

*High Availability Cluster
Multi-Processing for AIX*

Concepts and Facilities Guide

Version 5.2

Third Edition (July 2004)

Before using the information in this book, read the general information in [Notices for HACMP Concepts and Facilities Guide](#).

This edition applies to HACMP for AIX, version 5.2 and to all subsequent releases of this product until otherwise indicated in new editions.

© Copyright International Business Machines Corporation 1998, 2004. All rights reserved.

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

Contents

About This Guide	9
Chapter 1: HACMP for AIX	13
High Availability Cluster Multi-Processing for AIX	13
High Availability and Hardware Availability	13
High Availability vs. Fault Tolerance	14
Role of HACMP	14
Cluster Multi-Processing	15
The Availability Costs and Benefits Continuum	15
Enhancing Availability with the AIX Software	16
Journaled File System and Enhanced Journaled File System	17
Disk Mirroring	17
Process Control	17
Error Notification	17
High Availability with HACMP	18
HACMP and HACMP/ES	18
Physical Components of an HACMP Cluster	18
Nodes	19
Shared External Disk Devices	20
Networks	20
Clients	21
Goal of HACMP: Eliminating Scheduled Downtime	21
Chapter 2: HACMP Cluster Nodes, Networks and Heartbeating	23
Cluster Nodes and Cluster Sites	23
Nodes	23
Sites	24
Cluster Networks	25
Physical and Logical Networks	25
HACMP Communication Interfaces	27
HACMP Communication Devices	28
Subnet Routing Requirements in HACMP	28
Service IP Address/Label	29
IP Alias	29
IP Address Takeover	29
IPAT and Service IP Labels	30
IP Address Takeover via IP Aliases	30

IP Address Takeover via IP Replacement	31
Heartbeating Over Networks and Disks	31
Heartbeating in HACMP: Overview	31
Heartbeating over TCP/IP Networks	32
Heartbeating over Point-to-Point Networks	32

Chapter 3: HACMP Resources and Resource Groups 35

Cluster Resources: Identifying and Keeping Available	35
Identifying Cluster Resources	35
Types of Cluster Resources	36
Volume Groups	36
Logical Volumes	36
Filesystems	36
Applications	37
Service IP Labels/Addresses	38
Tape Resources	38
Communication Links	39
Cluster Resource Groups	39
Participating Nodelist	40
Default Node Priority	40
Home Node	41
Startup, Fallover and Fallback	41
Resource Group Policies and Attributes	42
Overview	43
Resource Group Startup, Fallover and Fallback	44
Settling Time, Dynamic Node Priority and Fallback Timer	44
Distribution Policy	45
Upgrading from Previous Releases	46
Cluster Networks and Resource Groups	46
Sites and Resource Groups	47
Resource Group Dependencies	47
Child and Parent Resource Groups	48

Chapter 4: HACMP Cluster Hardware and Software 51

Enhancing Availability with IBM Hardware	51
IBM pSeries	51
RS/6000 SP System	52
Disk Subsystems	52
HACMP Required and Supported Hardware	54
HACMP Cluster Software	54
Software Components of an HACMP Node	55
HACMP Software Components	56
Cluster Manager	56
Cluster Secure Communication Subsystem	57
IBM Reliable Scalable Cluster Technology Availability Services ..	58
Cluster SMUX Peer and SNMP Monitoring Programs	59
Cluster Information Program	60

	Highly Available NFS Server	61
	Shared External Disk Access	61
	Concurrent Resource Manager	64
	Complementary Cluster Software	65
Chapter 5:	Ensuring Application Availability	67
	Overview	67
	Eliminating Single Points of Failure in an HACMP Cluster	68
	Potential Single Points of Failure in an HACMP Cluster	68
	Eliminating Nodes as a Single Point of Failure	68
	Eliminating Applications as a Single Point of Failure	73
	Eliminating Communication Interfaces as a Single Point of Failure	74
	Eliminating Networks as a Single Point of Failure	76
	Eliminating Disks and Disk Adapters as a Single Point of Failure .	78
	Minimizing Scheduled Downtime with HACMP	79
	Dynamic Automatic Reconfiguration (DARE)	79
	Resource Group Management	82
	Cluster Single Point of Control (C-SPOC)	84
	Dynamic Adapter Swap	84
	Minimizing Unscheduled Downtime	84
	Recovering Resource Groups on Node Startup	85
	Fast Recovery	85
	Delayed Fallback Timer for Resource Groups	85
	Minimizing Takeover Time: Fast Disk Takeover	86
	Maximizing Disaster Recovery	86
	Cross-Site LVM Mirroring	86
	Cluster Events	87
	Processing Cluster Events	88
	Emulating Cluster Events	89
	Customizing Event Processing	89
	Customizing Event Duration	90
Chapter 6:	HACMP Cluster Configurations	91
	Sample Cluster Configurations	91
	Standby Configurations	91
	Standby Configurations: Example 1	92
	Standby Configurations: Example 2	94
	Takeover Configurations	95
	One-Sided Takeover	95
	Mutual Takeover	96
	Two-Node Mutual Takeover Configuration for Concurrent Access	97
	Eight-Node Mutual Takeover Configuration for Concurrent Access	97
	Cluster Configurations with Multi-Tiered Applications	98
	Cross-Site LVM Mirror Configurations for Disaster Recovery ..	99
	Cluster Configurations with Dynamic LPARs	100

Chapter 7:	HACMP Configuration Process and Facilities	103
	Information You Provide to HACMP	103
	Information on Physical Configuration of a Cluster	103
	AIX Configuration Information	104
	Establishing the Initial Communication Path	104
	Information Discovered by HACMP	105
	Cluster Configuration Options: Standard and Extended	105
	Configuring an HACMP Cluster Using the Standard Configuration Path	105
	Configuring an HACMP Cluster Using the Extended Configuration Path	106
	Overview: HACMP Administrative Facilities	106
	Cluster Security	106
	Installation, Configuration and Management Tools	106
	Two-Node Cluster Configuration Assistant	107
	Planning Worksheets	107
	SMIT Interface	107
	Web-Based SMIT Interface	108
	HACMP System Management (C-SPOC)	108
	Cluster Snapshot Utility	109
	Customized Event Processing	110
	Resource Group Management Utility	110
	HACMP File Collection Management	110
	Monitoring Tools	111
	Cluster Status Utility (clstat)	112
	HAView Cluster Monitoring Utility	112
	Cluster Monitoring and Administration with Tivoli Framework	113
	Application Availability Analysis Tool	113
	Persistent Node IP Labels	113
	HACMP Verification and Synchronization	114
	Troubleshooting Tools	115
	Cluster Diagnostic Utility	116
	Log Files	117
	Resetting HACMP Tunable Values	118
	Cluster Status Information File	119
	Automatic Error Notification	119
	Custom Pager Notification	120
	User-Defined Events	120
	Event Summaries	121
	Trace Facility	121
	Test Tool	121
	Emulation Tools	121
	HACMP Event Emulator	121
	Emulation of Error Log Driven Events	123
Chapter 8:	HACMP 5.2: Summary of Changes	125
	List of New Features	125

New Features that Enhance Ease of Use 125
New Features that Enhance Performance 129
New Features that Enhance Geographic Distance Capability 131
Discontinued Features 132
Where You Go From Here 132

Index

135

Contents

About This Guide

This guide introduces High Availability Cluster Multi-Processing for AIX (HACMP™), Version 5.2.

Who Should Use This Guide

System administrators, system engineers, and other information systems professionals who want to learn about features and functionality provided by the HACMP software should read this guide.

Highlighting

This guide uses the following highlighting conventions:

<i>Italic</i>	Identifies new terms or concepts, or indicates emphasis.
Bold	Identifies routines, commands, keywords, files, directories, menu items, and other items whose actual names are predefined by the system.
Monospace	Identifies examples of specific data values, examples of text similar to what you might see displayed, examples of program code similar to what you might write as a programmer, messages from the system, or information that you should actually type.

ISO 9000

ISO 9000 registered quality systems were used in the development and manufacturing of this product.

Related Publications

The following publications come with your HACMP system. They provide additional information about HACMP:

- *Release Notes* in `/usr/es/lpp/cluster/doc/release_notes` describe hardware and software requirements and last-minute information about installation, product usage, and known issues.
- *HACMP for AIX: Planning and Installation Guide - SC23-4861*
- *HACMP for AIX: Administration and Troubleshooting Guide - SC23-4862*
- *HACMP for AIX: Programming Client Applications - SC23-4865*
- *HACMP for AIX: Glossary - SC23-4867*
- *IBM International Program License Agreement*

The following publications offer more information on IBM technology related to or used by HACMP for AIX:

- *RS/6000 SP High Availability Infrastructure, SG24-4838*
- *IBM Reliable Scalable Cluster Technology for AIX 5L and Linux: Group Services Programming Guide and Reference, SA22-7888*
- *IBM Reliable Scalable Cluster Technology for AIX 5L and Linux: Administration Guide, SA22-7889*
- *IBM Reliable Scalable Cluster Technology for AIX 5L: Technical Reference, SA22-7890*
- *IBM Reliable Scalable Cluster Technology for AIX 5L: Messages, GA22-7891*

HACMP/XD

The HACMP/XD feature provides software solutions for disaster recovery. Added to the base HACMP software, they enable a cluster to operate over extended distances at two sites:

- **HACMP/XD for ESS PPRC** increases data availability for IBM TotalStorage Enterprise Storage Server (ESS) volumes that use Peer-to-Peer Remote Copy (PPRC) to copy data to a remote site for disaster recovery purposes. HACMP/XD for ESS PPRC takes advantage of the PPRC fallover/fallback functions and HACMP cluster management to reduce downtime and recovery time during disaster recovery.
- **HACMP/XD for HAGEO Technology** provides unlimited distance data mirroring. It is based on the IBM High Availability Geographic Cluster for AIX (HAGEO) v 2.4 product. HACMP/XD for HAGEO Technology extends an HACMP cluster to encompass two physically separate data centers. Data entered at one site is sent across a point-to-point TCP/IP network and mirrored at a second, geographically distant location.

Documentation for HACMP/XD

Documentation for HACMP/XD includes:

- *HACMP/XD for ESS PPRC: Planning and Administration Guide*
- *HACMP/XD for HAGEO Technology: Concepts and Facilities Guide*
- *HACMP/XD for HAGEO Technology: Planning and Administration Guide*

Accessing Publications

On the World Wide Web, use the following URL to access an online library of documentation covering AIX, IBM eServer pSeries, and related products:

<http://www.ibm.com/servers/aix/library/>

Trademarks

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

- AFS
- AIX

- AIX 5L
- DFS
- **@server**
- Enterprise Storage Server
- HACMP
- IBM
- NetView
- pSeries
- RS/6000
- Scalable POWERParallel Systems
- Shark
- SP
- xSeries

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

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

Chapter 1: HACMP for AIX

This chapter discusses the concepts of high availability and cluster multi-processing, presents the HACMP cluster diagram, and describes an HACMP cluster from a functional perspective.

High Availability Cluster Multi-Processing for AIX

The IBM tool for building UNIX-based mission-critical computing platforms is the HACMP software. The HACMP software ensures that critical resources, such as applications, are available for processing. HACMP has two major components: high availability (HA) and cluster multi-processing (CMP).

The primary reason to create HACMP clusters is to provide a highly available environment for mission-critical applications. For example, an HACMP cluster could run a database server program which services client applications. The clients send queries to the server program which responds to their requests by accessing a database, stored on a shared external disk.

In an HACMP cluster, to ensure the availability of these applications, the applications are put under HACMP control. HACMP takes measures to ensure that the applications *remain available* to client processes even if a component in a cluster fails. To ensure availability, in case of a component failure, HACMP moves the application (along with resources that ensure access to the application) to another node in the cluster.

High Availability and Hardware Availability

High availability is sometimes confused with simple hardware availability. Fault tolerant, redundant systems (such as RAID), and dynamic switching technologies (such as DLPAR) provide recovery of certain hardware failures, but do not provide the full scope of error detection and recovery required to keep a complex application highly available.

A modern, complex application requires access to all of these elements:

- Nodes (CPU, memory)
- Network interfaces (including external devices in the network topology)
- Disk or storage devices.

Recent surveys of the causes of downtime show that actual hardware failures account for only a small percentage of unplanned outages. Other contributing factors include:

- Operator errors
- Environmental problems
- Application and operating system errors.

Reliable and recoverable hardware simply cannot protect against failures of all these different aspects of the configuration. Keeping these varied elements—and therefore the application—highly available requires:

- Thorough and complete planning of the physical and logical procedures for access and operation of the resources on which the application depends. These procedures help to avoid failures in the first place.
- A monitoring and recovery package which automates the detection and recovery from errors.
- A well-controlled process for maintaining the hardware and software aspects of the cluster configuration while keeping the application available.

High Availability vs. Fault Tolerance

Fault tolerance relies on specialized hardware to detect a hardware fault and instantaneously switch to a redundant hardware component—whether the failed component is a processor, memory board, power supply, I/O subsystem, or storage subsystem. Although this cutover is apparently seamless and offers non-stop service, a high premium is paid in both hardware cost and performance because the redundant components do no processing. More importantly, the fault tolerant model does not address software failures, by far the most common reason for downtime.

High availability views availability not as a series of replicated physical components, but rather as a set of system-wide, shared resources that cooperate to guarantee essential services. High availability combines software with industry-standard hardware to minimize downtime by quickly restoring essential services when a system, component, or application fails. While not instantaneous, services are restored rapidly, often in less than a minute.

The difference between fault tolerance and high availability, then, is this: a fault tolerant environment has no service interruption, while a highly available environment has a minimal service interruption. Many sites are willing to absorb a small amount of downtime with high availability rather than pay the much higher cost of providing fault tolerance. Additionally, in most highly available configurations, the backup processors are available for use during normal operation.

High availability systems are an excellent solution for applications that can withstand a short interruption should a failure occur, but which must be restored quickly. Some industries have applications so time-critical that they cannot withstand even a few seconds of downtime. Many other industries, however, can withstand small periods of time when their database is unavailable. For those industries, HACMP can provide the necessary continuity of service without total redundancy.

Role of HACMP

HACMP helps you with the following:

- The HACMP planning process and documentation include tips and advice on the best practices for installing and maintaining a highly available HACMP cluster.
- Once the cluster is operational, HACMP provides the automated monitoring and recovery for all the resources on which the application depends.
- HACMP provides a full set of tools for maintaining the cluster while keeping the application available to clients.

HACMP lets you:

- Quickly and easily set up a basic two-node HACMP cluster by using the Two-Node Cluster Configuration Assistant.
- Set up an HACMP environment using online planning worksheets to simplify the initial planning and setup.
- Test your HACMP configuration by using the Cluster Test Tool. You can evaluate how a cluster behaves under a set of specified circumstances, such as when a node becomes inaccessible, a network becomes inaccessible, and so forth.
- Ensure high availability of applications by eliminating single points of failure in an HACMP environment.
- Leverage high availability features available in AIX.
- Manage how a cluster handles component failures.
- Secure cluster communications.
- Set up fast disk takeover for volume groups managed by the Logical Volume Manager (LVM).
- Monitor HACMP components and diagnose problems that may occur.

For a general overview of *all* HACMP features, see the IBM website:

http://www.ibm.com/servers/aix/products/ibmsw/high_avail_network/hacmp.html

For a list of *new features* in HACMP 5.2, see [Chapter 8: HACMP 5.2: Summary of Changes](#).

Cluster Multi-Processing

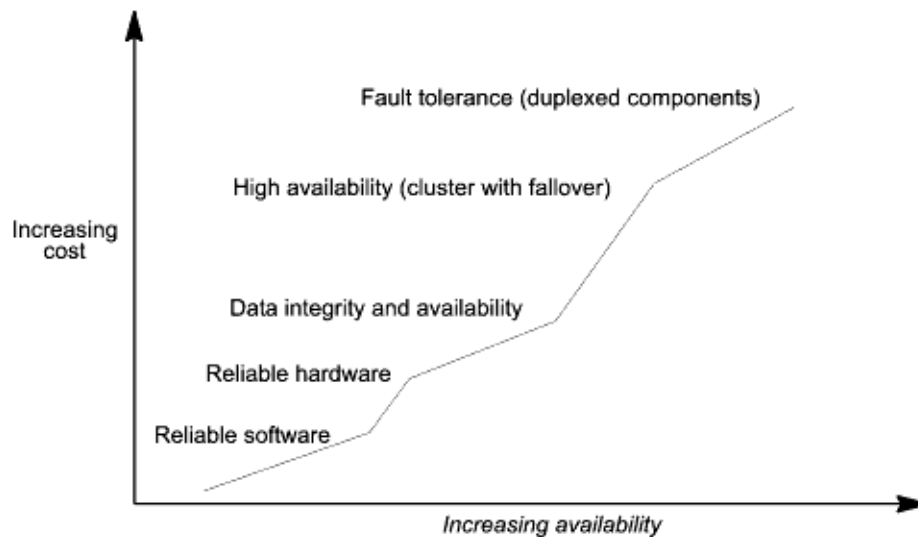
Cluster multi-processing is a group of loosely coupled machines networked together, sharing disk resources. In a cluster, multiple server machines cooperate to provide a set of services or resources to clients.

Clustering two or more servers to back up critical applications is a cost-effective high availability option. You can use more of your site's computing power while ensuring that critical applications resume operations after a minimal interruption caused by a hardware or software failure.

Cluster multi-processing also provides a gradual, scalable growth path. It is easy to add a processor to the cluster to share the growing workload. You can also upgrade one or more of the processors in the cluster to a more powerful model. If you were using a fault tolerant strategy, you must add *two* processors, one as a redundant backup that does no processing during normal operations.

The Availability Costs and Benefits Continuum

The following figure shows the costs and benefits of availability technologies.



Cost and Benefits of Availability Technologies

As you can see, availability is not an all-or-nothing proposition. Think of availability as a continuum. Reliable hardware and software provide the base level of availability. Advanced features such as RAID devices provide an enhanced level of availability. High availability software provides near continuous access to data and applications. Fault tolerant systems ensure the constant availability of the entire system, but at a higher cost.

Enhancing Availability with the AIX Software

HACMP takes advantage of the features in AIX—the high-performance UNIX operating system.

AIX 5.1 and 5.2 add new functionality to further improve security and system availability. This includes improved availability of mirrored data and enhancements to Workload Manager that help solve problems of mixed workloads by dynamically providing resource availability to critical applications. Used with the IBM eserver pSeries, HACMP can provide both horizontal and vertical scalability without downtime.

The AIX operating system provides numerous features designed to increase system availability by lessening the impact of both planned (data backup, system administration) and unplanned (hardware or software failure) downtime. These features include:

- Journaled File System and Enhanced Journaled File System
- Disk mirroring
- Process control
- Error notification.

Journalized File System and Enhanced Journalized File System

The AIX native filesystem, the Journalized File System (JFS), uses database journaling techniques to maintain its structural integrity. System or software failures do not leave the filesystem in an unmanageable condition. When rebuilding the filesystem after a major failure, AIX uses the JFS log to restore the filesystem to its last consistent state. Journaling thus provides faster recovery than the standard UNIX filesystem consistency check (fsck) utility.

In addition, the Enhanced Journalized File System (JFS2) is available in AIX. For more information, see the section [Journalized Filesystem and Enhanced Journalized Filesystem](#) in [Chapter 3: HACMP Resources and Resource Groups](#).

Disk Mirroring

Disk mirroring software provides data integrity and online backup capability. It prevents data loss due to disk failure by maintaining up to three copies of data on separate disks so that data is still accessible after any single disk fails. Disk mirroring is transparent to the application. No application modification is necessary because mirrored and conventional disks appear the same to the application.

Process Control

The AIX System Resource Controller (SRC) monitors and controls key processes. The SRC can detect when a process terminates abnormally, log the termination, pass messages to a notification program, and restart the process on a backup processor.

Error Notification

The AIX Error Notification facility detects errors, such as network and disk adapter failures, and triggers a predefined response to the failure.

HACMP builds on this AIX feature by providing:

- Automatically created error notification methods for volume groups that you configure in HACMP. These error notification methods let HACMP react to certain volume group failures and provide recovery.

HACMP also configures automatic error notification methods for those volume groups that do not themselves are configured to HACMP but that contain the corresponding filesystems configured to be kept highly available. This allows to automatically respond to the volume group failures by recovering the filesystems.

For more information on error notification and how it is used in HACMP, see [Eliminating Disks and Disk Adapters as a Single Point of Failure](#) in [Chapter 5: Ensuring Application Availability](#).

- Error emulation function. It allows you to test the predefined response without causing the error to occur. You also have an option to automatically configure notification methods for a set of device errors in one step.

High Availability with HACMP

The IBM HACMP software provides a low-cost commercial computing environment that ensures that mission-critical applications can recover quickly from hardware and software failures. The HACMP software is a high availability system that ensures that critical resources are available for processing. High availability combines custom software with industry-standard hardware to minimize downtime by quickly restoring services when a system, component, or application fails. While not instantaneous, the restoration of service is rapid, usually 30 to 300 seconds.

HACMP and HACMP/ES

HACMP 5.2 includes all the features of the product that was previously referred to as HACMP/ES (Enhanced Scalability). The product previously referred to as HACMP (Classic) is no longer available.

Note: Prior to version 5.1, the HACMP for AIX software included four features: HAS and CRM with core filesets named `cluster.base*`; and ES and ESCRM with core filesets named `cluster.es*`. Starting with HACMP 5.1, the HAS, CRM and ES features are no longer available, and the ESCRM feature is now called HACMP.

To summarize, HACMP 5.2 has the following characteristics:

- Includes all the features of ESCRM 4.5 in addition to the new functionality added in HACMP versions 5.1 and 5.2.
- Takes advantage of the enhanced Reliable Scalable Cluster Technology (RSCT). RSCT provides facilities for monitoring node membership; network interface and communication interface health; and event notification, synchronization and coordination via reliable messaging.
- HACMP clusters with both non-concurrent and concurrent access can have up to 32 nodes.
- Is supported on AIX 5.1 and AIX 5.2.

Physical Components of an HACMP Cluster

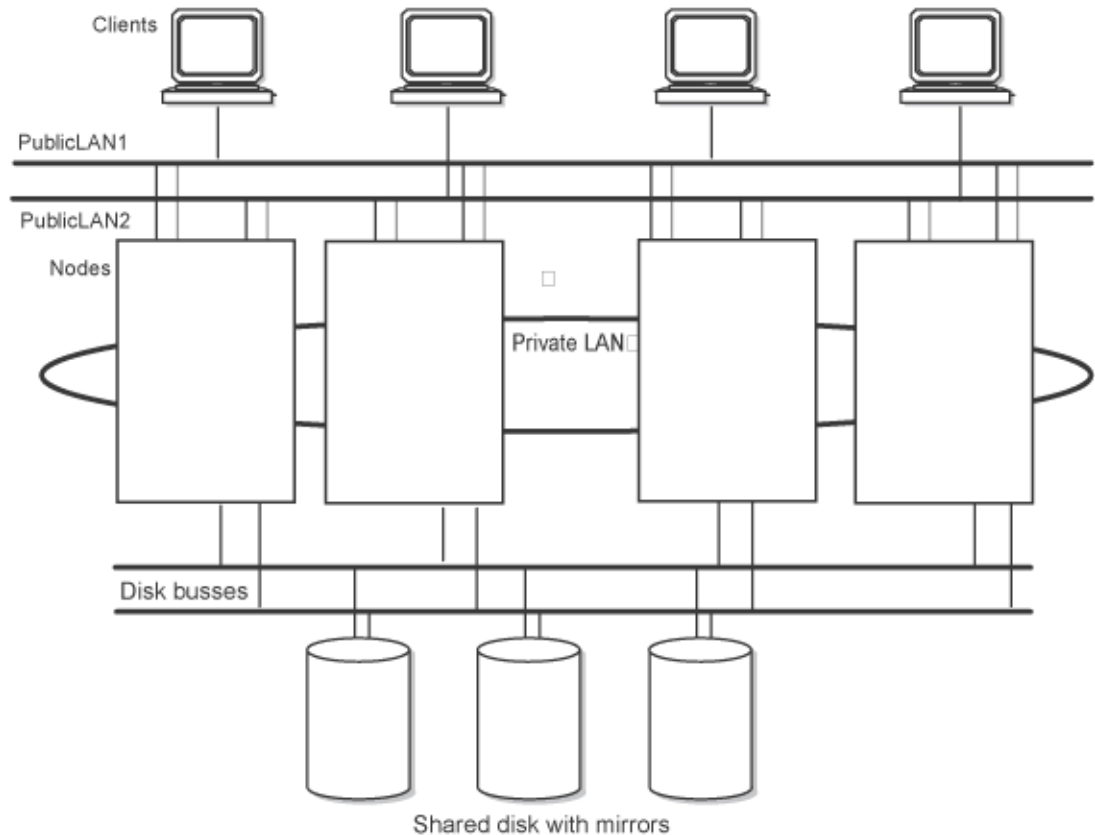
HACMP provides a highly available environment by identifying a set of resources essential to uninterrupted processing, and by defining a protocol that nodes use to collaborate to ensure that these resources are available. HACMP extends the clustering model by defining relationships among cooperating processors where one processor provides the service offered by a peer should the peer be unable to do so.

As shown in the following figure, an HACMP cluster is made up of the following physical components:

- Nodes
- Shared external disk devices
- Networks

- Network interfaces
- Clients.

The HACMP software allows you to combine physical components into a wide range of cluster configurations, providing you with flexibility in building a cluster that meets your processing requirements. The following figure shows one example of an HACMP cluster. Other HACMP clusters could look very different—depending on the number of processors, the choice of networking and disk technologies, and so on.



An Example of an HACMP Cluster

[Chapter 6: HACMP Cluster Configurations](#), describes examples of some types of cluster configurations supported by the HACMP software.

Nodes

Nodes form the core of an HACMP cluster. A node is a processor that runs AIX, the HACMP software, and the application software.

In an HACMP cluster, up to 32 pSeries divided into LPARS, RS/6000 or pSeries standalone systems, systems that run Parallel System Support Program (PSSP), or a combination of these cooperate to provide a set of services or resources to other entities. Clustering these servers to back up critical applications is a cost-effective high availability option. A business can use more of its computing power while ensuring that its critical applications resume running after a short interruption caused by a hardware or software failure.

In an HACMP cluster, each node is identified by a unique name. A node may own a set of resources—disks, volume groups, filesystems, networks, network addresses, and applications. Cluster resources are discussed in detail in [Chapter 3: HACMP Resources and Resource Groups](#). Typically, a node runs a server or a back end application that accesses data on the shared external disks. Applications are discussed in [Chapter 5: Ensuring Application Availability](#).

The HACMP software supports from two to thirty-two nodes in a cluster, depending on the disk technology used for the shared external disks.

A node in an HACMP cluster has several layers of software components. For the detailed description of these components, see the section [Software Components of an HACMP Node](#) in [Chapter 4: HACMP Cluster Hardware and Software](#).

Shared External Disk Devices

Each node has access to one or more shared external disk devices. A *shared external disk device* is a disk physically connected to multiple nodes. The shared disk stores mission-critical data, typically mirrored or RAID-configured for data redundancy. A node in an HACMP cluster must also have internal disks that store the operating system and application binaries, but these disks are not shared.

Depending on the type of disk used, the HACMP software supports the following types of access to shared external disk devices—non-concurrent access and concurrent access.

- *In non-concurrent access environments*, only one connection is active at any given time, and the node with the active connection owns the disk. When a node fails, disk takeover occurs when the node that currently owns the disk leaves the cluster and a surviving node assumes ownership of the shared disk.
- *In concurrent access environments*, the shared disks are actively connected to more than one node simultaneously. Therefore, when a node fails, disk takeover is not required.

Note that in such environments, either *all* nodes defined in the cluster can be part of the concurrent access, or *only a subset* of cluster nodes can have access to shared disks. In the second case, you configure resource groups only on those nodes that have shared disk access.

The differences between these methods are explained more fully in the section [Shared External Disk Access](#) in [Chapter 4: HACMP Cluster Hardware and Software](#).

Networks

As an independent, layered component of AIX, the HACMP software is designed to work with any TCP/IP-based network. Nodes in an HACMP cluster use the network to:

- Allow clients to access the cluster nodes
- Enable cluster nodes to exchange heartbeat messages
- Serialize access to data (in concurrent access environments).

The HACMP software has been tested with Ethernet, Token-Ring, ATM, and other networks.

Types of Networks

The HACMP software defines two types of communication networks, characterized by whether these networks use communication interfaces based on the TCP/IP subsystem (TCP/IP-based) or communication devices based on non-TCP/IP subsystems (device-based).

TCP/IP is a communications subsystem that lets you set up local area and wide area networks. TCP/IP provides facilities that make the computer system serve as an Internet host which can attach to a network and communicate with other Internet hosts.

- *TCP/IP-based network.* Connects two or more server nodes, and optionally allows client access to these cluster nodes, using the TCP/IP protocol. Ethernet, Token-Ring, ATM, HP Switch and SP Switch networks are defined as TCP/IP-based networks.
- *Device-based network.* Provides a point-to-point connection between two cluster nodes for HACMP control messages and *heartbeat* traffic. Device-based networks do not use the TCP/IP protocol, and therefore continue to provide communications between nodes in the event the TCP/IP subsystem on a server node fails.

Target mode SCSI devices, Target Mode SSA devices, *disk heartbeat* devices, or RS232 point-to-point devices are defined as device-based networks.

Clients

A client is a processor that can access the nodes in a cluster over a local area network. Clients each run a “front end” or client application that queries the server application running on the cluster node.

The HACMP software provides a highly available environment for critical data and applications on cluster nodes. *The HACMP software does not make the clients themselves highly available.* AIX clients can use the Client Information (Cinfo) services to receive notice of cluster events. Cinfo provides an API that displays cluster status information.

HACMP provides a cluster status utility, the `/usr/es/sbin/cluster/clstat`. It is based on Cinfo and reports the status of key cluster components—the cluster itself, the nodes in the cluster, the network interfaces connected to the nodes, and the resource groups on each node. The cluster status interface of the `clstat` utility includes web-based, Motif-based and ASCII-based versions.

See [Cluster Information Program](#) in [Chapter 4: HACMP Cluster Hardware and Software](#), for more information about how Cinfo obtains cluster status information.

Goal of HACMP: Eliminating Scheduled Downtime

The primary goal of high availability clustering software is to minimize, or ideally, eliminate, the need to take your resources out of service during maintenance and reconfiguration activities.

HACMP software optimizes availability by allowing for the *dynamic reconfiguration* of running clusters. Most routine cluster maintenance tasks, such as adding or removing a node or changing the priority of nodes participating in a resource group, can be applied to an active cluster without stopping and restarting cluster services.

1

HACMP for AIX

Goal of HACMP: Eliminating Scheduled Downtime

In addition, you can keep an HACMP cluster online while making configuration changes by using the *Cluster Single Point of Control (C-SPOC)* facility. C-SPOC makes cluster management easier, as it allows you to make changes to shared volume groups, users, and groups across the cluster from a single node. The changes are propagated transparently to other cluster nodes.

Chapter 2: HACMP Cluster Nodes, Networks and Heartbeating

This chapter introduces major cluster topology-related *concepts* and definitions that are used throughout the documentation and in the HACMP user interface.

The information in this chapter is organized as follows:

- [Cluster Nodes and Cluster Sites](#)
- [Cluster Networks](#)
- [Subnet Routing Requirements in HACMP](#)
- [IP Address Takeover](#)
- [IP Address Takeover via IP Aliases](#)
- [IP Address Takeover via IP Replacement](#)
- [Heartbeating Over Networks and Disks](#)

Cluster Nodes and Cluster Sites

A typical HACMP cluster environment consists of nodes which can serve as clients or servers. If you are using the HACMP/XD software or LVM cross-site mirroring, sites or groups of nodes become part of the cluster topology.

Nodes

A *node* is a processor that runs both AIX and the HACMP software. Nodes may share a set of resources—disks, volume groups, filesystems, networks, network IP addresses, and applications.

The HACMP software supports from two to thirty-two nodes in a cluster. In an HACMP cluster, each node is identified by a unique name. In HACMP, a node name and a hostname can usually be the same.

Nodes serve as core physical components of an HACMP cluster. For more information on nodes and hardware, see the section [Nodes](#) in [Chapter 1: HACMP for AIX](#).

Two types of nodes are defined:

- *Server nodes* form the core of an HACMP cluster. Server nodes run services or back end applications that access data on the shared external disks.
- *Client nodes* run “front end” applications that retrieve data from the services provided by the server nodes. Client nodes can run HACMP software to monitor the health of the nodes, and to react to failures.

Server Nodes

A cluster *server node* usually runs an application that accesses data on the shared external disks. Server nodes run HACMP daemons and keep resources highly available. Typically, applications are run, storage is shared between these nodes, and *clients* connect to the server nodes through a *Service IP Address*.

Client Nodes

A full high availability solution typically includes the client machine that uses services provided by the servers. Client nodes can be divided into two categories: naive and intelligent.

- A naive client views the cluster as a single entity. If a server fails, the client must be restarted, or at least must reconnect to the server.
- An intelligent client is cluster-aware. A cluster-aware client reacts appropriately in the face of a server failure, connecting to an alternate server, perhaps masking the failure from the user. Such an intelligent client must have knowledge of the cluster state.

HACMP extends the cluster paradigm to clients by providing both dynamic cluster configuration reporting and notification of cluster state changes, such as changes in subsystems or node failure.

Sites

You can define a group of one or more server nodes as belonging to a *site*. The site becomes a component, like a node or a network, that is known to the HACMP software. HACMP supports clusters divided into two sites.

Using sites, you can configure the cross-site LVM mirroring. You configure logical volume mirrors between physical volumes in separate storage arrays, and specify to HACMP which physical volumes are located at each site. Later when you use C-SPOC to create new logical volumes, HACMP automatically displays the site location of each defined physical volume, making it easier to select volumes from different sites for LVM mirrors. For more information on cross-site LVM mirroring, see the *Planning and Installation Guide*.

In addition, the HACMP/XD feature provides two distinct software solutions for disaster recovery. These solutions enable an HACMP cluster to operate over extended distances at two sites.

- **HACMP/XD for ESS PPRC** increases data availability for IBM TotalStorage Enterprise Storage Server (ESS) volumes that use Peer-to-Peer Remote Copy (PPRC) to copy data to a remote site for disaster recovery purposes. HACMP/XD for ESS PPRC takes advantage of the PPRC fallover/fallback functions and HACMP cluster management to reduce downtime and recovery time during disaster recovery.

When PPRC is used for data mirroring between sites, the physical distance between sites is limited to the capabilities of the ESS hardware.

- **HACMP/XD for HAGEO Technology** uses the TCP/IP network to enable *unlimited distance* for data mirroring between sites. (Note that although the distance is unlimited, practical restrictions exist on the bandwidth and throughput capabilities of the network).

This technology is based on the IBM High Availability Geographic Cluster for AIX (HAGEO) v 2.4 product. HACMP/XD for HAGEO Technology extends an HACMP cluster to encompass two physically separate data centers. Data entered at one site is sent across a point-to-point IP network and mirrored at a second, geographically distant location.

