System logic-generated checkstop (Model 250 only).

- 187 BIST was unable to identify the chip release level in the checkstop logout data.
- 195 Progress indicator. BIST checkstop logout completed.

## 200 - 299, 2e6-2e7

- 200 Key mode switch is in the secure position.
- 201 Checkstop occurred during system restart. If a 299 LED was shown before, recreate the boot logical volume (bosboot).
- 202 Unexpected machine check interrupt. System halts.
- 203 Unexpected data storage interrupt. System halts.
- 204 Unexpected instruction storage interrupt. System halts.
- 205 Unexpected external interrupt. System halts.
- 206 Unexpected alignment interrupt. System halts.
- 207 Unexpected program interrupt. System halts.
- 208 machine check due to an L2 uncorrectable ECC. System halts.
- 209 Reserved. System halts.
- 210 Unexpected switched virtual circuit (SVC) 1000 interrupt. System halts.
- 211 IPL ROM CRC miscompare occurred during system restart. System halts.
- 212 POST found processor to be bad. System halts.
- 213 POST failed. No good memory could be detected. System halts.
- 214 I/O planar failure has been detected. The power status register, the time-of-day clock, or NVRAM on the I/O planar failed. System halts.
- 215 Progress indicator. Level of voltage supplied to the system is too low to continue a system restart.
- 216 Progress indicator. IPL ROM code is being uncompressed into memory for execution.
- 217 Progress indicator. System has encountered the end of the boot devices list. System continues to loop through the boot devices list.
- 218 Progress indicator. POST is testing for 1MB of good memory.
- 219 Progress indicator. POST bit map is being generated.
- 21c L2 cache not detected as part of systems configuration (when LED persists for 2 seconds).
- 220 Progress indicator. IPL control block is being initialized.

- 221 NVRAM CRC miscompare occurred while loading the operating system with the key mode switch in Normal position. System halts.
- 222 Progress indicator. Attempting a Normal-mode system restart from the standard I/O planar-attached devices. System retries.
- 223 Progress indicator. Attempting a Normal-mode system restart from the SCSI-attached devices specified in the NVRAM list.
- 224 Progress indicator. Attempting a Normal-mode system restart from the 9333 High-Performance Disk-Drive Subsystem.
- 225 Progress indicator. Attempting a Normal-mode system restart from the bus-attached internal disk.
- 226 Progress indicator. Attempting a Normal-mode system restart from Ethernet.
- 227 Progress indicator. Attempting a Normal-mode system restart from Token-Ring.
- 228 Progress indicator. Attempting a Normal-mode system restart using the expansion code devices list, but cannot restart from any of the devices in the list.
- 229 Progress indicator. Attempting a Normal-mode system restart from devices in NVRAM boot devices list, but cannot restart from any of the devices in the list. System retries.
- 22c Progress indicator. Attempting a Normal-mode IPL from FDDI specified in the NVRAM device list.
- 230 Progress indicator. Attempting a Normal-mode system restart from Family 2 Feature ROM specified in the IPL ROM default devices list.
- 231 Progress indicator. Attempting a Normal-mode system restart from Ethernet specified by selection from ROM menus.
- 232 Progress indicator. Attempting a Normal-mode system restart from the standard I/O planar-attached devices specified in the IPL ROM default device list.
- 233 Progress indicator. Attempting a Normal-mode system restart from the SCSI-attached devices specified in the IPL ROM default device list.
- 234 Progress indicator. Attempting a Normal-mode system restart from the 9333 High-Performance Disk Drive Subsystem specified in the IPL ROM default device list.
- 234 Progress indicator. Attempting a Normal-mode system restart from the 9333 High-Performance Disk Drive Subsystem specified in the IPL ROM default device list.



- 235 Progress indicator. Attempting a Normal-mode system restart from the bus-attached internal disk specified in the IPL ROM default device list.
- 236 Progress indicator. Attempting a Normal-mode system restart from the ethernet specified in the IPL ROM default device list.
- 237 Progress indicator. Attempting a Normal-mode system restart from the token-ring specified in the IPL ROM default device list.
- 238 Progress indicator. Attempting a Normal-mode system restart from the token-ring specified by selection from ROM menus.
- 239 Progress indicator. A Normal-mode menu selection failed to boot.
- 23c Progress indicator. Attempting a Normal-mode IPL form FDDI in IPL ROM device list.
- 240 Progress indicator. Attempting a Service-mode system restart from the Family 2 Feature ROM specified in the NVRAM boot devices list.
- 241 Attempting a Normal-mode system restart from devices specified in NVRAM boot list.
- 242 Progress indicator. Attempting a Service-mode system restart from the standard I/O planar-attached devices specified in the NVRAM boot devices list.
- 243 Progress indicator. Attempting a Service-mode system restart from the SCSI-attached devices specified in the NVRAM boot devices list.
- 244 Progress indicator. Attempting a Service-mode system restart from the 9333 High-Performance Disk Drive Subsystem specified in the NVRAM boot devices list.
- 245 Progress indicator. Attempting a Service-mode system restart from the bus-attached internal disk specified in the NVRAM boot devices list.
- 246 Progress indicator. Attempting a Service-mode system restart from the Ethernet specified in the NVRAM boot devices list.
- 247 Progress indicator. Attempting a Service-mode system restart from the Token-Ring specified in the NVRAM boot devices list.
- 248 Progress indicator. Attempting a Service-mode system restart using the expansion code specified in the NVRAM boot devices list.
- 249 Progress indicator. Attempting a Service-mode system restart from devices in NVRAM boot devices list, but cannot restart from any of the devices in the list.
- 250 Progress indicator. Attempting a Service-mode system restart from the Family 2 Feature ROM specified in the IPL ROM default devices list.
- 251 Progress indicator. Attempting a Service-mode system restart from Ethernet by selection from ROM menus.

- 252 Progress indicator. Attempting a Service-mode system restart from the standard I/O planar-attached devices specified in the IPL ROM default devices list.
- 253 Progress indicator. Attempting a Service-mode system restart from the SCSI-attached devices specified in the IPL ROM default devices list.
- 254 Progress indicator. Attempting a Service-mode system restart from the 9333 High-Performance Subsystem devices specified in the IPL ROM default devices list.
- 255 Progress indicator. Attempting a Service-mode system restart from the bus-attached internal disk specified in the IPL ROM default devices list.
- 256 Progress indicator. Attempting a Service-mode system restart from the Ethernet specified in the IPL ROM default devices list.
- 257 Progress indicator. Attempting a Service-mode system restart from the Token-Ring specified in the IPL ROM default devices list.
- 258 Progress indicator. Attempting a Service-mode system restart from the Token-Ring specified by selection from ROM menus.
- 259 Progress indicator. Attempting a Service-mode system restart from FDDI specified by the operator.
- 260 Progress indicator. Menus are being displayed on the local display or terminal connected to your system. The system waits for input from the terminal.
- 261 No supported local system display adapter was found. The system waits for a response from an asynchronous terminal on serial port 1.
- 262 No local system keyboard was found.
- 263 Progress indicator. Attempting a Normal-mode system restart from the Family 2 Feature ROM specified in the NVRAM boot devices list.
- 269 Progress indicator. Cannot boot system, end of boot list reached.
- 270 Progress indicator. Ethernet/FDX 10 Mbps MC adapter POST is running.
- 271 Progress indicator. Mouse and mouse port POST is running.
- 272 Progress indicator. Tablet port POST is running.
- 276 Progress indicator. A 10/100 Mbps Ethernet MC adapter POST is running.
- 277 Progress indicator. Auto Token-Ring LAN streamer MC 32 adapter POST is running.
- 278 Progress indicator. Video ROM scan POST is running.
- 279 Progress indicator. FDDI POST is running

---

|     |                                                                                     |
|-----|-------------------------------------------------------------------------------------|
| 280 | Progress indicator. 3Com Ethernet POST is running.                                  |
| 281 | Progress indicator. Keyboard POST is running.                                       |
| 282 | Progress indicator. Parallel port POST is running.                                  |
| 283 | Progress indicator. Serial port POST is running.                                    |
| 284 | Progress indicator. POWER Gt1 graphics adapter POST is running.                     |
| 285 | Progress indicator. POWER Gt3 graphics adapter POST is running.                     |
| 286 | Progress indicator. Token-Ring adapter POST is running.                             |
| 287 | Progress indicator. Ethernet adapter POST is running.                               |
| 288 | Progress indicator. Adapter slot cards are being queried.                           |
| 289 | Progress indicator. POWER Gt0 graphics adapter POST is running.                     |
| 290 | Progress indicator. I/O planar test started.                                        |
| 291 | Progress indicator. Standard I/O planar POST is running.                            |
| 292 | Progress indicator. SCSI POST is running.                                           |
| 293 | Progress indicator. Bus-attached internal disk POST is running.                     |
| 294 | Progress indicator. TCW SIMM in slot J is bad.                                      |
| 295 | Progress indicator. Color Graphics Display POST is running.                         |
| 296 | Progress indicator. Family 2 Feature ROM POST is running.                           |
| 297 | Progress indicator. System model number could not be determined. System halts.      |
| 298 | Progress indicator. Attempting a warm system restart.                               |
| 299 | Progress indicator. IPL ROM passed control to loaded code.                          |
| 2e6 | Progress indicator. A PCI Ultra/Wide differential SCSI adapter is being configured. |
| 2e7 | An undetermined PCI SCSI adapter is being configured.                               |

## 500 - 599, 5c0 - 5c6

|     |                                                 |
|-----|-------------------------------------------------|
| 500 | Progress indicator. Querying standard I/O slot. |
| 501 | Progress indicator. Querying card in slot 1.    |
| 502 | Progress indicator. Querying card in slot 2.    |
| 503 | Progress indicator. Querying card in slot 3.    |
| 504 | Progress indicator. Querying card in slot 4.    |
| 505 | Progress indicator. Querying card in slot 5.    |

