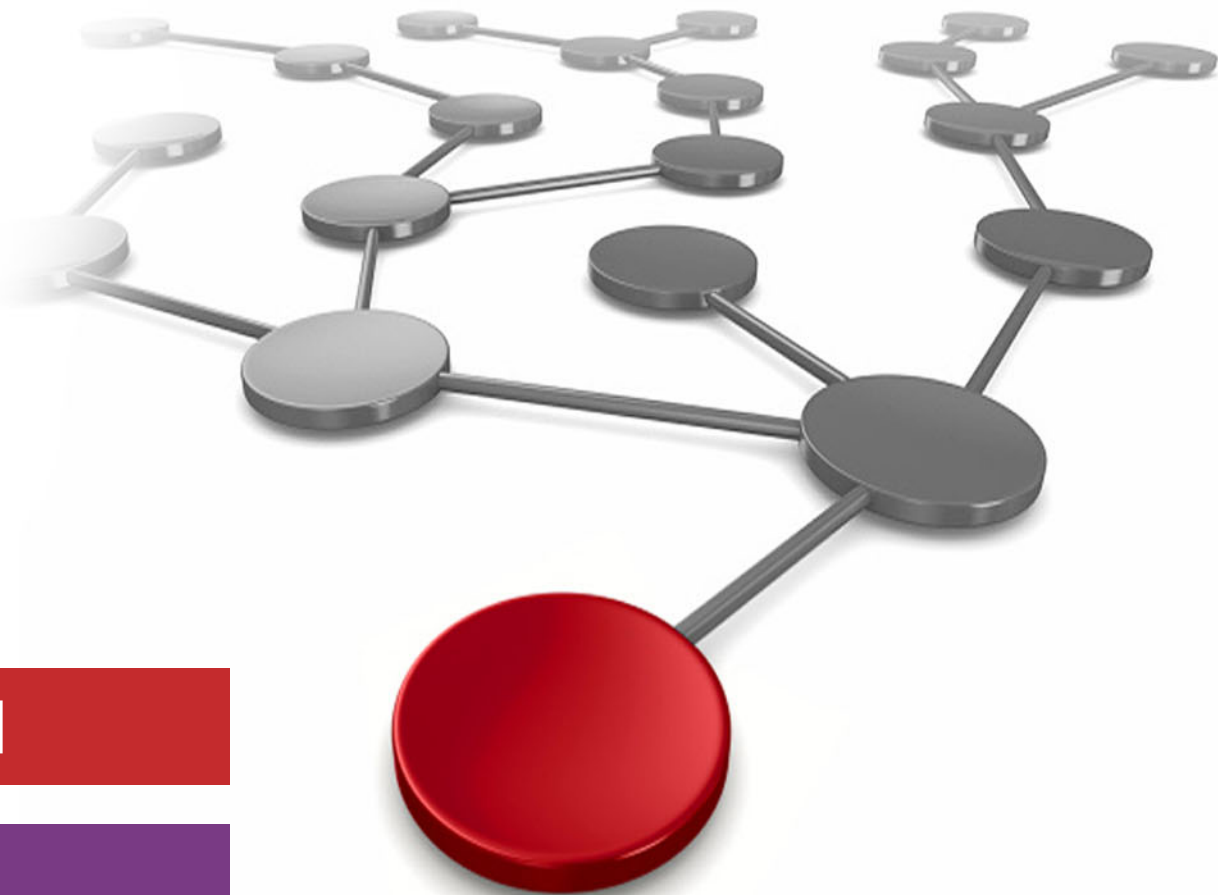


Using Pacemaker to Create Highly Available Linux Solutions on IBM Power

Tim Simon
Dishant Doriwala
Hemantha Gunasinghe
Munshi Hafizul Haque
Santhosh Joshi
Adhish Kapoor
Pravin Kedia
Youssef Largou
Henry Vo



 Cloud

Linux



IBM Redbooks

**Using Pacemaker to Create Highly Available Linux
Solutions on IBM Power**

December 2023

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

First Edition (December 2023)

This edition applies to the following software and hardware versions:

Red Hat Enterprise Linux 9.0 for IBM Power Little Endian

Red Hat Enterprise Linux 8.2 for IBM Power Little Endian

SUSE Linux Enterprise Server 15 SP1

DB2 Version 11.5 FP8

IBM MQ 9.2 LTS

IBM MQ 9.2 LT3

IBM VM Recovery Manager Version 1.7

IBM PowerVC 2.1.0

IBM Power 10

IBM Power 9

This document was created or updated on December 20, 2023.

© Copyright International Business Machines Corporation 2023. All rights reserved.

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

Contents

Notices	ix
Trademarks	x
Preface	xi
Authors	xi
Now you can become a published author, too!	xiii
Comments welcome	xiv
Stay connected to IBM Redbooks	xiv
Chapter 1. Introduction to High Availability	1
1.1 Overview of High Availability Concepts	2
1.2 High availability versus disaster recovery	2
1.2.1 High availability	3
1.2.2 Disaster recovery	3
1.2.3 Service Level Agreements	4
1.2.4 HA/DR Strategies – Eliminating Single Points of Failure	5
1.3 High availability and disaster recovery solutions for Power systems	6
1.4 High availability and disaster recovery options for Linux on IBM Power	8
1.4.1 VM recovery tools	8
1.4.2 Cluster solutions	10
1.5 Clustering implementations	11
1.5.1 Active/Passive clustering	12
1.5.2 Active/Active clustering	14
1.5.3 N +1 clustering	16
1.5.4 N-to-1 Clustering	17
1.5.5 N+N clustering	18
1.5.6 Three site clusters	19
1.5.7 Linux Pacemaker	22
1.5.8 High Availability with Red Hat Enterprise Linux Add-On	22
1.5.9 SUSE Linux Enterprise High Availability Extension	23
1.6 Linux Pacemaker vs VM Recovery Manager HA	23
1.7 Software licensing challenges in HA setups	24
Chapter 2. Pacemaker Overview	25
2.1 Architecture	26
2.1.1 Cluster Architecture	27
2.1.2 Pacemaker architecture and components	28
2.2 Pacemaker clustering options	29
2.2.1 Single site	30
2.2.2 Two site	31
2.2.3 Multi site	32
2.3 Failover policies	33
Chapter 3. Planning and Installing Your Pacemaker Environment	35
3.1 Prerequisites and requirements	36
3.2 Installing and Configuring Pacemaker environment	38
Chapter 4. High Availability for Your IBM Power Virtual Server Environment	43
4.1 Introduction to Power Virtual Server offering	44

4.1.1	IBM Power Virtual Server Compute Hardware Building Blocks	45
4.1.2	IBM Power Virtual Server Network Hardware Building Blocks	45
4.1.3	IBM Power Virtual Server Storage Hardware Building Blocks	46
4.2	IBM Power Virtual Server High Availability	46
4.3	IBM Power Virtual Server Networking Options	48
4.3.1	IBM Power Systems Virtualization	48
4.3.2	Cisco ACI Networking	55
4.3.3	External Networking	55
4.4	IBM Power Virtual Server storage options	58
4.4.1	GFS2	58
4.4.2	IBM Storage Scale	58
4.4.3	Global Replication Services (GRS)	59
4.4.4	Configuration of High Availability Storage with Pacemaker	59
Chapter 5. High Availability for Database Servers		63
5.1	Db2 database introduction	64
5.1.1	Db2 database concepts	64
5.1.2	Db2 high availability and disaster recovery	64
5.1.3	Db2 databases in High Availability deployment with HADR	65
5.2	Db2 High Level deployment Architecture with Pacemaker	66
5.3	Using Pacemaker to create highly available Db2 databases	67
5.3.1	Pacemaker core: foundational components & key elements	67
5.3.2	The Quorum device support on Pacemaker	69
5.4	Prerequisites for an integrated solution using Pacemaker	70
5.5	Overview of Installing Db2 on Red Hat Enterprise Linux and configuration of Pacemaker	71
5.6	Using the Db2 cluster manager utility for high availability	72
5.6.1	The HADR option for Pacemaker-managed Db2 instances	72
5.6.2	The Mutual Failover option for Pacemaker-managed Db2 instances	76
Chapter 6. Using Pacemaker for Highly Available SAP Landscapes		81
6.1	Introduction to SAP systems	82
6.1.1	High availability concepts for SAP solutions	83
6.2	Using Pacemaker for S/4HANA or NetWeaver application resources	85
6.2.1	HA Solution for NetWeaver or S/4 Based on ABAP Platform 1709 or older	86
6.2.2	HA Solution for S/4HANA based on ABAP Platform 1809 or newer	87
6.2.3	Configuring high availability for SAP S/4HANA (ASCS and ERS) in a RHEL HA Add-On cluster	88
6.3	Using Pacemaker to create highly available SAP HANA database instances	109
6.3.1	Performance Optimized	111
6.3.2	Cost Optimized	111
6.3.3	Active/Active (Read Enabled)	112
6.3.4	Multi-tier System Replication	113
6.3.5	Multi-target System Replication	113
6.3.6	Configuring SAP HANA Scale-Up System Replication in a RHEL HA Add-On cluster	114
Chapter 7. Additional Use Cases for Pacemaker High Availability		125
7.1	IBM MQ Replicated Data Queue Manager	126
7.2	Using VM Recovery Manager to Replicate Pacemaker cluster	132
7.2.1	VM Recovery Manager DR Overview and Architecture	132
7.2.2	Installation and Configuration requirements	134
7.2.3	Disaster recovery mechanism for the VM recovery manager DR solution	135
7.3	NSF Server High Availability	138

7.4 IBM MQ active/passive architecture using Pacemaker and NFS on IBM Power Virtual Server.....	139
7.5 IBM PowerVC three-node architecture using Pacemaker.....	140
Related publications	141
IBM Redbooks	141
Online resources	141
Help from IBM	142

Notices

This information was developed for products and services offered in the US. This material might be available from IBM in other languages. However, you may be required to own a copy of the product or product version in that language in order to access it.

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

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

IBM Director of Licensing, IBM Corporation, North Castle Drive, MD-NC119, Armonk, NY 10504-1785, US

INTERNATIONAL BUSINESS MACHINES CORPORATION PROVIDES THIS PUBLICATION "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some jurisdictions do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

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

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

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

The performance data and client examples cited are presented for illustrative purposes only. Actual performance results may vary depending on specific configurations and operating conditions.

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.

Statements regarding IBM's future direction or intent are subject to change or withdrawal without notice, and represent goals and objectives only.

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 actual people or business enterprises is entirely coincidental.


COPYRIGHT LICENSE:

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

Trademarks

IBM, the IBM logo, and ibm.com are trademarks or registered trademarks of International Business Machines Corporation, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the web at “Copyright and trademark information” at <https://www.ibm.com/legal/copytrade.shtml>

The following terms are trademarks or registered trademarks of International Business Machines Corporation, and might also be trademarks or registered trademarks in other countries.

AIX®	IBM FlashSystem®	pureScale®
Cognos®	IBM Spectrum®	Redbooks®
DB2®	Parallel Sysplex®	Redbooks (logo)  ®
Db2®	Passport Advantage®	Storwize®
FlashCopy®	POWER®	System z®
GDPS®	Power8®	SystemMirror®
IBM®	Power9®	WebSphere®
IBM Cloud®	PowerHA®	XIV®
IBM Cloud Pak®	PowerVM®	

The following terms are trademarks of other companies:

The registered trademark Linux® is used pursuant to a sublicense from the Linux Foundation, the exclusive licensee of Linus Torvalds, owner of the mark on a worldwide basis.

Java, and all Java-based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates.

OpenShift, Red Hat, are trademarks or registered trademarks of Red Hat, Inc. or its subsidiaries in the United States and other countries.

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

Preface

High availability and business continuity continue to be focus themes for businesses as they work hard to provide infrastructure components to support the “always available” requirements for their critical business functions. In addition, many enterprises are starting to build critical business functions on Linux servers to take advantage of the openness and portability of those applications across both on-premises and cloud-based infrastructures. These combined requirements lead to the need for a solution that allows enterprises to support highly available Linux environments for critical business applications.

This IBM Redbooks publication will help you install, tailor and configure a highly available cluster of IBM Power servers utilizing Pacemaker, which is an open source cluster manager that supports Linux environments. We describe basic concepts required for planning and designing an infrastructure to meet the requirements for high availability and business continuity and we compare and contrast different options and tools. While the publication is mostly targeted at skilled professionals who are designing and installing highly available Linux environments, it starts with a basic understanding of the technology to help less technical readers understand the concepts.

Authors

This book was produced by a team of specialists from around the world working at IBM Redbooks, Poughkeepsie Center.

Tim Simon is an IBM Redbooks® Project Leader in Tulsa, Oklahoma, USA. He has over 40 years of experience with IBM® primarily in a technical sales role working with customers to help them create IBM solutions to solve their business problems. He holds a BS degree in Math from Towson University in Maryland. He has worked with many IBM products and has extensive experience creating customer solutions using IBM Power, IBM Storage, and IBM System z® throughout his career. He currently lives in Tulsa, Oklahoma where he enjoys spending time with his grandchildren.

Dishant Doriwala is a Senior Staff Software engineer, Test Lead for VM Recovery Manager (HA and DR solution for Power Systems). He works in IBM Systems Development Labs (ISDL), Hyderabad, India and has 11 years of experience in industry with expertise in testing for HA-DR products. He has experience on working various HA and DR solutions: IBM PowerHA System Mirror, Reliable Scalable Cluster Technology, VM Recovery Manager and GLVM. He has expertise on enterprise storage: IBM San Volume controller/Storwize®, Hitachi, Dell EMC SRDF and Unity XT, IBM XIV®. He has authored technical white paper, red paper and conference articles in HA-DR domains. He holds a Master's degree in Software technologies from International Institute of information technology, Pune and Bachelors in Electronics and Communications from VTU, Bangalore.