Each site can be a backup data center for the other, maintaining an updated copy of essential data and running key applications. If a disaster disables one site, the data is available within minutes at the other site. The HACMP/XD software solutions thus increase the level of availability provided by the HACMP software by enabling it to recognize and handle a site failure, to continue processing even though one of the sites has failed, and to reintegrate the failed site back into the cluster.

If sites are configured in the cluster, you can use the Resource Group Management utility to bring a resource online, take it offline, or move it to another node *only within the boundaries of a site*. You cannot manually move resource groups between sites.

For information on HACMP/XD for ESS PPRC, and on HACMP/XD for HAGEO, see the documentation for each of those solutions.

Cluster Networks

Cluster nodes communicate with each other over communication networks. If one of the physical network interface cards on a node on a network fails, HACMP preserves the communication to the node by transferring the traffic to another physical network interface card on the same node. If a “connection” to the node fails, HACMP transfers resources to another node to which it has access.

In addition, RSCT sends heartbeats between the nodes over the cluster networks to periodically check on the health of the cluster nodes themselves. If HACMP detects no heartbeats from a node, a node is considered as failed and resources are automatically transferred to another node. For more information, see [Heartbeating Over Networks and Disks](#) in this chapter.

We highly recommend to configuring multiple communication paths between the nodes in the cluster. Having multiple communication networks prevents *cluster partitioning*, in which the nodes within each partition form their own entity. In a partitioned cluster, it is possible that nodes in each partition could allow simultaneous non-synchronized access to the same data. This can potentially lead to different views of data from different nodes.

Physical and Logical Networks

A physical network connects two or more physical network interfaces. There are many types of physical networks, and HACMP broadly categorizes them as those that use the TCP/IP protocol, and those that do not:

- *TCP/IP-based*, such as Ethernet or Token Ring
- *Device-based*, such as RS-232 or TM-SSA.

As stated in the previous section, configuring multiple TCP-IP-based networks helps to prevent cluster partitioning. Multiple device-based networks also help to prevent partitioned clusters by providing additional communications paths in cases when the TCP/IP-based network connections become congested or severed between cluster nodes.

Note: If you are considering a cluster where the physical networks use *external* networking devices to route packets from one network to another, consider the following: When you configure an HACMP cluster, HACMP verifies the connectivity and access to all interfaces defined on a particular physical network. However, HACMP cannot determine the presence of external network devices such as bridges and routers in the network path between cluster nodes. If the networks have external networking devices, ensure that you are using devices that are highly available and redundant so that they do not create a *single point of failure* in the HACMP cluster.

A *logical network* is a portion of a physical network that connects two or more logical network interfaces/devices. A logical network interface/device is the software entity that is known by an operating system. There is a one-to-one mapping between a physical network interface/device and a logical network interface/device. Each logical network interface can exchange packets with each other logical network interface on the same logical network.

If a subset of logical network interfaces on the logical network needs to communicate with each other (but with no one else) while sharing the same physical network, *subnets* are used. A *subnet mask* defines the part of the IP address that determines whether one logical network interface can send packets to another logical network interface on the same logical network.

Logical Networks in HACMP

HACMP has its own, similar concept of a logical network. All logical network interfaces in an HACMP network can communicate HACMP packets with each other directly. Each logical network is identified by a unique name. If you use an automatic discovery function for HACMP cluster configuration, HACMP assigns a name to each HACMP logical network it discovers, such as `net_ether_01`.

An HACMP logical network may contain one or more subnets. RSCT takes care of routing packets between logical subnets.

For more information on RSCT, see [Chapter 4: HACMP Cluster Hardware and Software](#).

Global Networks

A *global network* is a combination of multiple HACMP networks. The HACMP networks may be composed of any combination of physically different networks, and/or different logical networks (subnets), as long as they share the same network type, (for example, ethernet). HACMP treats the combined global network as a single network. RSCT handles the routing between the networks defined in a global network.

Global networks cannot be defined for all IP-based networks but only for those IP-based networks that are used for heartbeating.

Having multiple heartbeat paths between cluster nodes reduces the chance that the loss of any single network will result in a partitioned cluster. An example of when this would be useful is a typical configuration of the SP Administrative Ethernet on two separate SP systems.

Local and Global Network Failures

When a failure occurs on a cluster network, HACMP uses *network failure events* to manage such cases. HACMP watches for and distinguishes between two types of network failures: local network failure and global network failure events.

Local Network Failure

A *local network failure* is an HACMP event that describes the case in which packets cannot be sent or received by *one* node over an HACMP logical network. This may occur, for instance, if all of the node's network interface cards participating in the particular HACMP logical network fail. Note that in the case of a local network failure, the network is still in use by other nodes.

To handle local network failures, HACMP selectively moves the resources (on that network) from one node to another. This operation is referred to as *selective failover*.

Global Network Failure

A *global network failure* is an HACMP event that describes the case in which packets cannot be sent or received by *any* node over an HACMP logical network. This may occur, for instance, if the physical network is damaged.

Note: It is important to distinguish between these two terms in HACMP: a “global network”, and a “global network failure event.” A global network is a combination of HACMP networks; a global network failure event refers to a failure that affects all nodes connected to any logical HACMP network, not necessarily a global network.

HACMP Communication Interfaces

An *HACMP communication interface* is a grouping of a logical network interface, a service IP address and a service IP label *that you defined to HACMP*. HACMP communication interfaces combine to create IP-based networks.

An HACMP communication interface is a combination of:

- *A logical network interface* is the name to which AIX resolves a port (for example, en0) of a physical network interface card.
- *A service IP address* is an IP address (for example, 129.9.201.1) over which services, such as an application, are provided, and over which client nodes communicate.
- *A service IP label* is a label (for example, a hostname in the */etc/hosts* file, or a logical equivalent of a service IP address, for example, node_A_en_service) that maps to the Service IP address.

Communication interfaces in HACMP are used in the following ways:

- A communication interface refers to IP-based networks and NICs. The NICs that are connected to a common physical network are combined into logical networks that are used by HACMP.
- Each NIC is capable of hosting several TCP/IP addresses. When configuring a cluster, you define to HACMP the IP addresses that HACMP monitors (base or boot IP addresses), and the IP addresses that HACMP keeps highly available (the service IP addresses).

- Heartbeating in HACMP occurs over communication interfaces. HACMP uses the heartbeating facility of the RSCT subsystem to monitor its network interfaces and IP addresses. HACMP passes the network topology you create to RSCT, and RSCT provides failure notifications to HACMP.

HACMP Communication Devices

HACMP also monitors network devices which are not capable of IP communications. These devices include RS232 connections and Target Mode (disk-based) connections.

Device-based networks are point-to-point connections that are free of IP-related considerations such as subnets and routing—each device on a node communicates with only one other device on a remote node.

Communication devices make up device-based networks. The devices have names defined by the operating system (such as `tty0`). HACMP allows you to name them as well (such as `TTY1_Device1`).

For example, an RS232 or a point-to-point connection would use a device name of `/dev/tty2` as the device configured to HACMP on each end of the connection. Two such devices need to be defined—one on each node.

Note: The previous sections that described local and global network failures are true for TCP/IP-based HACMP logical networks. For device-based HACMP logical networks, these concepts do not apply. However, the heartbeating process occurs on device-based networks.

Subnet Routing Requirements in HACMP

A *subnet route* defines a path, defined by a subnet, for sending packets through the logical network to an address on another logical network. Beginning with AIX 5.1, you can add multiple routes for the same destination in the kernel routing table. If multiple matching routes have equal criteria, routing can be performed alternatively using one of the several subnet routes.

It is important to consider subnet routing in HACMP because of the following considerations:

- HACMP does not distinguish between logical network interfaces which share the same subnet route. If a logical network interface shares a route with another interface, HACMP has no means to determine its health. For more information on network routes, please see the AIX man page for the **route** command.
- Various constraints are often imposed on the IP-based networks by a network administrator or by TCP/IP requirements. The subnets and routes are also constraints within which HACMP must be configured for operation.

Note: We recommend that each communication interface on a node belongs to a unique subnet, so that HACMP can monitor each interface. This is not a strict requirement in all cases, and depends on several factors. In such cases where it is a requirement, HACMP enforces it. Also, ask your network administrator about the class and subnets used at your site.

Service IP Address/Label

A *Service IP label* is a label that maps to the service IP address and is used to establish communication between client nodes and the server node. Services, such as a database application, are provided using the connection made over the service IP label.

A service IP label can be placed in a *resource group* as a resource, which allows HACMP to monitor its health and keep it highly available, either within a node or, if *IP Address Takeover* is configured, between the cluster nodes by transferring it to another node in the event of a failure.

Note: A service IP label/address is configured as part of configuring cluster resources (not as part of topology, as in previous releases).

IP Alias

An *IP alias* is an IP label/address which is configured onto a network interface card *in addition* to the originally-configured IP label/address on the NIC. IP aliases are an AIX function which is supported by HACMP. AIX supports multiple IP aliases on a NIC. Each IP alias on a NIC can be configured on a separate subnet.

IP aliases are used in HACMP both as service and non-service addresses for *IP address takeover*, as well as for the configuration of the heartbeating method.

See the following sections for information on how HACMP binds a service IP label with a communication interface depending on which mechanism is used to recover a service IP label.

IP Address Takeover

If the physical network interface card on one node fails, and if there are no other accessible physical network interface cards on the same network on the same node (and, therefore, swapping IP labels of these NICs within the same node cannot be performed), HACMP may use the IP Address Takeover (IPAT) operation.

IP Address Takeover is a mechanism for recovering a service IP label by moving it to another NIC on another node, when the initial NIC fails. IPAT is useful because it ensures that an IP label over which services are provided to the client nodes remains available.

HACMP supports two methods for performing IPAT:

- *IPAT via IP Aliases* (this is the default)
- *IPAT via IP Replacement* (this method was known in previous releases as IPAT, or traditional IPAT).

Both methods are described in the sections that follow.

IPAT and Service IP Labels

The following list summarizes how IPAT manipulates the service IP label:

When IPAT via IP Aliases is used

The service IP address/label is aliased onto the same network interface as an existing communications interface.

That is, multiple IP addresses/labels are configured on the same network interface at the same time. In this configuration, all IP addresses/labels that you define must be configured on different subnets.

This method can save hardware, but requires additional subnets.

When IPAT via IP Replacement is used

The service IP address/label replaces the existing IP address/label on the network interface.

That is, only one IP address/label is configured on the same network interface at the same time. In this configuration, two IP addresses/labels on a node can share a subnet, while a backup IP address/label on the node must be on a different subnet.

This method can save subnets but requires additional hardware.

IP Address Takeover via IP Aliases

You can configure IP Address Takeover on certain types of networks using the IP aliasing network capabilities of AIX. Defining IP aliases to network interfaces allows creation of more than one IP label and address on the same network interface. IPAT via IP Aliases utilizes the gratuitous ARP capabilities available on many types of networks.

In a cluster with a network configured with *IPAT via IP Aliases*, when the resource group containing the service IP label falls over from the primary node to the target node, the initial IP labels that are used at boot time are added (and removed) as alias addresses on that NIC, or on other NICs that are available. Unlike in *IPAT via IP Replacement*, this allows a single NIC to support more than one service IP label placed on it as an alias. Therefore, the same node can host more than one *resource group* at the same time.

If the IP configuration mechanism for an HACMP network is via IP Aliases, the communication interfaces for that HACMP network must use routes that are different from the one used by the service IP address.

IPAT via IP Aliases provides the following advantages over the IPAT via IP Replacement scheme:

- Running IP Address Takeover via IP Aliases is faster than running IPAT via IP Replacement, because moving the IP address and the hardware address takes considerably longer than simply moving the IP address.
- IP aliasing allows co-existence of multiple service labels on the same network interface—you can use fewer physical network interface cards in your cluster. Note that upon failover, HACMP equally distributes aliases between available network interface cards.

IPAT via IP Aliases is the default mechanism for keeping a service IP label highly available.

IP Address Takeover via IP Replacement

The IP Address Takeover via IP Replacement facility moves the service IP label (along with the IP address associated with it) off a NIC on one node to a NIC on another node, should the NIC on the first node fail. IPAT via IP Replacement ensures that the service IP label which is included as a resource in a resource group in HACMP is accessible through its IP address, no matter onto which physical network interface card this service IP label is currently placed.

If the IP address configuration mechanism is IP Replacement, only one communication interface for that HACMP network must use a route that is the same as the one used by the Service IP Address.

In conjunction with IPAT via IP Replacement (also, previously known as *traditional IPAT*) you may also configure *Hardware Address Takeover (HWAT)* to ensure that the mappings in the ARP cache are correct on the target adapter.

Heartbeating Over Networks and Disks

A heartbeat is a type of a communication packet that is sent between nodes. Heartbeats are used to monitor the health of the nodes, networks and network interfaces, and to prevent cluster partitioning.

Heartbeating in HACMP: Overview

In order for an HACMP cluster to recognize and respond to failures, it must continually check the health of the cluster. Some of these checks are provided by the heartbeat function. Each cluster node sends heartbeat messages at specific intervals to other cluster nodes, and expects to receive heartbeat messages from the nodes at specific intervals. If messages stop being received, HACMP recognizes that a failure has occurred.

Heartbeats can be sent over:

- TCP/IP networks
- Point-to-point networks
- Shared disks.

The heartbeat function is configured to use specific paths between nodes. This allows heartbeats to monitor the health of all HACMP networks and network interfaces, as well as the cluster nodes themselves.

The TCP/IP heartbeat paths are set up automatically by RSCT; you have the option to configure point-to-point and disk paths as part of HACMP configuration.

HACMP passes the network topology you create to RSCT. RSCT Topology Services provides the actual heartbeat service, setting up the heartbeat paths, then sending and receiving the heartbeats over the defined paths. If heartbeats are not received within the specified time interval, Topology Services informs HACMP.

Heartbeating over TCP/IP Networks

RSCT Topology Services uses the HACMP network topology to dynamically create a set of heartbeat paths which provide coverage for all TCP/IP interfaces and networks. These paths form *heartbeat rings*, so that all components can be monitored without requiring excessive numbers of heartbeat packets.

In order for RSCT to reliably determine where a failure occurs, it must send and receive heartbeat packets over specific interfaces. This means that each NIC configured in HACMP must have an IP label on a separate subnet. There are two ways to accomplish this:

- Configure *heartbeating over IP interfaces*. If this method is used, you configure all service and non-service IP labels on separate subnets.
- Configure *heartbeating over IP Aliases*. If this method is used, you specify a base address for the heartbeat paths. HACMP then configures a set of IP addresses and subnets for heartbeating which are totally separate from those used as service and non-service addresses. With this heartbeating method, all service and non-service IP labels can be configured on the same subnet or on different subnets. Since HACMP automatically generates the proper addresses required for heartbeating, all other addresses are free of any constraints.

Heartbeating over IP Aliases provides the greatest flexibility for configuring boot (base) and service IP addresses at the cost of reserving a unique address and subnet range that is used specifically for heartbeating.

Note: Although heartbeating over IP Aliases bypasses the subnet requirements for HACMP to correctly perform the heartbeating function, the existence of multiple routes to the same subnet (outside of HACMP) may produce undesired results for your application. For information on subnet requirements, see [Subnet Routing Requirements in HACMP](#).

Heartbeating over Point-to-Point Networks

You can also configure non-IP point-to-point network connections that directly link cluster nodes. These connections can provide an alternate heartbeat path for a cluster that uses a single TCP/IP-based network. They also prevent the TCP/IP software itself from being a single point of failure in the cluster environment.

Point-to-point networks that you plan to use for heartbeating should be free of any other traffic for the exclusive use by HACMP.

You can configure non-IP point-to-point heartbeat paths over the following types of networks:

- Serial (RS232)
- Target Mode SSA
- Target Mode SCSI
- Disk heartbeating (over an enhanced concurrent mode disk).

Heartbeating over Disks

Heartbeating is supported on any shared disk that is part of an enhanced concurrent mode volume group.

Note: The volume group does not need to be configured as an HACMP resource.

Heartbeating over an enhanced concurrent mode disk operates with any type of disk—including those that are attached by fibre channel. This avoids the distance limitations (especially when using fibre channel connections) associated with RS232 links, making this solution more cost effective.

A single common disk serves as the heartbeat path between two cluster nodes. Enhanced concurrent mode supports concurrent read and write access to the non-data portion of the disk. Nodes use this part of the disk to periodically write heartbeat messages and read heartbeats written by the other node.

2 **HACMP Cluster Nodes, Networks and Heartbeating** Heartbeating Over Networks and Disks

Chapter 3: HACMP Resources and Resource Groups

This chapter introduces resource-related *concepts and definitions* that are used throughout the documentation, and also in the HACMP user interface.

The information in this chapter is organized as follows:

- [Cluster Resources: Identifying and Keeping Available](#)
- [Types of Cluster Resources](#)
- [Cluster Resource Groups](#)
- [Resource Group Policies and Attributes](#)
- [Resource Group Dependencies](#)

Cluster Resources: Identifying and Keeping Available

The HACMP software provides a highly available environment by:

- Identifying the set of *cluster resources* that are essential to processing.
- Defining the *resource group policies and attributes* that dictate how HACMP manages resources to keep them highly available at different stages of cluster operation (startup, fallover and fallback).

By identifying resources and defining resource group policies, the HACMP software makes numerous cluster configurations possible, providing tremendous flexibility in defining a cluster environment tailored to individual requirements.

Identifying Cluster Resources

Cluster resources can include both hardware and software:

- Disks
- [Volume Groups](#)
- [Logical Volumes](#)
- [Filesystems](#)
- [Service IP Labels/Addresses](#)
- [Applications](#)
- [Tape Resources](#)
- [Communication Links](#)
- Fast Connect Resources, and other resources.

A processor running HACMP owns a user-defined set of resources: disks, volume groups, filesystems, IP addresses, and applications. For the purpose of keeping resources highly available, sets of interdependent resources may be configured into *resource groups*.

Resource groups allow you to combine related resources into a single logical entity for easier configuration and management. The Cluster Manager handles the resource group as a unit, thus keeping the interdependent resources together on one node, *and* keeping them highly available.

Types of Cluster Resources

This section provides a brief overview of the resources that you can configure in HACMP and include into resource groups to let HACMP keep them highly available.

Volume Groups

A *volume group* is a set of physical volumes that AIX treats as a contiguous, addressable disk region. Volume groups are configured to AIX, and can be included in resource groups in HACMP. In the HACMP environment, a *shared volume group* is a volume group that resides entirely on the external disks that are shared by the cluster nodes. Shared disks are those that are physically attached to the cluster nodes and logically configured on all cluster nodes.

Logical Volumes

A *logical volume* is a set of *logical partitions* that AIX makes available as a single storage unit—that is, the logical volume is the “logical view” of a physical disk. Logical partitions may be mapped to one, two, or three physical partitions to implement mirroring.

In the HACMP environment, logical volumes can be used to support a journaled filesystem (non-concurrent access), or a raw device (concurrent access). Concurrent access does not support filesystems. Databases and applications in concurrent access environments must access raw logical volumes (for example, `/dev/rsharedlv`).

A *shared logical volume* must have a unique name within an HACMP cluster.

Filesystems

A filesystem is written to a single logical volume. Ordinarily, you organize a set of files as a filesystem for convenience and speed in managing data.

Shared Filesystems

In the HACMP system, a *shared filesystem* is a journaled filesystem that resides entirely in a shared logical volume.

For non-concurrent access, you want to plan shared filesystems so that they will be placed on external disks shared by cluster nodes. Data resides in filesystems on these external shared disks in order to be made highly available.

For concurrent access, you cannot use journaled filesystems. Instead, use raw logical volumes.

Journalized Filesystem and Enhanced Journalized Filesystem

An Enhanced Journalized Filesystem (JFS2) provides the capability to store much larger files than the Journalized Filesystem (JFS). JFS2 is the default filesystem for the 64-bit kernel. You can choose to implement either JFS which is the recommended filesystem for 32-bit environments, or JFS2 which offers 64-bit functionality.

JFS2 is more flexible than JFS because it allows to dynamically increase and decrease the number of files you can have in a filesystem. JFS2 also lets you include the filesystem log in the same logical volume as the data, instead of allocating a separate logical volume for logs for all filesystems in the volume group.

For more information on JFS2, see the *AIX 5.2 Differences Guide*:

<http://www.redbooks.ibm.com/pubs/pdfs/redbooks/sg245765.pdf>

Applications

The purpose of a highly available system is to ensure that critical services are accessible to users. Applications usually need no modification to run in the HACMP environment. Any application that can be successfully restarted after an unexpected shutdown is a candidate for HACMP.

For example, all commercial DBMS products checkpoint their state to disk in some sort of transaction journal. In the event of a server failure, the failover server restarts the DBMS, which reestablishes database consistency and then resumes processing.

If you use Fast Connect to share resources with non-AIX workstations, you can configure it as an HACMP resource, making it highly available in the event of node or network interface card failure, and making its correct configuration verifiable.

Applications are managed by defining the application to HACMP as an *application server* resource. The application server includes an application start and stop scripts. HACMP uses these scripts when the application needs to be brought online or offline on a particular node, to keep the application highly available.

Note: The start and stop scripts are the main points of control for HACMP over an application. It is very important that the scripts you specify operate correctly to start and stop all aspects of the application. If the scripts fail to properly control the application, other parts of the application recovery may be affected. For example, if the stop script you use fails to completely stop the application and a process continues to access a disk, HACMP will not be able to bring the volume group offline on the node that failed and to recover it on the backup node.

Add your application server to an HACMP resource group only after you have thoroughly tested your application start and stop scripts.

The resource group that will contain the application server should also contain all the resources that the application depends on including service IP addresses, volume groups and filesystems. Once such a resource group is created, HACMP manages the entire resource group, and

therefore all the interdependent resources in it as a single entity. (Note that HACMP coordinates the application recovery and manages the resources in the order that ensures activating all interdependent resources *before* other resources.)

In addition, HACMP includes application monitoring capability, whereby you can define a monitor to detect the unexpected termination of a process or to periodically poll the termination of an application and take automatic action upon detection of a problem.

In HACMP 5.2, you can configure multiple application monitors and associate them with one or more application servers. By supporting multiple monitors per application, HACMP can support more complex configurations. For example, you can configure one monitor for each instance of an Oracle parallel server in use. Or, you can configure a custom monitor to check the health of the database, and a process termination monitor to instantly detect termination of the database process.

You can also specify a mode for an application monitor: It can either track how the application is being run (running mode), or whether the application have started successfully (application startup mode). Using a monitor to watch the application startup is especially useful for complex cluster configurations.

Service IP Labels/Addresses

A *service IP label* is used to establish communication between client nodes and the server node. Services, such as a database application, are provided using the connection made over the service IP label.

A service IP label can be placed in a resource group as a resource which allows HACMP to monitor its health and keep it highly available, either within a node or, if IP address takeover is configured, between the cluster nodes by transferring it to another node in the event of a failure.

For more information about service IP labels, see [Service IP Address/Label](#) in [Chapter 2: HACMP Cluster Nodes, Networks and Heartbeating](#).

Note: Certain subnet requirements apply for configuring service IP labels as resources in different types of resource groups. For more information, see the *Planning and Installation Guide*.

Tape Resources

You can configure a SCSI or a Fibre Channel tape drive as a cluster resource in a non-concurrent resource group, making it highly available to two nodes in a cluster. Management of shared tape drives is simplified by the following HACMP functionality:

- Configuration of tape drives using SMIT
- Verification of proper configuration of tape drives
- Automatic management of tape drives during resource group start and stop operations
- Reallocation of tape drives on node failure and node recovery
- Controlled reallocation of tape drives on cluster shutdown
- Controlled reallocation of tape drives during a dynamic reconfiguration of cluster resources.

Communication Links

You can define the following communication links as resources in HACMP resource groups:

- SNA configured over LAN network interface cards
- SNA configured over X.25
- Pure X.25.

By managing these links as resources in resource groups, HACMP ensures their high availability. Once defined as members of an HACMP resource group, communication links are protected in the same way other HACMP resources are. In the event of a LAN physical network interface or an X.25 link failure, or general node or network failures, a highly available communication link falls over to another available network interface card on the same node, or on a takeover node.

- *SNA configured over LAN.* To be highly available, “SNA configured over LAN” need to be included in those resource groups that contain the corresponding service IP labels in them. These service IP labels, in turn, are defined on LAN network interface cards, such as Ethernet and Token Ring. In other words, the availability of the “SNA configured over LAN” resources is dependent upon the availability of service IP labels included in the resource group. If the NIC being used by the service IP label fails, and the service IP label is taken over by another interface, this interface will also take control over an “SNA configured over LAN” resource configured in the resource group.
- *SNA configured over X.25 links and pure X-25 links.* These links are usually, although not always, used for WAN connections. They are used as a means of connecting dissimilar machines, from mainframes to dumb terminals. Because of the way X.25 networks are used, these physical network interface cards are really a different class of devices that are *not* included in the cluster topology and are *not* controlled by the standard HACMP topology management methods. This means that heartbeats are not used to monitor X.25 link status, and you do not define X.25-specific networks in HACMP. To summarize, you can include X.25 links as resources in resource groups, keeping in mind that the health and availability of these resources also relies on the health of X.25 networks themselves (which are not configured within HACMP.)

Cluster Resource Groups

To be made highly available by the HACMP software, each resource must be included in a resource group. Resource groups allow you to combine related resources into a single logical entity for easier management.

This section provides definitions, such as a home node, a participating nodelist and a default node priority, that help to describe different types of resource group behavior. These definitions are used later in this chapter to explain how HACMP uses resource groups to keep the resources and applications highly available.

Prior to HACMP 5.2, in addition to custom resource groups, you could configure pre-defined types of resource groups, where a type was a set of pre-defined behaviors. These included rotating, cascading and concurrent resource groups. In HACMP 5.2, these predefined types are no longer supported as they provided only a restricted set of behaviors.

In HACMP 5.2, the resource groups are defined in the same way as custom resource groups were defined before HACMP 5.2. From here onward, custom resource groups are referred to as resource groups. It is possible to configure a resource group that replicates the behavior of pre-5.2 resource groups by using the different policies available for resource group startup, fallover and fallback.

In addition, the following functionality is available for resource groups in HACMP 5.2 and is described later in this chapter:

- Distribution policy for resource group startup
- Support for IPAT via IP replacement networks and service IP labels.

This section includes the terms that are used in HACMP 5.2 for resource groups, and contains the following topics:

- [Participating Nodelist](#)
- [Default Node Priority](#)
- [Home Node](#)
- [Startup, Fallover and Fallback](#).

Participating Nodelist

The *participating nodelist* defines a list of nodes that can host a particular resource group. You define a nodelist when you configure a resource group. The participating nodelist for non-concurrent resource groups can contain some or all nodes in the cluster. The participating nodelist for concurrent resource groups should contain *all* nodes in the cluster. Typically, this list contains all nodes sharing the same data and disks.

Default Node Priority

Default node priority is identified by the position of a node in the nodelist for a particular resource group. The first node in the nodelist has the *highest node priority*; it is also called the *home node* for a resource group. The node that is listed before another node has a *higher node priority* than the current node.

Depending on a *fallback policy* for a resource group, when a node with a higher priority for a resource group (that is currently being controlled by a lower priority node) joins or reintegrates into the cluster, it takes control of the resource group. That is, the resource group moves from nodes with lower priorities to the higher priority node.

At any given time, the resource group can have a default node priority specified by the participating nodelist. However, various resource group policies you select can override the default node priority and “create” the actual node priority according to which a particular resource group would move in the cluster.

Depending on the policies you select for a resource group, it selects a node either according to the default nodelist, or according to the *dynamic node priority* (DNP) that you can define for a non-concurrent resource group. In this case, HACMP determines a node priority based on the predefined DNP.

Home Node

The *home node* (or the *highest priority node for this resource group*) is the first node that is listed in the participating nodelist for a non-concurrent resource group. The home node is a node that normally owns the resource group. A non-concurrent resource group may or may not have a home node—it depends on the startup, fallover and fallback behaviors of a resource group.

The home node is a node that normally owns the resource group. Note that due to different changes in the cluster, the group may not always start on the home node. It is important to differentiate between the *home node* for a resource group and the *node that currently owns it*.

The term home node is not used for concurrent resource groups as they are owned by multiple nodes.

Startup, Fallover and Fallback

HACMP ensures the availability of cluster resources by moving resource groups from one node to another when the conditions in the cluster change. HACMP manages resource groups by activating them on a particular node or multiple nodes at cluster startup, or by moving them to another node if the conditions in the cluster change. These are the stages in a cluster lifecycle that affect how HACMP manages a particular resource group:

- *Cluster startup.* Nodes are up and resource groups are distributed between them according to the resource group startup policy you selected.
- *Node failure.* Resource groups that are active on this node fall over to another node.
- *Node recovery.* A node reintegrates into the cluster and resource groups could be reacquired, depending on the resource group policies you select.
- *Resource failure and recovery.* A resource group may fall over to another node, and be reacquired, when the resource becomes available.
- *Cluster shutdown.* There are different ways of shutting down a cluster, one of which ensures that resource groups fall over to another node gracefully.

During each of these cluster stages, the behavior of resource groups in HACMP is defined by the following:

- Which node, or nodes, activate the resource group at cluster startup
- How many resource groups are allowed to be acquired on a node during cluster startup
- Which node takes over the resource group when the node that owned the resource group fails and HACMP needs to move a resource group to another node
- Whether a resource group falls back to a node that has just joined the cluster or stays on the node that currently owns it.

The resource group policies that you select determine which cluster node originally controls a resource group and which cluster nodes take over control of the resource group when the original node relinquishes control.

Each combination of these policies allows you to specify varying degrees of control over which node, or nodes, control a resource group.

Cascading, rotating and concurrent resource groups are no longer supported. The terms Cascading without Fallback and Inactive Takeover are not used. In HACMP 5.2, the resource groups are defined in the same way as custom resource groups were defined in the previous release. The term “custom” is not used for resource groups. The term “concurrent resource group” defines a resource group with Online on All Available Nodes startup policy.

To summarize, the focus of HACMP on resource group ownership makes numerous cluster configurations possible and provides tremendous flexibility in defining the cluster environment to fit the particular needs of the application. The combination of startup, fallover and fallback policies summarizes all the management policies available in previous releases without the requirement to specify the set of options that modified the behavior of “predefined” group types.

When defining resource group behaviors, keep in mind that a resource group can be taken over by one or more nodes in the cluster.

Startup, fallover and fallback are specific behaviors that describe how resource groups behave at different cluster stages. It is important to keep in mind the difference between fallover and fallback. These terms appear frequently in discussion of the various resource group policies.

Startup

Startup refers to the activation of a resource group on a node (or multiple nodes) on which it currently resides, or on the home node for this resource group. Resource group startup occurs during cluster startup, or initial acquisition of the group on a node.

Fallover

Fallover refers to the movement of a resource group from the node that currently owns the resource group to another active node after the current node experiences a failure. The new owner is not a reintegrating or joining node.

Fallover is valid only for *non-concurrent resource groups*.

Fallback

Fallback refers to the movement of a resource group from the node on which it currently resides (which is *not* a home node for this resource group) to a node that is joining or reintegrating into the cluster.

For example, when a node with a higher priority for that resource group joins or reintegrates into the cluster, it takes control of the resource group. That is, the resource group falls back from nodes with lesser priorities to the higher priority node.

Defining a fallback behavior is valid only for *non-concurrent resource groups*.

Resource Group Policies and Attributes

In HACMP 5.2, you can configure resource groups that use specific startup, fallover and fallback policies.

This section describes resource group attributes and scenarios, and helps you to decide which resource groups suit your cluster requirements.

This section contains the following topics:

- [Overview](#)

- [Resource Group Startup, Fallover and Fallback](#)
- [Settling Time, Dynamic Node Priority and Fallback Timer](#)
- [Distribution Policy](#)
- [Upgrading from Previous Releases](#)
- [Cluster Networks and Resource Groups](#)
- [Sites and Resource Groups](#).

Overview

Prior to HACMP 5.1, in addition to cascading, rotating and concurrent resource groups, you could configure *custom resource groups*. Starting with HACMP 5.1, you could use custom resource groups if you preferred to specify parameters that precisely describe the resource group's behavior at startup, fallover and fallback. In HACMP 5.2, only custom resource groups can be configured, although they are now referred to simply as resource groups.

Note: The regular cascading, rotating and concurrent resource groups had predefined startup, fallover and fallback behaviors. You could refine fallover and fallback behavior of cascading resource groups by specifying Cascading without Fallback and Inactive Takeover attributes. In HACMP 5.2, you can create resource groups that replicate these settings, if needed. For information on how the pre-5.2 attributes and policies convert to the new configuration scheme, see the chapter on upgrading to the next version in the *Planning and Installation Guide*.

In HACMP 5.2, the policies for resource groups offer a wider variety of choices than predefined policies of pre-5.2 resource groups. Resource group policies can now be tailored to your needs. This allows you to have a greater control of the resource group behavior, increase resource availability, and better plan node maintenance.

The process of configuring a resource group is two-fold. First, you configure startup, fallover and fallback policies for a resource group. Second, you add specific resources to it.

HACMP 5.2 prevents you from configuring invalid combinations of behavior and resources in a resource group.

In addition, using resource groups in HACMP 5.2 potentially increases availability of cluster resources:

- For instance, you can configure resource groups in a way which ensures that they are brought back online on reintegrating nodes during off-peak hours
- or*
- You can specify that a resource group that contains a certain application is the only one that will be given preference and be acquired during startup on a particular node. You do so by specifying the node distribution policy. This is relevant if multiple non-concurrent resource groups can potentially be acquired on a node, but a specific resource group owns an application that is more important to keep available.

For resource group planning considerations, see the chapter on planning resource groups in the *Planning and Installation Guide*.

Resource Group Startup, Fallover and Fallback

In HACMP 5.2, the following policies exist for resource groups:

<p>Startup</p>	<ul style="list-style-type: none"> • Online on Home Node Only. The resource group should be brought online only on its home (highest priority) node during the resource group startup. This requires the highest priority node to be available. • Online on First Available Node. The resource group activates on the first participating node that becomes available. • Online on All Available Nodes. The resource group is brought online on all nodes. • Online Using Distribution Policy. Only one resource group is brought online on a node, or on a node per network, depending on the distribution policy specified (node or network).
<p>Fallover</p>	<ul style="list-style-type: none"> • Fallover to Next Priority Node in the List. In the case of fallover, the resource group that is online on only one node at a time follows the default node priority order specified in the resource group's nodelist. • Fallover Using Dynamic Node Priority. Before selecting this option, configure a dynamic node priority policy that you want to use. Or you can select one of the three predefined dynamic node priority policies. • Bring Offline (on Error Node Only). Select this option to bring a resource group offline on a node during an error condition.
<p>Fallback</p>	<ul style="list-style-type: none"> • Fallback to Higher Priority Node in the List. A resource group falls back when a higher priority node joins the cluster. If you select this option, you can use the delayed fallback timer. If you do not configure a delayed fallback policy, the resource group falls back immediately when a higher priority node joins the cluster. • Never Fallback. A resource group does <i>not</i> fall back when a higher priority node joins the cluster.

For more information on each policy, see the *Administration and Troubleshooting Guide*.