- 506 Progress indicator. Querying card in slot 6.
- 507 Progress indicator. Querying card in slot 7.
- 508 Progress indicator. Querying card in slot 8.
- 510 Progress indicator. Starting device configuration.
- 511 Progress indicator. Device configuration completed.
- 512 Progress indicator. Restoring device configuration from media.
- 513 Progress indicator. Restoring BOS installation files from media.
- 516 Progress indicator. Contacting server during network boot.
- 517 Progress indicator. The / (root) and /usr file systems are being mounted.
- 518 Mount of the /usr file system was not successful. System Halts.
- 520 Progress indicator. BOS configuration is running.
- 521 The /etc/inittab file has been incorrectly modified or is damaged. The configuration manager was started from the /etc/inittab file with invalid options. System halts.
- 522 The /etc/inittab file has been incorrectly modified or is damaged. The configuration manager was started from the /etc/inittab file with conflicting options. System halts.
- 523 The /etc/objrepos file is missing or inaccessible.
- 524 The /etc/objrepos/Config\_Rules file is missing or inaccessible.
- 525 The /etc/objrepos/CuDv file is missing or inaccessible.
- 526 The /etc/objrepos/CuDvDr file is missing or inaccessible.
- 527 You cannot run Phase 1 at this point. The /sbin/rc.boot file has probably been incorrectly modified or is damaged.
- 528 The /etc/objrepos/Config\_Rules file has been incorrectly modified or is damaged, or a program specified in the file is missing.
- 529 There is a problem with the device containing the ODM database or the root file system is full.
- 530 The savebase command was unable to save information about the base customized devices onto the boot device during Phase 1 of system boot. System halts.
- 531 The /usr/lib/objrepos/PdAt file is missing or inaccessible. System halts.
- 532 There is not enough memory for the configuration manager to continue. System halts.

- 533 The /usr/lib/objrepos/PdDv file has been incorrectly modified or is damaged, or a program specified in the file is missing.
- 534 The configuration manager is unable to acquire a database lock. System halts.
- 535 A HIPPI diagnostics interface driver is being configured.
- 536 The /etc/objrepos/Config\_Rules file has been incorrectly modified or is damaged. System halts.
- 537 The /etc/objrepos/Config\_Rules file has been incorrectly modified or is damaged. System halts.
- 538 Progress indicator. The configuration manager is passing control to a configuration method.
- 539 Progress indicator. The configuration method has ended and control has returned to the configuration manager.
- 540 Progress indicator. Configuring child of IEEE-1284 parallel port.
- 544 Progress indicator. An ECP peripheral configure method is executing.
- 545 Progress indicator. A parallel port ECP device driver is being configured.
- 546 IPL cannot continue due to an error in the customized database.
- 547 Rebooting after error recovery (LED 546 precedes this LED).
- 548 restbase failure.
- 549 Console could not be configured for the "Copy a System Dump" menu.
- 550 Progress indicator. ATM LAN emulation device driver is being configured.
- 551 Progress indicator. A varyon operation of the rootvg is in progress.
- 552 The ipl\_varyon command failed with a return code not equal to 4, 7, 8 or 9 (ODM or malloc failure). System is unable to vary on the rootvg.
- 553 The /etc/inittab file has been incorrectly modified or is damaged. Phase 1 boot is completed and the init command started.
- 554 The IPL device could not be opened or a read failed (hardware not configured or missing).
- 555 The fsck -fp /dev/hd4 command on the root file system failed with a non-zero return code.
- 556 LVM subroutine error from ipl\_varyon.
- 557 The root file system could not be mounted. The problem is usually due to bad information on the log logical volume (/dev/hd8) or the boot logical volume (hd5) has been damaged.

- 558 Not enough memory is available to continue system restart.
- 559 Less than 2 MB of good memory are left for loading the AIX kernel. System halts.
- 560 Unsupported monitor is attached to the display adapter.
- 561 Progress indicator. The TMSSA device is being identified or configured.
- 565 Configuring the MWAVE subsystem.
- 566 Progress indicator. Configuring Namkan twinaxx commo card.
- 567 Progress indicator. Configuring High-Performance Parallel Interface (HIPPI) device driver (fpdev).
- 568 Progress indicator. Configuring High-Performance Parallel Interface (HIPPI) device driver (fhip).
- 569 Progress indicator. FCS SCSI protocol device is being configured.
- 570 Progress indicator. A SCSI protocol device is being configured.
- 571 HIPPI common functions driver is being configured.
- 572 HIPPI IPI-3 master mode driver is being configured.
- 573 HIPPI IPI-3 slave mode driver is being configured.
- 574 HIPPI IPI-3 user-level interface is being configured.
- 575 A 9570 disk-array driver is being configured.
- 576 Generic async device driver is being configured.
- 577 Generic SCSI device driver is being configured.
- 578 Generic common device driver is being configured.
- 579 Device driver is being configured for a generic device.
- 580 Progress indicator. A HIPPI-LE interface (IP) layer is being configured.
- 581 Progress indicator. TCP/IP is being configured. The configuration method for TCP/IP is being run.
- 582 Progress indicator. Token-Ring data link control (DLC) is being configured.
- 583 Progress indicator. Ethernet data link control (DLC) is being configured.
- 584 Progress indicator. IEEE Ethernet (802.3) data link control (DLC) is being configured.
- 585 Progress indicator. SDLC data link control (DLC) is being configured.
- 586 Progress indicator. X.25 data link control (DLC) is being configured.

---

|     |                                                                                  |
|-----|----------------------------------------------------------------------------------|
| 587 | Progress indicator. Netbios is being configured.                                 |
| 588 | Progress indicator. Bisync read-write (BSCRW) is being configured.               |
| 589 | Progress indicator. SCSI target mode device is being configured.                 |
| 590 | Progress indicator. Diskless remote paging device is being configured.           |
| 591 | Progress indicator. Logical Volume Manager device driver is being configured.    |
| 592 | Progress indicator. An HFT device is being configured.                           |
| 593 | Progress indicator. SNA device driver is being configured.                       |
| 594 | Progress indicator. Asynchronous I/O is being defined or configured.             |
| 595 | Progress indicator. X.31 pseudo device is being configured.                      |
| 596 | Progress indicator. SNA DLC/LAPE pseudo device is being configured.              |
| 597 | Progress indicator. Outboard communication server (OCS) is being configured.     |
| 598 | Progress indicator. OCS hosts is being configured during system reboot.          |
| 599 | Progress indicator. FDDI data link control (DLC) is being configured.            |
| 5c0 | Progress indicator. Streams-based hardware driver being configured.              |
| 5c1 | Progress indicator. Streams-based X.25 protocol stack being configured.          |
| 5c2 | Progress indicator. Streams-based X.25 COMIO emulator driver being configured.   |
| 5c3 | Progress indicator. Streams-based X.25 TCP/IP interface driver being configured. |
| 5c4 | Progress indicator. FCS adapter device driver being configured.                  |
| 5c5 | Progress indicator. SCB network device driver for FCS is being configured.       |
| 5c6 | Progress indicator. AIX SNA channel being configured.                            |

## **c00 - c99**

|     |                                              |
|-----|----------------------------------------------|
| c00 | AIX Install/Maintenance loaded successfully. |
| c01 | Insert the AIX Install/Maintenance diskette. |
| c02 | Diskettes inserted out of sequence.          |
| c03 | Wrong diskette inserted.                     |
| c04 | Irrecoverable error occurred.                |



|     |                                                                                                                                                                                                 |
|-----|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| c05 | Diskette error occurred.                                                                                                                                                                        |
| c06 | The rc.boot script is unable to determine the type of boot.                                                                                                                                     |
| c07 | Insert next diskette.                                                                                                                                                                           |
| c08 | RAM file system started incorrectly.                                                                                                                                                            |
| c09 | Progress indicator. Writing to or reading from diskette.                                                                                                                                        |
| c10 | Platform-specific bootinfo is not in boot image.                                                                                                                                                |
| c20 | Unexpected system halt occurred. System is configured to enter the kernel debug program instead of performing a system dump. Enter bosboot -D for information about kernel debugger enablement. |
| c21 | The if config command was unable to configure the network for the client network host.                                                                                                          |
| c25 | Client did not mount remote mini root during network install.                                                                                                                                   |
| c26 | Client did not mount the /usr file system during the network boot.                                                                                                                              |
| c29 | System was unable to configure the network device.                                                                                                                                              |
| c31 | If a console has not been configured, the system pauses with this value and then displays instructions for choosing a console.                                                                  |
| c32 | Progress indicator. Console is a high-function terminal.                                                                                                                                        |
| c33 | Progress indicator. Console is a tty.                                                                                                                                                           |
| c34 | Progress indicator. Console is a file.                                                                                                                                                          |
| c40 | Extracting data files from media.                                                                                                                                                               |
| c41 | Could not determine the boot type or device.                                                                                                                                                    |
| c42 | Extracting data files from diskette.                                                                                                                                                            |
| c43 | Could not access the boot or installation tape.                                                                                                                                                 |
| c44 | Initializing installation database with target disk information.                                                                                                                                |
| c45 | Cannot configure the console. The cfgcon command failed.                                                                                                                                        |
| c46 | Normal installation processing.                                                                                                                                                                 |
| c47 | Could not create a PVID on a disk. The chgdisk command failed.                                                                                                                                  |
| c48 | Prompting you for input. BosMenus is being run.                                                                                                                                                 |
| c49 | Could not create or form the JFS log.                                                                                                                                                           |
| c50 | Creating rootvg on target disk.                                                                                                                                                                 |
| c51 | No paging devices were found.                                                                                                                                                                   |
| c52 | Changing from RAM environment to disk environment.                                                                                                                                              |



|     |                                                                                                       |
|-----|-------------------------------------------------------------------------------------------------------|
| c53 | Not enough space in /tmp to do a preservation installation. Make /tmp larger.                         |
| c54 | Installing either BOS or additional packages.                                                         |
| c55 | Could not remove the specified logical volume in a preservation installation.                         |
| c56 | Running user-defined customization.                                                                   |
| c57 | Failure to restore BOS.                                                                               |
| c58 | Displaying message to turn the key.                                                                   |
| c59 | Could not copy either device special files, device ODM, or volume group information from RAM to disk. |
| c61 | Failed to create the boot image.                                                                      |
| c70 | Problem mounting diagnostic CD-ROM disk in stand-alone mode.                                          |
| c99 | Progress indicator. The diagnostic programs have completed.                                           |



# Appendix D. PCI Firmware Checkpoints and Error Codes

This appendix shows firmware checkpoints and error codes for a 43P Model 140.