Hemantha Gunasinghe is an IT Specialist in IBM Systems Expert Labs within IBM Australia. He brings a unique experience and perspective with over 32 years of experience in IT industry in Sri Lanka and Australia. Before migrating to Australia in 1999, Hemantha was working as a Systems Engineer at IBM Sri Lanka. Hemantha joined IBM Australia as a Service Delivery administrator in IBM Global Services for IBM POWER® systems where he provided many managed infrastructure services to IBM Internal systems as well as for commercial customers at various capacities in the area of IBM AIX®, IBM PowerHA®, GPFS, Linux and IBM PowerVM® Virtualization technologies on IBM POWER systems. He joined IBM LAB Services in 2013 as an

IT Management Consultant and has delivered many IT projects and services to many IBM clients on IBM POWER systems. Also, he has provided implementation and migration services to IBM Power VS in IBM Cloud®. In recent times, he has delivered some of the largest SAP HANA systems running on Linux on Power using Red Hat Operating systems and also delivered OpenShift clusters deployed on Power Systems. His services expand beyond Australia into other countries such as New Zealand and US where some of the largest government, defense, commercial and financial institutions are involved with implementation and build services for such clients. He holds a degree in Mathematics, statistics and Computing from University of Kelaniya in Sri Lanka with a 1st class honours and further holds a Master of Computer Science from University of Colombo in Sri Lanka.

Munshi Hafizul Haque is a Senior Platform Consultant at Red Hat in Kuala Lumpur, Malaysia. Munshi is an experienced technologist in engineering, design, and architecture of PaaS and cloud infrastructures. He has over 16 years of experience in post-delivery technical support, designing full infrastructure solutions, performing installation, implementation and integration services of systems, storages and software solutions, and conducting post implementation assessment, presentations at customer engagements and corporate events. He is currently part of the Red Hat Consulting Services team where he helps organizations adopt automation, container technology and DevOps practices. Before that, he worked for IBM as a senior consultant with IBM Systems Lab Services in Petaling Jaya, Malaysia, where he took part in various projects with different people in different ASEAN countries, and as a specialist in IBM Power Systems and associated enterprise edition technology.

Santhosh Joshi is a Senior Staff Software Engineer in IBM India Systems Development Lab, IBM India. He has over 18 years of experience in the Information Technology field. Currently he works for IBM PowerVM Live Partition Mobility development team and prior to this he worked for IBM VM Recovery manager HA and DR for Power Systems solution development. He holds Bachelor of Engineering degree in Electronics and Communication from Visvesvaraya Technological University, Belagavi India. His areas of expertise include PowerVM virtualization, High availability and disaster recovery solutions, clustering technologies, SAN.

Adhish Kapoor serves as a Team Lead for VMRM Recovery Manager, overseeing the HA/DR solution for Power Systems. With over 20 years of experience in the system domain, he specializes in High Availability/Disaster Recovery solutions for intricate cluster environments, as well as in the development of backup applications and SAN/NAS storage appliances. Adhish has authored technical white papers and conference articles in the HA-DR domains. He earned his Bachelor of Computer Science Engineering degree from Punjab Technical University.

Pravin Kedia is a CTO for Data & AI Expert Labs running the Centre of Excellence in India. Pravin has 25+ years of experience in CP4D Data Fabric, Data Science and DataOps Solutions, Big Data, Data Warehouse, Replication, Data Science and AI Architect with international work experience in Banking, Telecommunication, Insurance and Financial Markets domain across WW customers. He holds a degree in MBA (Finance) from University of Kansas and Bachelors (Electronics and Telecommunications) from VESIT, Mumbai university. He has written extensively on Database and Replication technologies. Pravin regularly conducts Data Strategy, Data Discovery and Design workshops with customer C-level executives. As CTO and solutions Architect for IBM Data and AI solutions, Pravin is responsible for technical Delivery for IBM Expert Lab Services Projects across the Hybrid Cloud Product portfolio. Pravin is a PMP since 2005 and is a Master Inventor (Plateau 4) and on the IDT for Data and Governance teams. Pravin is an Open Group Certified Executive IT Specialist Thought Leader and Certified ITS Profession Champion. Pravin has received numerous rewards and led/participated in AOT and Cloud Studies on Data and Governance.

Youssef Largou is the founding director of PowerM, a platinum IBM Business Partner in Morocco. He has 22 years of experience in systems, HPC, middleware, and hybrid cloud, including IBM Power, IBM Storage, IBM Spectrum®, IBM WebSphere®, IBM Db2®, IBM Cognos®, IBM WebSphere Portal, IBM MQ, ESB, IBM Cloud Pak®, SAP HANA and Red Hat OpenShift. He has worked within numerous industries with many technologies. Youssef is an IBM Champion 2020, 2021, 2022 and 2023, an IBM Redbooks Platinum Author, and has designed many reference architectures. He has been recognized as an IBM Beacon Award Finalist in Storage, Software-Defined Storage, and LinuxONE five times. He holds an engineer degree in Computer Science from the Ecole Nationale Supérieure des Mines de Rabat and Excecutive MBA from EMLyon.

Henry Vo is an IBM Redbooks Project Leader with 10 years experience in IBM. He has technical expertise in business problem solving, risk/root-cause analysis, and writing technical plans for business. He has held multiple roles at IBM including Project management, ST/FT/ETE Test, Back End Developer, DOL agent for NY. He is a certified IBM zOS Mainframe Practitioner including IBM Z System programming, Agile, and Telecommunication Development Jumpstart. Henry holds a Master of MIS (Management Information System) from the University of Texas in Dallas.

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

Sanjeev Koul, Data, Performance & Cloud Architect Expert Labs
IBM Bangalore, India

Now you can become a published author, too!

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

Find out more about the residency program, browse the residency index, and apply online at:

ibm.com/redbooks/residencies.html

Comments welcome

Your comments are important to us!

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

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

ibm.com/redbooks

- ▶ Send your comments in an email to:

redbooks@us.ibm.com

- ▶ Mail your comments to:

IBM Corporation, IBM Redbooks
Dept. HYTD Mail Station P099
2455 South Road
Poughkeepsie, NY 12601-5400

Stay connected to IBM Redbooks

- ▶ Find us on LinkedIn:

<https://www.linkedin.com/groups/2130806>

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

<https://www.redbooks.ibm.com/subscribe>

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

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



Introduction to High Availability

In an age where digital technologies underpin every facet of our lives, from critical business processes to personal communications, one aspect stands as paramount: uninterrupted access to data, applications, and services. The modern world demands nothing less than a seamless, always-on experience, and any interruption – be it due to hardware failures, software glitches, or unforeseen disasters – is met with frustration and loss. It is in this landscape that the concept of High Availability (HA) emerges as a linchpin of resilience and continuity.

In this chapter, we describe the realm of High Availability, exploring its fundamental principles, implementation strategies, and the diverse technologies that empower it. We delve into the critical considerations that guide organizations in their pursuit of uninterrupted operations, from the tolerance for downtime to budget constraints.

In this chapter, the following topics are described:

- ▶ 1.1, “Overview of High Availability Concepts” on page 2
- ▶ 1.2, “High availability versus disaster recovery” on page 2
- ▶ 1.3, “High availability and disaster recovery solutions for Power systems” on page 6
- ▶ 1.4, “High availability and disaster recovery options for Linux on IBM Power” on page 8
- ▶ 1.5, “Clustering implementations” on page 11
- ▶ 1.6, “Linux Pacemaker vs VM Recovery Manager HA” on page 23
- ▶ 1.7, “Software licensing challenges in HA setups” on page 24

1.1 Overview of High Availability Concepts

The cost of IT downtime within the industry is influenced by various factors. Financial losses can fluctuate depending on factors such as revenue, the sector of operation, the duration of the outage, the number of individuals affected, the time of day, and more. To illustrate, businesses engaged in high-level data transactions, such as banks and online retailers, experience significantly greater losses per hour. If one of these businesses encounters an unexpected outage during peak traffic periods, the resulting damage will naturally be more substantial.

High Availability (HA) is the elimination of single points of failure to enable applications or workloads to operate continuously even if one or more software or hardware components, such as server, or network components fails.

When high availability measures fall short, disaster recovery steps in. Typically disaster recovery is primarily reserved for large-scale failures, such as a complete site failure or data corruption due to cyber attacks. Often disaster recovery involves some amount of manual processes, due to the high stakes involved. It may involve straightforward actions like restoring from a backup or executing a site failover, but it can also become highly intricate; involving actions like analyzing the database log to ensure logical consistency and making individual decisions on whether to apply specific database operations or not.

1.2 High availability versus disaster recovery

While the terms high availability and disaster recovery are often used interchangeably, they are actually two different concepts, with significant differences. This section describes both concepts and shows how they are different in scope and implementation.

Disaster recovery typically involves an unforeseen incident necessitating the restoration of either an application or the entire computing infrastructure at a remote location. In disaster recovery scenarios, the expected recovery times are typically measured in hours or even days, rather than minutes.

In contrast, high availability (HA) can address both planned and unplanned downtime. HA primarily involves automated, rapid intra-site failover procedures. Its core objective is the elimination of single points of failure, which necessitates redundant components. Most systems implement three types of redundancy: hardware, software, and environmental.

For planned downtime, such as hardware or software maintenance, coordination with end users is performed, and the downtime may last for minutes or hours. With an HA setup which features a primary and secondary database server, for instance, any downtime caused by a firmware upgrade on the primary database server can be minimized by switching over to the secondary server.

Near Zero Downtime planning requires a deliberately designed high availability infrastructure strategy. In cases of unplanned downtime – such as a malfunctioning hardware component or a corrupted database – an automatic failover to a secondary server takes place with minimal impact to end users.

The design of high availability within a server's components, the infrastructure architecture, and software all play crucial roles in reducing unplanned downtime to a minimum. Implementing HA in an infrastructure comes with a cost. It is essential to weigh this cost against the expenses associated with downtime to justify the expense of preventing or minimizing such disruptions.

1.2.1 High availability

A high availability system or component is specifically designed to be continuously operational for long periods of time. The HA system is built to eliminate any single point of failure by providing redundant components in the software, hardware, and applications.

Hardware redundancy means that the hardware design ensures that all critical hardware components in the system have one or more backups, ensuring that a failure of an individual component does not interrupt the operation of the system. This includes server hardware, networking components and storage. Hardware redundancy is critical for modern business data systems.

Highly available hardware is only one part of an HA system, software redundancy is also required to ensure that the business applications can tolerate the failure of the system where they are running. This is usually done using clustering technologies that spread the database and application code across multiple redundant servers. If an outage of short duration can be tolerated or if your software does not support clustering capability, a simpler infrastructure utilizing VM Restart based technology can be used to quickly restart the business applications on backup hardware in a quick and efficient manner.

1.2.2 Disaster recovery

There are some situations that can not be addressed by an HA system. These situations are generally the result of some incident that causes a full site failure such as fire, flooding or other natural disaster that would require the movement of your business application to another physical location, usually to another data center in a different geographic region. Disaster recovery (DR) is the process which is used when your high availability solution cannot continue to operate.

Disaster recovery involves a set of policies and procedure to enable recovery of your business applications on an infrastructure in a different location which is not affected by whatever disaster forces you to move operations. This must include planning for not only servers, storage, and local networking components in the new location, but also needs to consider user access to the new systems and personnel to operate the systems in the new location. This is critical to be able to continue your business operations in case of a site failure, natural disaster, or planned outage to your data center.

When you are planning and setting up your disaster recovery location it is important to understand your business requirements. The cost and complexity of your disaster recovery plan is dependent on a couple of important objectives that you need to be able to assign to your different business applications. These two objectives are Recovery Time Objective (RTO) and Recovery Point Objective (RPO).

Recovery Time Objective

The Recovery Time Objective describes the maximum amount of time that a system can be down. This is the length of time from when the outage starts until the application or function is completely restored to operation.

The RTO should be determined for each business function, based on the importance of that function to your business. For some business functions, RTO can be measured in hours or even days, but for critical functions the RTO could be measured in minutes or even seconds.

Recovery Point Objective

The Recovery Point Objective (RPO) defines the amount of data loss – measured in time – your business application can sustain. The RPO requirement is extremely critical to understand for each business function, as it directly affects the design and cost of your disaster recovery solution.

Disaster recovery solutions by definition, require replication of data between two locations that are geographically separated, sometime this separation is thousand of miles to ensure that both locations are not affected by the same event. The smaller your RTO, the more complex your disaster recovery solution will need to be in order to achieve that goal, and this extra complexity leads to increased costs.

High availability and disaster recovery are targeted to solve the same problem – ensuring your system is up and running and limiting loss of data caused by any failure. However you can characterize the differences between HA and DR as:

- ▶ HA is intended to handle issues while system is running:
 - Redundant components are used to recover from specific hardware or software failures.
 - The time required to recover your application depends on your RPO/RTO requirements. Using cluster technology, recovery can occur with little or no user impact. For those applications with less restrictive RPO/RTO requirements VM recovery technologies can restart the application in minutes on other available hardware.
 - Software licensing costs need to be considered – you may need to acquire additional licenses depending on the specific supplier.
- ▶ DR is used to handle problems caused by system failure or infrastructure failure that can not be addressed by your HA solution.
 - Recovery is made to a different physical location, often geographically distant from your primary location. There will almost always be an outage required – the length of the outage is dependent on many factors.
 - The DR location hardware can often be used for other purposes (testing and development for example) during normal operations and during a DR situation those functions will be displaced by production workloads until the primary site failure can be resolved.
 - Recovery back to production generally requires an additional outage.