Settling Time, Dynamic Node Priority and Fallback Timer

You can configure some additional parameters for resource groups that dictate how the resource group behaves at startup, fallover and fallback. They are:

- *Settling Time.* You can configure a startup behavior of a resource group by specifying the settling time for a resource group that is currently offline. When the settling time is not configured, the resource group starts on the first available higher priority node that joins the

cluster. If the settling time is configured, HACMP waits for the duration of the settling time period for a higher priority node to join the cluster before it activates a resource group. Specifying the settling time enables a resource group to be acquired on a node that has a higher priority, when multiple nodes are joining simultaneously. The settling time is a cluster-wide attribute that, if configured, affects the startup behavior of *all* resource groups in the cluster for which you selected Online on First Available Node startup behavior.

- *Distribution Policy.* You can configure a startup behavior of a resource group to use the distribution policy during startup. This policy ensures that during startup, only one resource group is acquired on a node (*node distribution*), or on a node per network (*network distribution*). See the following section for more information.
- *Dynamic Node Priority.* You can configure a failover behavior of a resource group to use the dynamic node priority.
- *Delayed Fallback Timer.* You can configure a fallback behavior of a resource group to occur at one of the predefined recurring times: daily, weekly, monthly and yearly, or on a specific date and time, by specifying and assigning a delayed fallback timer. This is useful, for instance, for scheduling the resource group fallbacks to occur during off-peak business hours.

Distribution Policy

In HACMP 5.2, on cluster startup, you can use a *distribution policy* for a resource group. You can select one of the two types of distribution policies:

- *A node distribution policy.* It prevents HACMP from acquiring more than one resource group on a node when the node joins the cluster.
- *A network distribution policy.* It prevents HACMP from acquiring more resource groups than the number of networks configured on a node, when the node joins the cluster.

The distribution policy is useful in the following scenario: When a node joins a cluster, for all the resource groups in which it participates, it determines which groups to acquire, based on the node priority, the node's available hardware, the resource group's management policy, and other parameters that you can configure.

In HACMP 5.2, for each resource group in the cluster, you can specify a startup policy Online Using Distribution Policy. The distribution policy that HACMP uses in this case, can be one of the two types:

- *Distribute Resource Groups per Node.* The distribution policy is a cluster-wide attribute that causes the resource groups to distribute themselves in a way that only *one* resource group is acquired on a node during startup. This policy is used by default.

Using this type of distribution lets you distribute your CPU-intensive applications on different nodes.

If a resource group startup uses this distribution policy, only this resource group is acquired on a node. If two resource groups use this policy, and both resource groups are offline at the time when a particular node joins, then only one of the two resource groups is acquired on a node. HACMP gives preference to those resource groups that have fewer nodes in the nodelist, and then sorts the list of resource groups alphabetically.

- *Distribute Resource Groups per Network.* This distribution policy is a cluster-wide attribute that causes the resource groups to distribute themselves in a way that only *one* resource group is acquired on *a node per cluster network*. The remaining resource groups

stay offline until you manually activate them on the nodes. However, if multiple IP labels exist on one node, a node may acquire multiple groups, one per network. If this policy is specified, it is recommended that the number of resource groups configured to use this policy is equal to (or less than) the number of configured networks per node. This policy requires that each resource group includes a service IP label.

Using this type of distribution lets you distribute your network-intensive applications on different nodes and on different networks.

Note: The network distribution policy is retained in HACMP 5.2 to ensure backward compatibility with the previous releases and may not be offered in future releases. It is recommended to change the distribution policy to node distribution after the migration to HACMP 5.2 is complete.

If, prior to HACMP 5.2, you had one rotating resource group configured on each node per network, only one resource group per network would be acquired during node startup. Upon migration to HACMP 5.2, such resource groups by default use the startup Online Using Distribution Policy, with the distribution policy set to “network”.

For more information on resource group distribution policy and how it is handled during migration from previous releases, see the *Planning and Installation Guide*.

For configuration information, and for information on resource group management, see the *Administration and Troubleshooting Guide*.

Each resource group has specific behaviors during startup, as well as during fallover and fallback. Each of these are discussed in detail in the chapter for planning resource groups in the *Planning and Installation Guide*.

Upgrading from Previous Releases

If you are upgrading from previous releases, the following issues about resource groups are important:

- The pre-5.2 resource groups are automatically converted to resource groups with particular startup, fallover and fallback policies. *All* functionality of pre-5.2 resource groups is retained in HACMP 5.2.
- If prior to migration, the cluster configuration included rotating resource groups, in HACMP 5.2 the network distribution policy is used as a startup policy for such resource groups. The network distribution ensures that only one resource group is brought online on a node per cluster network, during a node or cluster startup.

For more information, see the *Planning and Installation Guide*.

Cluster Networks and Resource Groups

In HACMP 5.1, some resource groups were supported only on IPAT via IP alias networks. In HACMP 5.2, all resource groups support service IP labels configured on IPAT via IP replacement networks, and on aliased networks.

A service IP label can be included in any non-concurrent resource group—that resource group could have any of the allowed startup policies except Online on All Available Nodes.

Sites and Resource Groups

In HACMP 5.2, the following options exist for configuring resource group inter-site policies:

- Prefer Primary Site
- Online On Either Site
- Online On Both Sites.

For information on inter-site policies, see the *Administration and Troubleshooting Guide*.

Resource Group Dependencies

In HACMP 5.1, the process of configuring simple clusters was made significantly easier. In HACMP 5.2, it is easier to set up more complex clusters by specifying dependencies between resource groups.

In previous releases, support for resource group ordering and customized serial processing of resources let you accommodate cluster configurations where a dependency existed between applications residing in different resource groups. With customized serial processing, you can specify that a resource group is processed before another resource group, on a local node. However, it is not guaranteed that a resource group will be processed in the order specified, as it depends on other resource group policies and conditions.

Also, in previous releases, you could use pre- and post-event scripts to configure a customized dependency mechanism between different applications defined to HACMP. Although you can continue using these scripts, HACMP 5.2 offers an easy way to configure a dependency between resource groups (and applications that belong to them).

In this release, the dependency between resource groups that you configure is:

- Explicitly specified using the SMIT interface
- Established cluster-wide, not just on the local node
- Guaranteed to be honored in the cluster.

Configuring a resource group dependency allows for easier cluster configuration and control for clusters with multi-tier applications where one application depends on the successful startup of another application, and both applications are required to be kept highly available with HACMP.

The following example illustrates the dependency behavior you can configure in HACMP 5.2:

- If resource group A depends on resource group B, upon node startup, resource group B must be brought online before resource group A is acquired on any node in the cluster. Upon failover, the order is reversed: Resource group A must be taken offline before resource group B is taken offline.
- In addition, if resource group A depends on resource group B, during a node startup or node reintegration, resource group A cannot be taken online before resource group B is brought online. If resource group B is taken offline, resource group A will be taken offline too, since it depends on resource group B.

Dependencies between resource groups offer a predictable and reliable way of building clusters with multi-tier applications. For more information on typical cluster environments that can use dependent resource groups, see [Cluster Configurations with Multi-Tiered Applications](#) in [Chapter 6: HACMP Cluster Configurations](#).

Child and Parent Resource Groups

These terms are defined in HACMP 5.2 to describe dependencies between resource groups:

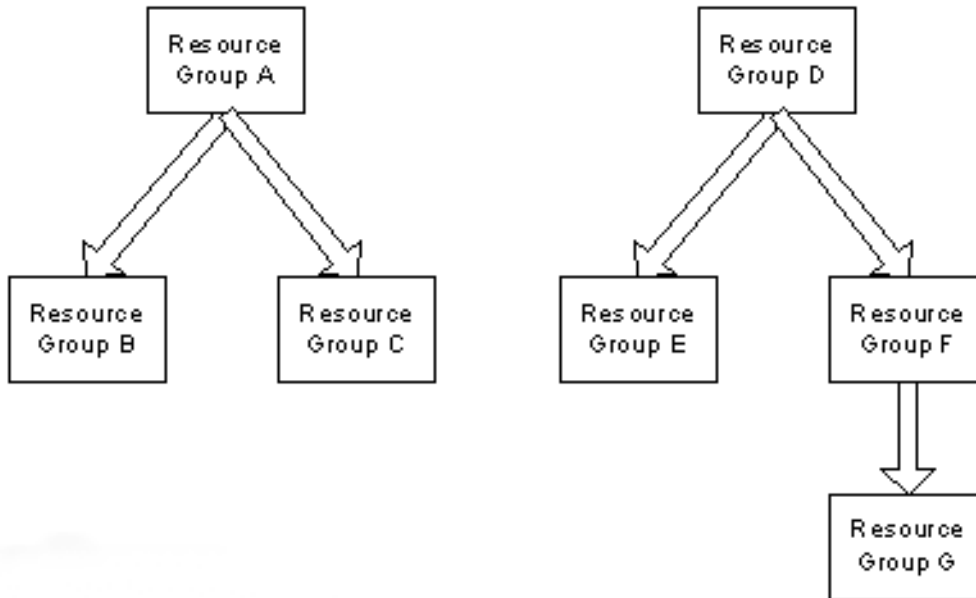
- *A parent resource group* has to be in an ONLINE state before the resource group that depends on it (*child*) can be started.
- *A child resource group* depends on a parent resource group. It will get activated on any node in the cluster *only after* the parent resource group has been activated. Typically, the child resource group depends on some application services that the parent resource group provides.

Upon resource group release (during failover or stopping cluster services, for instance) HACMP brings offline a child resource group before a parent resource group is taken offline.

The following statements describe the dependency relationship between resource groups in HACMP 5.2:

- You can configure a type of dependency where a parent resource group must be online on any node in the cluster before a child (dependent) resource group can be activated on a node.
- A resource group can serve as both a parent and a child resource group, depending on which end of a given dependency link it is placed.
- You can specify three levels of dependencies for resource groups.
- You cannot specify circular dependencies between resource groups.
- You can add, change or delete a dependency between resource groups, while the cluster services are running.
- When you delete a dependency between two resource groups, only the link between these resource groups is removed from the HACMP Configuration Database. The resource groups are not deleted.

The following graphic illustrates the statements:



Examples of Two and Three Levels of Dependencies between Resource Groups

During failover of a parent resource group, a child resource group containing the application temporarily goes offline and then online on any available node. The application that belongs to the child resource group is also stopped and restarted.

For more information on planning dependencies between resource groups, the application behavior in dependent resource groups and configuring dependencies that work successfully, see the *Planning and Installation Guide* and *Administration and Troubleshooting Guide*.

3 **HACMP Resources and Resource Groups**

Resource Group Dependencies

Chapter 4: HACMP Cluster Hardware and Software

This chapter describes:

- The IBM hardware that is used with HACMP and implements the base level of a highly available environment. The following IBM hardware is used with HACMP:

Nodes	<ul style="list-style-type: none">• IBM pSeries• RS/6000 SP System
Disks subsystems	<ul style="list-style-type: none">• IBM Serial Storage Architecture Disk Subsystem• IBM 2105 Enterprise Storage Server• IBM 2104 Expandable Storage Plus• IBM FASTT Storage Servers• SCSI Disks• OEM Disks. For an OEM disks overview, and for information on installing OEM disks, see <i>Appendix D: OEM Disk Accommodation</i> in the <i>Planning and Installation Guide</i>.
Networks	Ethernet, Token-Ring, ATM, SP Switch, and other networks. For information on administering particular types of networks, see the <i>Planning and Installation Guide</i> .

For detailed information, follow the links in the preceding table, or see the section [Enhancing Availability with IBM Hardware](#).

Other sections in this chapter are:

- [HACMP Required and Supported Hardware](#)
- [HACMP Cluster Software](#)

Enhancing Availability with IBM Hardware

Building a highly available cluster begins with reliable hardware. Within the AIX environment, the IBM EServer pSeries family of machines as well as the SP and its supported disk subsystems provide a robust, stable platform for building highly available clusters.

IBM pSeries

Some IBM pSeries servers, such as the 690 (Regatta), let you configure multiple logical partitions that can be used as separate nodes. The pSeries 690 delivers true logical partitioning (LPAR). Each system can be divided into as many as 16 virtual servers, each with its own set of system resources such as processors, memory and I/O. Unlike partitioning techniques

available on other UNIX servers, LPAR provides greater flexibility in matching resources to workloads. Based on business requirements, these resources can be allocated or combined for business-critical applications resulting in more efficient use of the system.

With the IBM eServer Cluster 1600 and AIX operating system, you can mix or match up to 128 units (512 via special order) including up to 32 pSeries 690 systems. An LPAR of a pSeries 690 is viewed by a Cluster 1600 as just another node or server in the cluster. Up to 16 LPARs per system and up to 128 LPARs per cluster are supported on pSeries 690. Up to 4 LPARs per system are supported on pSeries 650, and up to 2 LPARs are supported on pSeries 630.

RS/6000 SP System

The SP is a parallel processing machine that includes from two to 128 processors connected by a high-performance switch. The SP leverages the outstanding reliability provided by the RS/6000 series by including many standard RS/6000 hardware components in its design. The SP's architecture then extends this reliability by enabling processing to continue following the failure of certain components. This architecture allows a failed node to be repaired while processing continues on the healthy nodes. You can even plan and make hardware and software changes to an individual node while other nodes continue processing.

Disk Subsystems

The disk subsystems most often shared as external disk storage in cluster configurations are:

- 7133 Serial Storage Architecture (SSA) serial disk subsystems
- IBM 2105 Enterprise Storage Server (ESS)
- IBM 2104 Expandable Storage Plus
- IBM FASTT family storage servers
- SCSI disks.

See Appendix D on OEM disks in the *Planning and Installation Guide* for information on installing and configuring OEM disks.

IBM Serial Storage Architecture Disk Subsystem

You can use IBM 7133 and 7131-405 SSA disk subsystems as shared external disk storage devices in an HACMP cluster configuration.

If you include SSA disks in a volume group that uses LVM mirroring, you can replace a failed drive without powering off the entire subsystem.

IBM 2105 Enterprise Storage Server

IBM 2105 Enterprise Storage Server provides multiple concurrent attachment and sharing of disk storage for a variety of open systems servers. IBM EServer pSeries processors can be attached, as well as other UNIX and non-UNIX platforms. Attachment methods include Fibre Channel and SCSI.

The ESS uses IBM SSA disk technology (internally). ESS provides many availability features. RAID technology protects storage. RAID-5 techniques can be used to distribute parity across all disks in the array. *Sparing* is a function which allows you to assign a disk drive as a spare

for availability. Predictive Failure Analysis techniques are utilized to predict errors *before* they affect data availability. Failover Protection enables one partition, or *storage cluster*, of the ESS to takeover for the other so that data access can continue.

The ESS includes other features such as a web-based management interface, dynamic storage allocation, and remote services support. For more information on ESS planning, general reference material, and attachment diagrams, see this URL:

<http://www.storage.ibm.com/disk/ess>

IBM 2104 Expandable Storage Plus

The IBM 2104 Expandable Storage Plus (EXP Plus) provides flexible, scalable, and low-cost disk storage for RS/6000 and pSeries servers in a compact package. This new disk enclosure is ideal for enterprises--such as Internet or application service providers-- that need high-performance external disk storage in a small footprint.

- Scales from up to 2055 GB of capacity per drawer or tower to more than 28 TB per rack
- Shared storage for all major types of servers
- Single or split-bus configuration flexibility to one or two servers
- Incorporates high-performance Ultra3 SCSI disk storage with 160 MB/sec throughput
- Features up to fourteen 10,000 RPM disk drives, with capacities of 9.1 GB, 18.2GB, 36.4 GB and 73.4GB and 146.8GB
- High availability to safeguard data access
- Scalability for fast-growing environments.

See the following URL for complete information on IBM Storage Solutions:

<http://www.storage.ibm.com>

IBM FASt Storage Servers

IBM FASt700 Storage Server delivers superior performance with 2 Gbps Fibre Channel technology. The FASt700 is designed to offer investment protection with advanced functions and flexible features. Scales from 36GB to over 32TB to support growing storage requirements created by e-business applications. FASt700 offers advanced replication services to support business continuance. The FASt700 is an effective storage server for any enterprise seeking performance without borders.

IBM FASt500 Storage Server is the storage server of choice for medium range storage consolidation and data sharing on multiple or heterogeneous server platforms. The FASt500 supports rapid universal access to vast quantities of data through many advanced functions and features, making it a workhorse for business-critical applications.

In general, FASt200, 500 600, 700 and 900 series servers are supported with HACMP.

The FastT Storage Server family offers:

- Data protection with dual redundant components, multiple RAID levels, LUN masking and enhanced management options
- Storage Consolidation for SAN, NAS and direct-attach environments
- Investment protection throughout the FASt family of storage systems

- Heterogeneous support for IBM AIX[®], Microsoft[®] Windows[®] 2000, Windows NT[®], Sun Solaris, HP-UX, Red Hat Linux.
- Scales over 32 terabytes (TB) of capacity using flexible combinations of 18.2, 36.4, 73.4 and 146.8GB drives.
- Software upgrade to support advanced copy services and remote mirroring.

For complete information on IBM Storage Solutions, see the following URL:

<http://www.storage.ibm.com>

HACMP Required and Supported Hardware

For a complete list of required and supported hardware, see the sales guide for the product. You can locate this document from the following URL:

<http://www.ibm.com/common/ssi>

After selecting your country and language, select HW and SW Desc (SalesManual, RPQ) for a Specific Information Search.

HACMP Cluster Software

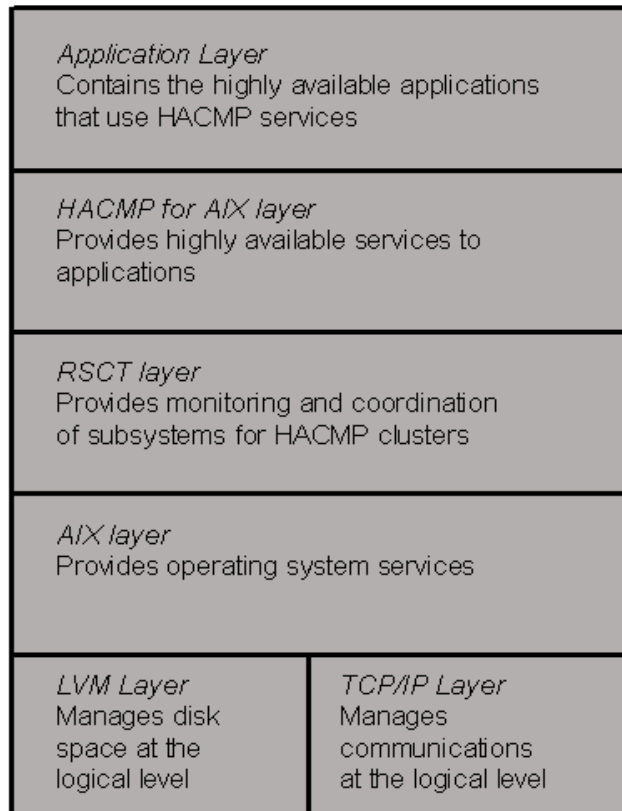
This section describes the HACMP software that implements a highly available environment.

It contains the following subsections:

- [Software Components of an HACMP Node](#)
- [HACMP Software Components](#)
- [Complementary Cluster Software](#)

Software Components of an HACMP Node

The following figure shows the layers of software on a node in an HACMP cluster:



A Model of an HACMP Cluster Node

The following list describes each layer:

- Application layer. Any applications made highly available through services provided by HACMP for AIX.
- HACMP for AIX layer. Software that recognizes changes within a cluster and coordinates the use of AIX features to create a highly available environment for critical data and applications. HACMP complements lower layers by providing additional services to enable applications with their associated communications, and data storage services to be highly available.
- RSCT layer. The IBM Reliable Scalable Cluster Technology services previously packaged with HACMP/ES (prior to HACMP 5.1), are now packaged with AIX 5.1 and 5.2. The RSCT layer provides facilities for monitoring node membership; network interface and communication interface health; event notification, synchronization and coordination via reliable messaging, in distributed or cluster environments. RSCT includes the Resource Monitoring and Control (RMC), Group Services, and Topology Services components. For more information, see the section [IBM Reliable Scalable Cluster Technology Availability Services](#) in this chapter and the RSCT documentation.
- AIX layer. Software that provides the underlying support for an HACMP cluster, including:

- Logical Volume Manager (LVM) subsystem layer, which manages data storage at the logical level.
- TCP/IP subsystem layer, which provides communications support for an HACMP cluster.

HACMP Software Components

The HACMP software has the following components:

- [Cluster Manager](#)
- [Cluster Secure Communication Subsystem](#)
- [IBM Reliable Scalable Cluster Technology Availability Services](#)
- [Cluster SMUX Peer and SNMP Monitoring Programs](#)
- [Cluster Information Program](#)
- [Highly Available NFS Server](#)
- [Shared External Disk Access](#)
- [Concurrent Resource Manager](#).

Cluster Manager

The Cluster Manager is a daemon that runs on each cluster node. The main task of the Cluster Manager is to respond to unplanned events, such as recovering from software and hardware failures, or user-initiated events, such as a joining node event. The RSCT subsystem informs the Cluster Manager about node and network-related events.

Note: In HACMP clusters, the RSCT software components—Group Services, Resource Monitoring and Control (RMC), and Topology Services—are responsible for most of the cluster monitoring tasks. For more information, see the RSCT documentation. For information on the architecture of the HACMP 5.2 (which includes the Enhanced Scalability functionality) product system, see the diagram in the section [IBM Reliable Scalable Cluster Technology Availability Services](#).

Changes in the state of the cluster are referred to as *cluster events*. The Cluster Manager runs scripts in response to cluster events. On each node, the Cluster Manager monitors local hardware and software subsystems for events, such as an “application failure” event. In response to such events, the Cluster Manager runs one or more event scripts, such as a “restart application” script. Cluster Managers on all nodes exchange messages to coordinate any actions required in response to an event.

Cluster Manager Connection to Other HACMP Daemons

The Cluster Manager maintains a connection to the Cluster SMUX Peer daemon, **clsmuxpd**, which gathers cluster information from the Cluster Manager relative to cluster state changes of nodes and interfaces. The Cluster Information Program (Cinfo) gets this information from

clsmuxpd and allows clients communicating with this program to be aware of changes in a cluster's state. This cluster state information is stored in the HACMP Management Information Base (MIB).

If your system is running TME 10 NetView, the Cluster Manager's connection to the local **clsmuxpd** also allows the HAView utility to obtain cluster state information and to display it graphically through the NetView map. See [Chapter 7: HACMP Configuration Process and Facilities](#), for information about how HAView communicates with **clsmuxpd**.

Cluster Secure Communication Subsystem

In HACMP 5.1, a common communication infrastructure was introduced to increase the security of intersystem communication. Cluster utilities use the Cluster Communications daemon that runs on each node for communication between the nodes. Because there is only one common communications path, all communications are reliably secured.

Although most components communicate through the Cluster Communications daemon, the following components use another mechanism for inter-node communications:

Component	Communication Method
Cluster Manager	RSCT
Cluster SMUX Peer (clsmuxpd)	SNMP
Cluster Information Program (Clinfo)	SNMP

For users who require additional security, HACMP 5.2 provides message authentication and message encryption for messages sent between cluster nodes.

Connection Authentication

HACMP provides two modes for connection authentication:

- **Standard.** Standard security mode checks the source IP address against an access list, checks that the value of the source port is between 571 and 1023, and uses the principle of least-privilege for remote command execution. Standard security is the default security mode.
- **Kerberos.** Kerberos security mode uses Kerberos security for authentication. Kerberos security is available only on systems running the PSSP software (SP or IBM eServer Cluster 1600).

For added security, you can set up a VPN for connections between nodes for HACMP inter-node communications.

Message Authentication and Encryption

Message authentication and message encryption provide additional security for HACMP messages sent between cluster nodes. Message authentication ensures the origination and integrity of a message. Message encryption changes the appearance of the data as it is transmitted and translates it to its original form when received by a node that authenticates the message.

You can configure the security options and options for distributing encryption keys using the SMIT interface.

IBM Reliable Scalable Cluster Technology Availability Services

The IBM Reliable Scalable Cluster Technology (RSCT) high availability services provide greater scalability, notify distributed subsystems of software failure, and coordinate recovery and synchronization among all subsystems in the software stack.

RSCT handles the heartbeats and network failure detection. The HACMP and RSCT software stack runs on each cluster node.

The HACMP Cluster Manager obtains indications of possible failures from several sources:

- RSCT monitors the state of the network devices
- AIX LVM monitors the state of the volume groups and disks
- Application monitors monitor the state of the applications.

The HACMP Cluster Manager drives the cluster recovery actions in the event of a component failure. RSCT running on each node exchanges a heartbeat with its peers so that it can monitor the availability of the other nodes in the cluster. If the heartbeat stops, the peer systems drive the recovery process. The peers take the necessary actions to get the critical applications running and to ensure that data has not been corrupted or lost.

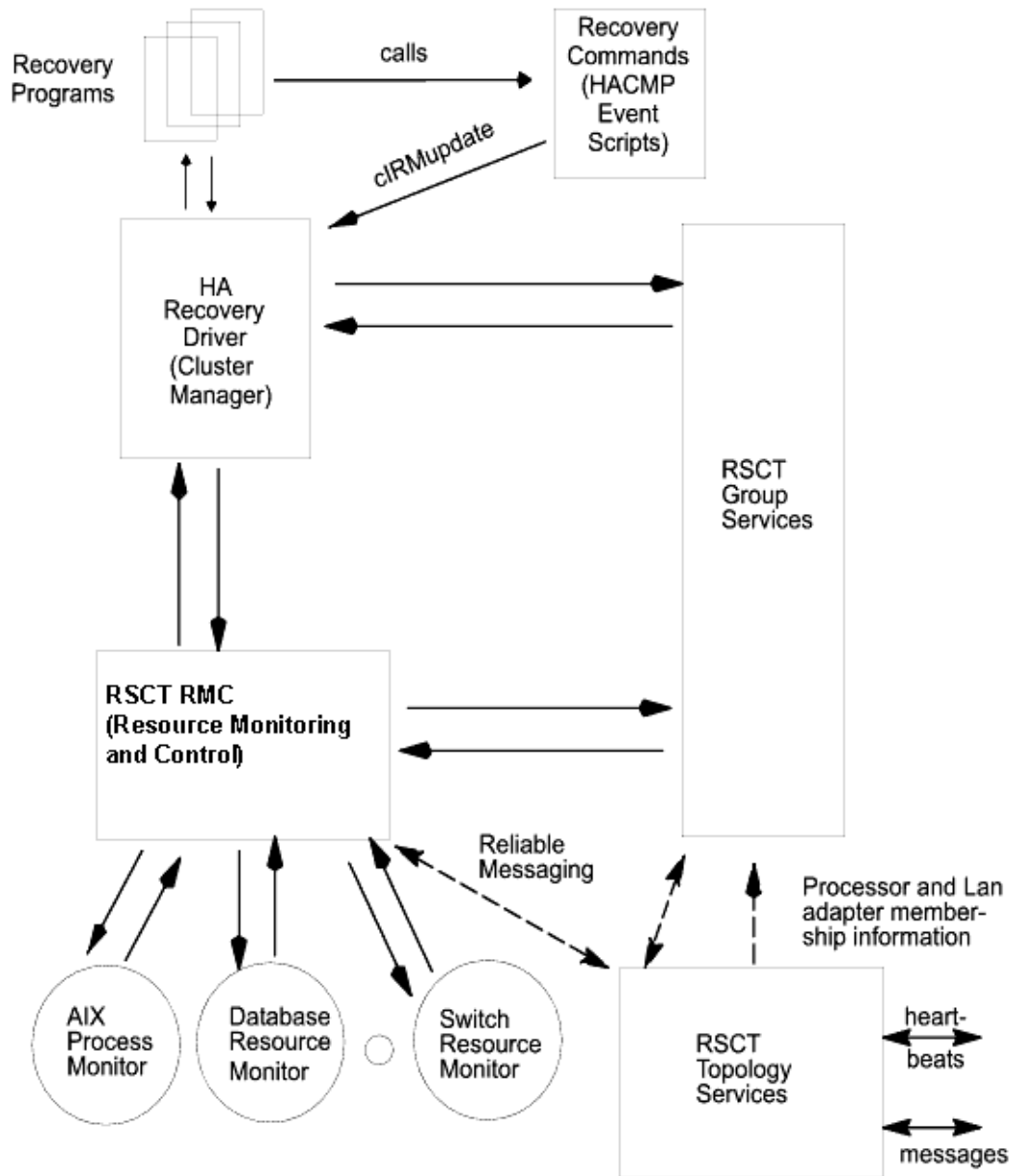
RSCT services include the following components:

- Resource Monitoring and Control (previous versions of HACMP use the Event Management subsystem). A distributed subsystem providing a set of high availability services. It creates events by matching information about the state of system resources with information about resource conditions of interest to client programs. Client programs in turn can use event notifications to trigger recovery from system failures.
- Group Services. A system-wide, highly available facility for coordinating and monitoring changes to the state of an application running on a set of nodes. Group Services helps both in the design and implementation of highly available applications and in the consistent recovery of multiple applications. It accomplishes these two distinct tasks in an integrated framework.
- Topology Services. A facility for generating heartbeats over multiple networks and for providing information about network interface membership, node membership, and routing. Network interface and node membership provide indication of NIC and node failures respectively. Reliable Messaging uses the routing information to route messages between nodes around adapter failures.

For more information on these services, see the following URL:

<http://www.ibm.com/servers/eserver/pseries/library/clusters/rsct.html>

The following figure shows the main components that make up the HACMP architecture:



HACMP Comprises IBM RSCT Availability Services and the HA Recovery Driver

Cluster SMUX Peer and SNMP Monitoring Programs

An HACMP cluster is dynamic and can undergo various transitions in its state over time. For example, a node can join or leave the cluster, or another IP label can replace a service IP label on a physical network interface card. Each of these changes affects the composition of the cluster, especially when highly available clients and applications must use services provided by cluster nodes.

SNMP Support

The Cluster SNMP Multiplexing (SMUX) Peer provides Simple Network Management Protocol (SNMP) support to client applications. SNMP is an industry-standard specification for monitoring and managing TCP/IP-based networks. SNMP includes a protocol, a database specification, and a set of data objects. A set of data objects forms a Management Information Base (MIB). SNMP provides a standard MIB that includes information such as IP addresses and the number of active TCP connections. The standard SNMP agent is the **snmpd** daemon.

SNMP can be extended through the use of the SMUX protocol to include *enterprise-specific* MIBs that contain information relating to a discrete environment or application. A SMUX peer daemon maintains information about the objects defined in its MIB and passes this information on to a specialized network monitoring or network management station.

HACMP MIB

The Cluster SMUX Peer daemon, **clsmuxpd**, maintains cluster status information in a special HACMP MIB. When **clsmuxpd** starts on a cluster node, it registers with the SNMP daemon, **snmpd**, and then continually gathers cluster information from the Cluster Manager daemon. The Cluster SMUX Peer daemon maintains an updated topology map of the cluster in the HACMP MIB as it tracks events and resulting states of the cluster.

For more information on the HACMP MIB, see the *Programming Client Applications Guide*.

Cluster Information Program

The Cluster Information Program (Cinfo), the **clinfo** daemon, is an SNMP-based monitor. Cinfo, running on a client machine or on a cluster node, queries the Cluster SMUX Peer for updated cluster information. Through Cinfo, information about the state of an HACMP cluster, nodes, and networks can be made available to clients and applications.

Clients can be divided into two categories: naive and intelligent.

- A *naive* client views the cluster complex as a single entity. If a cluster node fails, the client must be restarted (or at least must reconnect to the node), if IP address takeover (IPAT) is not enabled.
- An *intelligent* client, on the other hand, is cluster-aware—it reacts appropriately to node failure, connecting to an alternate node and perhaps masking the failure from the user. Such an intelligent client must have knowledge of the cluster state.

Note: For conceptual information about IP address takeover, see [Chapter 2: HACMP Cluster Nodes, Networks and Heartbeating](#).

The HACMP software extends the benefits of highly available servers, data, and applications to clients by providing notification of cluster state changes to clients through the **clsmuxpd** and Cinfo API functions.

Responding to Cluster Changes

Cinfo calls the `/usr/es/sbin/cluster/etc/clinfo.rc` script whenever a cluster, network, or node event occurs. By default, the **clinfo.rc** script flushes the system's ARP cache to reflect changes to network IP addresses, and it does not update the cache until another address responds to a **ping** request. Flushing the ARP cache typically is not necessary if the HACMP *hardware*

address swapping facility is enabled because hardware address swapping maintains the relationship between an IP address and a hardware address. Hardware address swapping is described in more detail in [Chapter 5: Ensuring Application Availability](#).

In a switched Ethernet network, you may need to flush the ARP cache to ensure that the new MAC address is communicated to the switch, or use the procedure, “MAC Address Is Not Communicated to the Ethernet Switch,” described in the *Administration and Troubleshooting Guide* to ensure that the hardware address is communicated correctly.

You can add logic to the `clinfo.rc` script if further action is desired.

Clinfo APIs

The Clinfo APIs provide application developers with both a C and a C++ language interface for accessing cluster status information. The HACMP software includes two versions of the Clinfo APIs: one for single-threaded applications and one for multi-threaded applications.

Clinfo and its associated APIs enable developers to write applications that recognize and respond to changes in a cluster. For more information, see the *Programming Client Applications Guide*.

Highly Available NFS Server

The highly available NFS server functionality is included in the HACMP product subsystem. A highly available NFS server allows a backup processor to recover current NFS activity should the primary NFS server fail. The NFS server special functionality includes highly available modifications and locks on network filesystems (NFS). You can do the following:

- Use the reliable NFS server capability that preserves locks and dupcache (2-node clusters only)
- Specify a network for NFS cross-mounting
- Define NFS exports and cross-mounts at the directory level
- Specify export options for NFS-exported directories and filesystems
- Configure two nodes to use NFS.

Note: While HACMP clusters can contain up to 32 nodes, clusters that use NFS can have a maximum of two nodes.

Shared External Disk Access

The HACMP software supports two methods of shared external disk access: non-concurrent and concurrent. Both methods of shared external disk access are described in the following sections.

Non-Concurrent Shared External Disk Access

In a non-concurrent environment, only one node has access to a shared external disk at a given time. If this node fails, one of the peer nodes acquires the disk, mounts filesystems defined as resources, and restarts applications to restore critical services. Typically, this takes from 30 to 300 seconds, depending on the number and size of the filesystems.

Supported Shared External Disk Types

A non-concurrent configuration can use:

- SCSI disks
- SCSI disk arrays
- serial disks
- SSA disks as shared external disks
- Fibre Channel direct-attached disks
- Fibre Channel SAN-attached disks.

For more information about supported devices, see the section [Disk Subsystems](#) in this chapter.

Mirroring

To prevent a failed disk from becoming a single point of failure, each logical volume in a shared volume group should be mirrored using the AIX LVM facility. If you are using an IBM Enterprise Storage System or other supported RAID array, do not use LVM mirroring. RAID devices provide their own data redundancy.

Applications

Most software that can run in single-machine mode can be managed by the HACMP software without modification.