## Firmware Checkpoints

|     |                                                                                                                              |
|-----|------------------------------------------------------------------------------------------------------------------------------|
| F01 | Performing system memory test                                                                                                |
| F05 | Transfer control to operating system (normal boot)                                                                           |
| F22 | No memory detected.<br><b>Note:</b> The disk drive light is on.                                                              |
| F2C | Processor card mismatch                                                                                                      |
| F4D | Loading boot image                                                                                                           |
| F4F | NVRAM initialization                                                                                                         |
| F51 | Probing primary PCI bus                                                                                                      |
| F52 | Probing for adapter FCODE, evaluate if present                                                                               |
| F55 | Probing PCI bridge secondary bus                                                                                             |
| F5B | Transferring control to operating system (service boot)                                                                      |
| F5F | Probing for adapter FCODE, evaluate if present                                                                               |
| F74 | Establishing host connection                                                                                                 |
| F75 | Bootp request                                                                                                                |
| F9E | Real-time clock (RTC) initialization                                                                                         |
| FDC | Dynamic console selection                                                                                                    |
| FDD | Processor exception                                                                                                          |
| FDE | Alternating pattern of FDE and FAD. Indicates a processor execution has been detected.                                       |
| FEA | Firmware flash corrupted, load from diskette                                                                                 |
| FEB | Firmware flush corrupted, load from diskette                                                                                 |
| FF2 | Power-On Password Prompt                                                                                                     |
| FF3 | Privileged-Access Password Prompt                                                                                            |
| FFB | SCSI bus initialization                                                                                                      |
| FFD | The operator panel alternates between the code FFD and another Fxx code, where Fxx is the point at which the error occurred. |

## Firmware Error Codes

|          |                                                       |
|----------|-------------------------------------------------------|
| 20100xxx | Power Supply                                          |
| 20A80xxx | Remote initial program load (RIPL) error              |
| 20D00xxx | Unknown/Unrecognized device                           |
| 20E00000 | Power on password entry error                         |
| 20E00001 | Privileged-access password entry error                |
| 20E00002 | Privileged-access password jumper not enabled         |
| 20E00003 | Power on password must be set for unattended mode     |
| 20E00004 | Battery drained or needs replacement                  |
| 20E00005 | EEPROM locked. Turn off, then turn on the system unit |
| 20E00008 | CMOS corrupted. Replace battery                       |
| 20E00009 | Invalid password entered. System locked               |
| 20E0000A | EEPROM lock problem. Check jumper position            |
| 20E0000B | EEPROM write problem. Turn off, turn on system unit   |
| 20E0000C | EEPROM read problem. Turn off, turn on system unit    |
| 20E00017 | Cold boot needed for password entry                   |
| 20EE0003 | SMS: Invalid RIPL address (3 dots needed)             |
| 20EE0004 | SMS: Invalid RIPL address                             |
| 20EE0005 | SMS: Invalid portion of RIPL IP address (> 255)       |
| 20EE0006 | SMS: No SCSI controllers present                      |
| 20EE0007 | Console selection: Keyboard not found                 |
| 20EE0008 | No configurable adapters found in the system          |
| 21A00xxx | SCSI disk driver errors                               |
| 21E00xxx | SCSI tape error                                       |
| 21ED0xxx | SCSI changer error                                    |
| 21EE0xxx | Other SCSI device type                                |
| 21F00xxx | SCSI CD-ROM error                                     |
| 21F20xxx | SCSI Read/Write Optical error                         |
| 25010xxx | Flash update                                          |
| 25A0xxy0 | Cache: L2 controller failure                          |
| 25A1xxy0 | Cache: L2 SRAM failure                                |

|          |                                            |
|----------|--------------------------------------------|
| 25A80xxx | NVRAM error                                |
| 25AA0xxx | EEPROM error                               |
| 25Cyyxxx | Memory error (DIMM fails or invalid)       |
| 28030xxx | Real-time clock (RTC) error                |
| 29000002 | Keyboard/Mouse controller failed self-test |
| 29A00003 | Keyboard not detected                      |
| 29A00004 | Mouse not detected                         |
| 2B2xxyrr | Processor or CPU error                     |



## Appendix E. Location Codes

The location code is a way of identifying physical devices in a RS/6000 system. It shows a path from the system unit (or a CPU drawer) through the adapter to the device itself.

### PCI Location Codes

| Device Name:                                    | Location Code: |
|-------------------------------------------------|----------------|
| Processor                                       | 00-00          |
| Motherboard                                     | 00-00          |
| PCI bus                                         | 00-00          |
| Diskette adapter                                | 01-A0          |
| Diskette drive                                  | 01-A0-00-00    |
| Parallel Port Adapter                           | 01-B0          |
| Parallel Printer                                | 01-B0-00-00    |
| Serial Port 1                                   | 01-C0          |
| Terminal attached to port 1                     | 01-C0-00-00    |
| Keyboard adapter                                | 01-E0          |
| PS2-Keyboard                                    | 01-E0-00-00    |
| ISA bus                                         | 04-A0          |
| Second PCI bus                                  | 04-D0          |
| On-board SCSI controller                        | 04-C0          |
| CD-ROM attached to on-board SCSI controller     | 04-C0-00-4,0   |
| Disk drive attached to on-board SCSI controller | 04-C0-00-8,0   |
| SCSI controller, not on-board                   | 04-01          |
| Graphics adapter                                | 04-02          |
| Token-ring Adapter, not on-board                | 04-03          |

The general format of a PCI location code is:

#### **AB-CD-EF-GH**

AB = Type of bus  
 CD = Slot  
 EF = Connector  
 GH = Port

The first two characters (AB) specify the type of bus where the device is located.

- **00** specifies a device that is located on the **processor bus**, for example the processor, a memory card or the L2 cache.
- **01** specifies a device that is attached to an ISA-bus. The term ISA (ISA = Industrial Standard Architecture) comes from the PC world and has a transfer rate of 8 MByte per second. Those devices are attached to the ISA-bus which does not need a high-speed connection, for example terminals or printers.
- **04** specifies a device that is attached to a PCI-bus. All location codes 04-A0, 04-B0, 04-C0, 04-D0 specify devices that are integrated on the standard I/O board. They can not be exchanged, because their electronic resides on the board.

Location codes 04-01, 04-02, 04-03, 04-04 specify devices that are not integrated into the motherboard. These cards can be replaced if newer adapters are available.



## Appendix F. Challenge Exercise

You will be presented with a series of problems to solve. The scenarios give several real-life problems that you may face as a system administrator. In some scenarios, you'll be given clear information about the problem but in some scenarios may not be given as much information as you would like. This is part of the troubleshooting process.

Like the other class exercises, the solutions are available but try to work through the scenarios without referring to the solutions. Try to solve the problems as if this were real. There's no solution section in the real world. Use your student notes, Web-based documentation, and the experience that you have gained from other exercises to troubleshoot and solve the problems.

### Day 1

Run the script: **/home/workshop/day1prob**

#### Scenario

You have just arrived at work and there are three trouble tickets waiting for you. Review the trouble tickets and solve the problem.

Trouble Ticket #1 - Several users have reported trying to create files in the **/home/data** file system but they keep receiving the error "There is not enough space in the file system."

Trouble ticket #2 - Several users have reported that some of the files in **/home/data/status** are missing and they need access to them right away. The missing files are **stat3** and **stat4**. The users accidentally removed the files and submitted a trouble ticket yesterday asking to have the files restored. They talked to the other administrator yesterday afternoon and were promised that the files would be restored overnight, but they are still missing.

Trouble ticket #3 - Users are complaining that the files in the **/home/project** directory are missing. There should be three files: **proj1**, **proj2** and **proj3**.

Extra: After talking to the other system administrator, he said he didn't do anything that would effect the **/home/data** file system. But, he did say he restored the **/home/data/status** file system from the backup (by inode) file **/home/workshop/status.bk** to the **/home/data/status** directory overnight.

### Day 2

Run the script **/home/workshop/day2prob**

Trouble ticket #4 - Users are complaining that the files in the **/home/project** directory are missing again. This is the fifth time in as many days. You check through the past week's trouble tickets and discover that there have been trouble tickets for this problem for the last 4 days. What might be the root cause of this recurring problem?.

Extra: You talk to the other administrator to determine if he did anything to impact the **/home/project** file system. He says he implemented a new backup script for the

/home/project file system. He's not sure exactly when he installed it. It was about 4 to 5 or maybe 6 days ago. He said he didn't document the date of the installation, but he tested the script five times and it worked perfectly all five times. He forgot the name of the script. He meant to write it in the system logbook but he forgot. After all, why document it when it works! He set the script up to run nightly.

### Day 3

Run the script: **/home/workshop/day3prob**

Power on the system and read the scenario.

You arrive at work. The other administrator looks very worried. He informs you he was cleaning up files, file systems and logical volumes. He said he deleted anything that looked like it wasn't in use. When he tried to reboot the system this morning, the machine wouldn't reboot. He is absolutely sure he didn't delete anything important... well, he is pretty sure that he didn't delete anything important... well, he might have deleted something important but he didn't know it was important. Of course, he didn't keep records of what he removed. But he did remember that he removed a logical volume. He knows it was a closed logical volume because he wouldn't attempt to remove an active logical volume. When he removed it, it prompted him to run another command... `chpv` something? He can't quite remember the command, but he did run the command just like it told him to.

What did he remove and can you fix it?

### Day 4

Run the script: **/home/workshop/day4prob**

Today, the administrator explains he did some more clean up last night. He was quite pleased with himself as he explained that this time when he removed **hd5**, he was not tempted to run the **chpv -c hdiskx** command because you made it clear to him that this was not a good thing. Next time, you need to make it clear not to remove a logical volume just because it is closed - especially **hd5**. He said that he rebooted the machine and it rebooted just fine. However, now you see some strange looking output from several commands.