1.2.3 Service Level Agreements

Service Level Agreements (SLAs) represent contractual agreements between the service provider and the service user, typically entailing penalties for the provider if SLA terms are breached, often in the form of financial compensation. SLAs primarily pertain to instances of unplanned service disruptions, as opposed to planned ones.

The higher the guaranteed percentage of uptime specified in the SLA, the less time the service is allowed to be unavailable. Each additional “9” following the decimal point signifies a substantial reduction in potential downtime. For example, transitioning from 3 “9s” (99.9%) to 5 “9s” (99.999%) represents a variance of 8 hours, 40 minutes, and 41 seconds in downtime. For service providers, adhering to higher SLA availability levels necessitates significant investments in and architectural enhancements to their IT infrastructure to enhance resilience against failures. Resilience manifests in various forms.

For instance, servers like IBM Power servers incorporate considerable built-in redundancy, enabling them to claim a hardware availability of five nines (99.999%). Similarly, IBM FlashSystem® storage boasts six nines (99.9999%) availability. However, if the switch responsible for connecting the fiber cables between the server and storage only guarantees

three nines (99.9%) availability, it becomes the weakest link, determining the overall infrastructure's availability. To address this, many clients opt to install a secondary switch along an entirely distinct fiber channel path connecting the server and storage.

Availability is defined by the SLA, data needs to be gathered on the systems to show when a system is available and the length of any unplanned outages to determine if the SLA has been met. Table 1-1 maps the amount of outage time allowed per year for a system to meet each availability level (from one nine to six nines). As can be seen by the values in the table, achieving availability numbers of four nines or higher requires a high level of automation and planning.

Table 1-1 Typical service-level agreements

Number of nines	% Uptime	Maximum Annual downtime
Six (6)	99.9999	32 seconds
Five (5)	99.999	5 minutes 16 seconds
Four (4)	99.99	52 minutes 36 seconds
Three (3)	99.9	8 hours 45 minutes 57 seconds
Two (2)	99.0	3 days 15 hours 39 minutes 29 seconds
One (1)	90.0	36.5 days

1.2.4 HA/DR Strategies – Eliminating Single Points of Failure

Every element within an IT infrastructure carries a certain, albeit low, probability of experiencing a failure. Consequently, the key to achieving high availability revolves around redundancy – encompassing hardware redundancy, network redundancy, and redundancy at the data center level.

Regarding data center redundancy, the depicted configuration in Figure 1-1 showcases a three-site architecture. It comprises two on-premise sites managed by the client and an additional failover site housing servers and storage systems within the IBM Cloud. Beyond replicating the data centers, it's crucial to address potential single points of failure at the site level, such as power supplies or network connections.

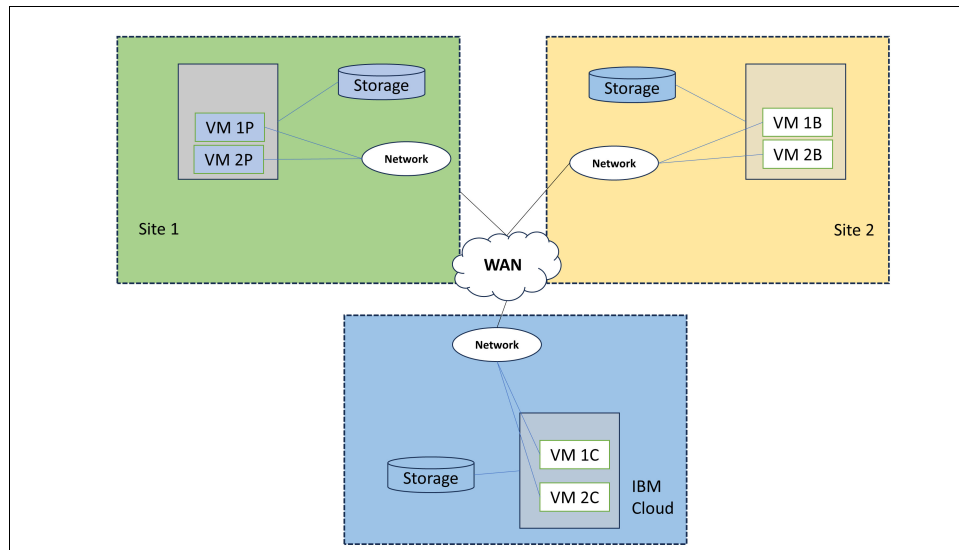


Figure 1-1 HA and DR configuration example

In addition to redundant hardware and redundant network connections, you also need to plan for redundant power. For example, using an Uninterrupted Power Supply device (UPS) can effectively balance power supplied from two different vendors, backed by an onsite battery backup and emergency generators.

In terms of hardware redundancy, IBM Power offers multiple layers of redundant server components. These include redundant power supplies, fans, enterprise-grade error-correcting memory, fully redundant network adapters, and spare CPUs that can be activated in the event of a failure or for load balancing purposes. Disk storage systems employ striping and mirroring techniques to ensure redundancy, facilitating automatic recovery from disk failures. All network components, including switches, are designed with full redundancy in mind.

In scenarios where an entire server or storage system encounters a failure, dedicated standby hosts are employed for seamless failover. Leveraging the capabilities of PowerVM, these standby hosts don't remain idle but can serve as development or test LPARs while simultaneously replicating production data from the primary server. In the event of a failover, the development and test LPARs are paused, allowing the production workload to seamlessly continue running on the standby hosts, virtually unaffected by the incident.

1.3 High availability and disaster recovery solutions for Power systems

IBM Power systems are one of the most reliable platforms in the industry and are designed to meet the requirements of the most critical data-intensive workloads. For the 15th straight year, IBM Power servers achieved the highest server reliability ranking in the ITIC 2023 Global Server Hardware and Server OS Reliability survey¹.

IBM Power and IBM Power Virtual Server – an IBM Cloud offering based on IBM Power systems – provide an entire product line of HA and DR solutions. Figure 1-2 provides an overview of the IBM Power HADR built in capabilities on IBM Power systems.

On the left half of the family tree shown in the figure are the IBM PowerHA SystemMirror® shared storage clustering solutions, and on the right side of the family tree is the IBM VM Recovery Manager (VMMR).

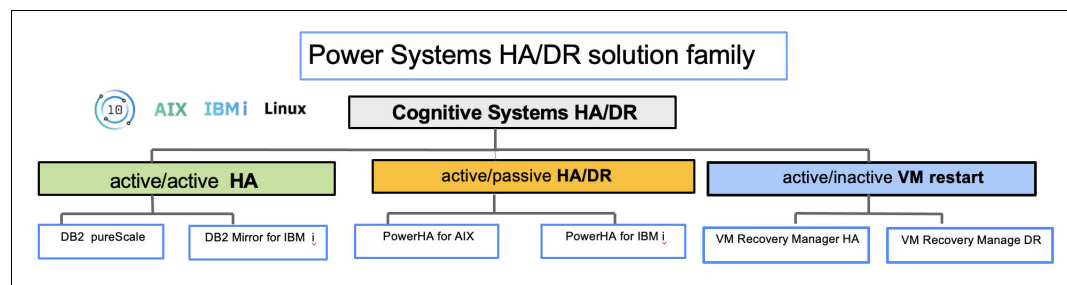


Figure 1-2 IBM Power Systems virtual server HADR solution family.

The PowerHA SystemMirror family of solutions is optimized for mission-critical applications where the total annual downtime for both planned and unplanned outages is zero or near-zero. The PowerHA System mirror product line covers all outage types, both software and hardware. There is at least one active OS on each of the nodes in the cluster, which enables the capability of doing software updates on a system other than the production node.

¹ <https://itic-corp.com/itic-reports-surveys/2023-itic-reports-surveys/>

PowerHA SystemMirror covers both data center and multi-site configurations. If you are looking to drive total outage time for both planned and unplanned events to near zero, this solution is the one that you should deploy. For more information refer to the [IBM PowerHA SystemMirror documentation](#).

VMRM solutions can be best understood by first understanding Live Partition Mobility (LPM). A set of logical partitions (LPARs) is virtualized by using IBM PowerVM and Virtual I/O Server (VIOS) to enable a partition to be moved for a firmware or hardware maintenance event by using LPM. If that virtual machine (VM) fails, it can be restarted on another server in the cluster. For DR operations, those same VMs are replicated by using storage area network (SAN) storage at the secondary location. For more information refer to this [IBM document](#).

The active-active solutions are IBM DB2® pureScale®, which supports AIX, and IBM DB2 Mirror for IBM i. In both cases, the solutions are classified as active-active, but they are achieved by different approaches. Db2 pureScale provides an active/active solution using a shared DB2 cluster configuration with distributed lock management to enable multiple application servers to simultaneously access the shared database. In Db2 pureScale, the active-active capability is first at the DB2 layer, where there are two copies that are replicated bi-directionally and synchronously. The application servers can be configured active-active or active-passive. Table 1-2 highlights differences between IBM Power Systems HA/DR solutions.

You can find more details in *IBM Power Systems High Availability and Disaster Recovery Updates: Planning for a Multicloud Environment*, REDP-5663.

Table 1-2 High Availability Topology Classification

Technology	Active/Active Clustering	Active/Passive Clustering	Active/Inactive Clustering
Definition	Application clustering; applications in the cluster have simultaneous access to the production data therefore no app restart upon an app node outage. Certain types enable read-only access from secondary nodes	OS clustering; one OS in the cluster has access to the production data, multiple active OS instances on all nodes in the cluster. Application is restarted on a secondary node upon outage of a production node.	VM Clustering, One VM in a cluster pair has access to the data, one logical OS, two physical copies. OS and applications must be restarted on a secondary node upon a primary node outage event. LPM enables the VM to be moved non-disruptively for a planned outage event.
Outage Types	SW,HW,HA, planned, unplanned RTO 0, limited distance	SW,HW,HA,DR, planned, unplanned, RTO>0, multi-site	HW,HA,DR, planned, unplanned, RTO>0, multi-site
OS integration	Inside the OS	Inside the OS	OS agnostic
RPO	Sync mode only	Sync/Async	Sync/Async
RTO	Zero	Fast (minutes)	Fast Enough (VM Reboot)
Licensing	N+N ^a	N+1 licensing	N+0 licensing
IBM Solution	DB2 pureScale, DB2 Mirror	PowerHA, Red Hat HA, Linux HA	VMR HA, LPM, VMR DR

a. N = number of licensed processor cores on each system in the cluster.

1.4 High availability and disaster recovery options for Linux on IBM Power

In today's digital landscape, where business operations are increasingly reliant on technology, ensuring the availability and resilience of critical systems is paramount. For organizations leveraging Linux on IBM Power systems, the pursuit of high availability and disaster recovery options takes center stage. This section explores the vital strategies and tools available to safeguard Linux-based workloads on IBM Power, delving into the key concepts, components, and real-world applications that underpin these essential practices.

1.4.1 VM recovery tools

The following are the key VM recovery tools available for Linux on IBM Power Systems:

- ▶ **PowerVM** is IBM's virtualization technology for IBM Power Systems. It provides the foundation for creating and managing virtualized environments, including virtual machines (LPARs - Logical Partitions).
- ▶ **Live Partition mobility** is a feature included in the PowerVM Enterprise Edition hardware offering, enables the relocation of logical partitions (LPARs) running AIX, IBM i, and Linux from one system to another. This migration process encompasses the entirety of the system environment, encompassing processor status, memory, attached virtual devices, and active user connections.
 - With active partition migration, also known as Live Partition Mobility, you have the capability to move AIX, IBM i, and Linux logical partitions, including their running operating systems and applications, from one system to another. Importantly, this migration can be executed without the need to halt or shut down the logical partition or the applications it hosts.
 - In the case of inactive partition migration, also known as cold partition mobility, you can transfer a powered-off AIX, IBM i, or Linux logical partition from one system to another.
 - To facilitate these migrations, the Hardware Management Console (HMC) serves as a valuable tool, allowing you to efficiently move both active and inactive logical partitions from one server to another.
 - The HMC is programmed to migrate the most recently activated profile. Consequently, an inactive logical partition that has never been activated cannot undergo migration. For cases of inactive partition mobility, you have the option to either choose the partition's state as defined in the hypervisor or opt for the configuration data specified in the last activated profile on the source server.
 - Partition mobility offers enhanced systems management flexibility and is explicitly engineered to enhance system uptime with the following advantages:
 - Planned outages for hardware or firmware maintenance can be circumvented by migrating logical partitions to an alternative server before conducting the necessary maintenance. Partition mobility proves invaluable in these situations as it permits the seamless adjustment of scheduled maintenance activities.
 - Downtime related to server upgrades can be entirely avoided by migrating logical partitions to an alternate server before proceeding with the upgrade process. This allows for uninterrupted operations.
 - In anticipation of potential server failures, logical partitions can be preemptively migrated to another server, thus mitigating the risk of unscheduled outages. Partition mobility serves as a preemptive measure against such incidents.

- The consolidation of workloads, typically dispersed across several underutilized servers, onto a single, more robust server can be achieved through partition mobility.
- Workloads can be dynamically relocated between servers to optimize resource utilization and workload performance within the computing environment. With active partition mobility, this can be achieved with minimal disruption to ongoing operations.

Note: It is essential to note that partition mobility, while delivering numerous advantages, does not encompass the following functionalities:

- Automatic workload balancing is not a feature of partition mobility.
- Partition mobility does not serve as a conduit for introducing new functionalities. To harness new features, logical partitions may require restarts and, in some cases, re-installation.

- ▶ The VM Recovery Manager HA solution implements recovery of the virtual machines based on the VM restart technology. The VM restart technology relies on an out-of-band monitoring and management component that restarts the VMs on another server when the host infrastructure fails. The VM restart technology is different from the conventional cluster based technology that deploys redundant hardware and software components for a near real-time failover operation when a component fails. The VM Recovery Manager HA solution is ideal to ensure high availability for many VMs when the workload on those VMs can tolerate short interruptions of service while the VM is restarted. Because the VM Recovery Manager HA solution is based on the VM restart technology instead of clustering technology, it is easier to manage as clustering introducing additional complexities.