Non-concurrent access typically does not require any code changes to server programs (a database management system, for example), or to applications to provide a highly available solution. To end users, node failure looks like a very fast machine reboot. One of the surviving nodes takes ownership of the failed node's resource groups and restarts the highly available applications. The Journaled Filesystem, the native AIX filesystem, guarantees filesystem integrity. The server program guarantees transaction data integrity.

End users simply log onto one of the surviving nodes and restart the application. The logon and application restart procedures can be driven by the HACMP software. In some HACMP configurations, users can continue without having to take any action—they simply experience a delay during failover.

Concurrent Shared External Disk Access

Note: For information on enhanced concurrent access, see the section [Enhanced Concurrent Mode](#) in this chapter.

The concurrent access feature enhances the benefits provided by an HACMP cluster. *Concurrent access* allows simultaneous access to a volume group on a disk subsystem attached to multiple (up to 32) nodes. Using concurrent access, a cluster can offer nearly continuous availability of data that rivals fault tolerance, but at a much lower cost. Additionally, concurrent access provides higher performance, eases application development, and allows horizontal growth.

Since concurrent access provides simultaneous access to data from multiple nodes, additional tools may be required to prevent multiple nodes from modifying the same block of data in a conflicting way. The HACMP software provides the Clinfo program that prepares an

application to run in a concurrent access environment. The Clinfo API provides an API through which applications may become "cluster-aware". The Clinfo tool is described earlier in this chapter.

The benefits of concurrent shared external disk access include the following:

- *Transparent Recovery Increases Availability.* Concurrent access significantly reduces the time for a fallover—sometimes to just a few seconds—because the peer systems already have physical access to the shared disk and are running their own instances of the application.

In a concurrent access environment, fallover basically involves backing out in-flight transactions from the failed processor. The server software running on the surviving nodes is responsible for recovering any partial transactions caused by the crash.

Since all nodes have concurrent access to the data, a client/server application can immediately retry a failed request on the surviving nodes, which continue to process incoming transactions.

- *Harnessing Multiple Processors Increases Throughput.* Applications are no longer limited to the throughput of a single processor. Instead, multiple instances of an application can run simultaneously on multiple processors. As more processing power is required, more systems can be added to the cluster to increase throughput.
- *Single Database Image Eases Application Development and Maintenance.* In a non-concurrent environment, the only route to improving performance is to partition an application and its data. Breaking code and data into pieces makes both application development and maintenance more complex.

Splitting a database requires a high degree of expertise to make sure that the data and workload are evenly distributed among the processors.

Partitioning code and data is not necessary in a concurrent access environment. To increase throughput, multiple instances of the same application running on different processors can simultaneously access a database on a shared external disk.

Supported Shared External Disk Types

A concurrent configuration can use:

- SCSI disks
- SCSI disk arrays
- serial disks
- SSA disks as shared external disks
- Fibre Channel direct-attached disks
- Fibre Channel SAN-attached disks.

SCSI disk arrays, serial disks, and SSA disks as shared external disks. For more information about supported devices, see the section [Disk Subsystems](#) in this chapter.

Mirroring

When creating concurrent access logical volumes, use LVM mirroring to avoid having the disks be a single point of failure, except for RAID disk subsystems that supply their own mirroring.

Applications

Concurrent access does not support the use of the Journaled File System. Therefore, the database manager must write directly to the raw logical volumes or `hdisks` in the shared volume group.

An application must use some method to arbitrate all requests for shared data. Most commercial UNIX databases provide a locking model that makes them compatible with the HACMP software. Check with your database vendor to determine whether a specific application supports concurrent access processing.

Concurrent Resource Manager

The Concurrent Resource Manager of HACMP provides concurrent access to shared disks in a highly available cluster, allowing tailored actions to be taken during takeover to suit business needs.

Concurrent Resource Manager adds enhanced-concurrent support for shared volume groups on all types of disks, and concurrent shared-access management for supported RAID and SSA disk subsystems.

Enhanced Concurrent Mode

AIX 5L v.5.1 and up provides a new form of concurrent mode: *enhanced concurrent mode*. In enhanced concurrent mode, the instances of the Concurrent Logical Volume Manager (CLVM) coordinate changes between nodes through the Group Services component of the Reliable Scalable Cluster Technology (RSCT) facility in AIX. Group Services protocols flow over the communications links between the cluster nodes. Support for enhanced concurrent mode and prior concurrent mode capabilities has the following differences:

- Any disk supported by HACMP for attachment to multiple nodes can be in an enhanced concurrent mode volume group; the special facilities of SSA disks are not required.
- The same capabilities for online changes of volume group and logical volume structure that have always been supported by AIX for SSA concurrent mode volume groups are available for enhanced concurrent mode volume groups.

You should keep in mind the following when planning an HACMP environment:

- When concurrent volume groups are created on AIX 5.1 and up, they are created as enhanced concurrent mode volume groups by default.
- If one node in a concurrent resource group runs a 64-bit kernel, then it must use an enhanced concurrent mode volume group.
- SSA concurrent mode is not supported on 64-bit kernels.
- SSA disks with the 32-bit kernel can use SSA concurrent mode.
- The C-SPOC utility does not work with RAID concurrent volume groups. You need to convert them to enhanced concurrent mode (otherwise, AIX 5L v.5.1 sees them as non-concurrent).
- The C-SPOC utility does not allow to create *new* SSA concurrent mode volume groups, if you are running AIX 5L v.5.2. (If you upgraded from previous releases of HACMP, you can use existing volume groups in SSA concurrent mode, but C-SPOC does not allow to create new groups of this type.) You can convert these volume groups to enhanced concurrent mode.

- You can include enhanced concurrent mode volume groups into shared resource groups. HACMP lists them in volume group picklists in resource group configuration SMIT panels. When enhanced concurrent volume groups are used in a non-concurrent environment, the volume groups are *not* concurrently accessed, they are still accessed by only one node at any given time.

Fast Disk Takeover

Failed volume groups are taken over faster than in previous releases of HACMP due to the improved disk takeover mechanism. If you have installed AIX 5L v. 5.2 and HACMP 5.1 and greater, and if you include in your non-concurrent resource groups enhanced concurrent mode volume groups, HACMP automatically detects these volume groups, and ensures that the faster option for volume group takeover is launched in the event of a node failure.

This functionality is especially useful for fallover of volume groups made up of a large number of disks.

During fast disk takeover, HACMP skips the extra processing needed to break the disk reserves, or update and synchronize the LVM information by running lazy update. As a result, the disk takeover mechanism used for enhanced concurrent volume groups is faster than disk takeover used for standard volume groups.

In addition, enhanced concurrent volume groups are included as choices in picklists for shared resource groups in SMIT panels for adding/changing resource groups, and in C-SPOC SMIT panels.

Complementary Cluster Software

A broad range of additional tools aids you in efficiently building, managing and expanding high availability clusters in AIX environments. These include:

- General Parallel File System (GPFS) for AIX, a cluster-wide filesystem that allows users shared access to files that span multiple disk drives and multiple nodes.
- Workload Manager for AIX provides resource balancing between applications.
- High Availability Geographic Cluster (HAGEO) for AIX provides disaster recovery.
- HACMP/XD features provide software solutions for disaster recovery. For more information, see [HACMP/XD](#) in this guide.

4 **HACMP Cluster Hardware and Software**

Complementary Cluster Software

Chapter 5: Ensuring Application Availability

This chapter describes how the HACMP software ensures application availability by ensuring the availability of cluster components. HACMP eliminates single points of failure for all key system components, and eliminates the need for scheduled downtime for most routine cluster maintenance tasks.

This chapter covers the following topics:

- [Eliminating Single Points of Failure in an HACMP Cluster](#)
- [Minimizing Scheduled Downtime with HACMP](#)
- [Minimizing Unscheduled Downtime](#)
- [Minimizing Takeover Time: Fast Disk Takeover](#)
- [Cluster Events](#)

Overview

The key facet of a highly available cluster is its ability to detect and respond to changes that could interrupt the essential services it provides. The HACMP software allows a cluster to continue to provide application services critical to an installation even though a key system component—a network interface card, for example—is no longer available. When a component becomes unavailable, the HACMP software is able to detect the loss and shift the workload from that component to another component in the cluster. In planning a highly available cluster, you attempt to ensure that key components do not become *single points of failure*.

In addition, HACMP software allows a cluster to continue providing application services while routine maintenance tasks are performed using a process called *dynamic reconfiguration*. In dynamic reconfiguration, you can change components in a running cluster, such as adding or removing a node or network interface, without having to stop and restart cluster services. The changed configuration becomes the active configuration dynamically. You can also dynamically replace a failed disk.

The following sections describe conceptually how to use the HACMP software to:

- Eliminate single points of failure in a cluster.
- Minimize scheduled downtime in an HACMP cluster with the dynamic reconfiguration, resource group management, and cluster management (C-SPOC) utilities.
- Minimize unscheduled downtime with the fast recovery feature, and by specifying a delayed fallback timer policy for custom resource groups.
- Minimize the time it takes to perform disk takeover.
- How to interpret and emulate cluster events.

Also, this chapter describes how to interpret and emulate cluster events.

Note: You may need to monitor the cluster activity while a key component fails and the cluster continues providing availability of an application. For more information on which monitoring and diagnostic tools you can use, see [Chapter 7: HACMP Configuration Process and Facilities](#).

Eliminating Single Points of Failure in an HACMP Cluster

The HACMP software enables you to build clusters that are both highly available and scalable by eliminating single points of failure. A *single point of failure* exists when a critical cluster function is provided by a single component. If that component fails, the cluster has no other way to provide that function and essential services become unavailable.

For example, if all the data for a critical application resides on a single disk that is not mirrored, and that disk fails, the disk has become a single point of failure for the entire system. Client nodes cannot access that application until the data on the disk is restored.

Potential Single Points of Failure in an HACMP Cluster

HACMP provides recovery options for the following cluster components:

- Nodes
- Applications
- Networks and network interfaces
- Disks and disk adapters.

To be highly available, a cluster must have no single point of failure. While the goal is to eliminate all single points of failure, compromises may have to be made. There is usually a cost associated with eliminating a single point of failure. For example, redundant hardware increases cost. The cost of eliminating a single point of failure should be compared to the cost of losing services should that component fail. Again, the purpose of the HACMP software is to provide a cost-effective, highly available computing environment that can grow to meet future processing demands.

Eliminating Nodes as a Single Point of Failure

Nodes leave the cluster either through a planned transition (a node shutdown or stopping cluster services on a node), or because of a failure.

Node failure begins when a node monitoring a neighbor node ceases to receive heartbeat traffic for a defined period of time. If the other cluster nodes agree that the failure is a node failure, the failing node is removed from the cluster and its resources are taken over by the nodes configured to do so. An active node may, for example, take control of the shared disks configured on the failed node. Or, an active node may masquerade as the failed node (by acquiring its service IP address) and run the processes of the failed node while still maintaining its own processes. Thus, client applications can switch over to a surviving node for shared-disk and processor services.

The HACMP software provides the following facilities for processing node failure:

- Disk takeover

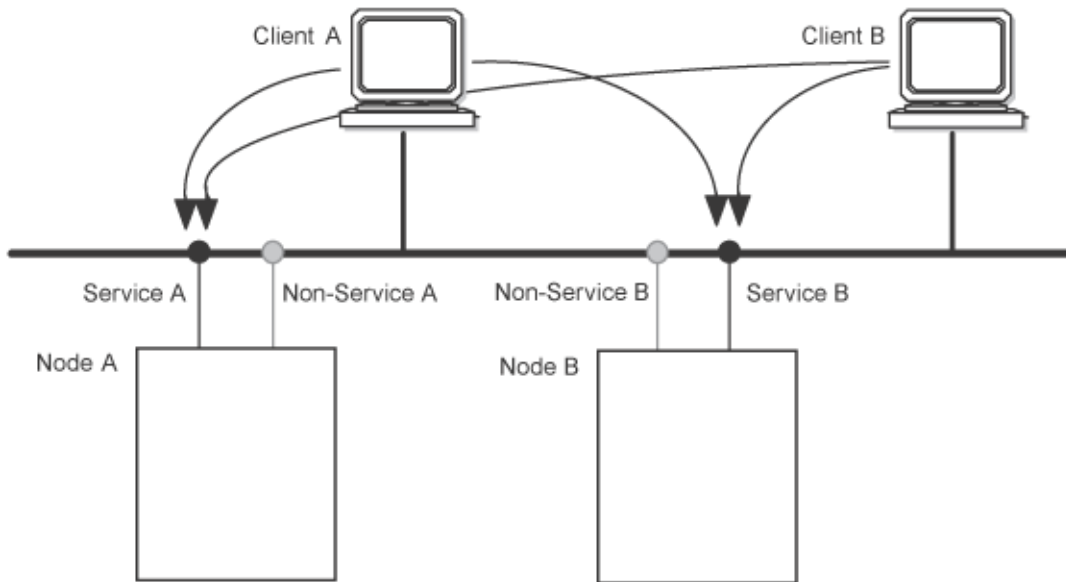
- IP Address Takeover via IP Aliases
- IP Address Takeover via IP Replacement (with or without Hardware Address Takeover).

Disk Takeover

In an HACMP environment, shared disks are physically connected to multiple nodes.

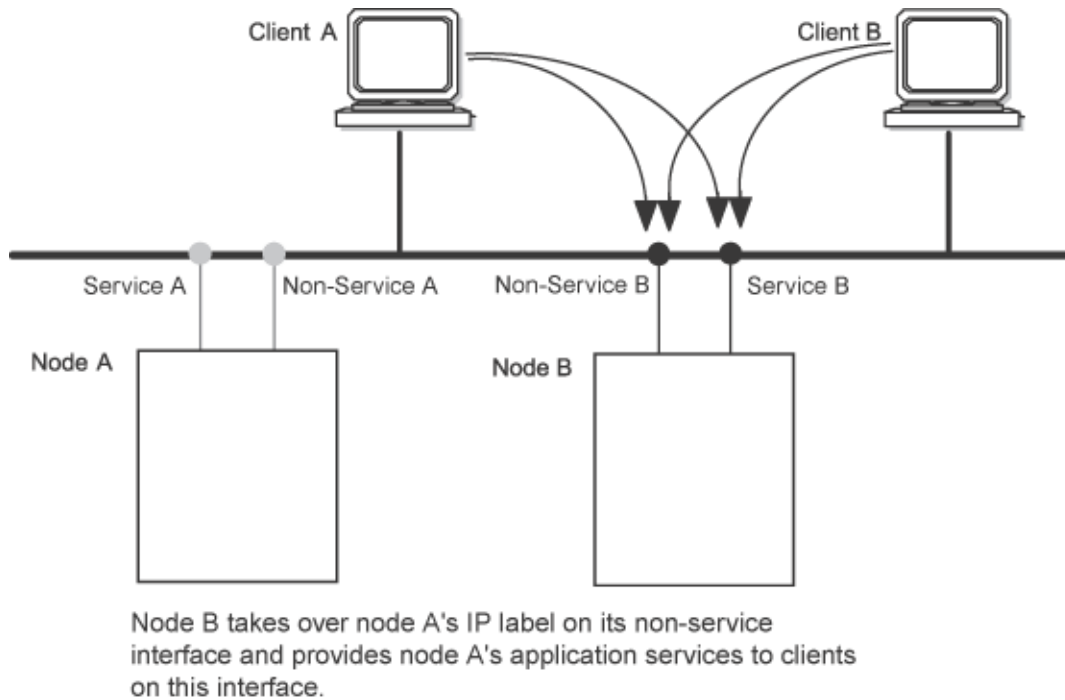
Disk Takeover in Concurrent Environments

In concurrent access configurations, the shared disks are actively connected to multiple nodes at the same time. Therefore, disk takeover is not required when a node leaves the cluster. The following figures illustrate disk takeover in concurrent environments.



Each node provides a separate network service.

The Concurrent Access Configuration Before Disk Takeover



The Concurrent Access Configuration After Disk Takeover

Fast Disk Takeover

In the case of a cluster failure, enhanced concurrent volume groups are taken over faster than in previous releases of HACMP due to the improved disk takeover mechanism. HACMP automatically detects enhanced concurrent volume groups and ensures that the faster option for volume group takeover is launched in the event of a node failure. For more information, see [Minimizing Takeover Time: Fast Disk Takeover](#) in this chapter.

Disk Takeover in Non-Concurrent Environments

In non-concurrent environments, only one connection is active at any given time, and the node with the active connection owns the disk. *Disk takeover* occurs when the node that currently owns the disk leaves the cluster and an active node assumes control of the shared disk so that it remains available. Note, however, that shared filesystems can be exported and NFS cross-mounted by other cluster nodes that are under the control of HACMP.

The `cl_export_fs` utility can use the optional `/usr/es/sbin/cluster/etc/exports` file instead of the standard `/etc/exports` file for determining export options.

IP Address Takeover

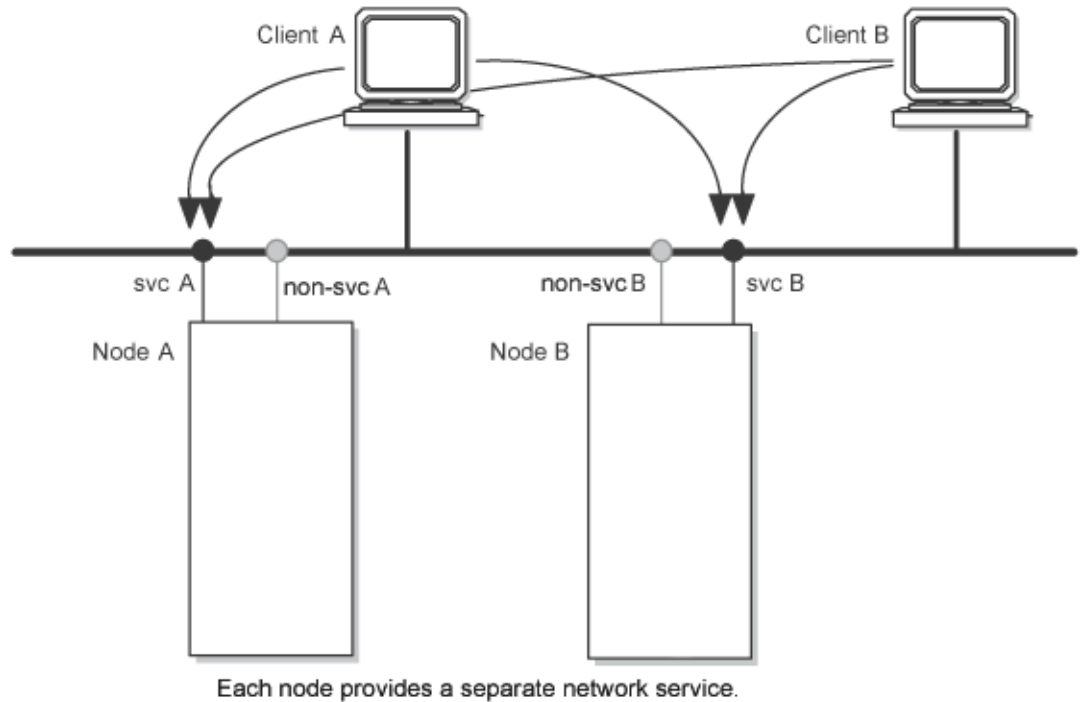
IP address takeover (IPAT) is a networking capability that allows a node to acquire the network address of a node that has left the cluster. IP address takeover is necessary in an HACMP cluster when a service being provided to clients is bound to a specific IP address, that is when a service IP label through which services are provided to the clients is included as a resource in a cluster resource group. If, instead of performing an IPAT, a surviving node simply did a disk and application takeover, clients would not be able to continue using the application at the specified server IP address.

HACMP uses two types of IPAT:

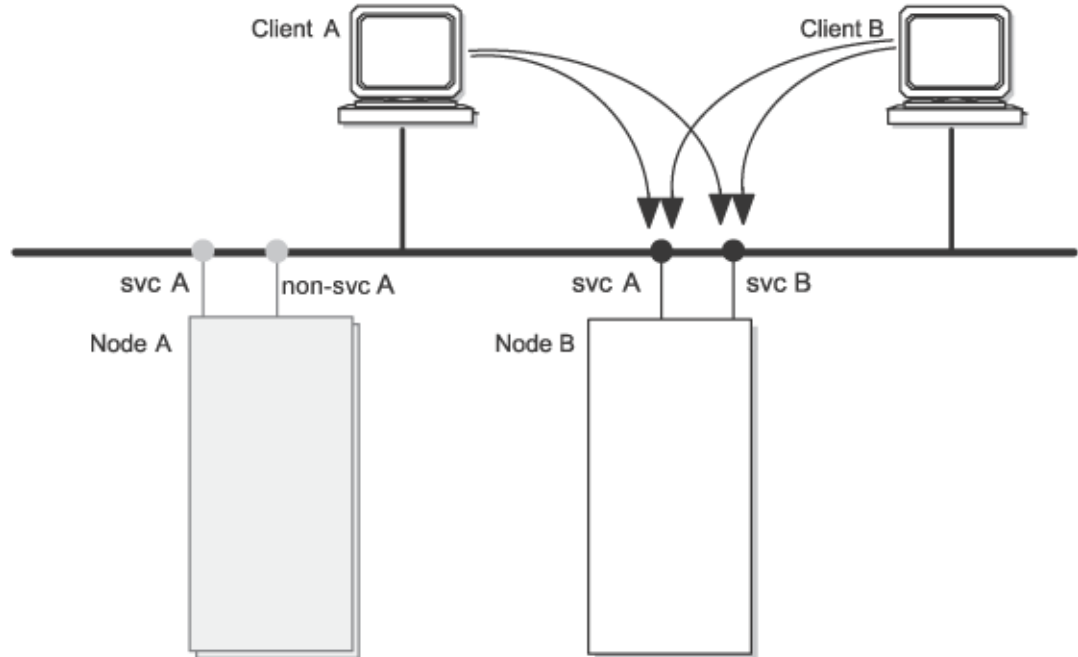
- IPAT via IP Aliases (this is the default)
- IPAT via IP Replacement.

For more information on each type, see [IP Address Takeover via IP Aliases](#) and [IP Address Takeover via IP Replacement](#) in [Chapter 2: HACMP Cluster Nodes, Networks and Heartbeating](#).

The following figures illustrate IP address takeover via IP Replacement.



The Configuration Before IP Address Takeover via IP Replacement



Node B assumes node A's IP address on its non-service interface and provides A's network service to clients.

The Configuration After IP Address Takeover via IP Replacement

Note: In an HACMP on the RS/6000 SP, special considerations apply to IP address takeover on the SP Switch network. For more information, see the *Planning and Installation Guide*.

Hardware Address Swapping and IP Address Takeover via IP Replacement

Hardware address swapping works in conjunction with IP address takeover via IP Replacement. With hardware address swapping enabled, a node also assumes the hardware network address (in addition to the IP address) of a node that has failed so that it can provide the service that the failed node was providing to the client nodes in the cluster. Hardware address swapping is also referred to as *hardware address takeover* (HWAT).

Without hardware address swapping, TCP/IP clients and routers which reside on the same subnet as the cluster nodes must have their Address Resolution Protocol (ARP) cache updated. The ARP cache contains a mapping of IP addresses to hardware addresses. The use of hardware address swapping is highly recommended for clients that cannot run the Clinfo daemon (machines not running AIX) or that cannot easily update their ARP cache.

Note: SP Switch networks do not support hardware address swapping. However, note that the SP switch networks can be configured so that their IP network interface cards update their ARP caches automatically when IP address takeover occurs. IP aliases are used in such cases.

Keep in mind that when an IP address takeover occurs, the netmask of the physical network interface card on which a service IP label is configured is obtained by the network interface card on another node; thus, the netmask follows the service IP address.

This means that with IPAT via IP Replacement, the netmask for all network interfaces in an HACMP network must be the same to avoid communication problems between network interfaces after an IP address takeover via IP Replacement, and during the subsequent release of the IP address acquired during takeover. The reasoning behind this requirement is as follows:

Communication problems occur when the network interface card (NIC) on another node releases the service IP address. This NIC assumes its original address, but retains the netmask of the service IP address. This address reassignment causes the NIC on another node to function on a different subnet from other backup NICs in the network. This netmask change can cause changes in the broadcast address and the routing information such that other backup NICs may now be unable to communicate on the same logical network.

Eliminating Applications as a Single Point of Failure

The primary reason to create HACMP clusters is to provide a highly available environment for mission-critical applications. For example, an HACMP cluster could run a database server program which services client applications. The clients send queries to the server program which responds to their requests by accessing a database, stored on a shared external disk.

In an HACMP cluster, these critical applications can be a single point of failure. To ensure the availability of these applications, the node configured to take over the resources of the node leaving the cluster should also restart these applications so that they remain available to client processes.

You can make an application highly available by using:

- An application server
- Cluster control
- Application monitors
- Application Availability Analysis Tool.

To put the application under HACMP control, you create an *application server* cluster resource that associates a user-defined name of the server with the names of user-provided written scripts to start and stop the application. By defining an application server, HACMP can start another instance of the application on the takeover node when a failover occurs.

Certain applications can be made highly available without application servers. You can place such applications under cluster control by configuring an aspect of the application as part of a resource group. For example, Fast Connect services can all be added as resources to a cluster resource group, making them highly available in the event of node or network interface failure.

Note: Application takeover is usually associated with IP address takeover. If the node restarting the application also acquires the IP service address on the failed node, the clients only need to reconnect to the same server IP address. If the IP address was not taken over, the client needs to connect to the new server to continue accessing the application.

Additionally, you can use the AIX System Resource Controller (SRC) to monitor for the presence or absence of an application daemon and to respond accordingly.

Application Monitors

You can also configure an *application monitor* to check for process failure or other application failures and automatically take action to restart the application.

In HACMP 5.2, you can configure multiple application monitors and associate them with one or more application servers. By supporting multiple monitors per application, HACMP can support more complex configurations. For example, you can configure one monitor for each instance of an Oracle parallel server in use. Or, you can configure a custom monitor to check the health of the database, and a process termination monitor to instantly detect termination of the database process.

Prior to HACMP 5.2, each application that is kept highly available could have only one of the two types of monitors configured for it. You could configure a monitor to:

- Check whether a specific process is terminated in the cluster
- or*
- Check the state of the application by the means of a customized script.

For example, you can supply a script to HACMP that sends a request to a database to check that it is functioning. A non-zero exit from the customized script will indicate a failure of the monitored application, and HACMP will respond by trying to recover the resource group that contains the application. In HACMP 5.2, you can use two monitors for one application.

In addition, in HACMP 5.1, the name of the monitor that you were configuring was required to be the same as the name of the application server. In HACMP 5.2, you can assign each monitor a unique name in SMIT.

Application Availability Analysis

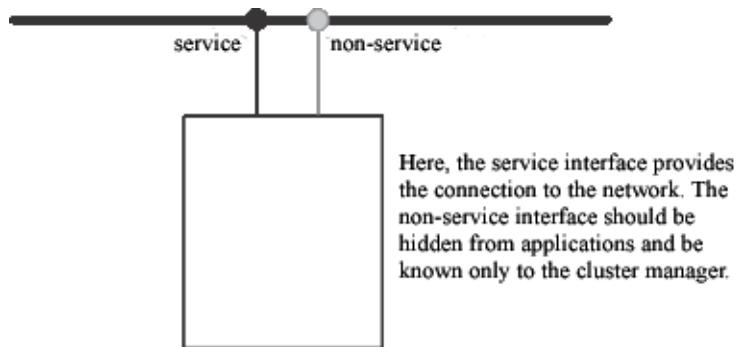
The Application Availability Analysis tool measures the exact amount of time that any of your applications has been available. The HACMP software collects, time-stamps, and logs extensive information about the applications you choose to monitor with this tool. Using SMIT, you can select a time period and the tool displays uptime and downtime statistics for a specific application during that period.

Eliminating Communication Interfaces as a Single Point of Failure

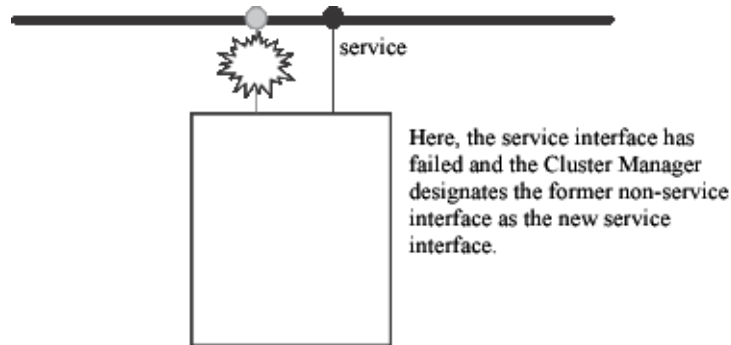
The HACMP software handles failures of network interfaces on which a service IP label is configured. Two types of such failures are:

- Out of two network interfaces configured on a node, the network interface with a service IP label fails, but an additional “backup” network interface card remains available on the *same* node. In this case, the Cluster Manager swaps the roles of these two interface cards on that node. Such a network interface failure is transparent to you except for a small delay while the system reconfigures the network interface on a node.
- Out of two network interfaces configured on a node, an additional, or a “backup” network interface fails, but the network interface with a service IP label configured on it remains available. In this case, the Cluster Manager detects a “backup” network interface failure, logs the event, and sends a message to the system console. If you want additional processing, you can customize the processing for this event.

The following figures illustrate network interface swapping that occurs on the *same* node:



The Configuration Before Network Adapter Swap



The Configuration After Network Adapter Swap

Hardware Address Swapping and Adapter Swapping

Hardware address swapping works in conjunction with adapter swapping (as well as IP address takeover via IP Replacement). With hardware address swapping enabled, the “backup” network interface assumes the hardware network address (in addition to the IP address) of the failed network interface that had the service IP label configured on it so that it can provide the service that the failed network interface was providing to the cluster clients.

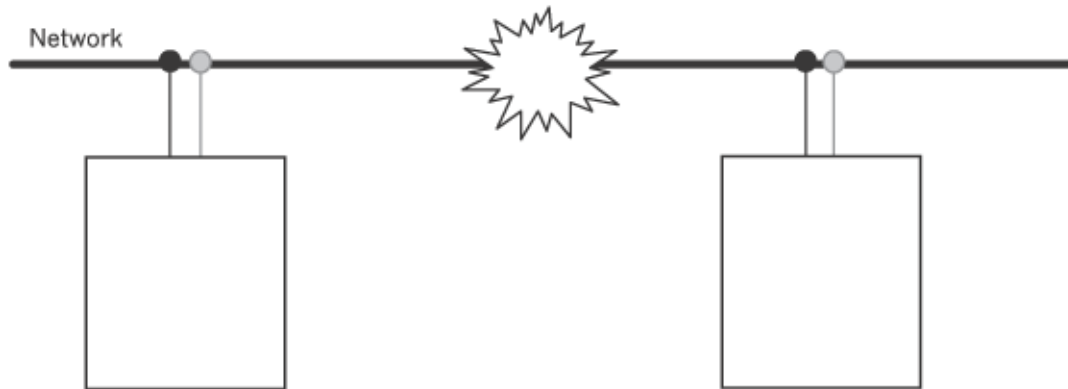
Without hardware address swapping, TCP/IP clients and routers which reside on the same subnet as the cluster nodes must have their Address Resolution Protocol (ARP) cache updated. The ARP cache contains a mapping of IP addresses to hardware addresses. The use of hardware address swapping is highly recommended for clients that cannot run the Clinfo daemon (machines not running AIX), or that cannot easily update their ARP cache.

Note: SP Switch networks do not support hardware address swapping. However, note that the SP switch networks can be configured such that their IP network interfaces update their ARP caches automatically when IP Address Takeover via IP Aliases occurs. For more information, see the *Administration and Troubleshooting Guide*.

Eliminating Networks as a Single Point of Failure

Network failure occurs when an HACMP network fails for all the nodes in a cluster. This type of failure occurs when none of the cluster nodes can access each other using any of the network interface cards configured for a given HACMP network.

The following figure illustrates a network failure:



Here, the network connecting the nodes has failed. The nodes are no longer able to communicate across this network.

Network Failure

The HACMP software's first line of defense against a network failure is to have the nodes in the cluster connected by multiple networks. If one network fails, the HACMP software uses a network that is still available for cluster traffic and for monitoring the status of the nodes.

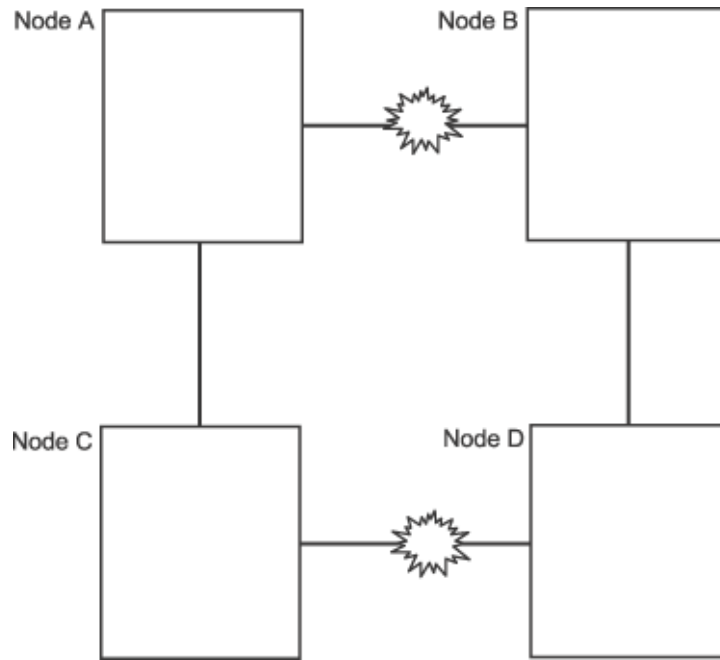
You can specify additional actions to process a network failure—for example, re-routing through an alternate network. Having at least two networks to guard against network failure is highly recommended.

When a local network failure event occurs, the Cluster Manager takes selective recovery actions for resource groups containing a service IP label connected to that network. The Cluster Manager attempts to move only the resource groups affected by the local network failure event, rather than all resource groups on a particular node.

Node Isolation and Partitioned Clusters

Node isolation occurs when all networks connecting two or more parts of the cluster fail. Each group (one or more) of nodes is completely isolated from the other groups. A cluster in which certain groups of nodes are unable to communicate with other groups of nodes is a *partitioned cluster*.

In the following illustration of a partitioned cluster, Node A and Node C are on one side of the partition and Node B and Node D are on the other side of the partition.



A Partitioned Cluster

The problem with a partitioned cluster is that the nodes on one side of the partition interpret the absence of heartbeats from the nodes on the other side of the partition to mean that those nodes have failed and then generate node failure events for those nodes. Once this occurs, nodes on each side of the cluster (if so configured) attempt to take over resources from a node that is still active and therefore still legitimately owns those resources. These attempted takeovers can cause unpredictable results in the cluster—for example, data corruption due to a disk being reset.

Using Device-Based Networks to Prevent Partitioning

To guard against the TCP/IP subsystem failure causing node isolation, each node in the cluster should be connected by a point-to-point non-IP-based network to its neighboring nodes, forming a logical “ring.” This logical ring of point-to-point networks reduces the chance of node isolation by allowing neighboring Cluster Managers to communicate even when all TCP/IP-based networks fail.