Try running:

**lslv hd5**

**lsvg -l rootvg**

Do you notice any problems with **hd5**? How are you going to fix it?

## Day 1 - Fix and Explanation

### Trouble Ticket 1 and 2 Fix:

The other administrator restored the files like he said except he did not do it correctly. He recovered the /home/data/status file system but did not mount the file system first. The result was the files were restored into the /home/data file system (instead of the /home/data/status file system) filling /home/data. The files **stat3** and **stat4** are missing because during the recovery, the file system ran out of space.

To correct the problem, the files from /home/data/status (directory) need to be removed. The /home/data/status file system needs to be mounted and the file need to be restored.

```
cd /home/data/status
rm -r *
cd ..
mount /home/data/status
cd status
restore -rqvf /home/workshop/status.bk
```

### Trouble Ticket 3 Fix and Explanation:

For some reason, the /home/project file system is umounted. Mounting the file system will resolve this problem.

```
mount /home/project
```

## Day 2 - Fix and Explanation

You know from checking the trouble tickets that this is a recurring problem. If it is a recurring problem, you should consider the crontab file as a possible source of the trouble.

View the crontab file for root: **crontab -l**

Every morning at 3 a.m. a script named **perfect.bkp** is executed.

Examine that file: **cat /home/workshop/perfect.bkp**

The file umounts **/home/project** and then backs up the file system. However, the file system is never re-mounted. The script performs the backup just fine but the file system is never made accessible after it finishes. Add a line to the script to make sure the file system is mounted when the backup is done.

```
mount /home/project
```

## Day 3 - Fix and Explanation

The administrator removed a closed logical volume that impacted the ability of the machine to reboot. This is certainly **hd5**. Anytime you move (or remove) **hd5**, you are prompted to run **chpv -c hdiskx** so that the boot record is cleared from that disk. Once that is run, the machine will not reboot until **bosboot** is run to recreate it.

To fix the problem, boot into maintenance mode from CD or tape. Activate the rootvg and mount all of the file systems. If you try to run a **bosboot** now, you will be informed that **hd5** does not exist. You must first recreate the missing logical volume. To do that, run: **mkiv -t boot -y hd5 rootvg 1 hdisk0**

Now you can run: **bosboot -ad /dev/hdisk0** Shutdown the system and reboot: **shutdown -Fr**

### **Day 4 - Fix and Explanation**

This situation is a little more challenging to fix. Since the boot record was never cleared, there was still a pointer to the physical area that was known as **hd5**. The data still existed on the physical disk and therefore the machine was still able to boot. However, **hd5** - the logical volume, doesn't exist so now you get some strange looking output.

Try to run **mkiv -t boot -y hd5 rootvg 1 hdisk0**. Why does it fail? Because the system thinks **hd5** exists. Where is this information coming from?

This requires some trouble shooting to find the missing pieces. Try running these commands to see what is missing:

**lqueryvg -Atp hdiskx** - This will confirm whether **hd5** is a part of the VGDA. It is not.

Run queries against the customized ODM object classes to see where **hd5** exists.

**odmget CuDv | grep hd5 odmget CuAt | grep hd5 odmget CuDvDr | grep hd5 odmget CuDep | grep hd5**

Entries for **hd5** exist in all of these.

You have entries in ODM but not the VGDA. The VGDA is accurate. What is the best way to clean up the ODM?

You can either run a series of **odmdelete's** to clean ODM manually or just run the **rvgrecover** script. This will clear the ODM for rootvg and rebuild it from the VGDA.

Then, you can finish the clean up by running: **mkiv -t boot -y hd5 rootvg 1 hdisk0 bosboot -ad /dev/hdisk0**

Run **lsvg -l rootvg** and **lslv hd5** to verify everything looks correct.

# Appendix G. Auditing Security Related Events

# Appendix Objectives

---

- Configure the Auditing Subsystem

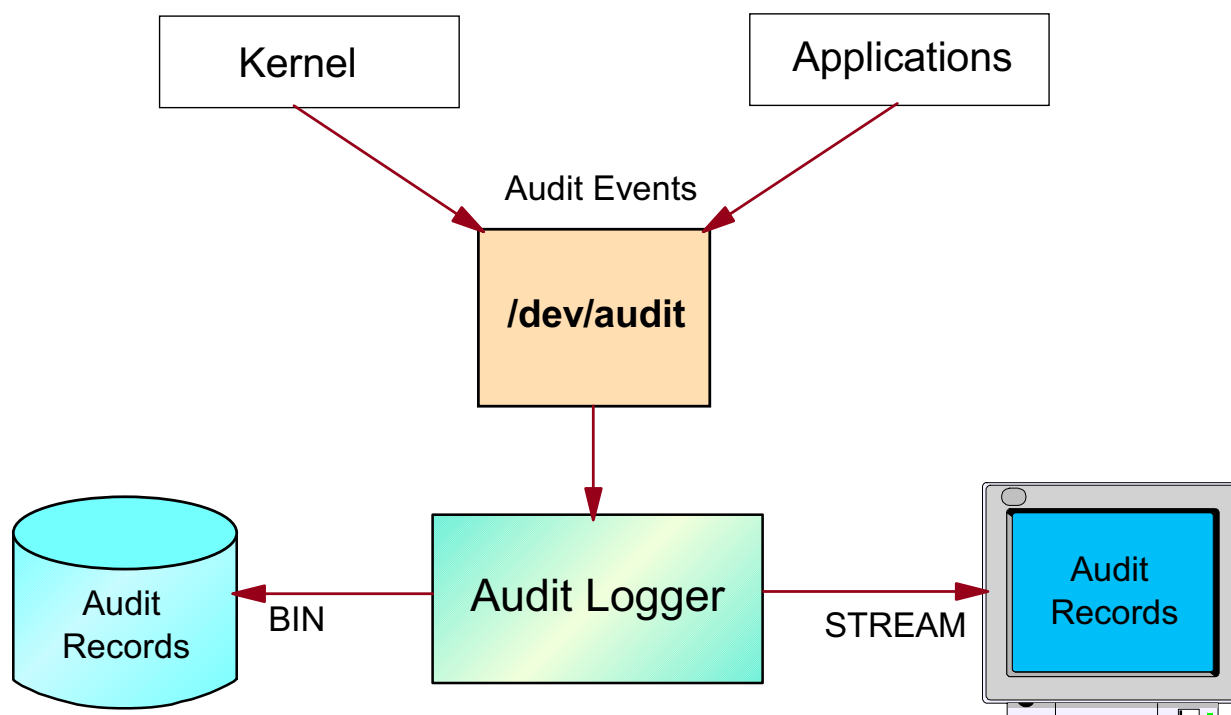
© Copyright IBM Corporation 2004

Figure G-1. Appendix Objectives

AU1612.0

## **Notes:**

## How the Auditing Subsystem Works



© Copyright IBM Corporation 2004

Figure G-2. How the Auditing Subsystem Works

AU1612.0

### Notes:

The AIX auditing subsystem provides a way to trace security-relevant events like **accessing an important system file** or the **execution of applications**, which might influence the security of your system.

The auditing subsystem works in the following way. The AIX Kernel or other security-related applications use a system call to process the security-related event in the auditing subsystem. This system call writes the auditing information to a special file **/dev/audit**. An **audit logger** reads the audit information from this device, formats it and writes the audit record either to files (in **BIN** mode) or to a specified device, for example a display, or a printer (in **STREAM** mode).

# Auditing Configuration Files

|                             |                                                                                                        |
|-----------------------------|--------------------------------------------------------------------------------------------------------|
| /etc/security/audit/objects | Contains the <b>audit events</b> triggered by file access                                              |
| /etc/security/audit/events  | Contains information about system <b>audit events</b> and <b>responses</b> to those events             |
| /etc/security/audit/config  | Contains <b>audit configuration</b> information:<br>- Start Mode<br>- Audit classes<br>- Audited Users |

© Copyright IBM Corporation 2004

Figure G-3. Audit Configuration Files

AU1612.0

## Notes:

All audit configuration files reside in directory **/etc/security/audit**. The following configuration files are used by the auditing subsystem:

- **objects**

This file describes all files and programs that are audited. For each file a unique audit event name is specified. These files are monitored by the AIX Kernel.

- **events**

This file contains one stanza called **auditpr**. Each audit event is named and the format of the output produced by each event is defined in this stanza. The **auditpr** command writes all audit output based on this information in this file.

- **config**

This file contains audit configuration information:

- The **start mode** for the audit logger (BIN or STREAM mode)



- **Audit classes**, which are groups of audit events. Each audit class name must be less than 16 characters and must be unique to the system. AIX supports up to 32 audit classes.
- **Audited users:** The users whose activities you wish to monitor are defined in the **users** stanza. A **users** stanza determines which combination of user and event class to monitor.

## Audit Configuration: objects

```
# vi /etc/security/audit/objects

/etc/security/user:
    w = "S_USER_WRITE"

...

/etc/filesystems:
    w = "MY_EVENT"

/usr/sbin/no:
    x = "MY_X_EVENT"
```

© Copyright IBM Corporation 2004

Figure G-4. Audit Configuration: objects

AU1612.0

### Notes:

To configure the auditing subsystem you first specify the **objects** (files or applications) that you want to audit in **/etc/security/audit/objects**. In this file you find predefined files, for example **/etc/security/user**.

To audit your own files you have to add stanzas for each file, in the following format:

```
file:
access_mode = "event_name"
```

An audit event name can be up to 15 bytes long. Valid access modes are read (r), write (w) and execute (x).

In the shown example we add two files. An event **MY\_EVENT** will be generated by the AIX Kernel, when somebody writes the file **/etc/filesystems**. Another event **MY\_X\_EVENT** will be generated when somebody executes the program **/usr/sbin/no**. After adding objects, you have to specify formatting information in the **events** file. That's shown on the next visual.

**Note:** Symbolic links cannot be monitored by the auditing subsystem.

## Audit Configuration: events

```
# vi /etc/security/audit/events
auditpr:

    USER_Login  = printf "user: %s tty: %s"
    USER_Logout = printf "%s"

    ...

    MY_EVENT = printf "%s"

    MY_X_EVENT = printf "%s"
```

© Copyright IBM Corporation 2004

Figure G-5. Audit Configuration: events

AU1612.0

### Notes:

All audit system events have a **format specification** that is used by the **auditpr** command, which prints the audit record. This format specification is defined in the **/etc/security/audit/events** file and specifies how the information will be printed when the audit data is analyzed.