Figure 1-3 shows the architecture of the VM Recovery Manager HA solution. A set of hosts is grouped to be backup for each other. When failures are detected, VMs are relocated and restarted on other healthy hosts within the group.

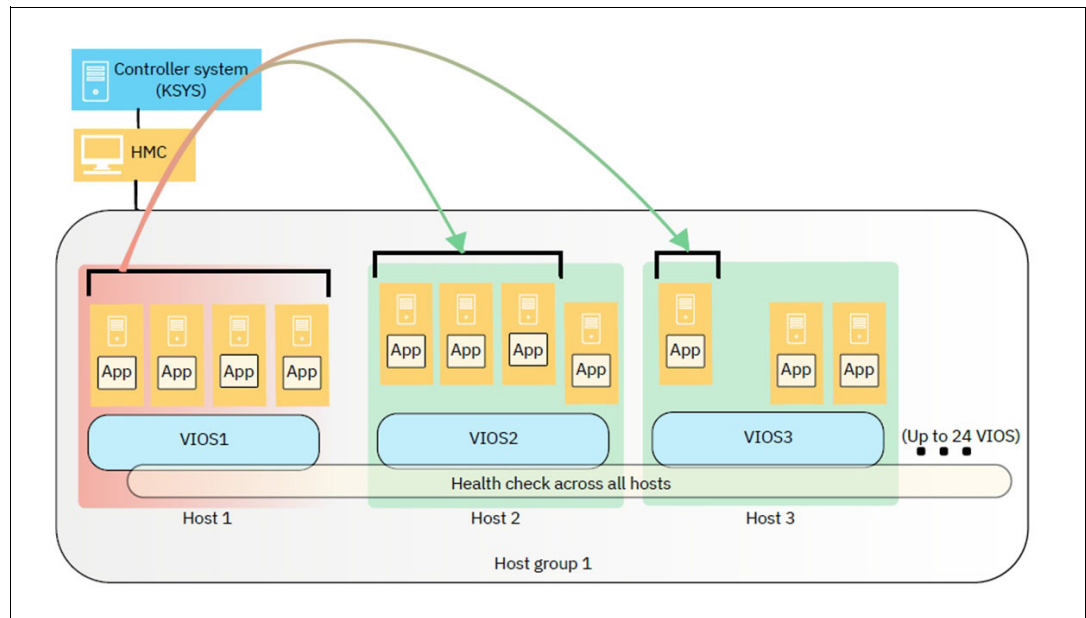


Figure 1-3 VM Recovery Manager HA solution architecture

1.4.2 Cluster solutions

Clustering is a strategic approach to ensuring the continuous availability of critical IT systems and services. It involves the use of multiple interconnected servers or nodes that work together to minimize downtime and protect against hardware failures, software glitches, or disasters.

Key concepts and components of clustering

The following are the key concepts and components of clustering:

- ▶ **Nodes:** Clustering typically involves multiple server nodes. These nodes can be physical servers or virtual machines. Each node is equipped to run the same applications or services.
- ▶ **Load Balancing:** Load balancing distributes incoming requests or workloads evenly across all nodes within the cluster. This ensures that no single node is overwhelmed with traffic, improving system performance and reliability.
- ▶ **Redundancy:** Redundancy provides replicated critical components, such as servers, storage, and network connections, to eliminate single points of failure.
- ▶ **Failover:** Failover is the ability to automatically and seamlessly switch from a failed or degraded node to a healthy one within the cluster. This ensures minimal downtime in case of a node failure.
- ▶ **Data Replication:** Data replication ensures that data is copied across multiple nodes or sites to ensure data availability and integrity in the event of a disaster.
- ▶ **Quorum:** Quorum is a mechanism to prevent “split-brain” scenarios, where multiple nodes in a cluster believe they are the active node. It ensures that there is consensus among nodes before taking certain actions, such as initiating a failover.

Types of clustering configurations

The following are common type of clustering configurations:

- ▶ **Active/Passive Clustering:** In this configuration, one node (the active node) actively serves the application or service, while the other node (the passive node) remains in standby mode. If the active node fails, the passive node takes over.
- ▶ **Active/Active Clustering:** In active/active clustering, both nodes actively serve the application or service simultaneously, sharing the workload. This configuration improves scalability and redundancy but requires careful load balancing.
- ▶ **N+1 Clustering:** N+1 clustering involves N active nodes and one standby node. If one of the active nodes fails, the standby node takes over. This configuration balances redundancy and resource utilization.
- ▶ **N+N Clustering:** N+N clustering has N active nodes running in parallel, with no standby nodes. All nodes are actively serving the application or service simultaneously. This configuration is highly redundant and scalable.
- ▶ **3-Site Clustering:** 3-site clustering extends clustering to multiple geographically separated sites (data centers). It provides disaster recovery capabilities by replicating data and services across different locations to withstand site-level failures or natural disasters.

Applications of clustering:

The following are some of the applications of clustering:

- ▶ **Server Failover:** Clustering ensures server failover in case of hardware or software failures, minimizing downtime and ensuring continuous service availability.

- ▶ **Load Balancing:** Clustering improves resource utilization and responsiveness by distributing incoming traffic or workloads across multiple nodes.
- ▶ **Disaster Recovery:** Clustering, especially in multi-site configurations, supports disaster recovery by replicating data and services across geographically dispersed locations.
- ▶ **Application Scaling:** Clustering allows for the dynamic scaling of applications and services to accommodate increased demand or resource requirements.
- ▶ **Highly Available Databases:** Database clustering provides redundancy and failover capabilities for critical database systems, ensuring data availability and integrity.

Clustering in the context of high availability and disaster recovery is a critical strategy for ensuring the resilience, redundancy, and continuous availability of IT systems and services, making it an essential component of modern business continuity planning.

1.5 Clustering implementations

Implementing clustering solutions, particularly in the context of High Availability (HA), is a strategic endeavor that requires careful consideration of various factors. One of the critical aspects to bear in mind is striking the right balance between the tolerance for downtime and the available budget. This balance ensures that the chosen HA solution aligns with the unique needs and constraints of each business function or project. When measuring or determining acceptable downtime, it's essential to take into account several parameters, as outlined below:

- ▶ **Loss of Data:** Assess the criticality of data loss to your organization. Some systems may require real-time data replication and failover to avoid even the slightest data loss, while others may tolerate minimal data loss.
- ▶ **Loss of an Application:** Consider the impact of losing access to a specific application. For some businesses, the unavailability of a critical application even for a short duration can have severe consequences.
- ▶ **Loss of a System:** Evaluate the consequences of a server or system failure. HA solutions should address the ability to recover from hardware or software failures swiftly.
- ▶ **Loss of Business Location:** For multi-site organizations, contemplate the implications of losing an entire business location due to a disaster. Clustering solutions that enable site failover can mitigate this risk.
- ▶ **Loss of Operations:** Gauge how a disruption in operations would affect your business. Some businesses can endure brief interruptions, while others demand continuous, uninterrupted operations.
- ▶ **Scheduled Maintenance:** Plan for scheduled maintenance activities. HA solutions should facilitate seamless failover or workload migration to minimize the impact of maintenance windows.
- ▶ **Impact of Downtime on User Experience:** Understand how downtime affects user experience and customer satisfaction. Ensuring minimal disruption to users is a priority for many organizations.
- ▶ **Loss of Business-Critical Applications:** Identify which applications are truly business-critical. HA solutions should prioritize these applications for rapid recovery and minimal downtime.

In practice, the choice of clustering implementation – whether it's active/passive, active/active, N+1, N+N, or a multi-site clustering setup like 3-site clustering – depends on the specific requirements of the organization. The aim is to align the HA solution with the organization's tolerance for downtime and budget constraints while safeguarding critical data, applications, and operations.

Ultimately, a robust clustering solution should be tailored to the unique needs of the business or project, taking into account the parameters mentioned above. Finding the right balance between HA requirements and budgetary considerations ensures that the chosen HA strategy optimally meets the organization's objectives for system availability and resilience.

1.5.1 Active/Passive clustering

Active/Passive clustering, also known as a failover cluster, is a high-availability configuration used to ensure the continuous availability of critical services or applications in the event of a failure. In this configuration, you have multiple servers (nodes) working together, but only one of them is actively serving the application or service at any given time, while the others remain in standby (passive) mode. The active node handles all requests and data processing, while the passive nodes are ready to take over in case the active node fails. This setup minimizes downtime and service disruption, making it an essential component of business continuity planning.

Components of Active/Passive clustering

The following are the key components of Active/Passive clustering configuration:

- ▶ **Nodes:** The servers participating in the cluster. In Active/Passive clustering, there is one active node and one or more passive nodes. The active node actively provides the service, while the passive nodes are on standby.
- ▶ **Quorum:** A mechanism to ensure that a majority of nodes agree on the cluster's state. This prevents “split-brain” scenarios where multiple nodes think they are the active node. A cluster must have a quorum before it can perform certain actions.
- ▶ **Cluster Resource Manager (CRM):** Software that manages the cluster resources and ensures they are running on the appropriate node. Pacemaker is an example of a CRM for Linux clusters.

Example of Active/Passive Clustering

Consider a scenario where a large payroll application needs to ensure high availability and handle internal user traffic as shown in Figure 1-4.

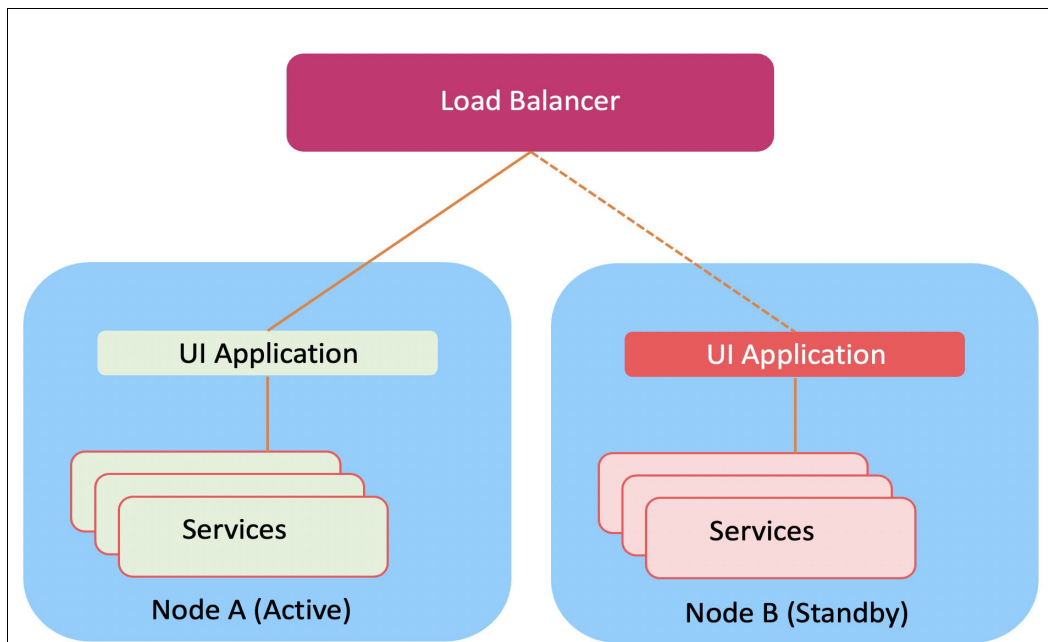


Figure 1-4 Active/Passive clustering

Scenario

Two data centers: Each data center is equipped with a server, Node A in data center 1 and Node B in data center 2.

- ▶ **Active Server (Node A):** This is the primary web server that actively serves web requests. It hosts the website, runs the web application, and handles incoming traffic.
- ▶ **Passive Server (Node B):** This is the secondary web server that remains in a standby or passive state. It mirrors the configuration and content of the active server but does not actively serve web traffic.
- ▶ **Load Balancer:** A load balancer sits in front of the active and passive servers. It distributes incoming web requests evenly between the two servers to ensure optimal resource utilization.

Configuration and Operation

- ▶ **Normal Operation (Active State):** In normal operation, Node A (the active server) receives and processes incoming web requests. It serves the website, runs the web application, and responds to user interactions.
- ▶ **Continuous Data Synchronization:** To maintain data consistency and redundancy, data, configurations, and database updates are continuously synchronized from Node A (active) to Node B (passive) using technologies like data replication or shared storage.
- ▶ **Failover Triggered by a Failure:** If Node A experiences a hardware failure, software crash, or any other issue that causes it to become unresponsive, the load balancer detects the failure.
- ▶ **Failover to Node B (Passive State):** The load balancer quickly reroutes incoming web traffic from Node A to Node B. Node B becomes the new active server, assuming the role of Node A. Users continue to access the website without interruption.
- ▶ **Notification and Remediation:** An alert or notification is sent to system administrators to investigate and address the issue with Node A. Once the issue is resolved, administrators can manually or automatically synchronize data and configurations back to Node A.
- ▶ **Load Balancer Adjustments:** When Node A is back in a healthy state, administrators may choose to manually or automatically switch back to Node A as the active server. The load balancer adjusts accordingly to maintain balanced traffic distribution.

Active/Passive clustering advantages

Active/Passive clustering offers several advantages for ensuring high availability of critical applications and services:

- ▶ **Reduced Downtime:** The primary advantage of Active/Passive clustering is the reduction of downtime. In the event of a failure on the active node, the passive node takes over almost immediately, minimizing service disruption. This ensures that critical applications remain available to users, even during hardware or software failures.
- ▶ **Improved Reliability:** Active/Passive clusters enhance the reliability of services by providing a backup node that can take over when the active node encounters problems. This redundancy increases system uptime and reduces the risk of data loss or service outages.
- ▶ **Simplified Maintenance:** Planned maintenance, such as software updates or hardware upgrades, can be performed on the passive node without impacting service availability. When maintenance is completed on the passive node, it can be switched back to the active role, and the previously active node can undergo maintenance in a controlled manner.