You can configure two kinds of point-to-point, non-IP-based networks in HACMP:

- Point-to-point networks, which use serial network interface cards and RS232 connections. Not all serial ports can be used for this function. For more information, see the *Planning and Installation Guide*.
- Disk networks, which use a shared disk and a disk bus as a point-to-point network. Any disk which is included in an HACMP enhanced concurrent volume group can be used. (You can also use TM SSA or SCSI disks which are not included in an enhanced concurrent volume group).

Point-to-point, device-based networks are especially important in concurrent access configurations so that data does not become corrupted when TCP/IP traffic among nodes is lost. Device-based networks do not carry TCP/IP communication between nodes; they only allow nodes to exchange heartbeats and control messages so that Cluster Managers have accurate information about the status of peer nodes.

Using Global Networks to Prevent Partitioning

You can also configure a “logical” global network that groups multiple networks of the same type. Global networks help to avoid node isolation when a network fails.

Eliminating Disks and Disk Adapters as a Single Point of Failure

The HACMP software does not itself directly handle disk and disk adapter failures. Rather, these failures are handled by AIX through LVM mirroring on disks and by internal data redundancy on the IBM 2105 ESS and SSA disks.

For example, by configuring the system with multiple SCSI-3 chains, serial adapters, and then mirroring the disks across these chains, any single component in the disk subsystem (adapter, cabling, disks) can fail without causing unavailability of data on the disk.

If you are using the IBM 2105 ESS and SSA disk arrays, the disk array itself is responsible for providing data redundancy.

The AIX Error Notification Facility

The AIX Error Notification facility allows you to detect an event not specifically monitored by the HACMP software—a disk adapter failure, for example—and to program a response to the event.

Permanent hardware errors on disk drives, controllers, or adapters can affect the fault resiliency of data. By monitoring these errors through error notification methods, you can assess the impact of a failure on the cluster’s ability to provide high availability. A simple implementation of error notification would be to send a mail message to the system administrator to investigate the problem further. A more complex implementation could include logic to analyze the failure and decide whether to continue processing, stop processing, or escalate the failure to a node failure and have the takeover node make the volume group resources available to clients.

It is strongly recommended that you implement an error notification method for all errors that affect the disk subsystem. Doing so ensures that degraded fault resiliency does not remain undetected.

AIX error notification methods are automatically used in HACMP to monitor certain recoverable LVM errors, such as volume group loss errors.

Automatic Error Notification

You can automatically configure error notification for certain cluster resources using a specific option in SMIT. If you select this option, error notification is automatically turned on on all nodes in the cluster for particular devices.

Certain non-recoverable error types are supported by automatic error notification: disk, disk adapter, and SP switch adapter errors. No media errors, recovered errors, or temporary errors are supported by this feature. One of two error notification methods is assigned for all error types supported by automatic error notification.

In addition, if you add a volume group to a resource group, HACMP creates an AIX Error Notification method for it. In the case a volume group loses quorum, HACMP uses this method to selectively move the affected resource group to another node. Do not edit or alter the error notification methods that are generated by HACMP.

Error Emulation

The Error Emulation utility allows you to test your error notification methods by simulating an error. When the emulation is complete, you can check whether your customized notification method was exercised as intended.

Minimizing Scheduled Downtime with HACMP

The HACMP software enables you to perform most routine maintenance tasks on an active cluster dynamically—without having to stop and then restart cluster services to make the changed configuration the active configuration. Several features contribute to this:

- Dynamic Reconfiguration (DARE)
- Resource Group Management (**clRGmove**)
- Cluster Single Point of Control (C-SPOC)
- Dynamic adapter swap for replacing hot-pluggable adapter cards.

Dynamic Automatic Reconfiguration (DARE)

This process, called *dynamic automatic reconfiguration* or *dynamic reconfiguration (DARE)*, is triggered when you synchronize the cluster configuration after making changes on an active cluster. Applying a cluster snapshot using SMIT also triggers a dynamic reconfiguration event.

For example, to add a node to a running cluster, you simply connect the node to the cluster, add the node to the cluster topology on any of the existing cluster nodes, and synchronize the cluster. The new node is added to the cluster topology definition on all cluster nodes and the changed configuration becomes the currently active configuration. After the dynamic reconfiguration event completes, you can start cluster services on the new node.

HACMP verifies the modified configuration before making it the currently active configuration to ensure that the changes you make result in a valid configuration.

How Dynamic Reconfiguration Works

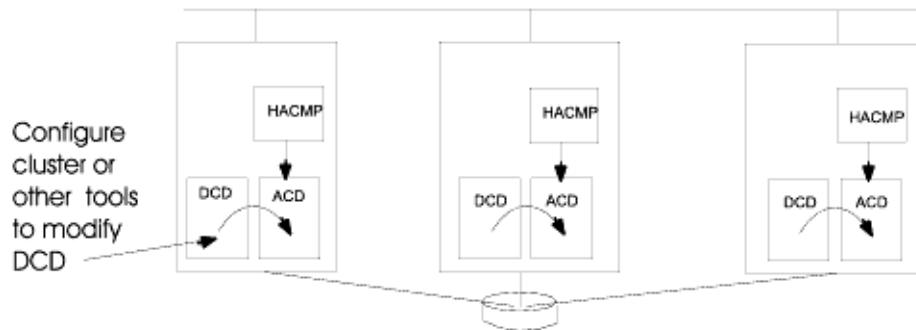
To allow for the dynamic reconfiguration of a running cluster, whenever HACMP starts, it creates a private copy of the HACMP-specific object classes stored in the system default Object Data Model (ODM). From now on, the HACMP ODM is referred to as the HACMP Configuration Database. All the HACMP daemons, scripts, and utilities on a running node reference the HACMP configuration database data stored in this private directory, called the Active Configuration Directory (ACD), instead of in the HACMP configuration database data stored in the system default HACMP configuration database, stored in the Default Configuration Directory (DCD).

Note: The operation of DARE is described here for completeness. No manual intervention is required to ensure that HACMP carries out these operations. HACMP correctly manages all dynamic reconfiguration operations in the cluster.

By default, the DCD is the directory named `/etc/objrepos`. This directory contains the default system object classes, such as the customized device database (CuDv) and the predefined device database (PdDv), as well as the HACMP-specific object classes. By default, the ACD is `/usr/es/sbin/cluster/etc/objrepos/active`.

Note: When you configure a cluster, you modify the HACMP configuration database data stored in the DCD—not data in the ACD. SMIT and other HACMP configuration utilities all modify the HACMP configuration database data in the DCD. In addition, all user commands that display HACMP configuration database data, such as the `cllsif` command, read data from the DCD.

The following figure illustrates how the HACMP daemons, scripts, and utilities all reference the ACD when accessing configuration information.



Relationship of HACMP to ACD at Cluster Start-Up

Reconfiguring a Cluster Dynamically

The HACMP software depends on the location of certain HACMP configuration database repositories to store configuration data. The presence or absence of these repositories are sometimes used to determine steps taken during cluster configuration and operation. The `ODMPATH` environment variable allows HACMP configuration database commands and subroutines to query locations other than the default location (held in the `ODMDIR` environment variable) if the queried object does not exist in the default location. You can set this variable, but it must not be set to include the `/etc/objrepos` directory or you will lose the integrity of the HACMP configuration information.

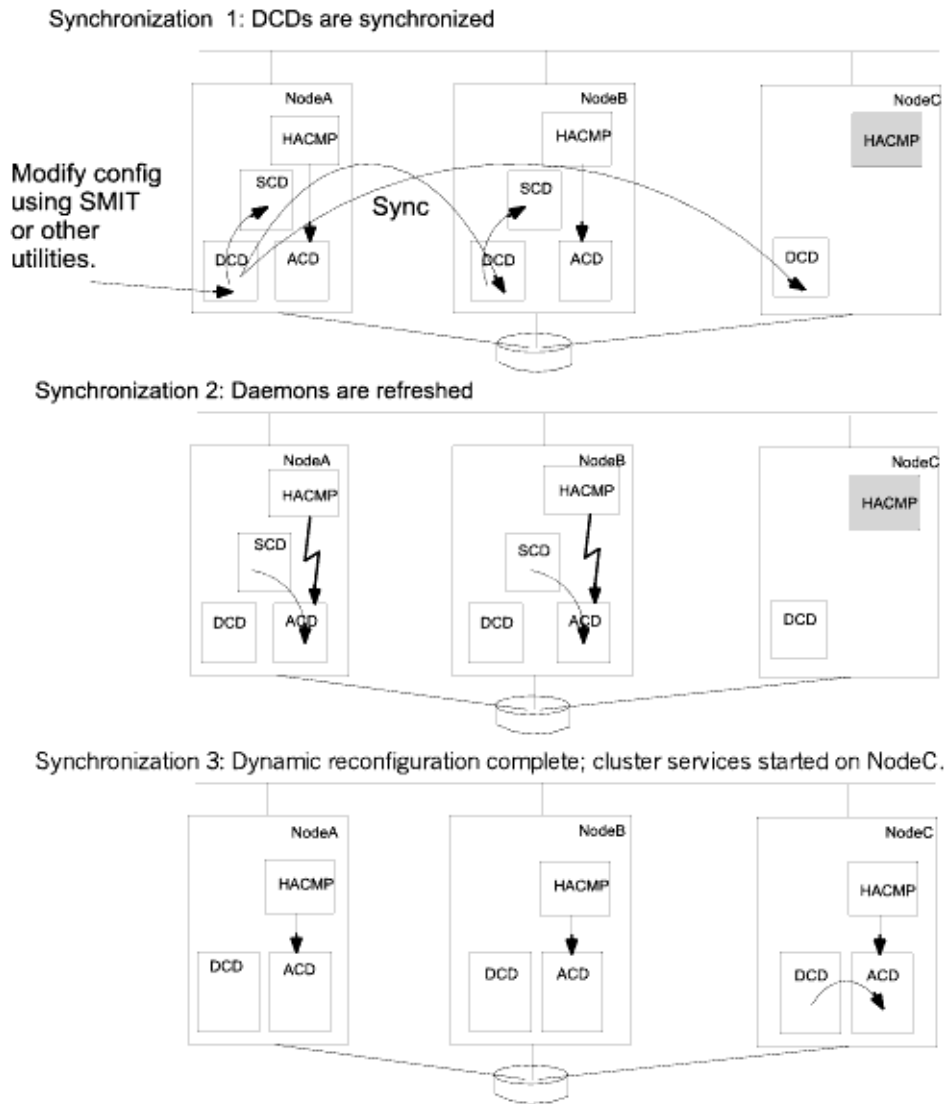
To change the configuration of an active cluster, you modify the cluster definition stored in the HACMP-specific HACMP configuration database classes stored in the DCD using SMIT. When you change the cluster configuration in an active cluster, you use the same SMIT paths to make the changes, but the changes do not take effect immediately. Therefore you can make several changes in one operation. When you synchronize your configuration across all cluster nodes, a cluster-wide dynamic reconfiguration event occurs. When HACMP processes a dynamic reconfiguration event, it updates the HACMP configuration database object classes

stored in the DCD on each cluster and replaces the HACMP configuration database data stored in the ACD with the new HACMP configuration database data in the DCD, in a coordinated, cluster-wide transition. It also refreshes the cluster daemons so that they reference the new configuration data.

After this processing, the cluster heartbeat is suspended briefly and the cluster is in an unstable state. The changed configuration becomes the active configuration. After cluster services are started on the newly added node, it can be integrated into the cluster.

The following figure illustrates the processing involved with adding a node to an active cluster using dynamic reconfiguration. The node to be added is connected to a running cluster but cluster services are inactive on this node. The configuration is redefined on NodeA. When the changes to the configuration are synchronized, the HACMP configuration database data stored in the DCD on NodeA is copied to the DCDs on other cluster nodes and a dynamic reconfiguration event is triggered. HACMP copies the new HACMP configuration database data in the DCD into a temporary location on each node, called the Staging Configuration Directory (SCD). The default location of the SCD is `/usr/es/sbin/cluster/etc/objrepos/stage`. By using this temporary location, HACMP allows you to start making additional configuration changes while a dynamic reconfiguration is in progress. Before copying the new HACMP configuration database data in the SCD over the current HACMP configuration database data in the ACD, HACMP verifies the new configuration.

Note: You can initiate a second reconfiguration while a dynamic reconfiguration is in progress, but you cannot synchronize it. The presence of an SCD on any cluster node acts as a lock, preventing the initiation of a new dynamic reconfiguration.



Dynamic Reconfiguration Processing

Resource Group Management

You can use the *Resource Group Management* utility to move resource groups to other cluster nodes.

Note: In previous releases, you could use the DARE Resource Migration utility for this purpose. In HACMP 5.1, The DARE Resource Migration is replaced with Resource Group Management (**cIRGmove**).

The Resource Group Management utility provides a means for managing resource groups in a cluster, and enhances failure recovery capabilities of HACMP. It allows you to change the status of any type of resource group, or the location of a non-concurrent resource group (along

with its resources—IP addresses, applications, and disks), without stopping cluster services. For instance, you can use this utility to free the node of any resource groups to perform system maintenance on a particular cluster node.

Use Resource Group Management to:

- Dynamically move a specified non-concurrent resource group from a node it currently resides on to the destination node that you have specified.
- Bring a resource group online or offline on one or all nodes in the cluster. This functionality is different for concurrent and non-concurrent resource groups. See the *Administration Guide* for detailed information on what kinds of online and offline operations you can perform on concurrent and non-concurrent resource groups.

Note: When you request HACMP to perform resource group management, it uses the **cIRGmove** utility, which moves resource groups by calling an **rg_move** event. It is important to distinguish between an **rg_move** event that in some cases is triggered automatically by HACMP, and an **rg_move** event that occurs when you request HACMP to manage resource groups for you. To track and identify the causes of operations performed on the resource groups in the cluster, look for the command output in SMIT and for the information in the **hacmp.out** file.

In addition, you can specify the *priority override location* for a migrated resource group. The priority override location is the destination node to which the resource group has been moved, or on which the resource group has been taken online or offline. The priority override location can refer to multiple nodes in the case of concurrent resource groups.

The priority override location can either be:

- *persistent* and remain in effect after HACMP cluster services are restarted on any or all nodes in the cluster

or

- *non-persistent*, that is remain in effect only until the cluster services are stopped on all nodes.

In releases prior to HACMP 5.1, you could specify the “sticky” migration attribute for a migrated resource group. In HACMP 5.1, the “sticky” attribute is removed. Instead, you can specify the priority override location attribute. The functionality of the “sticky” attribute in HACMP 5.1 is equivalent to that of the *persistent priority override location attribute*.

Note: The **clbare** command is not supported and is replaced with the **cIRGmove** command. To move resource groups to the specified node or state, use the **HACMP Resource Group Management SMIT** menu, or the **cIRGmove** command. See the man page for the **cIRGmove** command for more information.

For more information about resource group management, see the overview in [Chapter 7: HACMP Configuration Process and Facilities](#). Also see the chapter on changing resources and resource groups in the *Administration and Troubleshooting Guide* for complete information and instructions on performing resource group management through SMIT.

Cluster Single Point of Control (C-SPOC)

With the C-SPOC utility, you can make changes to the whole cluster from a single cluster node. Instead of performing administrative tasks on each cluster node, you can use the SMIT interface to issue a C-SPOC command once, on a single node, and the change is propagated across all cluster nodes.

For more information about C-SPOC, see the section [HACMP System Management \(C-SPOC\)](#) in [Chapter 7: HACMP Configuration Process and Facilities](#).

Dynamic Adapter Swap

The dynamic adapter swap functionality lets you swap the IP address of an active network interface card (NIC) with the IP address of a user-specified active, available “backup” network interface card on the *same* node and network. Cluster services do not have to be stopped to perform the swap.

This feature can be used to move an IP address off of a network interface card that is behaving erratically, to another NIC without shutting down the node. It can also be used if a hot pluggable NIC is being replaced on the node. Hot pluggable NICs can be physically removed and replaced without powering off the node. When the (hot pluggable) NIC to be replaced is pulled from the node, HACMP makes the NIC unavailable as a backup.

You can configure adapter swap using SMIT. The service IP address is moved from its current NIC to a user-specified NIC. The service IP address then becomes an available “backup” address. When the new card is placed in the node, the NIC is incorporated into the cluster as an available “backup” again. You can then swap the IP address from the backup NIC to the original NIC.

Note: The dynamic adapter swap feature is not supported on the SP switch network.

Note: This type of dynamic adapter swap can only be performed within a single node. You cannot swap the IP address with the IP address on a different node with this functionality. To move a service IP address to another node, move its resource group using the Resource Group Management utility.

Minimizing Unscheduled Downtime

Another important goal with HACMP is to minimize unscheduled downtime in response to unplanned cluster component failures. The HACMP software provides the following features to minimize unscheduled downtime:

- *Fast recovery* to speed up the failover in large clusters
- *A delayed fallback timer* to allow a custom resource group to fall back at a specified time
- *IPAT via IP Aliases* to speed up the processing during recovery of service IP labels
- *Automatic recovery of resource groups* that are in the ERROR state, whenever a cluster node comes up. For more information, see the following section.

Recovering Resource Groups on Node Startup

Prior to HACMP 5.2, when a node joined the cluster, it did not acquire any resource groups that had previously gone into an ERROR state on any other node. Such resource groups remained in the ERROR state and required use of the Resource Group Migration utility, **cIRGmove**, to manually bring them back online.

In HACMP 5.2, an attempt is made to bring the resource groups that are currently in the ERROR state into the online (active) state on the joining node. This further increases the chances of bringing the applications back online. When a node starts up, if a resource group is in the ERROR state on any node in the cluster, this node attempts to acquire the resource group. Note that the node must be included in the nodelist for the resource group.

The resource group recovery on node startup is different for non-concurrent and concurrent resource groups:

- If the starting node fails to activate a *non-concurrent resource group* that is in the ERROR state, the resource group continues to fall over to another node in the nodelist, if a node is available. The failover action continues until all available nodes in the nodelist have been tried.
- If the starting node fails to activate a *concurrent resource group* that is in the ERROR state on the node, the concurrent resource group is left in the ERROR state on that node. Note that the resource group might still remain online on other nodes.

Fast Recovery

The HACMP fast recovery feature speeds up failover in large clusters.

Fast recovery lets you select a filesystems consistency check and a filesystems recovery method:

- If you configure a filesystem to use a consistency check and a recovery method, it saves time by running **logredo** rather than **fsck** on each filesystem. If the subsequent **mount** fails, then it runs a full **fsck**.

If a filesystem suffers damage in a failure, but can still be mounted, **logredo** may not succeed in fixing the damage, producing an error during data access.

- In addition, it saves time by acquiring, releasing, and falling over all resource groups and filesystems in parallel, rather than serially.

Do not set the system to run these commands in parallel if you have shared, nested filesystems. These must be recovered sequentially. (Note that the cluster verification utility, **clverify**, does not report filesystem and fast recovery inconsistencies.)

The **varyonvg** and **varyoffvg** commands always run on volume groups in parallel, regardless of the setting of the recovery method.

Delayed Fallback Timer for Resource Groups

The Delayed Fallback Timer lets a resource group fall back to the higher priority node at a time that you specify. The resource group that has a delayed fallback timer configured and that currently resides on a non-home node falls back to the higher priority node at the recurring time (daily, weekly, monthly or yearly), or on a specified date.

The following example describes a case when configuring a delayed fallback timer would be beneficial: Suppose one node in the cluster has failed, and then has been repaired. You may want to integrate the node into a cluster during off-peak hours. Rather than writing a script or a cron job to do the work, which are both time-consuming and prone to error, you can set the delayed fallback timer for a particular custom resource group to the appropriate time. HACMP will automatically start the resource group failover at the time you specify.

For more information on the delayed fallback timer, see the *Planning and Installation Guide*.

Minimizing Takeover Time: Fast Disk Takeover

In the case of a cluster failure, enhanced concurrent volume groups are taken over faster than in previous releases of HACMP due to the improved disk takeover mechanism.

HACMP automatically detects enhanced concurrent volume groups and ensures that the faster option for volume group takeover is launched in the event of a node failure, if:

- You have installed AIX 5.2 and HACMP
- You include in your non-concurrent resource groups the enhanced concurrent mode volume groups (or convert the existing volume groups to enhanced concurrent volume groups).

This functionality is especially useful for failover of volume groups made up of a large number of disks.

During fast disk takeover, HACMP skips the extra processing needed to break the disk reserves, or update and synchronize the LVM information by running lazy update. As a result, the disk takeover mechanism of HACMP used for enhanced concurrent volume groups is faster than disk takeover used for standard volume groups included in non-concurrent resource groups.

Maximizing Disaster Recovery

HACMP can be an integral part of a comprehensive disaster recovery plan for your enterprise. Three possible ways to distribute backup copies of data to different sites, for possible disaster recovery operations, include:

- HACMP/XD - PPRC
- HACMP/XD - IP Mirroring
- Cross-Site LVM Mirroring.

For more information on the disaster recovery solutions included in HACMP/XD, see the Release Notes and documentation for the two solutions included in that package.

Cross-Site LVM Mirroring

In HACMP 5.2, you can set up disks located at two different sites for remote LVM mirroring, using a Storage Area Network (SAN), for example. Cross-site LVM mirroring replicates data between the disk subsystem at each site for disaster recovery.

A SAN is a high-speed network that allows the establishment of direct connections between storage devices and processors (servers) within the distance supported by Fibre Channel. Thus, two or more servers (nodes) located at different sites can access the same physical disks, which can be separated by some distance as well, through the common SAN. The disks can be combined into a volume group via the AIX Logical Volume Manager, and this volume group can be imported to the nodes located at different sites. The logical volumes in this volume group can have up to three mirrors. Thus, you can set up at least one mirror at each site. The information stored on this logical volume is kept highly available, and in case of certain failures, the remote mirror at another site will still have the latest information, so the operations can be continued on the other site.

HACMP 5.2 automatically synchronizes mirrors after a disk or node failure and subsequent reintegration. HACMP handles the automatic mirror synchronization even if one of the disks is in the PVREMOVED or PVMISSING state. Automatic synchronization is not possible for all cases, but you can use C-SPOC to manually synchronize the data from the surviving mirrors to stale mirrors after a disk or site failure and subsequent reintegration.

Cluster Events

This section describes how the HACMP software responds to changes in a cluster to maintain high availability.

The HACMP cluster software monitors all the different components that make up the highly available application including disks, network interfaces, nodes and the applications themselves. The Cluster Manager uses different methods for monitoring different resources:

- RSCT subsystem is responsible for monitoring networks and nodes
- The AIX LVM subsystem produces error notifications for volume group quorum loss
- The Cluster Manager itself dispatches application monitors.

An HACMP cluster environment is event-driven. An event is a change of status within a cluster that the Cluster Manager recognizes and processes. A cluster event can be triggered by a change affecting a network interface card, network, or node, or by the cluster reconfiguration process exceeding its time limit. When the Cluster Manager detects a change in cluster status, it executes a script designated to handle the event and its subevents.

Note: The logic of cluster events is described here for completeness. No manual intervention is required to ensure that HACMP carries out cluster events correctly.

The following examples show some events the Cluster Manager recognizes:

- **node_up** and **node_up_complete** events (a node joining the cluster)
- **node_down** and **node_down_complete** events (a node leaving the cluster)
- a local or global **network_down** event (a network has failed)
- **network_up** event (a network has connected)
- **swap_adapter** event (a network adapter failed and a new one has taken its place)
- dynamic reconfiguration events.

When a cluster event occurs, the Cluster Manager runs the corresponding event script for that event. As the event script is being processed, a series of subevent scripts may be executed. The HACMP software provides a script for each event and subevent. The default scripts are located in the `/usr/es/sbin/cluster/events` directory.

By default, the Cluster Manager calls the corresponding event script supplied with the HACMP software for a specific event. You can specify additional processing to customize event handling for your site if needed. For more information, see the section [Customizing Event Processing](#).

Processing Cluster Events

The two primary cluster events that HACMP software handles are fallover and reintegration:

- *Fallover* refers to the actions taken by the HACMP software when a cluster component fails or a node leaves the cluster.
- *Reintegration* refers to the actions that occur within the cluster when a component that had previously left the cluster returns to the cluster.

Event scripts control both types of actions. During event script processing, cluster-aware application programs see the state of the cluster as unstable.

Fallover

A fallover occurs when a resource group moves from its home node to another node because its home node leaves the cluster.

Nodes leave the cluster either by a planned transition (a node shutdown or stopping cluster services on a node), or by failure. In the former case, the Cluster Manager controls the release of resources held by the exiting node and the acquisition of these resources by nodes still active in the cluster. When necessary, you can override the release and acquisition of resources (for example, to perform system maintenance). You can also postpone the acquisition of the resources by integrating nodes (by setting the delayed fallback timer for custom resource groups).

Node failure begins when a node monitoring a neighboring node ceases to receive keepalive traffic for a defined period of time. If the other cluster nodes agree that the failure is a node failure, the failing node is removed from the cluster and its resources are taken over by the active nodes configured to do so.

If other components fail, such as a network interface card, the Cluster Manager runs an event script to switch network traffic to a backup network interface card (if present).

Reintegration

A reintegration, or a fallback occurs when a resource group moves to a node which has just joined the cluster.

When a node joins a running cluster, the cluster becomes temporarily unstable. The member nodes coordinate the beginning of the join process and then run event scripts to release any resources the joining node is configured to take over. The joining node then runs an event script to take over these resources. Finally, the joining node becomes a member of the cluster. At this point, the cluster is stable again.

Emulating Cluster Events

HACMP provides an emulation utility to test the effects of running a particular event without modifying the cluster state. The emulation runs on every active cluster node, and the output is stored in an output file on the node from which the emulation was launched.

For more information on the Event Emulator utility, see [Chapter 7: HACMP Configuration Process and Facilities](#).

Customizing Event Processing

The HACMP software has an event customization facility you can use to tailor event processing. The Cluster Manager's ability to recognize a specific series of events and subevents permits a very flexible customization scheme. Customizing event processing allows you to provide the most efficient path to critical resources should a failure occur.

You can define multiple pre- and post-events for a list of events that appears in the picklist in the **Change/Show Pre-Defined HACMP Events** SMIT panel.

Customization for an event could include notification to the system administrator before and after the event is processed, as well as user-defined commands or scripts before and after the event processing, as shown in the list:

- Notify system administrator of event to be processed
- Pre-event script or command
- HACMP for AIX event script
- Post-event script or command
- Notify system administrator event processing is complete.

Use this facility for the following types of customization:

- Pre- and post-event processing
- Event notification
- Event recovery and retry.

Note: In HACMP, the event customization information stored in the HACMP configuration database is synchronized across all cluster nodes when the cluster resources are synchronized. Thus, pre-, post-, notification, and recovery event script names must be the same on all nodes, although the actual processing done by these scripts can be different.

The **clverify** utility includes a function to automatically monitor cluster configuration with the means of a new event called **cluster_notify**. You can use this event to configure an HACMP Pager Notification Method to send out a page in case of detected errors in cluster configuration. The output of this event is also logged in **hacmp.out** throughout the cluster on each node that is running cluster services.

Defining New Events

In HACMP it is possible to define new events as well as to tailor the existing ones.

Pre- and Post-Event Processing

To tailor event processing to your environment, specify commands or user-defined scripts that execute before and after a specific event is generated by the Cluster Manager. For pre-processing, for example, you may want to send a message to specific users, informing them to stand by while a certain event occurs. For post-processing, you may want to disable login for a specific group of users if a particular network fails.

Event Notification

You can specify a command or user-defined script that provides notification (for example, mail) that an event is about to happen and that an event has just occurred, along with the success or failure of the event. You can also define a notification method through the SMIT interface to issue a customized page in response to a cluster event.

Event Recovery and Retry

You can specify a command that attempts to recover from an event command failure. If the retry count is greater than zero and the recovery command succeeds, the event script command is run again. You can also specify the number of times to attempt to execute the recovery command.

Customizing Event Duration

HACMP software issues a system warning each time a cluster event takes more time to complete than a specified timeout period.

Using the SMIT interface, you can customize the time period allowed for a cluster event to complete before HACMP issues a system warning for it.

Chapter 6: HACMP Cluster Configurations

This chapter provides examples of the types of cluster configurations supported by the HACMP software.

Sample Cluster Configurations

The following types of cluster configurations exist:

- **Standby Configurations**—These are the traditional redundant hardware configurations where one or more standby nodes stand idle, waiting for a server node to leave the cluster.
- **Takeover Configurations**—In these configuration, *all cluster nodes do useful work*, processing part of the cluster's workload. There are no standby nodes. Takeover configurations use hardware resources more efficiently than standby configurations since there is no idle processor. Performance can degrade after node detachment, however, since the load on remaining nodes increases.

Takeover configurations that use *concurrent access* use hardware efficiently and also minimize service interruption during fallover because there is no need for the takeover node to acquire the resources released by the failed node—the takeover node already shares ownership of the resources.

- **Cluster Configurations with Multi-Tiered Applications**—In these configurations, one application depends on another application. These cluster configurations use dependent resource groups.
- **Cross-Site LVM Mirror Configurations for Disaster Recovery**—In these geographically dispersed configurations, LVM mirroring replicates data between the disk subsystems at each of two sites for disaster recovery.
- **Cluster Configurations with Dynamic LPARs**—In these configurations, HACMP clusters use LPARs as cluster nodes. This lets you perform routine system upgrades through the dynamic allocation of system resources and redistribute CPU and memory resources to manage the application workload.

The sample cluster configurations shown in this chapter are by no means an exhaustive catalog of the possible configurations you can define using the HACMP software. Rather, use them as a starting point for thinking about the cluster configuration best suited to your environment.

Standby Configurations

The standby configuration is a traditional redundant hardware configuration, where one or more standby nodes stand idle, waiting for a server node to leave the cluster.

The sample standby configurations discussed in this chapter show how the configuration is defined for these types of resource groups:

- **Standby Configurations: Example 1** shows resource groups with the Online on Home Node Only startup policy, Fallover to Next Priority Node in the List fallover policy and Fallback to Higher Priority Node in the List fallback policy.
- **Standby Configurations: Example 2** shows resource groups with the startup Online Using Distribution Policy (network or node), fallover policy Next Priority Node in the List, and fallback policy Never Fallback.

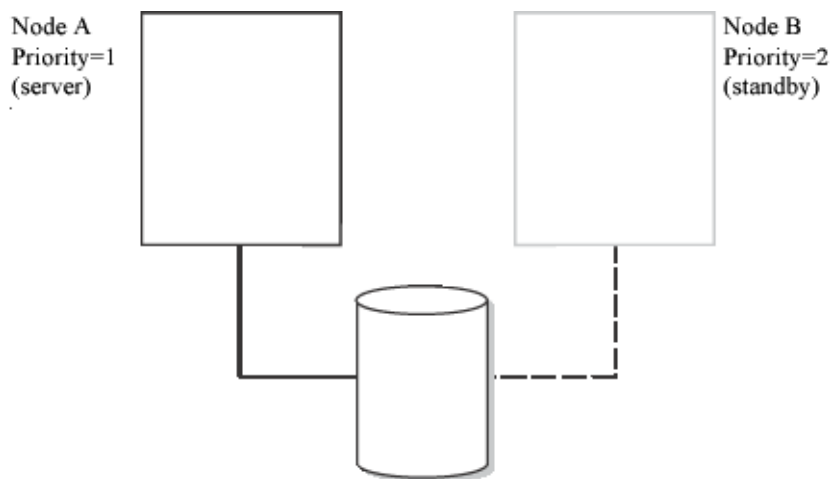
Concurrent resource groups require all nodes to have simultaneous access to the resource group and cannot be used in a standby configuration.

Standby Configurations: Example 1

The following figure shows a two-node standby configuration that uses resource groups with these policies:

- Startup policy: Online on Home Node Only
- Fallover policy: Fallover to Next Priority Node in the List
- Fallback policy: Fallback to Higher Priority Node in the List.

In the figure, a lower number indicates a higher priority:



One-for-One Standby Configuration Where IP Label Returns to the Home Node

In this setup, the cluster resources are defined as part of a single resource group. A nodelist is then defined as consisting of two nodes. The first node, Node A, is assigned a takeover (ownership) priority of 1. The second node, Node B, is assigned a takeover priority of 2.

At cluster startup, Node A (which has a priority of 1) assumes ownership of the resource group. Node A is the “server” node. Node B (which has a priority of 2) stands idle, ready should Node A fail or leave the cluster. Node B is, in effect, the “standby”.

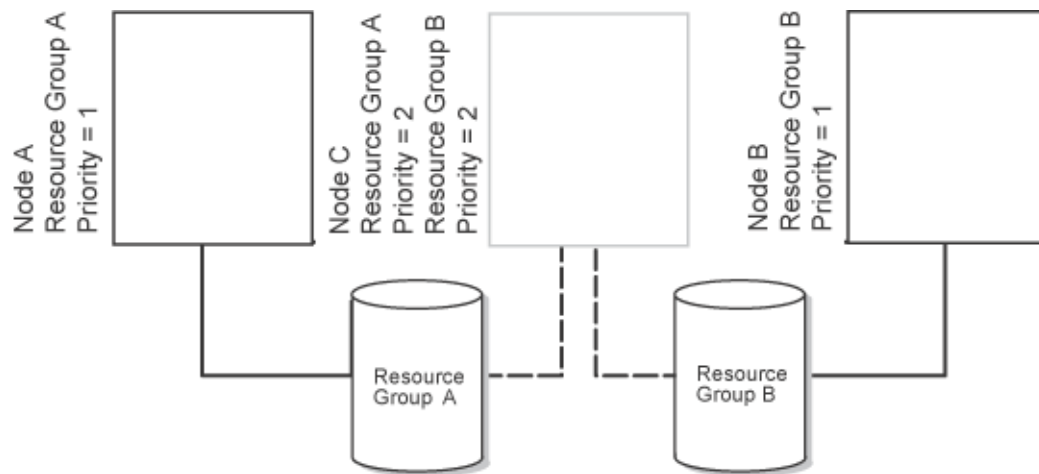
If the server node leaves the cluster, the standby node assumes control of the resource groups owned by the server, starts the highly available applications, and services clients. The standby node remains active until the node with the higher takeover priority rejoins the cluster. At that point, the standby node releases the resource groups it has taken over, and the server node reclaims them. The standby node then returns to an idle state.

Extending Standby Configurations From Example 1

The standby configuration from the previously described example can be easily extended to larger clusters. The advantage of this configuration is that it makes better use of the hardware. The disadvantage is that the cluster can suffer severe performance degradation if more than one server node leaves the cluster.

The following figure illustrates a three-node standby configuration using the resource groups with these policies:

- Startup policy: Online on Home Node Only
- Fallover policy: Fallover to Next Priority Node in the List
- Fallback policy: Fallback to Higher Priority Node in the List.



One-for-Two Standby with Three Resource Groups

In this configuration, two separate resource groups (A and B) and a separate nodelist for each resource group exist. The nodelist for Resource Group A consists of Node A and Node C. Node A has a takeover priority of 1, while Node C has a takeover priority of 2. The nodelist for Resource Group B consists of Node B and Node C. Node B has a takeover priority of 1; Node C again has a takeover priority of 2. (Remember, a resource group can be owned by only a single node in a non-concurrent configuration.)