Each attribute in the stanza is the **name of an audit event**, where the following formats are possible:

```
AuditEvent = printf "format-string"
AuditEvent = event_program arguments
```

To print out the audit record with all event arguments **printf** is used. Different format specifiers are used, depending on the audit event that occurs. If you want to trigger other applications that are called whenever an event occurs, you can specify an **event\_program**. If you do this, always use the full pathname of the **event\_program**.

If you specify your own events in the **objects** file, you need to add a format specification to the **events** file. For our self-defined events **MY\_EVENT** and **MY\_X\_EVENT** we use the

**printf** format command. Remember that the AIX Kernel monitors these objects and triggers the audit events.

## Audit Configuration: config

```
# vi /etc/security/audit/config

start:
    binmode = off
    streammode = on

...

classes:
    general = USER_SU, PASSWORD_Change, ...
    tcPIP = TCPIP_connect, TCPIP_data_in, ...
    ...
    init = USER_Login, USER_Logout

users:
    root = general
    michael = init
```

© Copyright IBM Corporation 2004

Figure G-6. Audit Configuration: config

AU1612.0

### Notes:

The `/etc/security/audit/config` file contains audit configuration information.

1. The stanza **start** specifies the start mode for the audit logger. If you work in **bin mode**, the audit records are stored in files. The **auditbin** daemon will be started. The **streammode** allows real-time processing of an audit event, for example to display the audit record on the system console or to print it on a printer.
2. The stanza **classes** groups audit events together to a class. These classes could then be assigned to users who are then audited for all events belonging to a class. Note that this is necessary for all events that are triggered by applications. Object events triggered by the kernel need not to be part of a class.

Note that the class name (for example **init**) must be less than 16 characters and must be unique on the system.

3. The stanza **users** assigns audit classes to a user. The username (for example **michael**) must be the login name of a system user, or the string **default** which stands for all system users.

In the example, the self-defined class **init** is assigned to the user **michael**. Whenever **michael** logs in or out from the system, an audit record will be written.

Note that you can also use the **chuser** command to establish an audit activity for a special user:

```
# chuser "auditclasses=init" michael
```

## Audit Configuration: bin Mode

```
# vi /etc/security/audit/config

start:
    binmode = on
    streammode = off

bin:
    trail = /audit/trail
    bin1 = /audit/bin1
    bin2 = /audit/bin2
    binsize = 10240
    cmds = /etc/security/audit/bincmds

...
```

- Use the **auditpr** command to display the audit records:

```
# auditpr -v < /audit/trail
```

© Copyright IBM Corporation 2004

Figure G-7. Audit Configuration: bin Mode

AU1612.0

### Notes:

To work in bin mode, specify **binmode = on** in the **start** stanza in **/etc/security/audit/config**. In this case, the **auditbin** daemon will be started.

The **bin** stanza specifies how the bin mode works: The audit records are stored in alternating files that have a fixed size (specified by **binsize**). The records are first written into the file specified by **bin1**. When this file fills, future records are written to **/audit/bin2** automatically and the content of **/audit/bin1** is written to **/audit/trail** to create the **permanent** record.

To display the audit records you must use the **auditpr** command:

```
# auditpr -v < /audit/trail
```

**In this example you display the audit records that are stored in /audit/trail.**

If you use bin-mode auditing, it's recommended that you do **not** specify bins that are in the **hd4** (root) file system.

## Audit Configuration: stream Mode

```
# vi /etc/security/audit/config


start:
    binmode = off
    streammode = on

stream:
    cmds = /etc/security/audit/streamcmds

...

# vi /etc/security/audit/streamcmds

/usr/sbin/auditstream | auditpr -v > /dev/console &
```



All audit records are displayed on the console

© Copyright IBM Corporation 2004

Figure G-8. Audit Configuration: stream Mode

AU1612.0

### Notes:

The **stream mode** allows real-time processing of the audit events. To configure **stream mode** auditing, you have to do two things in **/etc/security/audit/config**:

1. Specify **streammode = on** in the **start** stanza.
2. Specify the audit record destination in the stream mode backend file **/etc/security/audit/streamcmds**. In our example all records are displayed on the console, using the **auditpr** command. Note that you must specify the **&** sign after the command.

The **auditstream** command starts up an **auditstream** daemon. You can startup multiple daemons in **streamcmds** that monitors different classes, for example:

```
/usr/sbin/auditstream -c init | auditpr -v > /var/init.txt &
/usr/sbin/auditstream -c general | auditpr -v > /var/general.txt &
```

If you want to monitor selected events in these classes, use the **auditselect** command. See **man** pages for more information.



## The audit Command

|                  |   |                            |
|------------------|---|----------------------------|
| # audit start    | → | Start / Stop auditing      |
| # audit shutdown |   |                            |
| # audit query    | → | Display audit status       |
| # audit off      | → | Suspend / Restart auditing |
| # audit on       |   |                            |

© Copyright IBM Corporation 2004

Figure G-9. The audit Command

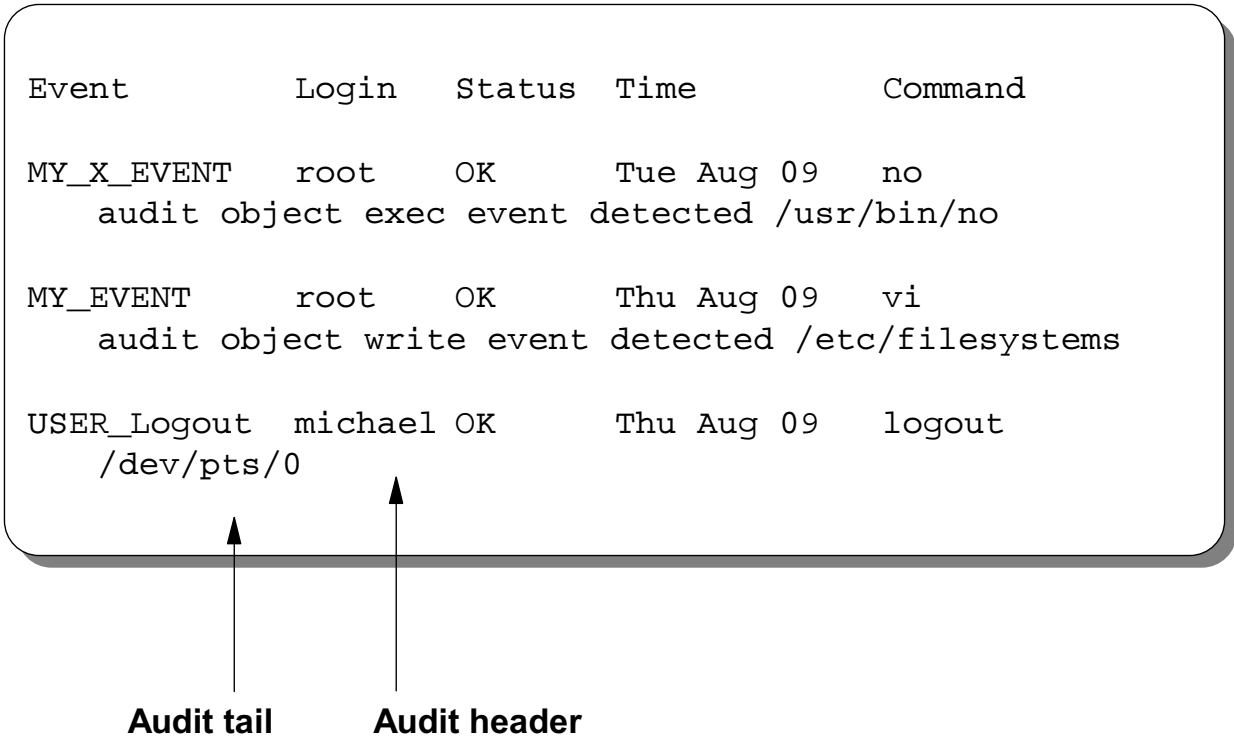
AU1612.0

### Notes:

The **audit** command controls system auditing. To start the auditing system use **audit start**, to stop auditing use **audit shutdown**. Note that you have to stop and restart auditing whenever you change a configuration file.

To query the current audit configuration, use **audit query**. If you want to suspend auditing, use **audit off** to restart it, use **audit on**.

# Example Audit Records



© Copyright IBM Corporation 2004

Figure G-10. Example Audit Records

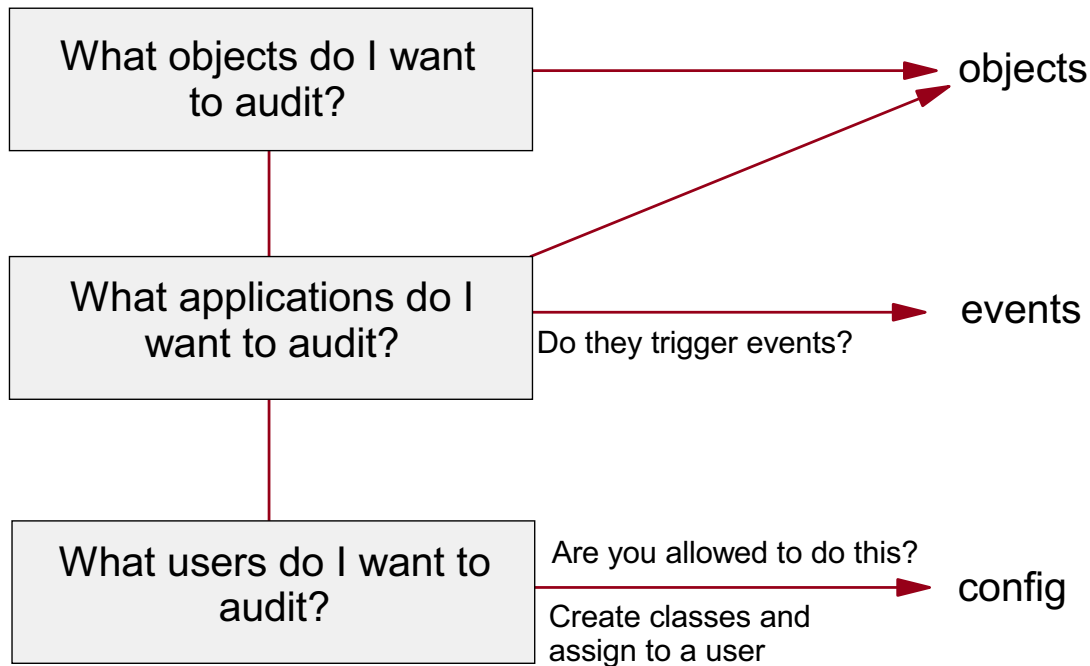
AU1612.0

## Notes:

Each audit record consists of two parts, an **audit header** and an **audit tail**. The tail is printed according to the format specification in `/etc/security/audit/events` and is only shown if you use the `-v` option in the `auditpr` command.