- ▶ **Highly Predictable Failover:** Active/Passive clusters provide predictable failover behavior. Administrators can configure the failover policies and resource constraints to ensure that the passive node only becomes active when specific conditions are met. This predictability is valuable for maintaining control over resource allocation and service availability.
- ▶ **Cost-Effective:** Active/Passive clustering is often more cost-effective than Active/Active clustering (where all nodes are active simultaneously) because fewer hardware resources are required for passive nodes. This makes it an attractive option for organizations with budget constraints.
- ▶ **Simplified Management:** Cluster resource managers like Pacemaker automate the failover process, making it easier for administrators to manage high-availability configurations. They provide tools for monitoring the cluster's health, resource status, and event logging.
- ▶ **Geographic Redundancy:** Active/Passive clusters can be set up across geographically separated data centers or locations. This geographic redundancy ensures that the service remains available even in the event of a site-wide failure, such as a natural disaster or power outage.
- ▶ **Compliance and SLA Fulfillment:** Active/Passive clustering helps organizations meet compliance requirements and service-level agreements (SLAs) by ensuring the availability of critical services. This is especially important for industries with strict uptime requirements, such as finance, healthcare, and e-commerce.
- ▶ **Risk Mitigation:** By providing a standby node that can take over in case of failures, Active/Passive clustering reduces the risk of costly business interruptions, reputations damage, and customer dissatisfaction associated with service downtime.

1.5.2 Active/Active clustering

Active/Active clustering is a high-availability configuration that allows multiple nodes or servers to actively serve applications or services simultaneously. In this setup, all nodes are actively processing user requests and workloads, sharing the load to improve performance, scalability, and redundancy. Active/Active clustering is commonly used in scenarios where high availability and resource optimization are critical.

Key Concepts and Components of Active/Active Clustering

The following are the key components of Active/Active clustering configuration:

- ▶ **Multiple Active Nodes:** Active/Active clustering involves two or more active nodes, often referred to as “members” or “instances.” These nodes are equally capable of running the same applications or services concurrently.
- ▶ **Load Balancing:** To evenly distribute incoming user requests and workloads among the active nodes, load balancing mechanisms are used. Load balancers direct traffic to the appropriate node based on factors like server load, response times, or other defined criteria.
- ▶ **Resource Sharing:** Active/Active clusters allow resources, such as CPU, memory, and storage, to be shared among all nodes. This sharing ensures efficient utilization of resources and prevents overloading of any single node.
- ▶ **Data Synchronization:** For stateful applications or services that require shared data, data synchronization mechanisms are implemented to keep data consistent across all active nodes. This can involve real-time replication or shared access to a centralized database.
- ▶ **Redundancy:** Active/Active clustering aims to eliminate single points of failure by providing multiple active nodes. If one node fails, the remaining nodes continue to provide service, minimizing downtime.

Example of Active/Active Clustering

Consider a scenario where a large e-commerce website needs to ensure high availability and handle a substantial volume of user traffic.

Scenario

Two Data Centers: Data Center A and Data Center B, each equipped with multiple servers.

E-commerce Application: The website's e-commerce application serves customers' online shopping needs.

This is illustrated in Figure 1-5.

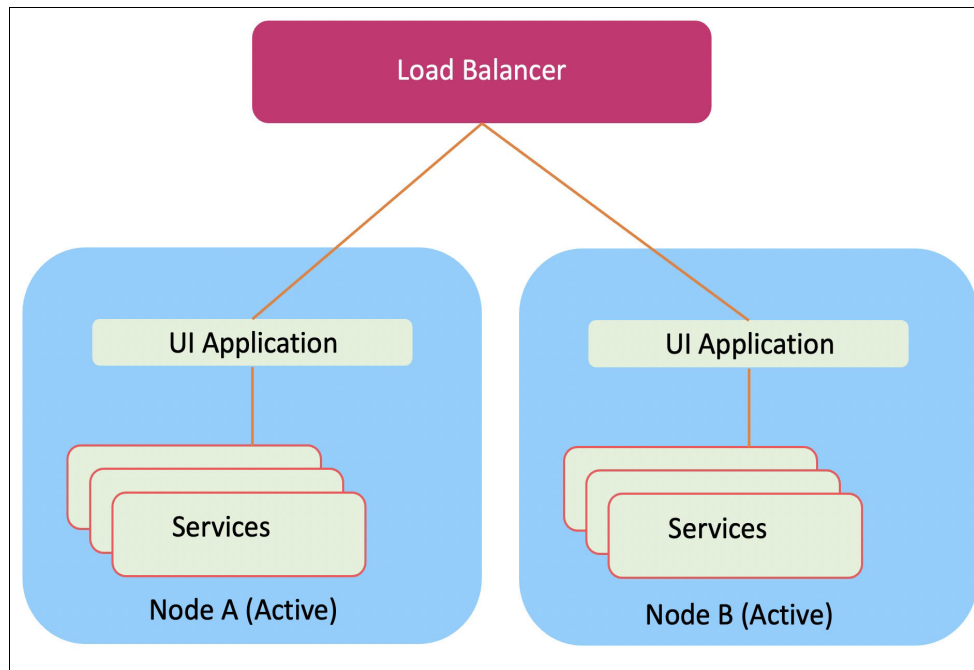


Figure 1-5 Active/Active Clustering

Configuration and Operation

- ▶ **Load Balancing:** Incoming user requests are directed to Data Center A or Data Center B based on factors such as server load and geographic proximity. Load balancers distribute the traffic evenly across all active nodes within each data center.
- ▶ **Active Nodes:** In each data center, there are multiple active nodes (servers) running the e-commerce application. All nodes are actively processing user requests and managing the website's database.
- ▶ **Data Replication:** To maintain data consistency, the website's database is replicated in real-time between Data Center A and Data Center B. Changes made in one data center are immediately reflected in the other.
- ▶ **Failover:** If a server in Data Center A fails due to hardware or software issues, the load balancer detects the failure and redirects traffic to the remaining servers in Data Center A and Data Center B. Customers experience minimal service interruption.
- ▶ **Resource Sharing:** During normal operation, the active nodes in both data centers share the processing load, ensuring efficient resource utilization.

Benefits of Active/Active clustering

The following are the key benefits of Active/Active clustering configuration:

- ▶ **High Availability:** Active/Active clustering offers excellent availability by allowing all nodes to actively serve traffic. Failures on one node do not result in service downtime.
- ▶ **Scalability:** It provides scalability and the ability to handle increased user traffic by adding more active nodes.
- ▶ **Resource Efficiency:** Resource sharing among nodes optimizes resource utilization and improves performance.
- ▶ **Load Balancing:** Load balancing ensures even distribution of workloads and responsive service delivery.
- ▶ **Data Consistency:** Real-time data synchronization guarantees data consistency and integrity.

Active/Active clustering is commonly used in scenarios where continuous availability, scalability, and resource optimization are essential, such as e-commerce websites, content delivery networks (CDNs), and large-scale cloud-based applications. It is a powerful approach to maintaining high availability while efficiently utilizing hardware and network resources.

1.5.3 N +1 clustering

N+1 clustering is a high-availability configuration that provides redundancy and fault tolerance for critical applications and services by having one extra (N+1) node in standby mode compared to the number of actively functioning nodes (N). In the event of a failure on one of the active nodes, the standby (extra) node takes over, ensuring continuous service availability. This setup is commonly used in data centers and critical infrastructure to minimize downtime and maintain uninterrupted operations.

Key Components of N+1 Clustering

The following are the key components of N+1 clustering configuration:

- ▶ **N Active Nodes:** These are the primary nodes that actively serve the application or service. They handle user requests, process data, and perform the core functions of the system.
- ▶ **1 Standby Node:** The standby node is kept in reserve, not actively participating in processing requests but ready to take over in case of a failure in any of the active nodes. This ensures redundancy and fault tolerance.
- ▶ **Failover Mechanism:** The cluster is equipped with a failover mechanism that automatically detects when an active node fails and initiates the transition of services to the standby node. This mechanism may involve monitoring for node health, resource availability, and network connectivity.

Example of N+1 Clustering

Let us consider an example of an N+1 cluster in the context of a web application hosted in a data center.

Scenario

- ▶ **N = 2 Active Nodes (Node A and Node B):** These nodes actively host the web application, serving user requests and processing data.
- ▶ **1 Standby Node (Node S):** This node is in standby mode and is not actively serving requests. It has the same configuration as Node A and Node B.

Configuration and Operation

- ▶ **Normal Operation:** During normal operation, Node A and Node B actively handle incoming web requests. Load balancing mechanisms distribute traffic between them, ensuring optimal performance. Node S remains in standby mode, monitoring the health of Node A and Node B.
- ▶ **Failure Detection:** If Node A encounters a hardware failure, becomes unresponsive, or experiences a critical software issue, the failover mechanism detects the problem.
- ▶ **Failover Process:** The failover mechanism triggers a failover process, directing incoming traffic away from the failed Node A. Node S, the standby node, is brought into active mode, assuming the responsibilities of Node A. It now starts processing web requests.
- ▶ **Service Continuity:** Users accessing the web application may experience a brief interruption or delay during the failover process, but the application remains accessible. Node B continues to operate as an active node, serving web requests alongside Node S.
- ▶ **Recovery and Maintenance:** After Node A's issue is resolved or maintenance is completed, it can be brought back into the cluster as a standby node. The cluster can be reconfigured to maintain N+1 redundancy, with Node A and Node B active, and Node S back in standby mode.

Benefits of N+1 Clustering

The following are the key benefits of N+1 clustering configuration:

- ▶ **High Availability:** N+1 clustering ensures continuous service availability even in the face of node failures, hardware issues, or software problems on active nodes.
- ▶ **Fault Tolerance:** It enhances fault tolerance by providing a ready-to-use standby node, reducing the risk of service disruptions.
- ▶ **Scalability:** Additional standby nodes can be added for further redundancy and capacity, making it suitable for scaling as needed.
- ▶ **Reliability:** N+1 clustering is a reliable solution for critical applications and services where downtime is costly or unacceptable.

N+1 clustering is widely used in various industries, including finance, healthcare, telecommunications, and e-commerce, to maintain the availability of mission-critical systems and services.

1.5.4 N-to-1 Clustering

The following are the key characteristics of N-to-1 clustering configuration:

- ▶ **Redundancy:** In an N-to-1 clustering configuration, “N” components (servers, nodes, or modules) work together as active units to handle the workload and serve applications.
- ▶ **1 Standby or Concentrator:** The “1” component, often referred to as the “concentrator” or “controller,” is a single unit that acts as a central point of management or control. It does not actively process workloads but coordinates and directs traffic among the “N” active components.
- ▶ **Failover:** If one of the active components fails, the standby concentrator takes over the management and routing of traffic. This ensures that even if one of the active units fails, the system can continue to operate with minimal disruption.
- ▶ **Examples:** N-to-1 clustering is commonly used in load balancing and traffic management scenarios, where multiple active servers distribute incoming requests, and a single concentrator handles traffic management and failover.

Note: The key distinction between N+1 and N-to-1 clustering lies in the role of the additional component(s). In N+1, you have “N” active components with a single standby or spare, while in N-to-1, you have “N” active components working in conjunction with a central concentrator that manages traffic and provides failover capabilities.

1.5.5 N+N clustering

N+N clustering is a high-availability configuration designed to provide redundancy and fault tolerance for critical applications and services by having an equal number of active nodes (N) working simultaneously. Unlike N+1 clustering, where there is one standby node for redundancy, N+N clustering has all N nodes actively serving the application or service at the same time. This approach aims to distribute the workload evenly across multiple nodes, improving performance and ensuring that if one node fails, the others can continue to provide services without interruption. N+N clustering is commonly used in scenarios where high availability and scalability are paramount.

Key Components of N+N clustering

The following are the key components of N+N clustering configuration:

- ▶ N Active Nodes: These are the primary nodes that actively serve the application or service.
- ▶ All N nodes are configured identically and work concurrently to process user requests and data.

Example of N+N clustering

Let us consider an example of N+N clustering in the context of a web application hosted in a data center as shown in Figure 1-6.

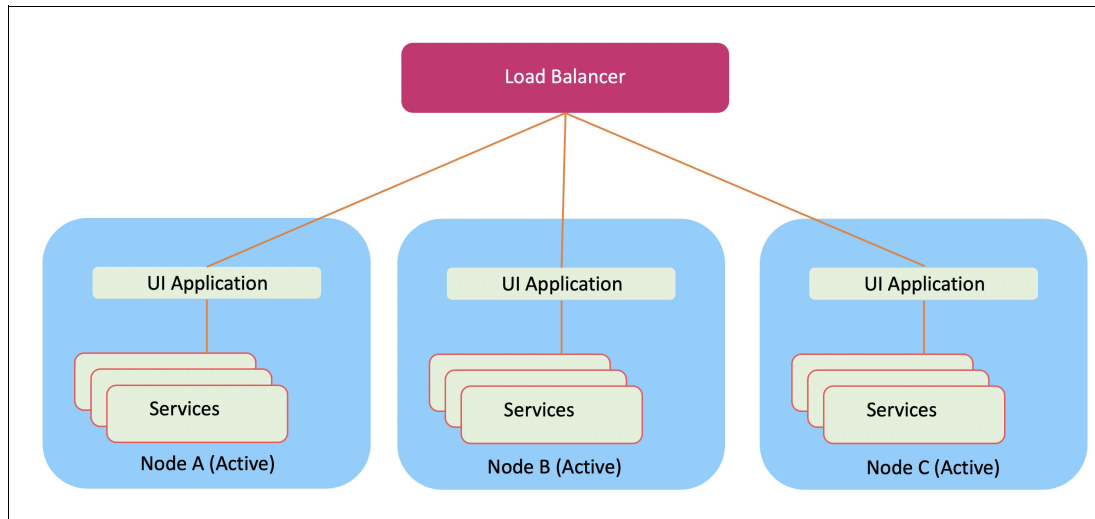


Figure 1-6 N+N Clustering

Scenario

- ▶ N = 3 Active Nodes (Node A, Node B, Node C): These nodes actively host the web application, serving user requests and processing data simultaneously.
- ▶ There is no standby node in this configuration.