Since each resource group has a different node at the head of its nodelist, the cluster's workload is divided, or partitioned, between these two resource groups. Both resource groups, however, have the same node as the standby in their nodelists. If either server node leaves the cluster, the standby node assumes control of that server node's resource group and functions as the departed node.

In this example, the standby node has three network interfaces (not shown) and separate physical connections to each server node's external disk. Therefore, the standby node can, if necessary, take over for both server nodes concurrently. The cluster's performance, however, would most likely degrade while the standby node was functioning as both server nodes.

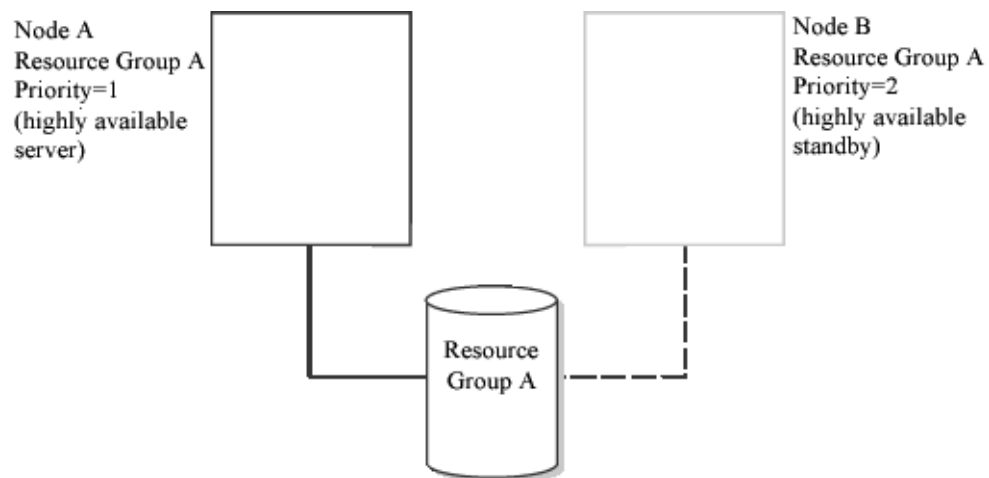
Standby Configurations: Example 2

In the following standby configuration, the resource groups have these policies:

- Startup policy: Online Using Distribution Policy (network-based or node-based)
- Fallover policy: Next Priority Node in the List
- Fallback policy: Never Fallback.

This configuration differs from a standby configuration in which the ownership priority of resource groups is not fixed. Rather, the resource group is associated with an IP address that can rotate among nodes. This makes the roles of server and standby fluid, changing over time.

The following figure shows illustrates the one-for-one standby configuration that is described in this section:



One-for-One Standby with Resource Groups Where IP Label Rotates

At system startup, the resource group attaches to the node that claims the shared IP address. This node “owns” the resource group for as long as it remains in the cluster. If this node leaves the cluster, the peer node assumes the shared IP address and claims ownership of that resource group. Now, the peer node “owns” the resource group for as long as it remains in the cluster.

When the node that initially claimed the resource group rejoins the cluster, it does not take the resource group back. Rather, it remains idle for as long as the node currently bound to the shared IP address is active in the cluster. Only if the peer node leaves the cluster does the node that initially “owned” the resource group claim it once again. Thus, ownership of resources rotates between nodes.

Extending Standby Configurations From Example 2

As with the first example of the standby configuration, configurations from example 2 can be easily extended to larger clusters. For example, in a one-for-two standby configuration from example 2, the cluster could have two separate resource groups, each of which includes a distinct shared IP address.

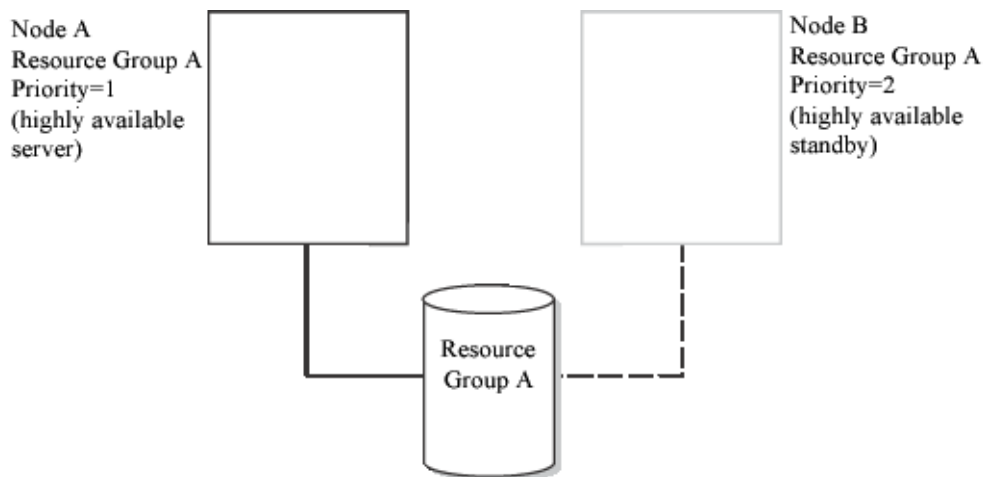
At cluster startup, the first two nodes each claim a shared IP address and assume ownership of the resource group associated with that shared IP address. The third node remains idle. If an active node leaves the cluster, the idle node claims that shared IP address and takes control of that resource group.

Takeover Configurations

All nodes in a takeover configuration process part of the cluster’s workload. There are no standby nodes. Takeover configurations use hardware resources more efficiently than standby configurations since there is no idle processor. Performance degrades after node detachment, however, since the load on remaining nodes increases.

One-Sided Takeover

The following figure illustrates a two-node, one-sided takeover configuration. In the figure, a lower number indicates a higher priority.



One-sided Takeover with Resource Groups in Which IP Label Returns to the Home Node

This configuration has two nodes actively processing work, but only one node providing highly available services to cluster clients. That is, although there are two sets of resources within the cluster (for example, two server applications that handle client requests), only one set of resources needs to be highly available. This set of resources is defined as an HACMP resource group and has a nodelist that includes both nodes. The second set of resources is not defined as a resource group and, therefore, is not highly available.

At cluster startup, Node A (which has a priority of 1) assumes ownership of Resource Group A. Node A, in effect, “owns” Resource Group A. Node B (which has a priority of 2 for Resource Group A) processes its own workload independently of this resource group.

If Node A leaves the cluster, Node B takes control of the shared resources. When Node A rejoins the cluster, Node B releases the shared resources.

If Node B leaves the cluster, however, Node A does not take over any of its resources, since Node B's resources are not defined as part of a highly available resource group in whose chain this node participates.

This configuration is appropriate when a single node is able to run all the critical applications that need to be highly available to cluster clients.

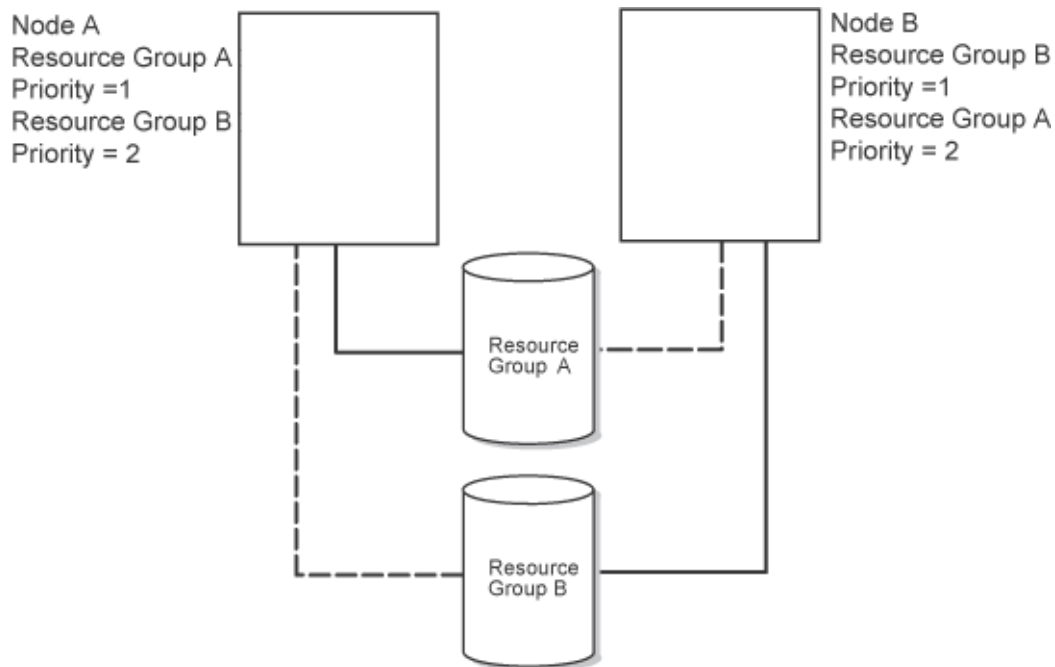
Mutual Takeover

The mutual takeover for non-concurrent access configuration has multiple nodes, each of which provides distinct highly available services to cluster clients. For example, each node might run its own instance of a database and access its own disk.

Furthermore, each node has takeover capacity. If a node leaves the cluster, a surviving node takes over the resource groups owned by the departed node.

The mutual takeover for non-concurrent access configuration is appropriate when each node in the cluster is running critical applications that need to be highly available and when each processor is able to handle the load of more than one node.

The following figure illustrates a two-node mutual takeover configuration for non-concurrent access. In the figure, a lower number indicates a higher priority.



Mutual Takeover for Non-Concurrent Access

The key feature of this configuration is that the cluster's workload is divided, or partitioned, between the nodes. Two resource groups exist, in addition to a separate resource chain for each resource group. The nodes that participate in the resource chains are the same. It is the differing priorities within the chains that designate this configuration as mutual takeover.

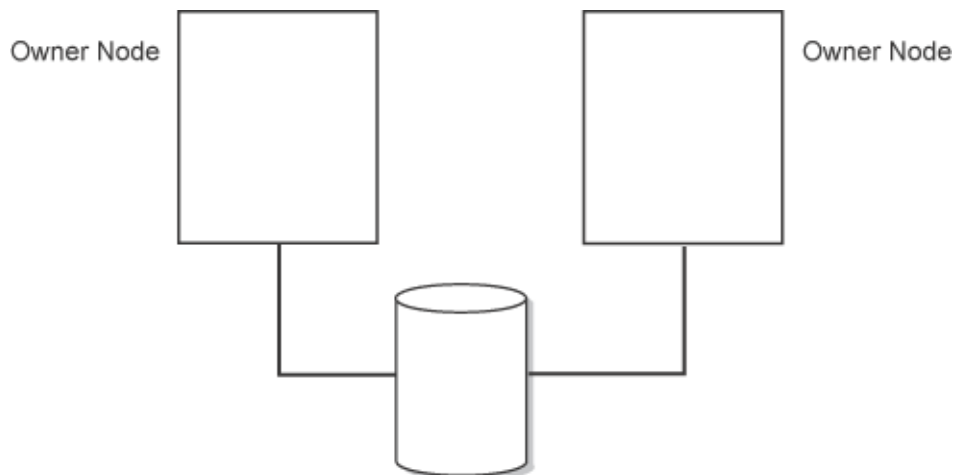
The chains for both resource groups consist of Node A and Node B. For Resource Group A, Node A has a takeover priority of 1 and Node B has a takeover priority of 2. For Resource Group B, the takeover priorities are reversed. Here, Node B has a takeover priority of 1 and Node A has a takeover priority of 2.

At cluster startup, Node A assumes ownership of the Resource Group A, while Node B assumes ownership of Resource Group B.

If either node leaves the cluster, its peer node takes control of the departed node’s resource group. When the “owner” node for that resource group rejoins the cluster, the takeover node relinquishes the associated resources; they are reacquired by the higher-priority, reintegrating node.

Two-Node Mutual Takeover Configuration for Concurrent Access

The following figure illustrates a two-node mutual takeover configuration for concurrent access:



Two-Node Mutual Takeover for Concurrent Access

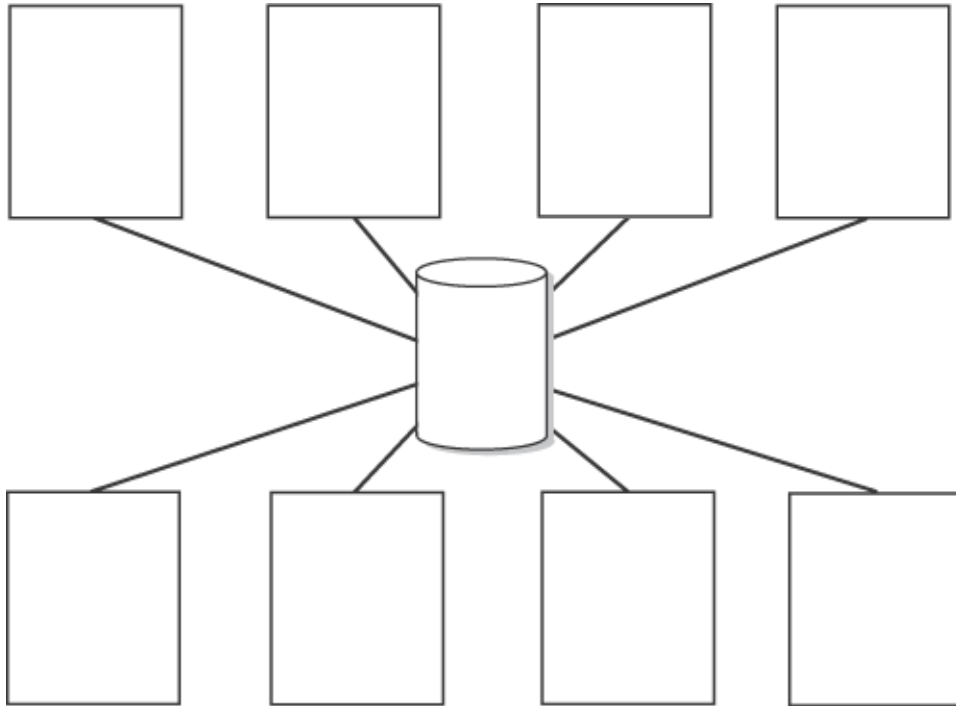
In this configuration, both nodes have simultaneous access to the shared disks and own the same disk resources. There is no “takeover” of shared disks if a node leaves the cluster, since the peer node already has the shared volume group varied on.

In this example both nodes are running an instance of a server application that accesses the database on the shared disk. The application’s proprietary locking model is used to arbitrate application requests for disk resources.

Running multiple instances of the same server application allows the cluster to distribute the processing load. As the load increases, additional nodes can be added to further distribute the load.

Eight-Node Mutual Takeover Configuration for Concurrent Access

The following figure illustrates an eight-node mutual takeover configuration for concurrent access:



Eight-Node Mutual Takeover for Concurrent Access

In this configuration, as in the previous configuration, all nodes have simultaneous—but not concurrent—access to the shared disks and own the same disk resources. Here, however, each node is running a different server application. Clients query a specific application at a specific IP address. Therefore, each application server and its associated IP address must be defined as part of a non-concurrent resource group, and all nodes that are potential owners of that resource group must be included in a corresponding nodelist.

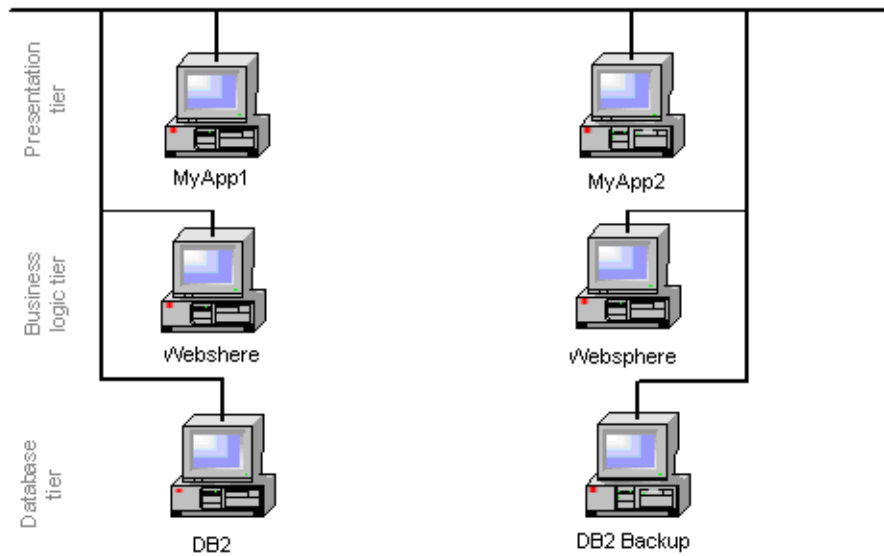
Concurrent access resource groups are supported in clusters with up to 32 nodes in HACMP.

Cluster Configurations with Multi-Tiered Applications

A typical cluster configuration that could utilize dependent resource groups is the environment in which an application such as WebSphere depends on another application such as DB2.

Note: It is important to distinguish the application server, such as WebSphere, from the HACMP application server that you configure in HACMP by specifying the application server start and stop scripts.

In order to satisfy business requirements, a cluster-wide dependency must be defined between two or more resource groups. The following figure illustrates the business scenario that utilizes dependencies between applications:



A Typical Multi-Tier Cluster Environment with Dependencies Between Applications

Multi-Tiered Applications

Business configurations that use layered, or multi-tiered applications can also utilize dependent resource groups. For example, the back end database must be online before the application server. In this case, if the database goes down and is moved to a different node, the resource group containing the application server would have to be brought down and back up on any node in the cluster.

Environments such as SAP require applications to be recycled anytime a database fails. There are many application services provided by an environment like SAP, and the individual application components often need to be controlled in a specific order.

Another area where establishing interdependencies between resource groups proves useful is when system services are required to support application environments. Services such as **cron** jobs for pruning log files or initiating backups need to move from node to node along with an application, but are typically not initiated until the application is established. These services can be built into application server start and stop scripts. When greater granularity is needed, they can be controlled through pre- and post- event processing. Dependent resource groups allow an easier way to configure system services to be dependent upon applications they serve.

For an overview of dependent resource groups, see the section [Resource Group Dependencies](#) in [Chapter 3: HACMP Resources and Resource Groups](#).

Cross-Site LVM Mirror Configurations for Disaster Recovery

In HACMP 5.2, you can set up disks located at two different sites for remote LVM mirroring, using a Storage Area Network (SAN). A SAN is a high-speed network that allows the establishment of direct connections between storage devices and processors (servers) within the distance supported by Fibre Channel. Thus, two or more distantly separated servers (nodes)

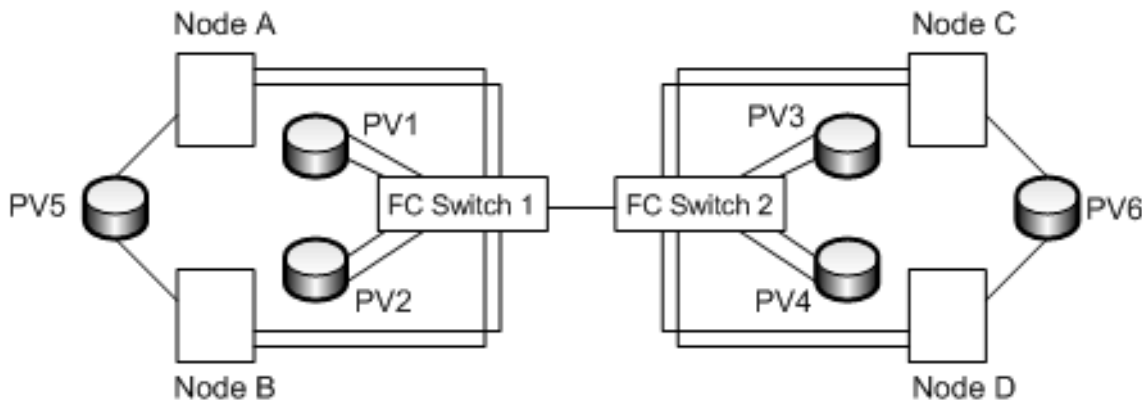
located at different sites can access the same physical disks, which may be distantly separated as well, via the common SAN. These remote disks can be combined into volume groups, using C-SPOC.

The logical volumes in a volume group can have up to three mirrors or copies, for example, one mirror at each site. Thus the information stored on this logical volume may be kept highly available, and in case of a certain failures—for example, all nodes at one site, including the disk subsystem at that site—the remote mirror at another site will still have the latest information and the operations can be continued on that site.

The primary intent of this feature is to support two-site clusters where LVM mirroring through a Storage Area Network (SAN) replicates data between the disk subsystem at each site for disaster recovery.

Another advantage of cross-site LVM mirroring is that after a site/disk failure and subsequent site reintegration, HACMP attempts to synchronize the data from the surviving disks to the joining disks automatically. The synchronization occurs in the background and does not significantly impact the reintegration time.

The following figure illustrates a cross-site LVM mirroring configuration using a SAN:



Cross-Site LVM Mirror Configuration for Disaster Recovery

The disks that are connected to at least one node at each of the two sites can be mirrored. In this example, PV4 is seen by nodes A and B on Site 1 via the Fibre Channel Switch 1-Fibre Channel Switch 2 connection, and is also seen on node C via Fibre Channel Switch 2. You could have a mirror of PV4 on Site 1. The disks that are connected to the nodes on one site only (PV5 and PV6) cannot be mirrored across sites.

The disk information is replicated from a local site to a remote site. The speed of this data transfer depends on the physical characteristics of the channel, the distance, and LVM mirroring performance.

Cluster Configurations with Dynamic LPARs

The advanced partitioning features of AIX 5L 5.2 provide the ability to dynamically allocate system CPU, memory, and I/O slot resources (*dynamic LPAR*).

Using HACMP in combination with LPARs lets you:

- **Perform routine system upgrades through the dynamic allocation of system resources.** When used with dynamic LPARs, HACMP can reduce the amount of downtime for well-planned systems upgrades by automating the transition of your application workload from one logical partition to another, so that the first logical partition may be upgraded without risk to the application.
- **Effectively redistribute CPU and memory resources to manage the workload.** Combining HACMP with dynamic LPAR lets you use customized application start and stop scripts to dynamically redistribute CPU and memory resources to logical partitions that are currently executing application workload, to further support application transition within a single frame. This way you maintain the processing power and resources necessary to support your applications, while minimal resources are devoted to upgrading, a less resource intensive task.

Note: Do not have all your cluster nodes configured as LPARs within the *same* physical server. This configuration could potentially be a significant single point of failure.

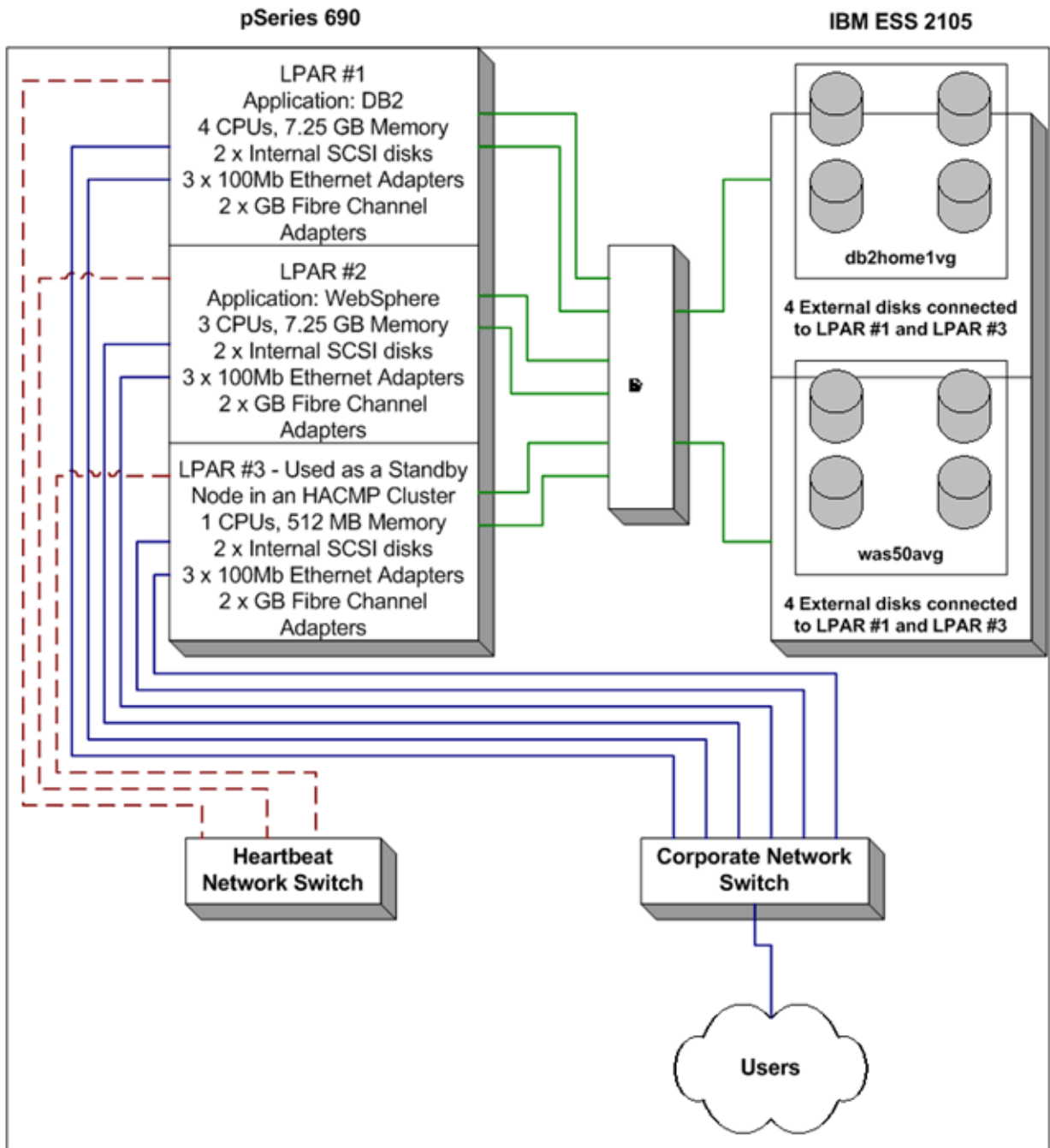
The following example illustrates a cluster configuration that uses three LPARs:

- LPAR #1 is running a back end database (DB2 UDB)
- LPAR #2 is running WebSphere Application Server (WAS)
- LPAR #3 is running as a backup (standby) for both the DB2 and WAS LPARs. This LPAR contains only minimal CPU and memory resources.

When it is time to move either the DB2 or WAS application to the third LPAR (due to a planned upgrade or a resource failure in these LPARs, for instance), you can use customized application start and stop scripts in HACMP to automate the dynamic reallocation of CPU and memory from the primary LPAR to the standby LPAR. This operation allows the third LPAR to acquire the CPU and memory resources necessary to meet business performance requirements. When HACMP moves the resource group containing the application back to its home LPAR, the CPU and memory resources automatically move with it.

Note: In general, dynamic LPARs allow dynamic allocation of CPU, memory and I/O slot resources. HACMP and dynamic LPAR I/O slot resources are not compatible (although you can dynamically allocate I/O slot resources outside of HACMP cluster).

The following figure illustrates this cluster environment:



A Cluster with Three LPARs

Chapter 7: HACMP Configuration Process and Facilities

This chapter provides an overview of the HACMP cluster configuration process, and covers the following topics:

- [Information You Provide to HACMP](#)
- [Information Discovered by HACMP](#)
- [Cluster Configuration Options: Standard and Extended.](#)

This chapter also provides an overview of the following administrative tools supplied with the HACMP software:

- [Cluster Security](#)
- [Installation, Configuration and Management Tools](#)
- [Monitoring Tools](#)
- [Troubleshooting Tools](#)
- [Emulation Tools.](#)

Information You Provide to HACMP

Prior to configuring a cluster, make sure the building blocks are planned and configured, and the initial communication path exists for HACMP to “reach” each node. This section covers the basic tasks you need to perform to configure a cluster.

Information on Physical Configuration of a Cluster

Physical configuration of a cluster consists of performing the following planning and configuration tasks:

- Ensure the TCP/IP network support for the cluster.
- Ensure the point-to-point network support for the cluster.
- Ensure the heartbeating support for the cluster.
- Configure the shared disk devices for the cluster.
- Configure the shared volume groups for the cluster.
- Consider the mission-critical applications for which you are using HACMP. Also, consider application server and what type of resource group management is best for each application.
- Examine issues relating to HACMP clients.
- Ensure physical redundancy by using multiple circuits or uninterruptable power supplies, redundant physical network interface cards, multiple networks to connect nodes and disk mirroring.

These tasks are described in detail in the *Planning and Installation Guide*.

AIX Configuration Information

Cluster components must be properly configured on the AIX level. For this task, ensure that:

- Basic communication to cluster nodes exists
- Volume groups, logical volumes, mirroring and filesystems are configured and set up. To ensure logical redundancy, consider different types of resource groups, and plan how you will group your resources in resource groups.

For the specifics of configuring volume groups, logical volumes and filesystems, refer to the AIX manuals and to the *Planning and Installation Guide*.

Establishing the Initial Communication Path

The initial communication path is a path to a node which you are adding to a cluster. To establish the initial communication path, you provide the name of the node, or other information that can serve as the name of the node.

In general, a node name and a hostname can be the same. When configuring a new node, you can enter any of the following denominations that will serve as an initial communication path to a node:

- An IP address of a physical network interface card (NIC) on that node, such as 1 . 2 . 3 . 4 . In this case that address is used as a communication path for contacting a node.
- An IP label associated with an IP address of a NIC on that node, such as `servername`. In this case that name is used to determine the communication path for contacting a node, based on the assumption that the local TCP/IP configuration (Domain Nameserver or Hosts Table) supplies domain qualifiers and resolves the IP label to an IP address.
- A *Fully Qualified Domain Name (FQDN)*, such as `"servername.thecompanyname.com"`. In this case the communication path is `"servername.thecompanyname.com"`, based on the assumption that the local TCP/IP configuration (Domain Nameserver or Hosts Table) supplies domain qualifiers and resolves the IP label to an IP address.

When you enter any of these names, HACMP ensures unique name resolution and uses the hostname as a node name, unless you explicitly specify otherwise.

Note: In HACMP, node names and hostnames have to be different in some cases where the application you are using requires that the AIX “hostname attribute” moves with the application in the case of a cluster component failure. This procedure is done through setting up special event scripts.

If the nodes and physical network interface cards have been properly configured to AIX, HACMP can use this information to assist you in the configuration process, by running the automatic discovery process discussed in the following section.

Information Discovered by HACMP

You can define the basic cluster components in just a few steps. To assist you in the cluster configuration, HACMP can automatically retrieve the information necessary for configuration from each node.

Note: For easier and faster cluster configuration, you can also use a cluster configuration assistant. For more information, see [Two-Node Cluster Configuration Assistant](#).

For the automatic discovery process to work, the following conditions should be met in HACMP:

- You have previously configured the physical components, and performed all the necessary AIX configurations.
- Working communications paths exist to each node. This information will be used to automatically configure the cluster TCP/IP topology when the *standard configuration path* is used.

Once these tasks are done, HACMP automatically discovers predefined physical components within the cluster, and selects default behaviors. In addition, HACMP performs discovery of cluster information if there are any changes made during the configuration process.

Running discovery retrieves current AIX configuration information from all cluster nodes. This information appears in picklists to help you make accurate selections of existing components.

The HACMP automatic discovery process is easy, fast, and does not place a "waiting" burden on you as the cluster administrator.

Cluster Configuration Options: Standard and Extended

The configuration process is significantly simplified. While the details of the configuration process are covered in the *Administration and Troubleshooting Guide*, this section provides a brief overview of two ways to configure an HACMP cluster.

Configuring an HACMP Cluster Using the Standard Configuration Path

You can add the basic components of a cluster to the HACMP configuration database *in a few steps*. The standard cluster configuration path simplifies and speeds up the configuration process, because HACMP automatically launches discovery to collect the information and to select default behaviors.

If you use this path:

- Automatic discovery of cluster information runs by default. Before starting the HACMP configuration process, you need to configure network interfaces/devices in AIX. In HACMP, you establish initial communication paths to other nodes. Once this is done, HACMP collects this information and automatically configures the cluster nodes and networks based on physical connectivity. All discovered networks are added to the cluster configuration.

- IP aliasing is used as the *default* mechanism for binding IP labels/addresses to network interfaces.
- You can configure the most common types of resources. However, customizing of resource group fallover and fallback behavior is limited.

Configuring an HACMP Cluster Using the Extended Configuration Path

In order to configure the less common cluster elements, or if connectivity to each of the cluster nodes is not established, you can manually enter the information in a way similar to previous releases of the HACMP software.

When using the HACMP extended configuration SMIT paths, if any components are on remote nodes, you must manually initiate the discovery of cluster information. That is, discovery is optional (rather than automatic, as it is when using the standard HACMP configuration SMIT path).

Using the options under the extended configuration menu, you can add the basic components of a cluster to the HACMP configuration database, as well as many additional types of resources. Use the extended configuration path to customize the cluster for all the components, policies, and options that are not included in the standard configuration menus.

Overview: HACMP Administrative Facilities

The HACMP software provides you with the following administrative facilities:

- [Cluster Security](#)
- [Installation, Configuration and Management Tools](#)
- [Monitoring Tools](#)
- [Troubleshooting Tools](#)
- [Emulation Tools](#).

Cluster Security

All communication between nodes is sent through the Cluster Communications daemon, **clcmd**, which runs on each node. The **clcmd** daemon manages the connection authentication between nodes and any message authentication or encryption configured.

In HACMP release 5.1 and higher, the Cluster Communications daemon uses the trusted **/usr/es/sbin/cluster/etc/rhosts** file, and removes reliance on an **.rhosts** file. In HACMP 5.2, the daemon provides support for message authentication and encryption.

Installation, Configuration and Management Tools

HACMP includes the tools described in the following sections for installing, configuring, and managing clusters.

Two-Node Cluster Configuration Assistant

HACMP provides the Two-Node Cluster Configuration Assistant to simplify the process for configuring a basic two-node cluster. The wizard-like application requires the minimum information to define an HACMP cluster and uses discovery to complete the cluster configuration. The application is designed for users with little knowledge of HACMP who want to quickly set up a basic HACMP configuration. The underlying AIX configuration must be in place before you run the Assistant.

Planning Worksheets

Along with your HACMP software and documentation set, you have two types of worksheets to aid in planning your cluster topology and resource configuration: online or paper.

Online Planning Worksheets

HACMP provides the Java-based Online Planning Worksheets application that lets you:

- Plan a cluster
- Create a cluster definition
- Examine the configuration for an HACMP cluster

You can review information about a cluster configuration in an easy-to-view format for use in testing and troubleshooting situations.

After you save an HACMP cluster definition to a cluster worksheets file, you can open that file in Online Planning Worksheets running on a node, a laptop, or other computer running the application. This lets you examine the cluster definition on a non-cluster node, or share the worksheets file with a colleague.

For more information on the requirements and instructions for using the online planning worksheets, see the *Planning and Installation Guide*.