The **audit header** specifies the event name, the user, the status, the time and the command that triggers the audit event. The **audit tail** shows additional information, for example the terminal where the user logged out, as shown on the visual.

# Set Up Auditing in Your Environment



© Copyright IBM Corporation 2004

Figure G-11. Set Up Auditing in Your Environment

AU1612.0

## Notes:

If used correctly, the auditing subsystem is a very good tool for auditing events. However, problems can arise if the auditing subsystem gathers too much data to be analyzed. To prevent this problem from occurring, careful planning is required when configuring auditing. This flowchart provides an aid to configure auditing in your environment so that the auditing data can be managed.

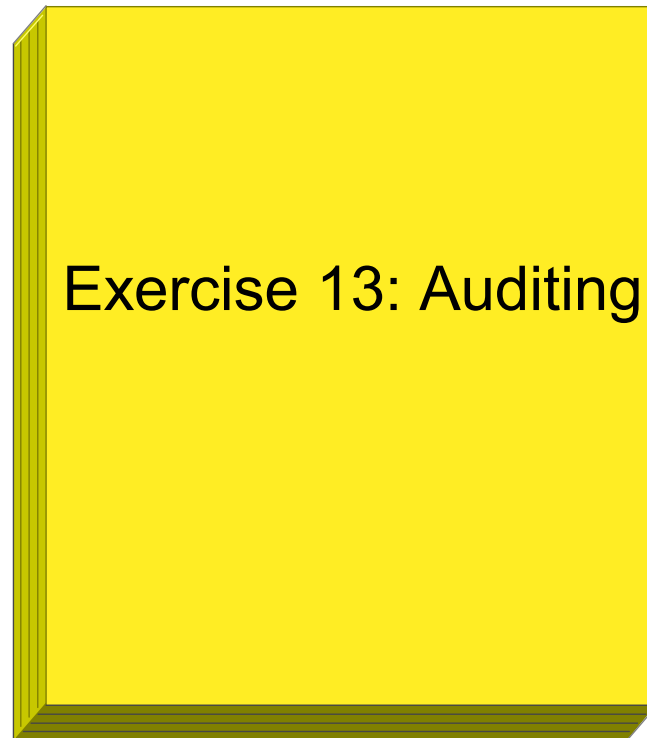
- Decide what **objects** you want to monitor. Objects are files that you can audit for read, write or execute actions. For example, files that make good candidates for monitoring are those in the **/etc** directory. Unfortunately, the audit subsystem can only monitor **existing** files. If you wanted to monitor files like **.rhosts**, you first need to create the files.
- Decide if you want to monitor special **applications**. This could be done by adding an execute event into the **objects** file. If you are interested in application events, you must determine if the application triggers audit events. For example, you might want to audit all TCP/IP-related events on a system where the transfer of data needs to be monitored. These events can be found in the **events** file.

- Decide if you want to trace **users**. Before doing this, confirm that there are no legal issues within your organization that would prohibit tracing users. To trace users, create audit classes and assign these classes to the users you want to audit.

---

## Next Step

---



© Copyright IBM Corporation 2004

Figure G-12. Next Step

AU1612.0

### **Notes:**

After the lab exercise, you should be able to:

- Audit objects and application events
- Create audit classes and audit users
- Set up auditing in bin and stream mode



# Glossary

## A

**access mode** A matrix of protection information stored with each file specifying who may do what to a file. Three classes of users (owner, group, all others) are allowed or denied three levels of access (read, write, execute).

**access permission** See **access mode**.

**access privilege** See **access mode**.

**address space** The address space of a process is the range of addresses available to it for code and data. The relationship between real and perceived space depends on the system and support hardware.

**AIX** Advanced Interactive Executive. IBM's implementation of the UNIX Operating System.

**AIX Family Definition** IBM's definition for the common operating system environment for all members of the AIX family. The AIX Family Definition includes specifications for the AIX Base System, User Interface, Programming Interface, Communications Support, Distributed Processing, and Applications.

**alias** The command and process of assigning a new name to a command.

**ANSI** American National Standards Institute. A standards organization. The United States liaison to the International Standards Organization (ISO).

**application program** A program used to perform an application or part of an application.

**argument** An item of information following a command. It may, for example, modify the command or identify a file to be affected.

**ASCII** American Standard Code for Information Interchange. A collection of public domain character sets considered standard throughout the computer industry.

**awk** An interpreter, included in most UNIX operating systems, that performs sophisticated text pattern matching. In combination with shell scripts, awk can be used to prototype or implement applications far more quickly than traditional programming methods.

## B

**background (process)** A process is "in the background" when it is running independently of the initiating terminal. It is specified by ending the ordinary command with an ampersand (&). The parent of the background process does not wait for its "death".

**backup diskette** A diskette containing information copied from another diskette. It is used in case the original information is unintentionally destroyed.

**Berkeley Software Distribution** Disseminating arm of the UNIX operating system community at the University of California at Berkeley; commonly

abbreviated "BSD". Complete versions of the UNIX operating system have been released by BSD for a number of years; the latest is numbered 4.3. The phrase "Berkeley extensions" refers to features and functions, such as the C shell, that originated or were refined at UC Berkeley and that are now considered a necessary part of any fully-configured version of the UNIX operating system.

**bit bucket** The AIX file "/dev/null" is a special file which will absorb all input written to it and return no data (null or end of file) when read.

**block** A group of records that is recorded or processed as a unit.

**block device** A device that transfers data in fixed size blocks. In AIX, normally 512 or 1024 bytes.

**block special file** An interface to a device capable of supporting a file system.

**booting** Starting the computer from scratch (power off or system reset).

**break key** The terminal key used to unequivocally interrupt the foreground process.

**BSD** Berkeley Software Distribution.

- BSD 2.x - PDP-11 Research
- BSD 4.x - VAX Research
- BSD 4.3 - Current popular VAX version of UNIX.

## button

1. A word, number, symbol, or picture on the screen that can be selected. A button may represent a command, file, window, or value, for example.
2. A key on a mouse that is used to select buttons on the display screen or to scroll the display image.

**byte** The amount of storage required to represent one character; a byte is 8 bits.

## C

**C** The programming language in which the UNIX operating system and most UNIX application programs are written. The portability attributed to UNIX operating systems is largely due to the fact that C, unlike other higher level languages, permits programmers to write systems-level code that will work on any computer with a standard C compiler.

**change mode** The chmod command will change the access rights to your own files only, for yourself, your group or all others.

**character I/O** The transfer of data byte by byte; normally used with slower, low-volume devices such as terminals or printers.

**character special file** An interface to devices not capable of supporting a file system; a byte-oriented device.

**child** The process emerging from a fork command with a zero return code, as distinguished from the parent which gets the process id of the child.

**client** User of a network service. In the client/server model, network elements are defined as either using (client) or providing (server) network resources.

**command** A request to perform an operation or run a program. When parameters, arguments, flags, or other operands are associated with a command, the resulting character string is a single command.

**command file** A data file containing shell commands. See **shell file**, or **shell script**.

**command interpreter** The part of the operating system that translates your commands into instructions that the operating system understands.

**concatenate** The process of forming one character string or file from several. The degenerate case is one file from one file just to display the result using the **cat** command.

**console** The only terminal known explicitly to the Kernel. It is used during booting and it is the destination of serious system messages.

**context** The hardware environment of a process, including:

- CPU registers
- Program address
- Stack
- I/O status

The entire context must be saved during a process swap.

**control character** Codes formed by pressing and holding the **control** key and then some other key; used to form special functions like **End Of File**.

**cooked input** Data from a character device from which backspace, line kill, and interrupt characters have been removed (processed). See **raw input**.

**current directory** The currently active directory. When you specify a file name without specifying a directory, the system assumes that the file is in your current directory.

**current subtree** Files or directories attached to the current directory.

**courses** A C subroutine library providing flexible screen handling. See **Termlib** and **Termcap**.

**cursor** A movable symbol (such as an underline) on a display, usually used to indicate to the operator where to type the next character.

**customize** To describe (to the system) the devices, programs, users, and user defaults for a particular data processing system.

## D

**DASD** Direct Access Storage Device. IBM's term for a hard disk.

**device driver** A program that operates a specific device, such as a printer, disk drive, or display.

**device special file** A file which passes data directly to/from the device.

**directory** A type of file containing the names and controlling information for other files or other directories.

**directory pathname** The complete and unique external description of a file giving the sequence of connection from the root directory to the specified directory or file.

**diskette** A thin, flexible magnetic plate that is permanently sealed in a protective cover. It can be used to store information copied from the disk.

**diskette drive** The mechanism used to read and write information on diskettes.

**display device** An output unit that gives a visual representation of data.

**display screen** The part of the display device that displays information visually.

## E

**echo** To simply report a stream of characters, either as a message to the operator or a debugging tool to see what the file name generation process is doing.

**editor** A program used to enter and modify programs, text, and other types of documents.

**environment** A collection of values passed either to a C program or a shell script file inherited from the invoking process.

**escape** The backslash "\" character specifies that the single next character in a command is ordinary text without special meaning.