Configuration and Operation

- ▶ **Normal Operation:** During normal operation, all N nodes (Node A, Node B, and Node C) actively handle incoming web requests. Load balancing mechanisms distribute traffic evenly across them, ensuring optimal performance and scalability.
- ▶ **Failure Detection:** If one of the nodes, say Node B, encounters a hardware failure, software issue, or becomes unresponsive, the system detects the problem.
- ▶ **Continued Service:** Despite the failure of Node B, the web application continues to operate without disruption. User requests are automatically routed to Nodes A and C, which are still active and operational.
- ▶ **Service Continuity:** Users accessing the web application may not even notice the failure of Node B. The N+N clustering configuration ensures that there is no single point of failure, and the system continues to provide services reliably.
- ▶ **Recovery and Maintenance:** After resolving the issue with Node B or performing necessary maintenance, it can be reintegrated into the cluster as an active node. This process can be repeated for any failed node.

Benefits of N+N clustering

The following are the key benefits of N+N clustering configuration:

- ▶ **High Availability:** N+N clustering offers high availability by distributing the workload across multiple active nodes. Even if one or more nodes fail, others continue to operate, minimizing service downtime.
- ▶ **Redundancy:** There is no single point of failure in N+N clustering, as multiple nodes are actively serving the application or service.
- ▶ **Scalability:** N+N clusters are highly scalable. Additional nodes can be added to increase capacity and further distribute the workload, accommodating growth.
- ▶ **Performance:** N+N clustering typically results in improved performance and response times since multiple nodes share the load, making it suitable for high-demand environments.

N+N clustering is often used in scenarios where service uptime is crucial and where handling fluctuations in workload is essential. It provides a robust solution for businesses and organizations that require continuous availability and scalability for their critical applications and services.

1.5.6 Three site clusters

A 3-site cluster, often referred to as a multi-site or three-node clustering configuration, is a high-availability setup designed to provide redundancy and disaster recovery capabilities across three geographically separate locations or data centers. This configuration is suitable for organizations that require extremely high levels of availability, data protection, and continuity of critical services, even in the face of site-level failures or natural disasters.

Key Components of 3-Site clustering

The following are the key components of 3-site clustering configuration:

- ▶ **Three Sites:** In this configuration, there are three distinct physical locations or data centers, each housing a portion of the cluster nodes and resources. These sites are geographically separated to reduce the risk of simultaneous failures due to local disasters.

- ▶ **Cluster Nodes:** Each site contains a set of cluster nodes (typically two or more) that actively serve the application or service. These nodes work together to provide redundancy and high availability.
- ▶ **Data Replication:** Data is actively replicated between the sites, ensuring that data is available at multiple locations. This replication can be synchronous (real-time) or asynchronous (with a delay) depending on the specific requirements and technologies used.
- ▶ **Load Balancing:** Load balancing mechanisms distribute user requests and workloads across the cluster nodes within and across the three sites.

Example of 3-Site Clustering

Let us consider an example of a 3-site clustering configuration in the context of a financial institution that relies on continuous availability and data protection:

Scenario

- ▶ **Three Sites:** Site A, Site B, and Site C, each located in different cities or regions.
- ▶ **N Active Nodes at Each Site (Node A1, Node A2 at Site A; Node B1, Node B2 at Site B; Node C1, Node C2 at Site C):** Each site has multiple active nodes that host the core banking application.

Configuration and Operation

- ▶ **Normal Operation:** During normal operation, all nodes at each site actively handle user transactions and data processing. Data is continuously replicated and synchronized between the sites in real-time (synchronous replication) to ensure data consistency and availability.
- ▶ **Load Balancing and Geographic Redundancy:** Load balancing mechanisms distribute incoming user requests across the nodes within each site to ensure optimal performance. In addition to load balancing within a site, DNS or Global Server Load Balancers can be used to distribute traffic evenly across the three geographically separated sites, providing geographic redundancy.
- ▶ **Failure and Disaster Recovery:** In the event of a site-level failure (e.g., a natural disaster or complete data center outage at Site A), the remaining two sites (B and C) continue to provide services. Failover mechanisms redirect user traffic and data access to the surviving sites, ensuring service continuity.
- ▶ **Data Resilience:** Data replication ensures that even if one site becomes unavailable, data is accessible from the remaining sites. This redundancy protects against data loss.
- ▶ **Failback and Recovery:** After the failed site (e.g., Site A) is restored, data is resynchronized, and it can rejoin the cluster, resuming normal operations.

Benefits of 3-Site clustering

The following are the key benefits of 3-site clustering configuration:

- ▶ **High Availability:** 3-site clustering offers exceptional availability by spreading services and data across three geographically separated locations.
- ▶ **Disaster Recovery:** It provides robust disaster recovery capabilities, protecting against site-level failures or natural disasters.
- ▶ **Data Protection:** Continuous data replication ensures data integrity and minimizes the risk of data loss.
- ▶ **Geographic Redundancy:** Geographic dispersion of sites reduces the risk of simultaneous failures due to localized events.

- **Compliance:** It helps organizations meet regulatory compliance requirements, particularly in industries like finance and healthcare.

3-site clustering is a complex and costly solution, typically reserved for mission-critical systems where even a brief downtime can result in significant financial loss or regulatory non-compliance. It is essential to carefully plan, implement, and test such configurations to ensure they meet the specific requirements of the organization.

Table on page 22 highlights the key differences between Active/Passive, N+1, N+N, and 3-site clustering configurations

Table 1-3 Clustering options

Aspect	Active/Passive	N+1	N+N	Active/Active	3-Site
Number of Active Nodes	1	N	N	N (all actively serving)	N (across multiple sites)
Number of Standby Nodes	1	1	None or optional N+1	None	None
Fault Tolerance	Limited	Moderate	High	High	Exceptionally High
Redundancy	Standby Node Redundancy	Standby Node Redundancy	Full Node Redundancy	Full Node Redundancy	Geographic Redundancy
Load Balancing	No	No (standby is idle)	Yes (evenly distributed)	Yes (evenly distributed)	Yes (load balanced within sites)
Scalability	Limited	Moderate	High	High	High
Geographic Redundancy	No	No	No (typically single-site)	No (typically single-site)	Yes (across geographically separate sites)
Complexity	Low	Moderate	Moderate to High	Moderate to High	High (complex to configure)
Data Replication	Optional	Optional	Optional	Optional	Real-time (synchronous or asynchronous)
Use Cases	Simple applications with low tolerance for downtime	Moderate tolerance for downtime, scalability needed	High availability with redundancy and scalability	High availability with efficient resource utilization	Mission-critical applications requiring the highest levels of availability, disaster recovery, and data protection
Cost	Lower (fewer hardware resources)	Moderate (additional standby node)	Moderate to High (multiple active nodes)	Moderate to High (additional hardware and load balancers)	High (multiple sites, extensive hardware, and network infrastructure)
Examples	Single server with a standby backup	A web server farm with one standby server	Database clusters with multiple active nodes	Large-scale e-commerce websites with efficient resource utilization	Large-scale data centers with geographically dispersed sites

Please note that the suitability of each clustering configuration depends on specific organizational requirements, budgets, and the level of availability needed for the applications or services in question. The choice between these configurations should be based on a thorough analysis of business needs and risk tolerance.

1.5.7 Linux Pacemaker

Pacemaker serves as a high-availability cluster resource manager, a software solution designed to operate on a group of hosts, known as a cluster of nodes, with the objective of maintaining system integrity and minimizing downtime for critical services and resources.²

Key features of Pacemaker include:

- ▶ Detecting and recovering from failures at both the node and service levels.
- ▶ Ensuring data integrity by isolating and addressing malfunctioning nodes.
- ▶ Supporting single or multiple nodes within a cluster.
- ▶ Providing compatibility with various resource interface standards, enabling the clustering of virtually any scriptable component.
- ▶ Offering optional support for shared storage.
- ▶ Adapting to a wide range of redundancy configurations, including active/passive and N+1 setups.
- ▶ Automatically replicating configurations, allowing updates from any cluster node.
- ▶ Enabling the specification of cluster-wide relationships between services, including ordering, collocation, and anti-collocation constraints.
- ▶ Supporting advanced service types such as clones (services active on multiple nodes), promotable clones (clones with dual roles), and containerized services.
- ▶ Offering a unified and scriptable set of cluster management tools.

Pacemaker is developed and maintained by the ClusterLabs community³. This collaborative effort unites numerous open-source projects related to High Availability into a comprehensive cluster solution suitable for deployments of varying sizes. Comprising components like Corosync, Pacemaker, DRBD, ScanCore, and others, the ClusterLabs stack has been instrumental in detecting and recovering from hardware and application-level failures in production clusters since 1999. Its flexibility extends to accommodating virtually any imaginable redundancy configuration.

For more information refer to Chapter 2, “Pacemaker Overview” on page 25.

1.5.8 High Availability with Red Hat Enterprise Linux Add-On

Red Hat Enterprise Linux offers the flexibility to establish managed, highly available clusters comprising RHEL servers⁴. This Add-On boasts extensive configurability, catering to a wide range of applications, whether they are off-the-shelf or custom-built. Key components of this Add-On include:

- ▶ The Pacemaker cluster resource manager.
- ▶ Network management via Corosync and Kronosnet, facilitating seamless communication among cluster nodes.
- ▶ Multisite and stretch cluster support through Booth, QDevice, and QNetd.

² https://clusterlabs.org/pacemaker/doc/2.1/Clusters_from_Scratch/html/intro.html#what-is-pacemaker

³ <https://www.clusterlabs.org>

⁴ https://access.redhat.com/documentation/fr-fr/red_hat_enterprise_linux/8/html/configuring_and_managing_high_availability_clusters/assembly_overview-of-high-availability-configuring-and-managing-high-availability-clusters

- ▶ Management options including a Pacemaker command line interface (pcs) and GUI (pcsd).
- ▶ A comprehensive library of resource agents for commonly used applications, alongside the ability to create custom resource agents.
- ▶ Fencing agents for bare metal servers, virtual machines, and cloud platforms.

Red Hat Enterprise Linux is well-suited for a diverse array of clustering workloads, encompassing configurations such as active/passive, active/active, and primary/secondary. Its automatic service restart and failover mechanisms contribute to maintaining high availability by swiftly transitioning failing services and their dependencies to functional nodes.

This high availability solution seamlessly integrates with popular applications like SAP HANA, Apache, PostgreSQL, and DB2. Furthermore, custom applications can be effortlessly managed within the RHEL High Availability framework, either by utilizing their systemd startup scripts or by creating straightforward custom resource agents.

1.5.9 SUSE Linux Enterprise High Availability Extension

SUSE Linux Enterprise High Availability Extension is a powerful solution designed to enhance the reliability and availability of critical services and applications in enterprise environments. At the heart of this extension lies Pacemaker, a robust open-source cluster resource manager that plays a pivotal role in orchestrating high availability clustering. Additional information can be found in the [SUSE HA Extension documentation](#).

1.6 Linux Pacemaker vs VM Recovery Manager HA

The choice between Linux Pacemaker and IBM Power VM Recovery Manager depends on the specific needs of your environment. Linux Pacemaker is a versatile open-source solution suitable for Linux-based HA clustering across different hardware platforms. In contrast, IBM Power VM Recovery Manager is tailored for IBM Power Systems, focusing on VM-level high availability and integration with PowerVM virtualization features. The decision should consider your infrastructure, budget, and the level of high availability required for your workloads.

Table 1-4 highlights comparison between Pacemaker and VM Recovery Manager HA.

Table 1-4 Linux Pacemaker vs VMRM HA.

Solution	Pacemaker for IBM Power Systems	VM Recovery Manager – HA for IBM Power Systems
Base technology	Linux OS	IBM PowerVM & shared storage
Automated failover	Yes	Yes
Ideal for DR	No	No
Ideal for HA	Yes	Yes
Ease of Management	Easy	Advanced
Reuse secondary for non-production workload	Partly	Partly
DR Rehearsal feature for non-disruptive DR testing	No	No
Cloud support	Yes	No
License cost	Part of Linux OS	Low cost

1.7 Software licensing challenges in HA setups

When it comes to setting up a cluster scenario for high availability, it's crucial to exercise caution, especially when it comes to software licensing considerations. While high availability clustering can provide resilience and minimize downtime, it may have implications for software licensing costs. Some software licenses are structured in a way that charges organizations based on the number of installations, even in an active/passive scenario. Here are a few examples to illustrate this point:

- ▶ Many DBMS vendors charge licensing fees based on the number of installations or nodes in a cluster, regardless of whether those nodes are actively processing transactions or in a standby (passive) state. For instance, if you have an active/passive cluster for a database server, you may be required to pay licensing fees for both the active and passive nodes, effectively doubling the licensing costs.
- ▶ Middleware and application server software, commonly used for hosting web applications and services, may also have licensing models that consider each node in a cluster as a separate installation. This means that even if one server is actively serving requests while the other is on standby, you could still be charged for both.
- ▶ Some operating system vendors may have licensing policies that require licensing for each node, irrespective of whether they are actively serving workloads or serving as a failover in an active/passive cluster. This can significantly impact licensing costs for organizations running high availability clusters.
- ▶ Various third-party software components, including monitoring tools, security solutions, and backup software, may have licensing models that aren't cluster-aware. As a result, organizations need to carefully review their licensing agreements to ensure compliance.