Paper Worksheets

The HACMP documentation includes a set of planning worksheets to guide your entire cluster planning process, from cluster topology to resource groups and application servers. You can use these worksheets as guidelines when installing and configuring your cluster. You may find these paper worksheets useful in the beginning stages of planning, when your team might be around a conference table discussing various configuration decisions. The planning worksheets are found in the *Planning and Installation Guide*.

SMIT Interface

You can use the SMIT panels supplied with the HACMP software to perform the following tasks:

- Configure clusters, nodes, networks, resources, and events
- Capture and restore snapshots of cluster configurations
- Read log files
- Diagnose cluster problems
- Manage a cluster using the C-SPOC utility

- Perform resource group management tasks
- Configure Automatic Error Notification
- Perform dynamic adapter swap
- Configure cluster performance tuning
- Configure custom disk methods.

Web-Based SMIT Interface

HACMP 5.2 includes a Web-enabled user interface (WebSMIT) that provides consolidated access to the SMIT functions of configuration and management as well as to interactive cluster status and the HACMP documentation.

The WebSMIT interface is similar to the ASCII SMIT interface. You do not need to learn a new user interface or terminology and can easily switch between ASCII SMIT and the Web version. Because WebSMIT runs in a Web browser, it can be accessed from any platform.

To use the WebSMIT interface, you must configure and run a Web server process on the cluster node(s) to be administered. The `/usr/es/sbin/cluster/wsm/README` file contains information on basic Web server configuration, the default security mechanisms in place when HACMP 5.2 is installed, and the configuration files available for customization.

Cluster Status Display Linked to Management Functions

When using the WebSMIT interface to see the cluster status display, you have links to the related WebSMIT management functions; thus HACMP 5.2 provides a consolidated user interface for cluster status with management capabilities.

For example, the node status display has a link to (among other options) the SMIT panels for starting and stopping Cluster Services. Now you can manipulate entities in the status display interactively rather than having to go to an ASCII SMIT interface on the node.

HACMP System Management (C-SPOC)

To facilitate management of a cluster, HACMP provides a way to run commands from one node and then verify and synchronize the changes to all the other nodes. You can use the HACMP System Management, or the Cluster Single Point of Control (C-SPOC) to automatically add users, files, and hardware without stopping mission-critical jobs.

C-SPOC lets you perform the following tasks:

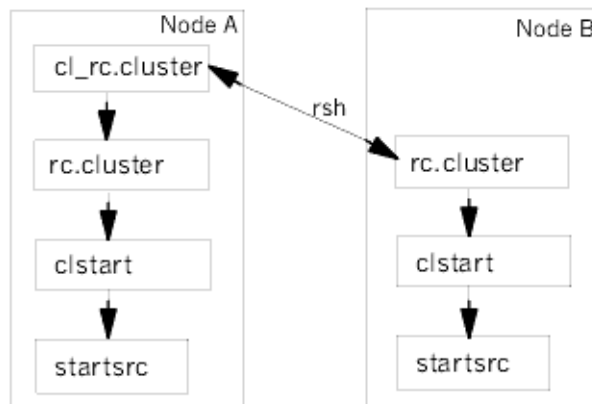
- Start/Stop HACMP Services
- HACMP Communication Interface Management
- HACMP Resource Group and Application Management
- HACMP File Collection Management
- HACMP Log Viewing and Management
- HACMP Security & Users Management
- HACMP Logical Volume Management
- HACMP Concurrent Logical Volume Management
- HACMP Physical Volume Management

- GPFS Filesystem Support
- Open a SMIT Session on a Node.

The C-SPOC utility simplifies maintenance of shared LVM components in clusters of up to 32 nodes. C-SPOC commands provide comparable functions in a cluster environment to the standard AIX commands that work on a single node. By automating repetitive tasks, C-SPOC eliminates a potential source of errors, and speeds up the process.

Without C-SPOC functionality, the system administrator must execute administrative tasks individually on each cluster node. For example, to add a user you usually must perform this task on each cluster node. Using the C-SPOC utility, a command executed on one node is also executed on other cluster nodes. Thus C-SPOC minimizes administrative overhead and reduces the possibility of inconsistent node states. Using C-SPOC, you issue a C-SPOC command once on a single node, and the user is added to all specified cluster nodes.

C-SPOC also makes managing logical volume components and controlling cluster services more efficient. You can use the C-SPOC utility to start or stop cluster services on nodes from a single node. The following figure illustrates a two-node configuration and the interaction of commands, scripts, and nodes when starting cluster services from a single cluster node. Note the prefix `cl_` begins all C-SPOC commands.



Flow of Commands Used at Cluster Startup by C-SPOC Utility

C-SPOC provides this functionality through its own set of cluster administration commands, accessible through SMIT menus and panels. To use C-SPOC, select the **Cluster System Management** option from the HACMP SMIT menu.

Cluster Snapshot Utility

The Cluster Snapshot utility allows you to save cluster configurations you would like to restore later. You also can save additional system and cluster information that can be useful for diagnosing system or cluster configuration problems. You can create your own custom snapshot methods to store additional information about your cluster.

A cluster snapshot lets you skip saving log files in the snapshot. Cluster snapshots are used for recording the cluster configuration information, whereas cluster logs only record the operation of the cluster and not the configuration information. By default, HACMP no longer collects

cluster log files when you create the cluster snapshot, although you can still specify collecting the logs in SMIT. Skipping the logs collection speeds up the running time of the snapshot utility, and reduces the size of the snapshot.

Customized Event Processing

You can define multiple pre- and post-events to tailor your event processing for your site's unique needs. For more information about writing your own scripts for pre- and post-events, see the *Administration and Troubleshooting Guide*.

Resource Group Management Utility

Prior to HACMP 5.1, you could use the DARE Resource Migration utility to move resource groups to other cluster nodes. In HACMP 5.1, the *Resource Group Management* utility replaced DARE Resource Migration.

The Resource Group Management utility provides a means for managing resource groups in the cluster, and enhances failure recovery capabilities of HACMP. It allows you to change the status or the location of any type of resource group (along with its resources—IP addresses, applications, and disks) without stopping cluster services.

Resource group management helps you manage your cluster more effectively, giving you better use of your cluster hardware resources. Resource group management also lets you perform selective maintenance without rebooting the cluster or disturbing operational nodes. For instance, you can use this utility to free the node of any resource groups to perform system maintenance on a particular cluster node.

Using the Resource Group Management utility does not affect other resource groups currently owned by a node. The node that currently owns the resource group to be moved releases it as it would during a “graceful shutdown with takeover,” and the node to which the resource group is being moved acquires the resource group just as it would during a node failover.

Use Resource Group Management to:

- Dynamically move a specified non-concurrent resource group from a node it currently resides on to the destination node that you have specified.
- Bring a resource group online or offline on one or all nodes in the cluster. This functionality is different for concurrent and non-concurrent resource groups. See the *Administration and Troubleshooting Guide* for detailed information on what kinds of online and offline operations you can perform on concurrent and non-concurrent resource groups.

HACMP File Collection Management

Like volume groups, certain files located on each cluster node need to be kept in sync in order for HACMP (and other applications) to behave correctly. Such files include event scripts, application scripts, and some AIX and HACMP configuration files.

HACMP File Collection management provides an easy way to request that a list of files be kept in sync across the cluster. Using HACMP file collection, you do not have to manually copy an updated file to every cluster node, verify that the file is properly copied, and confirm that each node has the same version of it.

Also, if one or more of these files is inadvertently deleted or damaged on one or more cluster nodes, it can take time and effort to determine the problem. Using HACMP file collection, this scenario is mitigated. HACMP detects when a file in a file collection is deleted or if the file size is changed to zero, and logs a message to inform the administrator.

Two predefined HACMP file collections are installed by default:

- **Configuration_Files**. A container for essential system files, such as **/etc/hosts** and **/etc/services**.
- **HACMP_Files**. A container for all the user-configurable files in the HACMP configuration. This is a special file collection that the underlying file collection propagation utility uses to reference all the user-configurable files in the HACMP configuration database (ODM) classes.

For a complete list of configuration files and user-configurable HACMP files, see the *Planning and Installation Guide*.

For information on configuring file collections in SMIT, see the *Administration and Troubleshooting Guide*.

Monitoring Tools

The following tools are used to monitor an HACMP cluster:

- **SNMP, or the Cluster SMUX Peer Daemon (clsmuxpd)** provides SNMP information and traps for SNMP clients. It gathers cluster information from the Cluster Manager relative to cluster state changes of nodes and interfaces. Cluster information can be retrieved using SNMP commands or by SNMP based client programs such as HATivoli. See also the section [Cluster SMUX Peer and SNMP Monitoring Programs](#) in [Chapter 4: HACMP Cluster Hardware and Software](#).
- **The Cluster Information Program (Clinfo)** gets the information from **clsmuxpd** and allows clients communicating with this program to be aware of changes in a cluster state. For information on Clinfo, see the section [Cluster Information Program](#) in [Chapter 4: HACMP Cluster Hardware and Software](#).
- The **clstat** utility reports the status of key cluster components: the cluster itself, the nodes in the cluster, the network interfaces connected to the nodes, and the resource groups on each node. It reports whether the cluster is up, down, or unstable. It also reports whether a node is up, down, joining, leaving, or reconfiguring, and the number of nodes in the cluster. The clstat utility provides ASCII, Motif, X Windows and HTML interfaces. You can run **clstat** from both ascii SMIT and WebSMIT.
- The **HAView** utility extends Tivoli NetView services so you can monitor HACMP clusters and cluster components across a network from a single node. Using HAView, you can also view the full cluster event history in the **/usr/es/sbin/cluster/history/cluster.mmddyyyy** file.
- **Cluster Monitoring with Tivoli** allows you to monitor clusters and cluster components and perform cluster administration tasks through your Tivoli Framework console.
- **Application Monitoring** allows you to configure multiple monitors for an application server to monitor specific applications and processes; and define action to take upon detection of an unexpected termination of a process or other application failures. See the section [Application Monitors](#) in [Chapter 5: Ensuring Application Availability](#).

- **Show Cluster Applications** SMIT option provides an application-centric view of the cluster configuration. This utility displays existing interfaces and information in an “application down” type of view. You can access it from both ASCII SMIT and WebSMIT.
- The **Application Availability Analysis** tool measures uptime statistics for applications with application servers defined to HACMP.
- **Assigning a persistent node IP label** is a useful administrative “tool” that lets you contact a node even if the HACMP cluster services are down on that node.
- **HACMP Verification and Synchronization** verifies that HACMP-specific modifications to AIX system files are correct, that the cluster and its resources are configured correctly, and that security, if set up, is configured correctly.

Each of these utilities is described in greater detail in the following sections. For more information, see the *Administration and Troubleshooting Guide*.

Cluster Status Utility (clstat)

The Cluster Status utility, `/usr/es/sbin/cluster/clstat`, monitors cluster status. For the cluster as a whole, **clstat** indicates the cluster state and the number of cluster nodes. For each node, **clstat** displays the IP label and address of each service network interface attached to the node, and whether that interface is up or down. **clstat** also displays resource group state.

You can view cluster status information in ASCII or X Window display mode or through a web browser.

Note: **clstat** uses the Clinfo API to retrieve information about the cluster. Therefore, have Clinfo running on the client machine to view the **clstat** display.

HAView Cluster Monitoring Utility

The HAView cluster monitoring utility makes use of the Tivoli TME 10 NetView for AIX graphical interface to provide a set of visual maps and submaps of HACMP clusters. HAView extends NetView services to allow you to monitor HACMP clusters and cluster components across a network from a single node. HAView creates symbols that reflect the state of all nodes, networks, and network interface objects associated in a cluster. You can also monitor resource groups and their resources through HAView.

HAView monitors cluster status using the Simple Network Management Protocol (SNMP). It combines periodic polling and event notification through traps to retrieve cluster topology and state changes from the HACMP Management Information Base (MIB). The MIB is maintained by the Cluster SMUX peer daemon (**clsmuxpd**), the HACMP management agent. HAView allows you to:

- View maps and submaps of cluster symbols showing the location and status of nodes, networks, and addresses, and monitor resource groups and resources.
- View detailed information in NetView dialog boxes about a cluster, network, IP address, and cluster events
- View cluster event history using the HACMP Event Browser
- View node event history using the Cluster Event Log

- Open a SMIT hacmp session for an active node and perform cluster administration functions from within HAView, using the HAView Cluster Administration facility.

Cluster Monitoring and Administration with Tivoli Framework

You can monitor the state of an HACMP cluster and its components through your Tivoli Framework enterprise management system. Using various windows of the Tivoli Desktop, you can monitor the following aspects of your cluster:

- Cluster state and substate
- Configured networks and network state
- Participating nodes and node state
- Configured resource groups and resource group state
- Resource group location.

In addition, you can perform the following cluster administration tasks through Tivoli:

- Start cluster services on specified nodes
- Stop cluster services on specified nodes
- Bring a resource group online
- Bring a resource group offline
- Move a resource group to another node.

For complete information about installing, configuring, and using the cluster monitoring through Tivoli functionality, see the *Administration and Troubleshooting Guide*.

Application Availability Analysis Tool

The Application Availability Analysis tool measures the exact amount of time that any of your applications have been available. The HACMP software collects, time-stamps, and logs extensive information about the applications you choose to monitor with this tool. Using SMIT, you can select a time period and the tool displays uptime and downtime statistics for a given application during that period.

Persistent Node IP Labels

A *persistent node IP label* is a useful administrative “tool” that lets you contact a node even if the HACMP cluster services are down on that node (in this case, HACMP attempts to put an IP address on the node). Assigning a persistent node IP label to a network on a node allows you to have a node-bound IP address on a cluster network that you can use for administrative purposes to access a specific node in the cluster.

A persistent node IP label is an IP alias that can be assigned to a specific node on a cluster network and that:

- Always stays on the same node (is *node-bound*)
- Co-exists on a network interface card that already has a service IP label defined
- Does not require installing an additional physical network interface card on that node
- Is *not* part of any *resource group*.

You can have one persistent node IP label per network per node.

HACMP Verification and Synchronization

After you configure, reconfigure, or update a cluster, you should run the cluster verification procedure on one node to check that all nodes agree on the cluster topology, network configuration, and the ownership and takeover of HACMP resources. If the verification is successful, the configuration is synchronized. Synchronization takes effect immediately on an active cluster.

The Cluster Verification utility, `/usr/es/sbin/cluster/diag/clverify`, verifies that HACMP-specific modifications to AIX system files are correct, that the cluster and its resources are configured correctly, and that security, if set up, is configured correctly. `clverify` also indicates whether custom cluster snapshot methods exist and whether they are executable on each cluster node.

The `clverify` utility keeps a detailed record of the information in the HACMP configuration database on each of the nodes after it runs. Subdirectories for each node contain information for the last successful verification (pass), the next-to-last successful verification (pass.prev), and the last unsuccessful verification (fail).

Messages output by the `clverify` utility follow a common, standardized format where feasible, indicating the node(s), devices, command, etc. in which the error occurred.

Verification with Automatic Cluster Configuration Monitoring

HACMP 5.2 provides automatic cluster configuration monitoring. By default, HACMP runs the `clverify` utility automatically on the node that is first in alphabetical order once every 24 hours at midnight. The cluster administrator is notified if the cluster configuration has become invalid.

When cluster verification completes on the selected cluster node, this node notifies the other cluster nodes. Every node stores the information about the date, time, which node performed the verification, and the results of the verification in the `/var/hacmp/log/clutils.log` file. If the selected node becomes unavailable or cannot complete cluster verification, you can detect this by the lack of a report in the `/var/hacmp/log/clutils.log` file.

In case cluster verification completes and detects some configuration errors, you are notified about the potential problems:

- The exit status of `clverify` is published across the cluster along with the information about cluster verification process completion.
- Broadcast messages are sent across the cluster and displayed on `stdout`. These messages inform you about detected configuration errors.
- A general_notification event runs on the cluster and is logged in `hacmp.out` (if cluster services is running).

Verification with Corrective Actions

Cluster verification consists of a series of checks performed against various HACMP user-configured components. Each check attempts to detect either a cluster consistency issue or a configuration error. Some error conditions result when information important to the operation of HACMP, but not part of the HACMP software itself, is not propagated properly to all cluster nodes.

In HACMP 5.2, when **clverify** detects any of the following conditions, you can authorize a corrective action before **clverify** continues error checking:

- HACMP shared volume group time stamps do not match on all nodes
- The **/etc/hosts** file on a node does not contain all HACMP-managed labels/IP addresses
- SSA concurrent volume groups need SSA node numbers
- A filesystem is not created on a node that is part of the resource group, although disks are available
- Disks are available, but a volume group has not been imported to a node
- Required **/etc/services** entries are missing on a node
- Required HACMP **snmpd** entries are missing on a node.

If an error found by **clverify** triggers any corrective actions, then the utility runs all checks again after it finishes the first pass. If the same check fails again and the original problem is an error, the error is logged and verification fails. If the original condition is a warning, verification succeeds.

Software and Cluster Verification

clverify has two categories of verification options. Running the *software* option ensures that the HACMP-specific modifications to the AIX configuration settings are correct. Verifying the *cluster* ensures that all resources used by the HACMP software are properly configured, and that ownership and takeover of those resources are assigned properly and are in agreement across all cluster nodes.

Custom Verification Methods

Through SMIT you also can add, change, or remove custom-defined verification methods that perform specific checks on your cluster configuration.

You can perform verification from the command line or through the SMIT interface to issue a customized page in response to a cluster event.

Troubleshooting Tools

Typically, a functioning HACMP cluster requires minimal intervention. If a problem occurs, however, diagnostic and recovery skills are essential. Thus, troubleshooting requires that you identify the problem quickly and apply your understanding of the HACMP software to restore the cluster to full operation.

You can use:

- **cldiag**, the Cluster Diagnostic utility. It provides an interface to several HACMP and AIX diagnostic tools you can use to troubleshoot an HACMP cluster.
- **Log files**. The **/usr/es/adm/cluster.log** file tracks cluster events; the **/tmp/hacmp.out** file records the output generated by configuration scripts as they execute; the **/usr/es/sbin/cluster/history/cluster.mmddyyyy** log file logs the daily cluster history; the **/tmp/cspoc.log** file logs the status of C-SPOC commands executed on cluster nodes. You should also check the RSCT log files.

HACMP lets you view, redirect, save and change parameters of the log files, so that you can tailor them to your particular needs.

- **Collecting log files for problem reporting.** If you encounter a problem with HACMP and report it to IBM support, you may be asked to collect log files pertaining to the problem. A SMIT panel under the **HACMP Problem Determination Tools** menu aids in this process. It is recommended to use this panel if requested by the IBM support personnel. If you use this utility without direction from IBM support, be careful to fully understand the actions and the potential consequences.
- The **Cluster Test Tool** lets you test your HACMP configuration to evaluate how a cluster behaves under a set of specified circumstances, such as when a node becomes inaccessible, a network becomes inaccessible, a resource group moves from one node to another and so forth. You can start the test, let it run unattended, and return later to evaluate the results of your testing.

Run the tool after you initially configure HACMP and before you put your cluster into a production environment; after you make cluster configuration changes while the cluster is out of service; or at regular intervals even though the cluster appears to be functioning well.

- **Resetting HACMP Tunable Values.** In HACMP 5.2, you can reset cluster tunable values using the SMIT interface. We recommend that you create a cluster snapshot, prior to resetting. After the values have been reset to defaults, if you want to return to customized cluster settings, you can apply the cluster snapshot.
- The **cluster status information** file, produced by the Cluster Snapshot utility. A cluster snapshot lets you skip saving log files in the snapshot. By default, HACMP no longer collects cluster log files when you create the cluster snapshot, although you can still specify collecting the logs in SMIT. Skipping the logs collection reduces the size of the snapshot and speeds up running the snapshot utility.
- **Automatic Error Notification** lets you configure AIX error notification for certain cluster resources using a specific option in SMIT. If you select option, error notification is automatically turned on or off on all nodes in the cluster for particular devices.
- **Custom Pager Notification** lets you define a notification method through SMIT.
- **User-Defined Events** let you define your own events for which HACMP can run recovery programs that you specify.
- **Event summaries** that appear at the end of each event's details make it easier to check the **hacmp.out** file for errors. The event summaries contain pointers back to the corresponding event, which allow you to easily locate the output for any event.
- **The trace facility** helps you to isolate a problem within an HACMP system by allowing you to monitor selected events. Using the trace facility, you can capture a sequential flow of time-stamped system events that provide a fine level of detail on the activity within an HACMP cluster.

These utilities are described in the following sections. For more detailed information on each of these utilities, see the *Administration and Troubleshooting Guide*.

Cluster Diagnostic Utility

The Cluster Diagnostic utility, **cldiag**, provides a common interface to several HACMP and AIX diagnostic tools you can use to troubleshoot an HACMP cluster. Use **cldiag** to perform the following tasks:

- View the cluster log files for error and status messages
- Check volume group definitions
- Activate and deactivate tracing of HACMP daemons.

Log Files

The HACMP software writes the messages it generates to the system console and to several log files. Because each log file contains a different level of detail, system administrators can focus on different aspects of HACMP processing by viewing different log files. The HACMP software writes messages into the log files described in the following sections.

Note that each log file has a default directory. You may redirect a log file to a storage location other than its default directory, if needed.

Note that in HACMP, you can also collect log files for problem reporting. For more information, see the *Administration and Troubleshooting Guide*.

/usr/es/adm/cluster.log File

The **/usr/es/adm/cluster.log** file contains time-stamped, formatted messages generated by HACMP scripts and daemons.

/tmp/hacmp.out File

The **/tmp/hacmp.out** file contains messages generated by HACMP scripts.

In verbose mode, this log file contains a line-by-line record of every command executed by these scripts, including the values of all arguments to these commands.

Event summaries appear after the verbose output for certain events (those initiated by the Cluster Manager), making it easier to scan the **hacmp.out** file for important information. In addition, event summaries provide HTML links to the corresponding events within the **hacmp.out** file.

Due to the recent changes in the way resource groups are handled and prioritized in fallover circumstances, the **hacmp.out** file and its event summaries are very important in tracking the activity and resulting location of your resource groups. Always check this log early on when investigating resource group movement after takeover activity.

System Error Log File

The system error log contains time-stamped, formatted messages from all AIX subsystems, including HACMP scripts and daemons.

/usr/es/sbin/cluster/history/cluster.mmddyyyy File

The **/usr/es/sbin/cluster/history/cluster.mmddyyyy** file contains time-stamped, formatted messages generated by HACMP scripts. The system creates a cluster history file every day, identifying each file by the file name extension, where *mm* indicates the month, *dd* indicates the day, and *yyyy* indicates the year.

This log reports specific events. Note that when resource groups are processed in parallel, certain steps that were formerly run as events are now processed differently and cannot be reported in this log. Resource group activity and location, especially when processed in parallel, is best tracked using the **hacmp.out** log.

/tmp/clstrmgr.debug File

Contains time-stamped, formatted messages generated by HACMP **clstrmgr** activity. This log is primarily for the use of IBM support personnel.

/var/hacmp/clcomd/clcomd.log

Contains time-stamped, formatted messages generated by **clcomd** Cluster Communication Daemon.

/var/hacmp/clverify/clverify.log

The **/var/hacmp/clverify/clverify.log** file contains the verbose messages output by the **clverify** utility. Cluster verification consists of a series of checks performed against various HACMP configurations. Each check attempts to detect either a cluster consistency issue or an error. The messages output by the **clverify** utility follow a common, standardized format, where feasible, indicating the node(s), devices, command, in which the error occurred.

/tmp/cspoc.log File

Contains time-stamped, formatted messages generated by C-SPOC commands.

/tmp/emuhacmp.out File

Contains time-stamped, formatted messages generated by the HACMP event emulator scripts.

Resetting HACMP Tunable Values

While configuring and testing a cluster, you may change a value for one of the HACMP tunable values that affects the cluster performance. Or, you may want to reset tunable values to their default settings without changing any other aspects of the configuration. A third-party cluster administrator or a consultant may be asked to take over the administration of a cluster that they did not configure and may need to reset the tunable values to their defaults.

Before HACMP 5.2, HACMP displayed the current cluster configuration, but it was sometimes difficult to obtain all the HACMP tunable values that had been changed from their defaults to troubleshoot the cluster. Also, although HACMP allowed you to remove the cluster definition altogether, removing the cluster did not reset customization settings that could have been made.

In HACMP 5.2, you can reset cluster tunable values using the SMIT interface. We recommend that you create a cluster snapshot, prior to resetting. After the values have been reset to defaults, if you want to return to customized cluster settings, you can apply the cluster snapshot.

Resetting the cluster tunable values resets information in the cluster configuration database. The information that is reset or removed comprises two categories:

- Information supplied by the users (for example, pre- and post- event scripts and network parameters, such as netmasks). Note that resetting cluster tunable values *does not* remove the pre- and post-event scripts that you already have configured. However, if you reset the

tunable values, HACMP's knowledge of pre- and post-event scripts is removed from the configuration, and these scripts are no longer used by HACMP to manage resources in your cluster. You can reconfigure HACMP to use these scripts again, if needed.

- Information automatically generated by HACMP during configuration and synchronization. This includes node and network IDs, and information discovered from the operating system, such as netmasks. Typically, users cannot see generated information.

You can reset the tunable values in HACMP 5.2 using either of the following methods:

- A SMIT panel that lets you restore the cluster to the installation-time values, and optionally create a snapshot, prior to resetting

or

- A command line interface that provides an interactive mode to display the results of all verification checks in a cluster and prompts you to restore them to their default values. This mode is useful if you want to reset specific tunable values, such as heartbeat tuning options, without resetting *all* tunable values. This command is intended for use by IBM support.

For a complete list of tunable values that you can restore to their default settings, see the *Planning and Installation Guide*.

For instructions on how to reset the tunable values using SMIT, see the *Administration and Troubleshooting Guide*.

Cluster Status Information File

When you use the HACMP Cluster Snapshot utility to save a record of a cluster configuration (as seen from each cluster node), you optionally cause the utility to run many standard AIX commands and HACMP commands to obtain status information about the cluster. This information is stored in a file, identified by the **.info** extension, in the snapshots directory. The snapshots directory is defined by the value of the SNAPSHOTPATH environment variable. By default, the cluster snapshot utility includes the output from the commands, such as **cllssif**, **cllssnw**, **df**, **ls**, and **netstat**. You can create custom snapshot methods to specify additional information you would like stored in the **.info** file.

A cluster snapshot lets you skip saving log files in the snapshot. Cluster snapshots are used for recording the cluster configuration information, whereas cluster logs only record the operation of the cluster and not the configuration information. By default, HACMP no longer collects cluster log files when you create the cluster snapshot, although you can still specify collecting the logs in SMIT. Skipping the logs collection reduces the size of the snapshot and speeds up running the snapshot utility. The size of the cluster snapshot depends on the configuration. For instance, a basic two-node configuration requires roughly 40KB.

Automatic Error Notification

You can use the AIX Error Notification facility to detect events not specifically monitored by the HACMP software—a disk adapter failure, for example—and specify a response to take place if the event occurs.

Normally, you define error notification methods manually, one by one. HACMP provides a set of pre-specified notification methods for important errors that you can automatically “turn on” in one step through the SMIT interface, saving considerable time and effort by not having to define each notification method manually.

Custom Pager Notification

You can define a notification method through the SMIT interface to issue a customized page in response to a cluster event.

After configuring a pager notification method, you can send a test message to confirm that the configuration is correct.

You can configure any number of notification methods, for different events and with different text messages and telephone numbers to dial. The same notification method can be used for several different events, as long as the associated text message conveys enough information to respond to all of the possible events that trigger the page.

User-Defined Events

You can define your own events for which HACMP can run your specified recovery programs. This adds a new dimension to the predefined HACMP pre- and post-event script customization facility.

Note: HACMP 5.2 interacts with the RSCT Resource Monitoring and Control (RMC) subsystem instead of with the Event Management subsystem. (the Event Management subsystem continues to be used for interaction with Oracle 9i). Only a subset of Event Management user-defined event definitions is automatically converted to the corresponding RMC event definitions, upon migration to HACMP 5.2. After migration is complete, all user-defined event definitions must be manually reconfigured with the exception of seven UDE definitions defined by DB2. For more information, see the *Administration and Troubleshooting Guide*.

You specify the mapping between events that you define and recovery programs defining the event recovery actions through the SMIT interface. This lets you control both the scope of each recovery action and the number of event steps synchronized across all nodes. For details about registering events, see the RSCT documentation.

You must put all the specified recovery programs on all nodes in the cluster, and make sure they are executable, before starting the Cluster Manager on any node.

- *AIX resource monitor.* This monitor generates events for OS-related events such as the percentage of CPU that is idle or percentage of disk space in use. The attribute names start with:
 - `IBM.Host.`
 - `IBM.Processor.`
 - `IBM.PhysicalVolume.`
- *Program resource monitor.* This monitor generates events for process-related occurrences such as unexpected termination of a process. It uses the resource attribute `IBM.Program.ProgramName.`

Note: You cannot use the Event Emulator to emulate a user-defined event.

Event Summaries

Details of cluster events are recorded in the **hacmp.out** file. The verbose output of this file contains many lines of event information; you see a concise summary at the end of each event's details. For a quick and efficient check of what has happened in the cluster lately, you can view a compilation of only the event summary portions of current and previous **hacmp.out** files, by using the **Display Event Summaries** panel in SMIT. You can also select to save the compiled event summaries to a file of your choice. Optionally, event summaries provide HTML links to the corresponding events in the **hacmp.out** file.

For details on viewing event summaries, see the troubleshooting chapter in the *Administration and Troubleshooting Guide*.

Trace Facility

If the log files have no relevant information and the component-by-component investigation does not yield concrete results, you may need to use the HACMP trace facility to attempt to diagnose the problem. The trace facility provides a detailed look at selected system events.

Note that both the HACMP and AIX software must be running in order to use HACMP tracing.

For details on using the trace facility, see the troubleshooting chapter in the *Administration and Troubleshooting Guide*.

Test Tool

The Cluster Test Tool is a utility that lets you test an HACMP cluster configuration to evaluate how a cluster behaves under a set of specified circumstances, such as when a node becomes inaccessible, a network becomes inaccessible, a resource group moves from one node to another and so forth. You can start the test, let it run unattended, and return later to evaluate the results of your testing.

If you want to run an automated suite of basic cluster tests for topology and resource group management, you can run the automated test suite from SMIT. If you are an experienced HACMP administrator and want to tailor cluster testing to your environment, you can also create custom tests that can be run from SMIT.

Emulation Tools

HACMP includes the Event Emulator for running cluster event emulations and the Error Emulation functionality for testing notification methods.

HACMP Event Emulator

The HACMP Event Emulator is a utility that emulates cluster events and dynamic reconfiguration events by running event scripts that produce output but that do not affect the cluster configuration or status. Emulation allows you to predict a cluster's reaction to a particular event just as though the event actually occurred.

The Event Emulator follows the same procedure used by the Cluster Manager given a particular event, but does not execute any commands that would change the status of the Cluster Manager. For descriptions of cluster events and how the Cluster Manager processes these events, see the *Administration and Troubleshooting Guide*.

You can run the Event Emulator through SMIT or from the command line. The Event Emulator runs the events scripts on every active node of a stable cluster, regardless of the cluster's size. The output from each node is stored in an output file on the node from which the event emulator is invoked. You can specify the name and location of the output file using the environment variable **EMUL_OUTPUT**, or use the default output file, **/tmp/emuhacmp.out**.

Note: The Event Emulator requires that both the Cluster SMUX peer daemon (**clsmuxpd**) and the Cluster Information Program (**clinfo**) be running on your cluster.

The events emulated are categorized in two groups:

- Cluster events
- Dynamic reconfiguration events.

Emulating Cluster Events

The cluster events that can be emulated are:

node_up	fail_standby
node_down	join_standby
network_up	swap_adapter
network_down	

Emulating Dynamic Reconfiguration Events

The dynamic reconfiguration event that can be emulated is Synchronize the HACMP Cluster.

Restrictions on Event Emulation

Note: If your current cluster does not meet any of the following restrictions, you can use the cluster test tool as an alternative to executing cluster events in emulation mode. The cluster test tool performs real cluster events and pre- and post-event customizations.

The Event Emulator has the following restrictions:

- You can run only one instance of the event emulator at a time. If you attempt to start a new emulation in a cluster while an emulation is already running, the integrity of the results cannot be guaranteed.
- **clinfo** must be running.
- You cannot run successive emulations. Each emulation is a standalone process; one emulation cannot be based on the results of a previous emulation.

- When you run an event emulation, the Emulator's outcome may be different from the cluster manager's reaction to the same event under certain conditions:
 - The Event Emulator will not change the configuration of a cluster device. Therefore, if your configuration contains a process that makes changes to the Cluster Manager (disk fencing, for example), the Event Emulator will not show these changes. This could lead to a different output, especially if the hardware devices cause a fallover.
 - The Event Emulator runs customized scripts (pre- and post-event scripts) associated with an event, but does not execute commands within these scripts. Therefore, if these customized scripts change the cluster configuration when actually run, the outcome may differ from the outcome of an emulation.
- When emulating an event that contains a customized script, the Event Emulator uses the **ksh** flags **-n** and **-v**. The **-n** flag reads commands and checks them for syntax errors, but does not execute them. The **-v** flag indicates verbose mode. When writing customized scripts that may be accessed during an emulation, be aware that the other **ksh** flags may not be compatible with the **-n** flag and may cause unpredictable results during the emulation. See the **ksh** man page for flag descriptions.

Emulation of Error Log Driven Events

Although the HACMP software does not monitor the status of disk resources, it does provide a SMIT interface to the AIX Error Notification facility.

HACMP uses the following utilities for monitoring purposes:

- RSCT
- AIX Error Notification
- RMC
- User-defined events
- Application monitoring.

The AIX Error Notification facility allows you to detect an event not specifically monitored by the HACMP software—a disk adapter failure, for example—and to program a response (notification method) to the event. In addition, if you add a volume group to a resource group, HACMP automatically creates an AIX Error Notification method for it. In the case the loss of quorum error occurs for a mirrored volume group, HACMP uses this method to selectively move the affected resource group to another node. Do not edit or alter the error notification methods that are generated by HACMP.

HACMP provides a utility for testing your error notification methods. After you add one or more error notification methods with the AIX Error Notification facility, you can test your methods by emulating an error. By inserting an error into the AIX error device file (**/dev/error**), you cause the AIX error daemon to run the appropriate pre-specified notification method. This allows you to determine whether your pre-defined action is carried through, without having to actually cause the error to occur.

When the emulation is complete, you can view the error log by typing the **errpt** command to be sure the emulation took place. The error log entry has either the resource name **EMULATOR**, or a name as specified by the user in the **Resource Name** field during the process of creating an error notification object.

You will then be able to determine whether the specified notification method was carried out.

Chapter 8: HACMP 5.2: Summary of Changes

This chapter lists all new or enhanced features in HACMP 5.2 and also notes discontinued features.

List of New Features

HACMP 5.2 includes the following new or enhanced features. The features are listed under a particular type of enhancement, although some improve both usability and performance, and some also offer improved interoperability with other IBM products.