**Ethernet** A baseband protocol, invented by the XEROX Corporation, in common use as the local area network for UNIX operating systems interconnected via TCP/IP.

**event** One of the previous lines of input from the terminal. Events are stored in the (Berkeley) History file.

**event identifier** A code used to identify a specific event.

**execution permission** For a file, the permission to execute (run) code in the file. A text file must have execute permission to be a shell script. For a directory, the permission to search the directory.



**F**

**field** A contiguous group of characters delimited by blanks. A field is the normal unit of text processed by text processes like sort.

**field separator** The character used to separate one field from the next; normally a blank or tab.

**FIFO** "First In, First Out". In AIX, a FIFO is a permanent, named pipe which allows two unrelated processes to communicate. Only related processes can use normal pipes.

**file** A collection of related data that is stored and retrieved by an assigned name. In AIX, files are grouped by directories.

**file index** Sixty-four bytes of information describing a file. Information such as the type and size of the file and the location on the physical device on which the data in the file is stored is kept in the file index. This index is the same as the AIX Operating System i-node.

**filename expansion or generation** A procedure used by the shell to generate a set of filenames based on a specification using metacharacters, which define a set of textual substitutions.

**file system** The collection of files and file management structures on a physical or logical mass storage device, such as a diskette or minidisk.

**filter** Data-manipulation commands (which, in UNIX operating systems, amount to small programs) that take input from one process and perform an operation yielding new output. Filters include editors, pattern-searchers, and commands that sort or differentiate files, among others.

**fixed disk** A storage device made of one or more flat, circular plates with magnetic surfaces on which information can be stored.

**fixed disk drive** The mechanism used to read and write information on a fixed disk.

**flag** See **Options**.

**foreground (process)** An AIX process which interacts with the terminal. Its invocation is not followed by an ampersand.

**formatting** The act of arranging text in a form suitable for reading. The publishing equivalent to compiling a program.

**fsck** A utility to check and repair a damaged file structure. This normally results from a power failure or hardware malfunction. It looks for blocks not assigned to a file or the free list and puts them in the free list. (The use of blocks not pointed at cannot be identified.)

**free list** The set of all blocks not assigned to a file.

**full path name** The name of any directory or file expressed as a string of directories and files beginning with the root directory.

**G**

**gateway** A device that acts as a connector between two physically separate networks. It has interfaces

to more than one network and can translate the packets of one network to another, possibly dissimilar network.

**global** Applying to all entities of a set. For example:

- A global search - look everywhere
- A global replace - replace all occurrences
- A global symbol - defined everywhere.

**grep** An AIX command which searches for strings specified by a regular expression. (Global Regular Expression and Print.)

**group** A collection of AIX users who share a set of files. Members of the group have access privileges exceeding those of other users.

**H**

**hardware** The equipment, as opposed to the programming, of a system.

**header** A record at the beginning of the file specifying internal details about the file.

**heterogeneous** Descriptor applied to networks composed of products from multiple vendors.

**hierarchy** A system of objects in which each object belongs to a group. Groups belong to other groups. Only the "head" does not belong to another group. In AIX this object is called the "Root Directory".

**highlight** To emphasize an area on the display screen by any of several methods, such as brightening the area or reversing the color of characters within the area.

**history** A list of recently executed commands.

**home (directory)**

1. A directory associated with an individual user.
2. Your current directory on login or after issuing the **cd** command with no argument.

**homogeneous** Descriptor applied to networks composed of products from a single vendor.

**hypertext** Term for on-line interactive documentation of computer software; to be included with AIX.

**I**

**IEEE** Institute of Electrical and Electronics Engineers. A professional society active in standards work, the IEEE is the official body for work on the POSIX (Portable Operating System for Computer Environments) open system interface definition.

**index** See **file index**.

**indirect block** A file element which points at data sectors or other indirect blocks.

**init** The initialization process of AIX. The ancestor of all processes.

**initial program load** The process of loading the system programs and preparing the system to run jobs.

**i-node** A collection of logical information about a file including owner, mode, type and location.

**i number** The internal index or identification of an i-node.

**input field** An area into which you can type data.

**input redirection** The accessing of input data from other than standard input (the keyboard or a pipe).

**interoperability** The ability of different kinds of computers to work well together.

**interpreter** A program which "interprets" program statements directly from a text (or equivalent) file. Distinguished from a compiler which creates computer instructions for later direct execution.

**interrupt** A signal that the operating system must reevaluate its selection of which process should be running. Usually to service I/O devices but also to signal from one process to another.

**IP** Internet Protocol.

**ipl** See initial program load.

**ISO** International Standards Organization. A United Nations agency that provides for creation and administration of worldwide standards.

## J

**job** A collection of activities.

**job number.** An identifying number for a collection of processes devolving from a terminal command.

## K

**kernel** The part of an operating system that contains programs that control how the computer does its work, such as input/output, management and control of hardware, and the scheduling of user tasks.

**keyboard** An input device consisting of various keys allowing the user to input data, control cursor and pointer locations, and to control the user/work station dialogue.

**kill** To prematurely terminate a process.

**kill character** The character which erases an entire line (usually @).

## L

**LAN** Local Area Network. A facility, usually a combination of wiring, transducers, adapter boards, and software protocols, which interconnects workstations and other computers located within a department, building, or neighborhood. Token-Ring and Ethernet are local area network products.

**libc** A basic set of C callable routines.

**library** In UNIX operating systems, a collection of existing subroutines that allows programmers to make use of work already done by other programmers. UNIX operating systems often include separate libraries for communications, window management, string handling, math, and so forth.

**line editor** An editor which processes one line at a time by the issuing of a command. Usually associated with sequential only terminals such as a teletype.

**link** An entry in an AIX directory specifying a data file or directory and its name. Note that files and directories are named solely by virtue of links. A name is not an intrinsic property of a file. A file is uniquely identified only by a system generated identification number.

**lint** A program for removing "fuzz" from C code. Stricter than most compilers. Helps former Pascal programmers sleep at night.

**Local Area Network (LAN)** A facility, usually a combination of wiring, transducers, adapter boards, and software protocols, which interconnects workstations and other computers located within a department, building, or neighborhood. Token-Ring and Ethernet are local area network products.

**log in** Identifying oneself to the system to gain access.

**login directory** See home directory.

**login name** The name by which a user is identified to the system.

**log out** Informing the system that you are through using it.

## M

**mail** The process of sending or receiving an electronically delivered message within an AIX system. The message or data so delivered.

**make** Programming tool included in most UNIX operating systems that helps "make" a new program out of a collection of existing subroutines and utilities, by controlling the order in which those programs are linked, compiled, and executed.

**map** The process of reassigning the meaning of a terminal key. In general, the process of reassigning the meaning of any key.

**memory** Storage on electronic memory such as random access memory, readonly memory, or registers. See **storage**.

**message** Information displayed about an error or system condition that may or may not require a user response.

**motd** "Message of the day". The login "billboard" message.

**Motif™** The graphical user interface for OSF, incorporating the X Window System. Behavior of this interface is compatible with the IBM/Microsoft Presentation Manager user interface for OS/2. Also called OSF/Motif.

**mount** A logical (that is, not physical) attachment of one file directory to another. "remote mounting" allows files and directories that reside on physically separate computer systems to be attached to a local system.

**mouse** A device that allows you to select objects and scroll the display screen by means of buttons.

**move** Relinking a file or directory to a different or additional directory. The data (if any) is not moved, only the links.

**multiprogramming** Allocation of computer resources among many programs. Used to allow many users to operate simultaneously and to keep the system busy during delays occasioned by I/O mechanical operations.

**multitasking** Capability of performing two or more computing tasks, such as interactive editing and complex numeric calculations, at the same time. AIX and OS/2 are multi-tasking operating systems; DOS, in contrast, is a single-tasking system.

**multiuser** A computer system which allows many people to run programs "simultaneously" using multiprogramming techniques.

## N

**named pipe** See **FIFO**.

**Network File System (NFST)** A program developed by SUN Microsystems, Inc. for sharing files among systems connected via TCP/IP. IBM's AIX, VM, and MVS operating systems support NFS.

**NFS™** See **Network File System**.

**NIST** National Institute of Science and Technology (formerly the National Bureau of Standards).

**node** An element within a communication network.

- Computer
- Terminal
- Control Unit

**null** A term denoting emptiness or nonexistence.

**null device** A device used to obtain empty files or dispose of unwanted data.

**null string** A character string containing zero characters.

## O

**object-oriented programming** Method of programming in which sections of program code and data are represented, used, and edited in the form of "objects", such as graphical elements, window components, and so forth, rather than as strict computer code. Through object-oriented programming techniques, toolkits can be designed that make programming much easier. Examples of object-oriented programming languages include Pareplace Systems, Inc.'s Smalltalk-80™, AT&T's C++™, and Stepstone Inc.'s Objective-C®.

**oem** original equipment manufacturer. In the context of AIX, OEM systems refer to the processors of a heterogeneous computer network that are not made or provided by IBM.

**Open Software Foundation™ (OSF)**. A non-profit consortium of private companies, universities, and research institutions formed to conduct open technological evaluations of available components of UNIX operating systems, for the purpose of assembling selected elements into a complete version of the UNIX operating system available to those who wish to license it. IBM is a founding sponsor and member of OSF.

**operating system** The programs and procedures designed to cause a computer to function, enabling the user to interact with the system.

**option** A command argument used to specify the details of an operation. In AIX an option is normally preceded by a hyphen.

**ordinary file** Files containing text, programs, or other data, but not directories.

**OSF** See Open Software Foundation.

**output redirection** Passing a programs standard output to a file.

**owner** The person who created the file or his subsequent designee.

## P

**packet switching** The transmission of data in small, discrete switching "packets" rather than in streams, for the purpose of making more efficient use of the physical data channels. Employed in some UNIX system communications.

**page** To move forward or backward on screen full of data through a file usually referring to an editor function.

**parallel processing** A computing strategy in which a single large task is separated into parts, each of which then runs in parallel on separate processors.

**parent** The process emerging from a Fork with a non-zero return code (the process ID of the child process). A directory which points at a specified directory.