To address these licensing challenges in high availability scenarios, organizations should take the following steps:

- ▶ Engage with software vendors to understand their licensing policies specifically regarding clustered environments. Seek clarification on whether passive nodes are subject to licensing fees.
- ▶ Thoroughly review software licensing agreements to identify any clauses or terms related to clustering and passive nodes. Seek legal counsel if necessary.
- ▶ Explore alternative software options that have more cluster-friendly licensing models, such as licensing based on active nodes only.
- ▶ Depending on licensing costs, consider optimizing your cluster configuration to minimize passive nodes or explore active/active clustering where both nodes are actively serving traffic to avoid additional licensing fees.



Pacemaker Overview

This chapter provides a more in depth view of Pacemaker, describing how Pacemaker can be used to provide high availability for your Linux workloads. We discuss the basic components of a Pacemaker installation and help you understand how they work together to ensure that your application environment is protected from component failures.

We also discuss the various implementation topologies that can be used from the simplest active failover cluster to multi site solutions that can provide application failover not only within a single location, but across multiple sites to provide a disaster recovery solution.

This chapter contains the following topics:

- ▶ 2.1, “Architecture” on page 26.
- ▶ 2.2, “Pacemaker clustering options” on page 29.
- ▶ 2.3, “Failover policies” on page 33.

2.1 Architecture

Pacemaker is a high-availability cluster resource manager – software that runs on a set of hosts to create a cluster in order to preserve integrity and minimize downtime of desired services or resources. Pacemaker achieves maximum availability for the cluster resources by detecting and recovering from node and resource level failures – making use of the messaging and membership capabilities provided for the preferred cluster infrastructure using CoroSycn.

CoroSycn

CoroSycn provides cluster infrastructure functionality. CoroSycn provides messaging and membership functionality and maintains the quorum information. This feature has been utilized by Pacemaker to provide a high availability solution.

What is Quorum?

In order to maintain cluster integrity and availability, cluster systems use a concept known as quorum to prevent data corruption and loss. A cluster has quorum when more than half of the cluster nodes are online. To mitigate the chance of data corruption due to failure, Pacemaker by default stops all resources if the cluster does not have quorum. The term "split-brain" generally describes a problematic condition that high availability cluster solutions must avoid or address when members of a cluster lose contact with each other. When cluster members are unable to communicate, then they may be at risk of performing conflicting activities because they cannot coordinate with each other. For more details refer to <https://access.redhat.com/articles/2824071>.

Quorum is established using a voting system. When a cluster node does not function as it should or loses communication with the rest of the cluster, the majority of working nodes can vote to isolate and if needed, fence the node for servicing.

What is Fencing?

Fencing is the mechanism that the cluster uses to resolve issues and failures. Fencing ensures that any problematic cluster node will be cut off quickly and the remaining nodes in the cluster can take over the cluster services, ensuring a more resilient and stable cluster.

For an example, if communication with a single node in the cluster fails, then other nodes in the cluster must be able to restrict or release access to resources that the failed cluster node may have access to. This cannot be accomplished by contacting the cluster node itself as the cluster node may not be responsive. Instead, you must provide an external method, which is called fencing with a fence agent. A fence device is an external device that can be used by the cluster to restrict access to shared resources by an errant node, often by issuing a hard reboot on the cluster node.

Fence Device Types

A fence agent (or device) is an external device that can be used by the cluster to restrict access to shared resources by an errant node (or hard reboot the cluster node). The two most common types of fencing are:

- ▶ **Power fence agents:** The cluster software logs in via telnet, ssh, or SNMP to the device such as an APC switch, Dell DRAC, HP iLO, IBM RSA, or similar device and turns off (and optionally on) the power for the cluster node. This method will execute a hard "off" action.
- ▶ **I/O fence agents:** The cluster software logs in to a fiber channel switch via telnet or ssh and disables the ports for that node, thereby cutting off its access to shared storage. This

method requires that an administrator manually reboot or shutdown the errant node to recover it, and log in to the switch interface to re-enable the appropriate ports.

For more details refer to <https://access.redhat.com/solutions/15575>.

2.1.1 Cluster Architecture

The following provides a high-level overview of the terminology used in high-availability cluster architecture:

- ▶ **Node:** One of a set of hosts or computers that work together to perform a task as a member of the cluster. The maximum number of nodes that can be members of a RHEL High Availability cluster is 32. For more details refer to <https://access.redhat.com/articles/3069031>.
- ▶ **Resources:** The services that need to be kept highly available. For example: IP Address, Filesystem, or Application Service.
- ▶ **Resource group:** A resource group is a collection of resources that function together and can fail over from one node to another as a single unit.
- ▶ **Resource agents:** These are scripts or operating system components that start, stop, and monitor resources, given a set of resource parameters.
- ▶ **Fence agents:** These are scripts that execute node fencing actions, given a target and fence device parameters.
- ▶ **Cluster membership layer:** This component provides reliable messaging, membership, and quorum information about the cluster.
- ▶ **Cluster resource manager:** Pacemaker provides the brain that processes and reacts to events that occur in the cluster. These events may include nodes joining or leaving the cluster; resource events caused by failures, maintenance, or scheduled activities; and other administrative actions. To achieve the desired availability, Pacemaker may start and stop resources and fence nodes.
- ▶ **Cluster tools:** These provide various command-line and graphical (GUI) interfaces for users to interact with the cluster.

Figure 2-1 illustrates Pacemaker's architecture and main components:

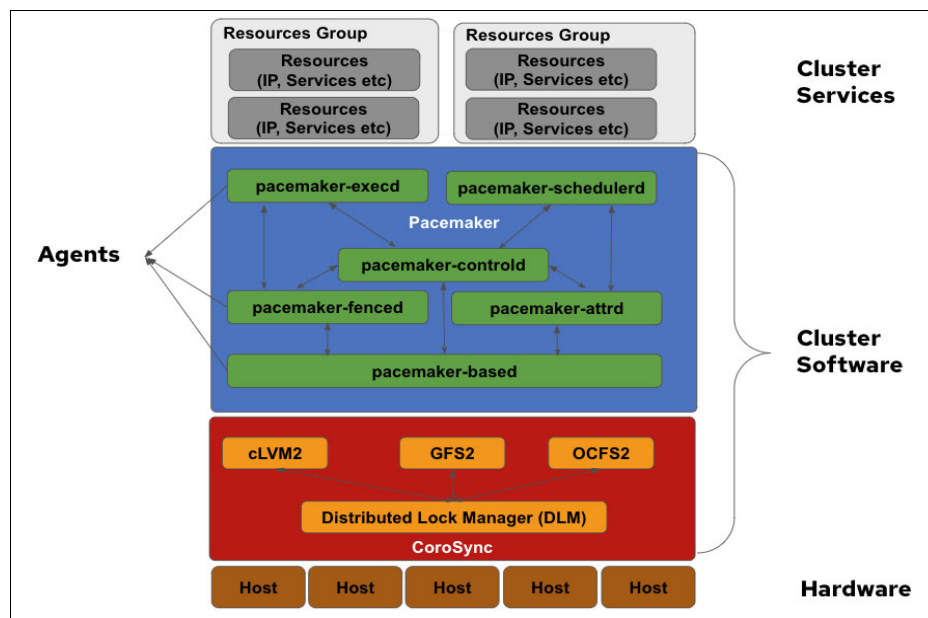


Figure 2-1 Pacemaker architecture and components

For more information refer to [Pacemaker Administration](#).

2.1.2 Pacemaker architecture and components

A cluster configured with Pacemaker comprises separate component daemons that monitor cluster membership, scripts that manage the services, and resource management subsystems that monitor the disparate resources. Pacemaker's main process (`pacemakerd`) spawns all the other daemons, and respawns them if they unexpectedly exit.

The following components form the Pacemaker architecture:

- ▶ Pacemaker centralizes cluster decision-making by electing one of the controller instances as the **Designated Controller (DC)**. Should the elected DC process (or the node it is on) fail, a new one is quickly established. The DC responds to cluster events by taking a current snapshot of the CIB, feeding it to the scheduler, then asking the executors (either directly on the local node, or via requests to controller peers on other nodes) and the fencer to execute any necessary actions.
- ▶ Cluster Information Base (CIB): The Pacemaker information daemon, which uses XML internally to distribute and synchronize current configuration and status information from the Designated Coordinator (DC) – a node assigned by Pacemaker to store and distribute cluster state and actions by means of the CIB to all other cluster nodes.
- ▶ The CIB manager (***pacemaker-based***) keeps the CIB synchronized across the cluster, and handles requests to modify it.
- ▶ The attribute manager (***pacemaker-attribd***) maintains a database of attributes for all nodes, keeps it synchronized across the cluster, and handles requests to modify them. These attributes are usually recorded in the CIB.
- ▶ The scheduler (***pacemaker-schedulerd***) determines what actions are necessary to achieve the desired state of the cluster by given a snapshot of the CIB as input.
- ▶ Cluster Resource Management Daemon (**CRMD**): Pacemaker cluster resource actions are routed through this daemon. Resources managed by CRMD can be queried by client systems, moved, instantiated, and changed when needed. Each cluster node also includes a local resource manager daemon (**LRMD**) that acts as an interface between CRMD and resources. LRMD passes commands from CRMD to agents, such as starting and stopping and relaying status information.
- ▶ The local executor (***pacemaker-execd***) handles requests to execute resource agents on the local cluster node, and returns the result.
- ▶ Shoot the Other Node in the Head (**STONITH**): STONITH is the Pacemaker fencing implementation. It acts as a cluster resource in Pacemaker that processes fence requests, forcefully shutting down nodes and removing them from the cluster to ensure data integrity. STONITH is configured in the CIB and can be monitored as a normal cluster resource.
- ▶ The fencer (***pacemaker-fenced***) handles requests to fence nodes. Given a target node, the fencer decides which cluster nodes should execute which fencing devices, and calls the necessary fencing agents (either directly, or via requests to the fencer peers on other nodes), and returns the result.
- ▶ The controller (***pacemaker-controld***) is Pacemaker's coordinator, maintaining a consistent view of the cluster membership and orchestrating all the other components.

For a cluster to operate, it requires an effective communications between the nodes and to manage membership of the cluster. Pacemaker utilizes CoroSyc (the function with daemon of the same name) to serve the core membership and member-communication needs for high availability clusters. It is required for any High Availability Add-On capabilities to function.

The CoroSync membership and messaging functions also:

- ▶ Manage quorum rules and determination.
- ▶ Provide messaging capabilities for applications that coordinate or operate across multiple members of the cluster and thus must communicate stateful or other information between instances.
- ▶ Use the kronosnet library as its network transport to provide multiple redundant links and automatic failover.

Most managed services (applications) are not themselves, *cluster-aware*. However, many popular open-source cluster filesystems make use of a common **Distributed Lock Manager (DLM)**, which makes direct use of CoroSync for its messaging and membership capabilities and Pacemaker for the ability to fence nodes.

For more information refer to this [Pacemaker document](#).

2.2 Pacemaker clustering options

Pacemaker supports practically any node redundancy configuration including Active/Active, Active/Passive, N+1, N+M, N-to-1, and N-to-N. The most common size for a high-availability cluster is a two-node cluster, since that is the minimum required to provide redundancy, but many clusters consist of many more, sometimes dozens of nodes.

A high-availability cluster categorized into one of the following models¹:

- Active/active: Traffic intended for the failed node is either passed onto an existing node or load balanced across the remaining nodes. This is usually only possible when the nodes use a homogeneous software configuration.
- Active/passive: Provides a fully redundant instance of each node, which is only brought online when its associated primary node fails. This configuration typically requires the most extra hardware.
- N+1: Provides a single extra node that is brought online to take over the role of the node that has failed. In the case of heterogeneous software configuration on each primary node, the extra node must be universally capable of assuming any of the roles of the primary nodes it is responsible for. This normally refers to clusters that have multiple services running simultaneously; in the single service case, this degenerates to active/passive.
- N+M: In cases where a single cluster is managing many services, having only one dedicated failover node might not offer sufficient redundancy. In such cases, more than one (M) standby server is included and available. The number of standby servers is a trade-off between cost and reliability requirements.
- N-to-1: Allows the failover standby node to become the active one temporarily, until the original node can be restored or brought back online, at which point the services or instances must be failed-back to it in order to restore high availability.
- N-to-N: A combination of active/active and N+M clusters, N to N clusters redistribute the services, instances or connections from the failed node among the remaining active nodes, thus eliminating (as with active/active) the need for a 'standby' node, but introducing a need for extra capacity on all active nodes.

All of these high availability cluster models consist of multiple physical resources (compute nodes, share storage, network, fence device) and logical resources (node, resources, resource group) that could be collocated at one geographical location or data center (DC), or could be distributed across multiple geographical locations or data centers (DCs).

¹ https://en.wikipedia.org/wiki/High-availability_cluster#Node_configurations

2.2.1 Single site

The single-site cluster is setup where all cluster members and the associated resources are in the same physical location or data center (DC), connected by a local area network.

Figure 2-2 show an example of a single site Pacemaker cluster architecture for database with replication. The database in this setup would be configured so that data is synchronized from the primary to secondary node.