New Features that Enhance Ease of Use

These features make the product easier to use:

- **All Resource Groups are Custom Resource Groups.** In HACMP 5.2, all resource groups are configured in the same way as you configured custom resource groups in the previous release. All groups are referred to as simply resource groups.

In HACMP 5.2, you can create resource groups that replicate pre-5.2 cascading, rotating and concurrent resource groups, as well as Inactive Takeover and Fallover without Fallback settings, if needed.

For a complete table of how pre-5.2 resource group characteristics are mapped to the combinations of startup, fallover and fallback policies available in HACMP 5.2, see *Chapter 10: Upgrading an HACMP Cluster* in the *Planning and Installation Guide*.

In addition, the following functionality is available for resource groups in HACMP 5.2:

- Distribution policy for resource group startup. This policy ensures that during a node or cluster startup, only one resource group is brought online on a node (node distribution), or on a node per network (network distribution).
- Extended support for IPAT via IP replacement networks and service IP labels.

HACMP 5.2 prevents you from configuring invalid combinations of behavior, resources and site policies in a resource group, and displays only valid choices during the configuration process.

- **Configuring Dependencies between Resource Groups.** In HACMP 5.2, it is easier to set up more complex clusters with multi-tier applications by specifying dependencies between resource groups that contain these applications. In prior releases, with customized serial processing, you could specify that a resource group is processed before another resource group, on a local node. However, it was not guaranteed that a resource group would be processed in the order specified, as this order is affected by other resource group processing information. Also, previously, you could use pre- and post-event scripts to configure a customized dependency mechanism between different applications defined to HACMP. Although you can continue using these scripts, HACMP 5.2 offers an easy way to configure a dependency between resource groups (and applications that belong to them).

In this release, the dependency that you configure is:

- Explicitly specified using the SMIT interface
- Established cluster-wide, not just on the local node
- Guaranteed to be honored in the cluster, that is, it is not affected by the current cluster conditions.

Configuring a resource group dependency allows for easier configuration and control for clusters with multi-tier applications where one application depends on the successful startup of another application, and both applications are required to be kept highly available with HACMP.

The following example illustrates the dependency behavior you can configure in HACMP 5.2:

- If resource group A (*child*) depends on resource group B (*parent*), upon node startup, resource group B must be brought online before resource group A is acquired on any node in the cluster. Upon failover, the order is reversed: Resource group A (*child*) must be taken offline before resource group B (*parent*) is taken offline.
- In addition, if resource group A depends on resource group B, then during a node startup or node reintegration, resource group A cannot be taken online before resource group B is brought online. If resource group B is taken offline, resource group A will be taken offline too, since it depends on resource group B.
- **Two-Node Cluster Configuration Assistant.** The Two-Node Cluster Configuration Assistant simplifies the process for configuring an HACMP cluster that includes two nodes and associated volume groups, networks, service IP labels and applications. The Assistant guides you through the configuration process and provides online user assistance to help you enter correct information for the five entry fields. The application is designed for users with little knowledge of HACMP who want to quickly set up a basic two-node cluster with one non-concurrent resource group, one application server, volume groups and one service label.
- **Cluster Test Tool.** In HACMP 5.2, you can use the Cluster Test Tool to test your HACMP configuration to evaluate how a cluster behaves under a set of specified circumstances, such as when a node becomes inaccessible, a network becomes inaccessible, a resource group moves from one node to another and so forth. You can start the test, let it run unattended, and return later to evaluate the results of your testing. You should run the tool:
 - After you initially configure HACMP and before you put your cluster into a production environment
 - After you make cluster configuration changes while the cluster is out of service
 - At regular intervals even though the cluster appears to be functioning well.
- **HACMP File Collections.** Similar to volume groups, certain files located on each cluster node need to be kept in sync in order for HACMP (and other applications) to behave correctly. Such files include event scripts, application scripts, and some AIX and HACMP configuration files.

HACMP File Collection provides an easy way to request that a list of files be kept in sync across the cluster. Using HACMP File Collection, you no longer have to manually copy an updated file to every cluster node, verify that the file is properly copied, and confirm that each node has the same version of it.

Also, if one or more of these files is inadvertently deleted or damaged on one or more cluster nodes, it can take time and effort to determine the problem. Using HACMP file collection, this scenario is mitigated. HACMP detects when a file in a file collection is deleted or if the file size is changed to zero, and logs a message to inform the administrator.

Two predefined HACMP file collections are installed by default:

- **Configuration_Files.** A container for essential system files, such as **/etc/hosts** and **/etc/services**.
- **HACMP_Files.** A container for all the user-configurable files in the HACMP configuration. This is a special file collection that the underlying file collection propagation utility uses to reference all the user-configurable files in the HACMP configuration database classes.

- **Automatic cluster configuration checking by the cluster verification utility.** By default, HACMP 5.2 runs the **clverify** utility automatically on the node that is first in alphabetical order once every 24 hours at midnight. HACMP notifies the cluster administrator if problems with the cluster configuration are detected.

When cluster verification completes on the selected cluster node, this node notifies the other cluster nodes. Every node stores the information about the date, time, which node performed the verification, and the results of the verification in the **/var/hacmp/log/clutils.log** file. If the selected node becomes unavailable or cannot complete cluster verification, you can detect this by the lack of a report in the log file.

In case cluster verification completes and detects some configuration errors, you are notified about the potential problems.

- **Web-Based Cluster Management.** HACMP 5.2 includes a Web-enabled user interface (WebSMIT) that provides consolidated access to the SMIT functions of configuration and management as well as to interactive cluster status and the HACMP documentation.

The WebSMIT interface is similar to the ASCII SMIT interface. You do not need to learn a new user interface or terminology and can easily switch between ASCII SMIT and the Web version. Because WebSMIT runs in a Web browser, it can be accessed from any platform.

- **Cluster Status Display (clstat) Linked to Management Functions.** When using the WebSMIT interface to see the cluster status display, you have links to the related WebSMIT management functions; thus HACMP 5.2 provides a consolidated user interface for cluster status with management capabilities.

For example, among other options, the node status display has a link to the SMIT panels for starting and stopping Cluster Services. Now you can manipulate entities in the status display interactively rather than having to access an ASCII SMIT interface on the node.

- **Show Cluster Applications Utility.** HACMP 5.2 provides a new **Show Cluster Applications** SMIT option that provides an application-centric view of the cluster configuration. This utility displays existing interfaces and information from an application point of view. You can access it from both SMIT and WebSMIT. Using the WebSMIT version lets you expand and collapse areas of the information. Colors reflect the state of individual items (for example, green indicates online).

- **Resetting HACMP Tunable Values.** In HACMP 5.2, you can change the settings for a list of tunable values that were changed during the cluster maintenance and reset them to their default settings, or installation-time cluster settings. Resetting the tunable values to their defaults helps to troubleshoot cluster performance and assists third-party administrators

(such as IBM support personnel) in case they take over the administration of a cluster. We recommend that you create a cluster snapshot, prior to resetting. After the values have been reset to defaults, if you want to return to customized cluster settings, you can apply the cluster snapshot. Resetting the tunable values *does not* change any aspects of the configuration other than tunable values. For a list of tunable values that can be reset, see the *Administration and Troubleshooting Guide*.

- **Cluster Snapshot without Logs and Collecting Cluster Log Files for Problem Reporting.** In HACMP 5.2, cluster snapshot lets you skip saving log files in the snapshot. Cluster snapshots are used for recording the cluster configuration information, whereas cluster logs only record the operation of the cluster and not the configuration information. By default, HACMP no longer collects cluster log files when you create the cluster snapshot, although you can still specify collecting the logs in SMIT. Skipping the logs collection speeds up the running time of the snapshot utility, and reduces the size of the snapshot.

In addition, you can use SMIT to collect cluster log files for problem reporting. This option is available under the **HACMP Problem Determination Tools > HACMP Log Viewing and Management** SMIT menu. It is recommended to use this option only if requested by the IBM support personnel.

- **Cluster-wide Password Change for Users.** Users can now change their password across cluster nodes when authorized to do so. A new Cluster Password utility links to the AIX password utility to support this change. System administrators enable the Cluster Password utility on all nodes in a cluster or nodes in specified resource groups, then give specified users permission to change their password on particular nodes. The utility also provides a new command, **clpasswd**. For users trying to change their own passwords, good message support gives them information about issues that may arise during this process.
- **Online Planning Worksheets.** The Online Planning Worksheets application now lets you:
 - View the cluster definition for a local, active HACMP cluster
 - Create a worksheets file from a local, active HACMP cluster running HACMP 5.2 or greater
 - Save a cluster definition in a worksheets file from SMIT or from Online Planning Worksheets.

A worksheets file that stores a cluster definition from an active HACMP cluster contains configuration data about the components in an HACMP cluster that are supported by the Online Planning Worksheets application. The application stores information you enter for applications, disks, volume groups, logical volumes, or NFS mounted drives, but it does not provide this data for HACMP to configure a cluster. Other AIX system components manage the configuration data for these components.

Note that the interface lets you save and export the cluster configuration information from the worksheets utility and then apply the configuration to a cluster, using the **cl_opsconfig** command.

Viewing an active cluster definition, or a worksheets file saved from a cluster definition, lets you review information about a cluster configuration in an easy-to-view format for use in testing and troubleshooting situations. You can also generate a report from this information to produce formal documentation of the cluster state.

In HACMP 5.2, the application uses a new format for the worksheets file, as identified by the **.haw** extension. The application lets you open a worksheets file with the **.ws** extension from HACMP 5.1, then save that file to the **.haw** format.

There are also changes to the user interface in Online Planning Worksheets to reflect changes in resource group and site configuration in HACMP. In addition, the application has the following new panels:

- File Collections
- File Collections Global Settings
- Cross-Site LVM Mirroring
- Cluster Verification.

New Features that Enhance Performance

These HACMP 5.2 features enhance performance of the product:

- **Recovering Resource Groups on Node Startup.** Prior to HACMP 5.2, when a node joined the cluster, it did not make an attempt to acquire any resource groups that had previously gone into an ERROR state on any other node. Such resource groups remained in the ERROR state and required use of the Resource Group Migration utility, **clRGmove**, to manually bring them back online.

In HACMP 5.2, an attempt is made to bring online the resource groups that are currently in the ERROR state. This further increases the chances of bringing the applications back online. When a node starts up, if a resource group is in the ERROR state on any node in the cluster, this node attempts to acquire the resource group. Note that the node must be included in the nodelist for the resource group.

The resource group recovery on node startup is different for non-concurrent and concurrent resource groups:

- If the starting node fails to activate a *non-concurrent resource group* that is in the ERROR state, the resource group continues to fall over to another node in the nodelist, if a node is available. The fallover action continues until all available nodes in the nodelist have been tried.
- If the starting node fails to activate a *concurrent resource group* that is in the ERROR state on the node, the concurrent resource group is left in the ERROR state on that node. Note that the resource group might still remain online on other nodes.
- **Cluster verification with auto-corrective actions.** In HACMP 5.2, you can authorize a corrective action before **clverify** continues error checking, when **clverify** detects any of the following conditions:
 - HACMP shared volume group time stamps do not match on all nodes
 - The **/etc/hosts** file on a node does not contain all HACMP-managed labels/IP addresses
 - SSA concurrent volume groups need SSA node numbers
 - A filesystem does not exist on a node that is part of the resource group, although disks are available
 - Disks are available, but a volume group has not been imported to a node
 - Required **/etc/services** entries are missing on a node
 - Required HACMP **snmpd** entries are missing on a node.

If an error found by **clverify** triggers any corrective actions, then the utility runs all checks again after it finishes the first pass. If the same check fails again and the original problem is an error, the error is logged and verification fails. If the original condition is a warning, verification succeeds.

- **Resource Monitoring and Control (RMC) Subsystem Replaces RSCT Event Management.** HACMP 5.2 interacts with the RSCT Resource Monitoring and Control (RMC) subsystem instead of with the Event Management subsystem. HACMP 5.2 uses the RMC subsystem for these instances:
 - Dynamic Node Priority
 - Application Monitoring
 - User Defined Events.

Note: The Event Management daemon is still present, but is only used by Oracle 9i to access network status through the EMAPI.

- **Multiple Application Monitors.** In HACMP 5.2, you can configure multiple application monitors and associate them with one or more application servers. You can assign each monitor a unique name in SMIT. Prior to HACMP 5.2, for each application that is kept highly available, you could configure only one of the two types of monitors: a monitor to check whether a specific process is terminated in the cluster, or a monitor to check the state of the application by the means of a customized script.

By supporting multiple monitors per application, HACMP can support more complex configurations. For example, you can configure one monitor for each instance of an Oracle parallel server in use. Or, you can configure a custom monitor to check the health of the database along with a process termination monitor to instantly detect termination of the database process.

- **New Application Monitoring Mode Added: Application Startup Monitoring.** You can configure application monitors to function in the application startup monitoring mode. The monitors in this mode monitor the startup of the application server start script within the specified stabilization interval. Specifying the application startup monitoring mode is strongly recommended for applications included in resource groups on which other resource group(s) depend. This ensures that the child resource group(s) can be started successfully, after the parent resource group has started.
- **Improved Security.** In addition to connection authentication to protect HACMP communications between cluster nodes, HACMP now provides additional security for HACMP messages sent between nodes:
 - Message authentication ensures the origination and integrity of a message.
 - Message encryption changes the appearance of the data as it is transmitted and returns it to its original form when received by a node that authenticates the message.

You can configure message authentication alone or with message encryption. You cannot configure message encryption by itself. Both of these mechanisms require careful management of encryption keys between cluster nodes.

HACMP supports the following types of encryption keys for message authentication and encryption:

- Message Digest 5 (MD5) with Data Encryption Standard (DES)
- MD5 with Triple DES

- MD5 with Advanced Encryption Standard (AES)

AIX 5.1 does not support the encryption modules.

- **Communications Performance Enhancements.** Versions previous to HACMP 5.2 use Remote Procedure Calls (RPC) to pass information between the event scripts and cluster utilities of HACMP and the Cluster Manager. In HACMP 5.2, a faster and more robust inter-process communication (IPC) implementation replaces the use of RPC in the Cluster Manager and its external utilities. Cluster events should take less time to complete than in previous releases.

All local IPC requests are unencrypted. All remote IPC requests go through **clcomd** and therefore use the security provided by **clcomd** for encryption between nodes.

The new IPC interface logs all detectable errors and either returns or processes error codes. All successful and unsuccessful IPC requests, including the name of the procedure to be called on the server end of the IPC request, are logged from the **clstrmgrES** and **clcomd** daemons. The logging location differs for each component:

- Cluster utilities log errors (failed IPC requests) to standard error.
- **clstrmgrES** daemon logs to **clstrmgr.debug**.
- **clcomd** daemon logs to **clcomd.log**.

New Features that Enhance Geographic Distance Capability

These features add to the capability for distributing the cluster over a geographic distance, for improved availability and disaster recovery:

- **Cross-Site LVM Mirroring.** In HACMP 5.2, you can set up disks located at two different sites for cross-site LVM mirroring (using a Storage Area Network (SAN), for example). Cross-Site LVM mirroring replicates data between the disk subsystems at each site for disaster recovery.

Two or more nodes located at different sites can access the same physical disks, which may be separated by some distance via the common SAN. These remote disks can be combined into a volume group via the AIX Logical Volume Manager and this volume group may be imported to the nodes located at different sites. The logical volumes in this volume group can have up to three mirrors. Thus you can set up at least one mirror at each site. The information stored on this logical volume is kept highly available, and in case of certain failures, the remote mirror at another site will still have the latest information and the operations can be continued on the other site.

HACMP 5.2 automatically synchronizes mirrors after a disk or node failure and subsequent reintegration. HACMP handles the automatic mirror synchronization even if one of the disks is in the PVREMOVED or PVMISSING state. The automatic synchronization is not possible for all cases, but then you can use C-SPOC to synchronize the data from the surviving mirrors to stale mirrors after a disk/site failure and subsequent reintegration.

- **High Availability Cluster Multi-Processing XD (Extended Distance) for ESS PPRC (HACMP/XD for ESS PPRC)** increases data availability for IBM TotalStorage Enterprise Storage Server (ESS) volumes that use Peer-to-Peer Remote Copy (PPRC) to copy data to a remote site for disaster recovery purposes. With HACMP 5.2, you can use the (*optional*) Enterprise Remote Copy Management Facility (eRCMF) with HACMP/XD for ESS PPRC.

Discontinued Features

The following features, utilities, or functions are discontinued in HACMP 5.2:

- The **cllockd** or **cllockdES** (the Cluster Lock Manager) **is no longer supported in HACMP 5.2**. During node-by-node migration, it is uninstalled. Installing HACMP 5.2 removes the Lock Manager binaries and definitions. Once a node is upgraded to HACMP 5.2, the Lock Manager state information in SNMP and **clinfo** shows the Lock Manager as being in the down state on all nodes, regardless of whether or not the Lock Manager is still running on a back-level node.

Before upgrading, make sure your applications use their proprietary locking mechanism. Check with your application's vendor about concurrent access support.

- **Cascading, rotating and predefined concurrent resource groups are not supported.** Also, Cascading without Fallback and Inactive Takeover settings are not used. In HACMP 5.2 you can continue using the groups migrated from previous releases. You can also configure these types of groups using the combinations of startup, fallover and fallback policies available for resource groups in HACMP 5.2. For information on how the predefined resource groups and their settings are mapped to the startup, fallover and fallback policies, see the chapter on upgrading to HACMP 5.2 in the *Planning and Installation Guide*.
- **Manual reconfiguration of user-defined events is required.** HACMP 5.2 interacts with the RSCT Resource Monitoring and Control (RMC) subsystem instead of with the Event Management subsystem. This affects the following utilities:
 - Dynamic Node Priority (DNP)
 - Application Monitoring
 - User-Defined Events (UDE).

You must manually reconfigure all user-defined event definitions with the exception of the several user-defined event definitions defined by DB2. The **clconvert** command only converts a subset of Event Management user-defined event definitions to the corresponding RMC event definitions. For complete information on the mapping of the Event Management resource variables to the RMC resource attributes, see *Appendix G: RSCT: Resource Monitoring and Control Subsystem* in the *Administration and Troubleshooting Guide*.

Where You Go From Here

For planning an HACMP cluster and installing HACMP, see the *Planning and Installation Guide*.

For configuring HACMP cluster components and troubleshooting HACMP clusters, see the *Administration and Troubleshooting Guide*.

Notices for HACMP Concepts and Facilities Guide

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

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

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

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

For license inquiries regarding double-byte (DBCS) information, contact the IBM Intellectual Property Department in your country or send inquiries, in writing, to:

IBM World Trade Asia Corporation
Licensing
2-31 Roppongi 3-chome, Minato-ku
Tokyo 106, Japan

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

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

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

Licensees of this program who wish to have information about it for the purpose of enabling: (i) the exchange of information between independently created programs and other programs (including this one) and (ii) the mutual use of the information which has been exchanged, should contact:

IBM Corporation
Dept. LRAS / Bldg. 003
11400 Burnet Road
Austin, TX 78758-3493
U.S.A.

Such information may be available, subject to appropriate terms and conditions, including in some cases, payment of a fee.

The licensed program described in this document and all licensed material available for it are provided by IBM under terms of the IBM Customer Agreement, IBM International Program License Agreement or any equivalent agreement between us.

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

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

Index

+-* /

- /.rhosts file 106
- /tmp/cspoc.log file 118
- /tmp/emuhacmp.out file 118
- /tmp/hacmp.out file 117
 - event summaries 121
- /usr/es/sbin/cluster/etc/objrepos/active directory 80
- /usr/es/sbin/cluster/wsm/README file
 - WebSMIT information 108
- /usr/sbin/cluster/etc/objrepos/stage directory 81
- /var/hacmp/clcomd/clcomd.log log file 118
- /var/hacmp/clverify/clverify.log file 118
- /var/hacmp/log/clutils.log 114, 127

0,1,2...

- 7133 SSA disk subsystem
 - cluster support 52

A

- ACD 80
- Active Configuration Directory (ACD) 79
- adapter swap
 - network 74
- adapters
 - swapping dynamically 84
- administrative facilities (overview) 106
- AIX
 - 5.1 and enhanced concurrent mode 64
 - error notification 78
 - overview 16
 - System Resource Controller (SRC) 73
- API
 - Clinfo 61
- applications
 - availability analysis tool 112, 113
 - eliminating as SPOF 73
 - keeping highly available 37
 - monitoring with multiple monitors 111
 - multiple monitors 74
 - multi-tiered example 99
 - server 37
 - starting and stopping 108
 - takeover 73
- ARP
 - gratuitous 30
- authentication of inter-node communications 130
- automated cluster testing 121
- automatic monitoring 114

- automatic synchronization of mirrors 87

C

- cascading resource groups
 - not supported in HACMP 5.2 39
- changing your password 128
- child resource group 48
- clcomd 106
- cldare command 82
- cldiag utility
 - overview 116
- clients
 - "cluster-aware" 60
 - defined 21
- Clinfo 57, 59, 60
 - APIs
 - C 61
 - C++ 61
 - clinfo daemon 60
 - clinfo.rc script 60
 - clRGmove utility 83
 - clsmuxpd daemon 60
 - clstat
 - enhancements 127
 - clstat utility 112
- cluster
 - components 18
 - nodes 19
 - shared disks 20
 - concurrent access
 - eight-node mutual takeover 97
 - two-node mutual takeover 97
 - cross-site LVM mirroring configuration 100
 - dynamic LPAR configuration 100
 - example configurations 91
 - high-level description 18
 - mutual takeover configurations
 - eight-node 97
 - two-node 97
 - networks 20
 - non-concurrent access configurations
 - standby 91
 - takeover 95
 - partitioned 76
- Cluster Communication Subsystem 57
- Cluster Communications daemon 106

Index

D – D

- cluster configuration
 - extended SMIT path 106
 - saving with snapshots 109
 - standard SMIT path 105
 - Two-Node Cluster Configuration Assistant 107
- cluster definition with Online Planning Worksheets 128
- cluster diagnostic utility
 - overview 116
- cluster events
 - event customization facility 89
 - events 89
 - notification script 90
 - pre- and post-event scripts 90
 - processing
 - failover 88
 - reintegration 88
 - recovery script 90
- Cluster Information Program
 - overview 60
- cluster lock manager
 - not supported in HACMP 5.2 132
- cluster management
 - WebSMIT 108, 127
- Cluster Manager
 - event customization 89
- cluster monitoring
 - application availability analysis 112
 - clstat utility
 - overview 112
 - HAVIEW utility
 - overview 112
 - Tivoli
 - overview 113
- cluster monitoring tools
 - clstat 111
 - haview 111
 - Tivoli 111
- cluster multi-processing
 - defined 15
- Cluster Password utility 128
- Cluster SMUX Peer and SNMP programs 59
- cluster snapshot
 - .info file 119
 - overview 109
- cluster software 54
 - GPFS 65
 - HAGEO 65
 - overview 54
 - PPRC 65
 - Workload Manager 65
- cluster status (clstat) utility 112
- Cluster Test Tool 116, 121, 126
- cluster testing 121
- cluster with dynamic LPAR 102
- cluster.log file 117
- cluster.mmdyyy file 117
- clverify 114
 - automatic monitoring 114
 - corrective action 115, 129
- clverify log file 118
- communication path to a node 104
- concurrent access mode
 - applications 64
 - defined 62
 - enhanced concurrent mode in AIX 5.1 64
 - mirroring 63
- concurrent resource groups 39
- concurrent volume groups
 - enhanced concurrent mode 64
- configuration path
 - extended 106
 - standard 105
- configuring clusters
 - tools 106
- corrective action
 - clverify 115, 129
- Cross-site LVM mirroring 86
- cross-site LVM mirroring
 - disaster recovery 131
 - illustration 100
- C-SPOC utility 108
 - /tmp/cspoc.log file 118
- custom cluster testing 121
- custom pager notification method 120
- custom resource groups and HACMP 5.2 40

D

- DARE Resource Migration
 - replaced by Resource Group Management 82, 110
- DB2 98
- DCD
 - creating 79
- Default Configuration Directory
 - DCD 79
- default node priority
 - overview for resource groups 40
- delayed fallback timer 45
 - example when to use it 85
 - overview 85
- dependent resource groups 47
- diagnostic information
 - cluster information file
 - overview 119
- diagnostic utility
 - cldiag utility
 - overview 116
- disaster recovery 86
 - sample configuration 100
- discovery process 105
- disk adapters
 - eliminating as SPOF 78

- disk heartbeating 32
- disk subsystems
 - supported for HACMP 52
- disk takeover 69
 - fast 86
- disks
 - 2105 Enterprise Storage Server 52
 - eliminating as SPOF 78
 - SCSI 20
 - shared 20
 - SSA subsystem 52
- distribution
 - node and network 45
- distribution policy 45
- dynamic adapter swap 84
- dynamic node priority 45
- dynamic node priority (DNP) 40
- dynamic reconfiguration
 - defined 79
 - description of processing 79

E

- eliminating single points of failure 67
- emulating
 - cluster events 121
 - dynamic reconfiguration events 121
 - error log entries 123
- encryption of inter-node communications 130
- enhanced concurrent mode 64
 - fast disk takeover 65
- Enterprise Storage Server (ESS)
 - overview 52
- error notification 78
 - automatic 78, 116
 - volume group loss 78
- error notification methods
 - testing 123
- Ethernet 20
- event customization facility 89
- event emulator 121
 - /tmp/emuhacmp.out file 118
- Event Management subsystem
 - replaced by RMC 130
- event summaries
 - displaying 121
- events
 - customizing duration before config_too_long 90
 - customizing event processing 89
 - emulation 89, 121
 - notification 90
 - pre- and post-event scripts 90
 - processing 88
 - fallover 88
 - reintegration 88
 - recovery and retry 90

- extended configuration path
 - overview 106

F

- facilities, administrative (overview) 106
- fallback of a resource group 41
- fallback options for resource groups 44
- fallover
 - defined 88
 - speeding up with fast recovery 85
- fallover of a resource group 41
- fallover options for resource groups 44
- fast disk takeover 65, 86
- fast recovery
 - configuring resource groups for 85
- file collection management 110
- files
 - /.rhosts 106
- filesystems
 - as shared LVM component 36
 - jfs2 17, 37

G

- global
 - network failure event 27
 - networks 27
- GPFS support 65, 109

H

- HACMP
 - new features in version 5.2 125
- HACMP and HACMP/ES
 - name of the product 18
- HACMP Configuration Database
 - ODM 79
- HACMP file collection 126
- HACMP for AIX
 - a cluster diagram 19
 - LPP software 56
 - new features in this version 125
- hacmp.out 115
- HACMP/XD
 - description of feature 10, 24, 65
- HACMP/XD for ESS PPRC
 - new features 131
- hardware address swapping 70
- HAView
 - overview 112
- heartbeat ring 32
- heartbeating
 - over disk busses 32
 - over IP Aliases 32
 - overview 31
- high availability
 - dynamic reconfiguration 79

Index

I – N

- home node
 - for non-concurrent resource groups 41
- hostname 104

I

- IBM 2104 Expandable Storage Plus 53
- IBM eServer Cluster 1600 52
- IBM FAStT500 Storage Server 53
- IBM FAStT700 Storage Server 53
- IBM Pseries 690 51
- IBM Reliable Scalable Cluster Technology 58
- initial communication path 104
- inter-node communications, authentication and encryption 130
- inter-process communication (IPC)
 - replaces RPC in HACMP 5.2 131
- IP address
 - as a cluster resource 38
 - persistent 113
- IP address takeover 29, 70
 - via IP Aliases 30
 - via IP Replacement 31
- IP alias 29
 - heartbeating 32
- IP service label
 - defined 29
 - requirements for IPAT 30
- IPC
 - replaces RPC in HACMP 5.2 131

J

- journalized file systems 37
 - enhanced 37

K

- keepalives
 - overview 31
- Kerberos security 57

L

- local network failure event
 - recovery actions 76
- lock manager
 - not supported in HACMP 5.2 132

- log files 117
 - /tmp/clstrmgr.debug 118
 - /tmp/cspoc.log 118
 - /tmp/emuhacmp.out 118
 - /tmp/hacmp.out file 117
 - /var/hacmp/log/clutils.log 114, 127
 - cluster.log file 117
 - cluster.mmdd file 117
 - hacmp.out 115
 - problem reporting 128
 - system error log 117
- logical volume manager (LVM) 56
- logical volumes
 - as shared LVM component 36
- logs
 - viewing 108
- LPAR
 - example of a cluster with dynamic LPARs 100
 - pSeries 690 logical partition 52
- LVM 56

M

- managing file collections 110
- message authentication and encryption 130
- message authentication and message 57
- MIB
 - SNMP 60
 - the HACMP MIB 60
- migrating resource groups
 - overview of Resource Group Management 82
- mirroring
 - shared disks 62
- mirrors
 - automatic synchronization 87
- monitoring
 - applications 111
 - cluster 111
 - tools (overview) 111
 - network interfaces 111
 - node status 111
 - nodes and network interfaces
 - with clstat 111
 - pager notification 116
- monitoring a cluster
 - list of tools 115
- multi-tiered applications 99
- mutual takeover configurations
 - eight-node 97
 - two-node 97

N

- network
 - point-to-point 21

- network adapters
 - monitoring 111
 - swapping 74
- network distribution 45
- network failure
 - defined 76
- network interfaces
 - eliminating as SPOF 74
- networks
 - ATM 20
 - eliminating as SPOF 76
 - Ethernet 20
 - global 26
 - IP-based and device-based 21, 25
 - logical 26
 - physical 25
 - Token-Ring 20
- new features 125
- new or enhanced features
 - in this version 125
- node distribution 45
- node isolation 76
 - preventing 77
- node name
 - hostname 104
- nodelist
 - establishing node priorities 40, 42
- nodes
 - defined 19
 - eliminating as SPOF 68
 - monitoring 111
- non-concurrent access
 - applications 62
 - defined 61
 - mirroring 62
- notification
 - event 90

O

- Online Planning Worksheets
 - overview 107
 - working with a cluster definition 128

P

- pager notification 120
- parent resource group 48
- participating nodelist
 - for a resource group 40
- partitioned clusters 76
- passwords, changing 128
- planning
 - shared disks
 - SSA disk subsystem 52
 - shared LVM components
 - filesystems 36

- logical volumes 36
- point-to-point network 21
- PPRC support 65
- pre- and post-event scripts 90
- priorities in nodelists 40, 42
- priority override location (POL) 83
- problem reporting 128
- Pseries 690 51

R

- RAID concurrent volume groups
 - convert to enhanced concurrent 64
- recovery
 - event 90
 - fast recovery 85
 - resource groups 84, 85
- recovery actions
 - for a local network failure event 76
- Regatta
 - IBM Pseries 690 51
- reintegration
 - defined 88
- Remote Procedure Calls
 - replaced by Inter-Process Communication 131
- resetting tunable values 118
- resource group
 - sample standby configuration (2) 94
- resource group management utility 110
- resource groups
 - activating on nodes 108
 - and cluster networks 46
 - and sites 47
 - automatic recovery 84
 - configuring for fast recovery 85
 - delayed fallback timer 45
 - dependent 47
 - distribution policy 45
 - dynamic node priority 40, 45
 - fallback options 44
 - home node 41
 - migrating dynamically

Index

S – S

- overview 110
- mutual takeover configurations 96
- node priority 40
- nodelist 40
- one-sided configurations 95
- one-sided takeover configuration 95
- options for fallover 44
- options for startup 44
- parameters 44
- parent and child definitions 48
- policies and attributes 42
- priority override location (POL) 83
- sample configuration 92
- sample standby configuration (1) 92
- selective fallover 27
- startup, fallover and fallback 41
- upgrading to HACMP 5.2 46
- what is new in HACMP 5.2 125
- Resource Monitoring and Control (RMC)
 - RSCT subsystem 130
- resources
 - cluster
 - introduction/overview 35
 - highly available 35
 - types 35
- retry
 - event 90
- rg_move
 - user-requested 83
- RMC
 - RSCT resource monitoring and control subsystem 130
- rotating resource groups
 - upgrading to HACMP 5.2 46
- rotating resource groups and HACMP 5.2 39
- routing
 - subnet requirements 28
- RPC
 - replaced by IPC 131
- RSCT
 - overview 58
 - services for high availability 58
- RSCT services
 - introduction/overview 55

S

- SCD 81
 - during dynamic reconfiguration 81
- SCSI devices 20
 - in non-concurrent access 62
- SCSI-3 disk arrays
 - in concurrent access 63
 - in non-concurrent access 62
- security 57
 - managing 108
 - message authentication and encryption 130

- serial disks
 - in non-concurrent access 62
- Serial Storage Architecture (SSA) 52
- settling time 44
- shared
 - disk access
 - concurrent access 62
 - non-concurrent access 61
 - disks
 - defined 20
 - supported by HACMP 20
 - LVM components
 - filesystems 36
 - logical volumes 36
- shared disk devices
 - ESS 52
- shared disks
 - SSA disk subsystems 52
- show cluster applications utility 127
- single points of failure
 - applications 73
 - disk adapters 78
 - disks 78
 - eliminating (overview) 67
 - in clusters with LPARs 101
 - network adapters 74
 - networks 76
 - nodes 68
- SMIT
 - extended configuration path 106
 - interface overview 107
 - opening a session on a node 109
 - standard configuration path 105
- SNA
 - configuring communication links 39
- SNMP
 - and the Cluster SMUX Peer 57, 59
 - overview 60
 - snmpd daemon 60
- snmpd daemon 60
- software
 - components of HACMP software 54
- SP Switch
 - and IP address takeover 72
- SSA
 - disk subsystems 52
- SSA disks
 - convert to enhanced concurrent 64
- Staging Configuration Directory (SCD) 81
- standard configuration path 105
- standard security mode 57
- starting and stopping HACMP services 108
- startup 41
- startup of a resource group 41
- startup options for resource groups 44
- Storage Area Network
 - cross-site LVM mirroring 86

- storage servers 52, 53
- subnet
 - logical networks 26
 - requirements for IPAT 30
 - routing 28
 - routing requirements 28
- swapping
 - hardware addresses 70
 - network adapters 74
- system error log file 117
- system management tasks in HACMP 108
- System Resource Controller (SRC) 73

T

- takeover
 - applications 73
 - disk 69
 - eight-node mutual takeover 97
 - hardware address 70
 - IP address 70
 - sample configuration 95
 - two-node mutual takeover 97
- tape drives
 - as cluster resource 38
- Tivoli, cluster monitoring/administration
 - overview 113
- tools in HACMP
 - configuration 106
 - emulation 121
 - installation 106
 - monitoring 111
- trace facility 121
- troubleshooting
 - event summaries 116
 - trace facility 116
- Two-Node Cluster Configuration Assistant 107, 126

U

- upgrading
 - user-defined events 120
- user-defined events
 - upgrading to HACMP 5.2 120

V

- verification
 - automatic monitoring 114
 - corrective action 115
- verification of cluster configuration 114
- volume group loss
 - error notification 78
- VPN for inter-node communications 57

W

- WebSMIT
 - cluster management 108, 127
- WebSphere 98
- worksheets
 - online worksheet program
 - overview 107
 - paper vs. online worksheets 107

XYZ

- X.25
 - configuring X.25 communication links 39