**password** A secret character string used to verify user identification during login.

**PATH** A variable which specifies which directories are to be searched for programs and shell files.

**path name** A complete file name specifying all directories leading to that file.

**pattern-matching character** Special characters such as \* or ? that can be used in a file specification to match one or more characters. For example, placing a ? in a file specification means that any character can be in that position.

**permission** The composite of all modes associated with a file.

**pipes** UNIX operating system routines that connect the standard output of one process with the standard input of another process. Pipes are central to the

function of UNIX operating systems, which generally consist of numerous small programs linked together into larger routines by pipes. The "piping" of the list directory command to the word count command is **ls | wc**. The passing of data by a pipe does not (necessarily) involve a file. When the first program generates enough data for the second program to process, it is suspended and the second program runs. When the second program runs out of data it is suspended and the first one runs.

**pipe fitting** Connecting two programs with a pipe.

**pipeline** A sequence of programs or commands connected with pipes.

**portability** Desirable feature of computer systems and applications, referring to users' freedom to run application programs on computers from many vendors without rewriting the program's code. Also known as "applications portability", "machine-independence", and "hardware-independence"; often cited as a cause of the recent surge in popularity of UNIX operating systems.

**port** A physical I/O interface into a computer.

**POSIX** "Portable Operating Systems for Computer Environments". A set of open standards for an operating system environment being developed under the aegis of the IEEE.

**preprocessor** The macro generator preceding the C compiler.

**process** A unit of activity known to the AIX system, usually a program.

**process 0 (zero)** The scheduler. Started by the "boot" and permanent. See **init**.

**process id** A unique number (at any given time) identifying a process to the system.

**process status** The process's current activity.

- Non existent
- Sleeping
- Waiting
- Running
- Intermediate
- Terminated
- Stopped.

**profile** A file in the users home directory which is executed at login to customize the environment. The name is **.profile**.

**prompt** A displayed request for information or operator action.

**protection** The opposite of permission, denying access to a file.

## Q

**quotation** Temporarily cancelling the meaning of a metacharacter to be used as an ordinary text character. A backslash (\) "quotes" the next character only.

## R

**raw I/O** I/O conducted at a "physical" level.

**read permission.** Allows reading (not execution or writing) of a file.

**recursive** A recursive program calls itself or is called by a subroutine which it calls.

**redirection** The use of other than standard input (keyboard or pipe output) or standard output (terminal display or pipe). Usually a file.

**regular expression** An expression which specifies a set of character strings using metacharacters.

**relative path name** The name of a directory or file expressed as a sequence of directories followed by a file name, beginning from the current directory.

**RISC** Reduced Instruction Set Computer. A class of computer architectures, pioneered by IBM's John Cocke, that improves price-performance by minimizing the number and complexity of the operations required in the instruction set of a computer. In this class of architecture, advanced compiler technology is used to provide operations, such as multiplication, that are infrequently used in practice.

**root directory** The directory that contains all other directories in the file system.

## S

**scalability** Desirable feature of computer systems and applications. Refers to the capability to use the same environment on many classes of computers, from personal computers to supercomputers, to accommodate growth or divergent environments, without rewriting code or losing functionality.

**SCCS** Source Code Control System. A set of programs for maintaining multiple versions of a file using only edit commands to specify alternate versions.

**scope** The field of an operation or definition. Global scope means all objects in a set. Local scope means a restriction to a subset of the objects.

**screen** See **display screen**.

**scroll** To move information vertically or horizontally to bring into view information that is outside the display screen or pane boundaries.

**search and replace** The act of finding a match to a given character string and replacing each occurrence with some other string.

**search string** The pattern used for matching in a search operation.

**sed** Non-interactive stream editor used to do "batch" editing. Often used as a tool within shell scripts.



**server** A provider of a service in a computer network; for example, a mainframe computer with large storage capacity may play the role of database server for interactive terminals. See **client**.

**setuid** A permission which allows the access rights of a program owner to control the access to a file. The program can act as a filter for user data requests.

**shell** The outermost (user interface) layer of UNIX operating systems. Shell commands start and control other processes, such as editors and compilers; shells can be textual or visual. A series of system commands can be collected together into a "shell script" that executes like a batch (.BAT) file in DOS.

**shell program** A program consisting of a sequence of shell commands stored in an ordinary text file which has execution permission. It is invoked by simply naming the file as a shell command.

**shell script** See **shell program**.

**single user (mode)** A temporary mode used during "booting" of the AIX system.

**signal** A software generated interrupt to another process. See **kill**.

**sockets** Destination points for communication in many versions of the UNIX operating system, much as electrical sockets are destination points for electrical plugs. Sockets, associated primarily with 4.3 BSD, can be customized to facilitate communication between separate processes or between UNIX operating systems.

**software** Programs.

**special character** See **metacharacter**.

**special file** A technique used to access I/O devices in which "pseudo files" are used as the interface for commands and data.

**standard error** The standard device at which errors are reported, normally the terminal. Error messages may be directed to a file.

**standard input** The source of data for a filter, which is by default obtained from the terminal, but which may be obtained from a file or the standard output of another filter through a pipe.

**standard output** The output of a filter which normally is by default directed to the terminal, but which may be sent to a file or the standard input of another filter through a pipe.

**stdio** A "Standard I/O" package of C routines.

**sticky bit** A flag which keeps commonly used programs "stick" to the swapping disk for performance.

**stopped job** A job that has been halted temporarily by the user and which can be resumed at his command.

**storage** In contrast to memory, the saving of information on physical devices such as fixed disk or tape. See **memory**.

**store** To place information in memory or onto a diskette, fixed disk, or tape so that it is available for retrieval and updating.

**streams** Similar to sockets, streams are destination points for communications in UNIX operating systems. Associated primarily with UNIX System V, streams are considered by some to be more elegant than sockets, particularly for interprocess communication.

**string** A linear collection of characters treated as a unit.

**subdirectory** A directory which is subordinate to another directory.

**subtree** That portion of an AIX file system accessible from a given directory below the root.

**suffix** A character string attached to a file name that helps identify its file type.

**superblock** Primary information repository of a file system (location of i-nodes, free list, and so forth).

**superuser** The system administration; a user with unique privileges such as upgrading execution priority and write access to all files and directories.

**superuser authority** The unrestricted ability to access and modify any part of the Operating System. This authority is associated with the user who manages the system.

**SVID** System V Interface Definition. An AT&T document defining the standard interfaces to be used by UNIX System V application programmers and users.

**swap space (disk)** That space on an I/O device used to store processes which have been swapping out to make room for other processes.

**swapping** The process of moving processes between main storage and the "swapping device", usually a disk.

**symbolic debugger** Program for debugging other programs at the source code level. Common symbolic debuggers include sdb, dbx, and xdbx.

**sync** A command which copies all modified blocks from RAM to the disk.

**system** The computer and its associated devices and programs.

**system unit** The part of the system that contains the processing unit, the disk drive and the disk, and the diskette drive.

**System V** AT&T's recent releases of its UNIX operating system are numbered as releases of "UNIX System V".

## T

**TCP** Transmission Control Protocol. A facility for the creation of reliable bytestreams (byte-by-byte, end-to-end transmission) on top of unreliable datagrams. The transmission layer of TCP/IP is used to interconnect applications, such as FTP, so that issues of re-transmission and blocking can be subordinated in a standard way. See **TCP/IP**.

**TCP/IP** Transmission Control Protocol/Internet Protocol. Pair of communications protocol considered defacto standard in UNIX operating system environments. IBM TCP/IP for VM and IBM TCP/IP for MVS are licensed programs that provide VM and MVS users with the capability of participating in networks using the TCP/IP protocol suite.

**termcap** A file containing the description of several hundred terminals. For use in determining communication protocol and available function.

**termlib** A set of C programs for using termcap.

**tools** Compact, well designed programs to perform specific tasks. More complex processes are performed by sequences of tools, often in the form of pipelines which avoid the need for temporary files.

**two-digit display.** Two seven-segment light-emitting diodes (LEDs) on the operating panel used to track the progress of power#on self-tests (POSTs).

## U

**UNIX® Operating System** A multi-user, multi-tasking interactive operating system created at AT&T Bell Laboratories that has been widely used and developed by universities, and that now is becoming increasingly popular in a wide range of commercial applications. See **Kernel, Shell, Library, Pipes, Filters**.

**user interface** The component of the AIX Family Definition that describes common user interface functions for the AIX PS/2, AIX/RT, and AIX/370 operating systems.

**/usr/grp®** One of the oldest, and still active, user groups for the UNIX operating systems. IBM is a member of /usr/grp.

**uucp** A set of AIX utilities allowing

- Autodial of remote systems
- Transfer of files
- Execution of commands on the remote system
- Reasonable security.

## V

**vi** Visual editor. A character editor with a very powerful collection of editing commands optimized for ASCII terminals; associated with BSD versions of the UNIX operating system.

**visual editor** An optional editor provided with AIX in which changes are made by modifying an image of the file on the screen, rather than through the exclusive use of commands.

## W

**wild card** A metacharacter used to specify a set of replacement characters and thus a set of file names. For example "\*" is any zero or more characters and "?" is any one character.

**window** A rectangular area of the screen in which the dialog between you and a given application is displayed.

**working directory** The directory from which file searches are begun if a complete pathname is not specified. Controlled by the cd (change directory) command.

**workstation** A device that includes a keyboard from which an operator can send information to the system, and a display screen on which an operator can see the information sent to or received from the computer.

**write** Sending data to an I/O device.

**write permission** Permission to modify a file or directory.

## X

**X/Open™** An international consortium, including many suppliers of computer systems, concerned with the selection and adoption of open system standards for computing applications. IBM is a corporate sponsor of X/Open. See **Common Application Environment**.

**X Windows.** IBM's implementation of the X Window System developed at the Massachusetts Institute of Technology with the support of IBM and DEC, that gives users "windows" into applications and processes not located only or specifically on their own console or computer system. X-Windows is a powerful vehicle for distributing applications among users on heterogeneous networks.

## Y

**yacc.** "Yet Another Compiler - Compiler". For producing new command interfaces.

## Z

**zeroeth argument** The command name; the argument before the first.