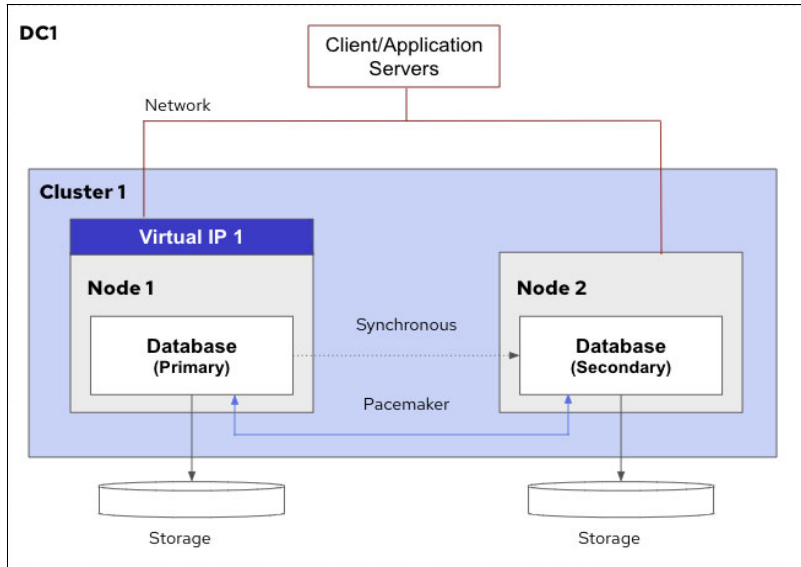


Figure 2-2 Example of single site Pacemaker cluster architecture for database with replication

Figure 2-3 is an example of a single site Pacemaker cluster architecture for general application with shared storage. In this scenario the application does not have any synchronization method to maintain identical data across the node and can leverage the shared storage configuration in the Pacemaker cluster.

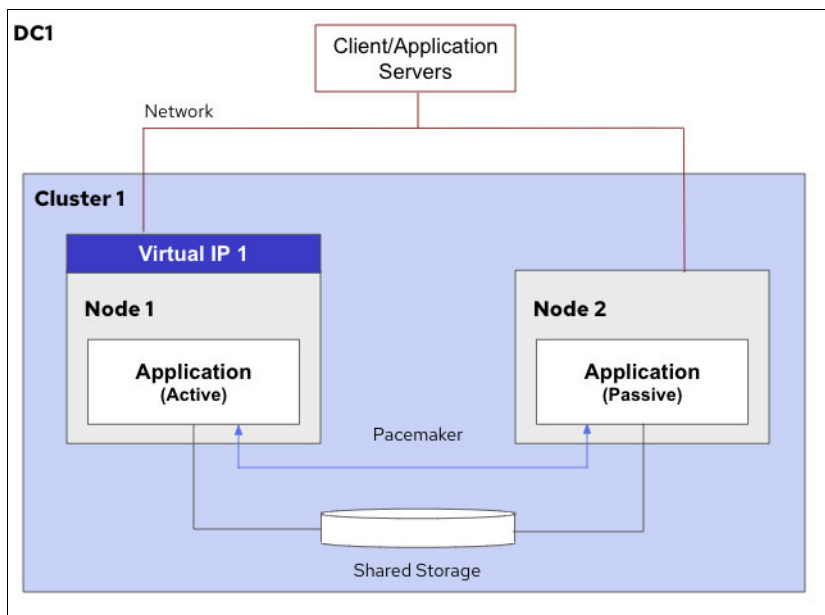


Figure 2-3 Example of single site Pacemaker cluster architecture for general application with shared storage

2.2.2 Two site

The two site clusters configuration consists of two clusters² – one active and one for disaster recovery. A method for creating a two-site cluster with disaster recovery for a high availability cluster is to configure two clusters; one cluster as the primary site cluster and the second cluster as the disaster recovery cluster. Alternatively, you could create a single cluster that is managing (stretching) cluster resources across two sites.

For example, a database running in the primary cluster in promoted mode and running in the disaster recovery cluster in demoted mode. The database is setup and configured in a way, so that data is synchronized from the primary site to the disaster recovery site.

When the primary cluster goes down, the Pacemaker command interface can be used to manually fail the resources over to the disaster recovery site. Once the primary cluster has recovered, the Pacemaker command interface can be used to manually move resources back to the primary site.

Figure 2-4 is an example of a two site Pacemaker cluster architecture for a disaster recovery scenario.

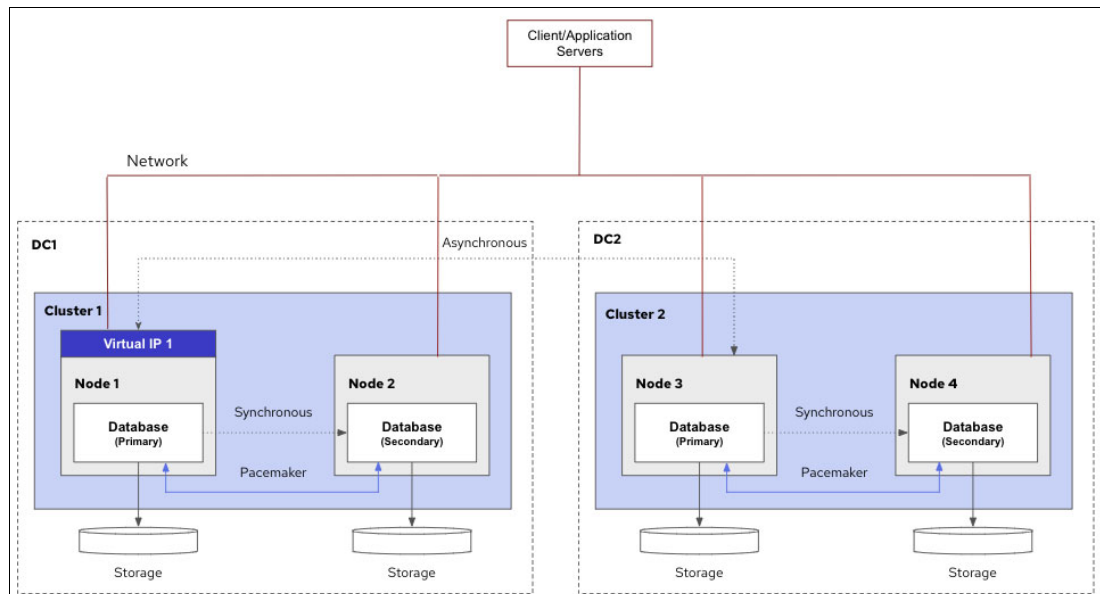


Figure 2-4 Example of two site Pacemaker cluster architecture for a disaster recovery scenario.

Figure 2-5 on page 32 shows a stretch cluster that is designed to withstand the loss or failure of all members at a given physical site. A stretch cluster comprises a single infrastructure and membership spanning all sites. Membership of the cluster is logically divided into two groups so that cluster services can continue with minimal disruption when an entire group fails or becomes unreachable.

If there is shared storage, it is replicated via either hardware or software replication mechanisms so that each group has access to a replica. The groups are typically, but not necessarily, at different physical locations, often with reduced communication inter-connectivity and increased delay compared to a single site. For more details refers to: <https://access.redhat.com/articles/27136>.

² https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux/9/html-single/configuring_and_managing_high_availability_clusters/index#assembly_configuring-disaster-recovery-configuring-and-managing-high-availability-clusters

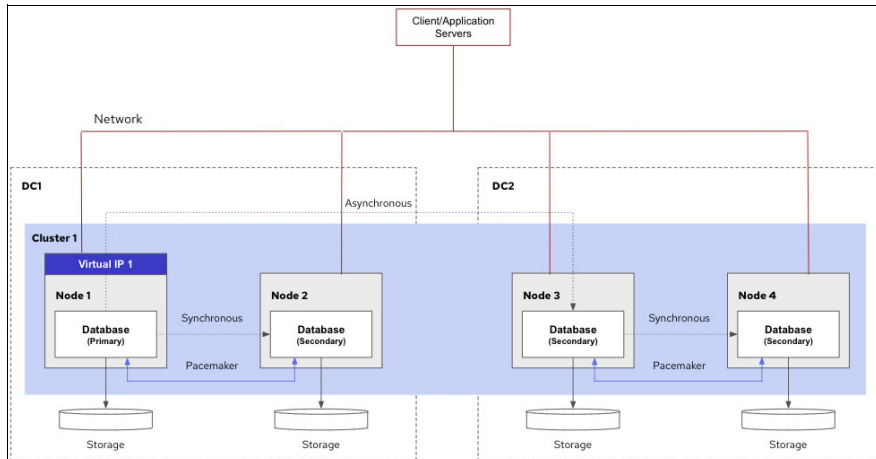


Figure 2-5 Example of two site Pacemaker cluster architecture for a stretch cluster

2.2.3 Multi site

Multisite or also known as geo clusters³, are clusters stretched out over multiple physical locations. Multi-site or disaster-tolerant clusters typically use SAN-based storage replication to replicate data and usually used in an active/passive manner for disaster recovery with manual failover of the active cluster to the passive cluster. This is shown in Figure 2-6.

In a two-site cluster, there will be potential issues with network connectivity between the sites that can lead to split-brain situations. When connectivity drops, there is no way for a node on one site to determine whether a node on another site has failed or is still functioning with a failed site interlink. In addition, it can be problematic to provide high availability services across two sites which are too far apart to keep synchronous. To address these issues, Pacemaker provides full support for the ability to configure high availability clusters that span multiple sites through the use of a Booth cluster ticket manager.

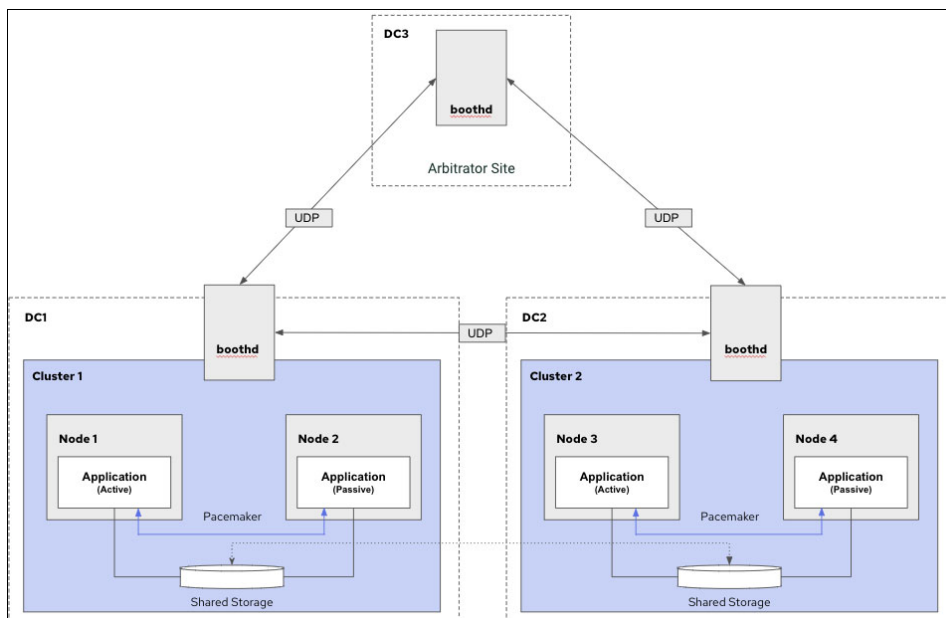


Figure 2-6 Example of Multisite or Geo Pacemaker cluster with two sites and a single arbitrator on a third site.

³ <https://documentation.suse.com/sle-ha/15-SP1/html/SLE-HA-all/cha-ha-geo-concept.html>

Multisite or Geo clusters can be considered overlay clusters where each cluster site corresponds to a cluster node in a traditional cluster. The overlay cluster is managed by the booth cluster ticket manager (in the following called booth).

Each of the parties involved in a Geo cluster runs a service, the boothd. It connects to the booth daemons running at the other sites and exchanges connectivity details. For making cluster resources highly available across sites, booths rely on cluster objects called tickets. A ticket grants the right to run certain resources on a specific cluster site. Booth guarantees that every ticket is granted to no more than one site at a time.

If the communication between two booth instances breaks down, it might be because of a network breakdown between the cluster sites or because of an outage of one cluster site. In this case, you need an additional instance (a third cluster site or an arbitrator) to reach consensus about decisions (such as failover of resources across sites). Arbitrators are single machines (outside of the clusters) that run a booth instance in a special mode. Each Geo cluster can have one or multiple arbitrators. The most common scenario probably is a Multisite or Geo cluster with two sites and a single arbitrator on a third site. This requires three booth instances.

2.3 Failover policies

Pacemaker offers several failover policies and constraints to manage resources within a cluster environment. These policies and constraints control how resources behave in the event of failures or changes in the cluster. Table 2-1 provides descriptions and examples of some Pacemaker failover policies.

Table 2-1 Common Pacemaker failover policies

Failover Policy	Description	Example
Priority-Based Failover	Resources are assigned priority levels, and Pacemaker attempts to start resources on nodes with the highest priority that are available and suitable for hosting the resources.	In a scenario with two nodes, Node A has a higher priority (e.g., 100) for running a critical service compared to Node B (e.g., 50). If Node A fails, Pacemaker will try to start the service on Node B.
Location Constraints	Defines preferred or forbidden nodes for specific resources, guiding Pacemaker on where to place or avoid relocating resources during failover.	An administrator sets a location constraint that the database service should run on Node X unless Node X fails. In case of a failure on Node X, Pacemaker relocates the service to Node Y, the next preferred node.
Collocation Constraints	Groups resources together or enforces separation to specify which resources can or cannot run simultaneously on the same node.	To ensure optimal performance, a web server and a database are collocated, meaning they run on the same node. However, a database and a certain critical service are anti-collocated to avoid resource contention.
Ordering Constraints	Defines the sequence in which resources should start or stop to maintain dependencies.	A scenario where a web server must start after a network interface is up. An ordering constraint ensures that the network interface starts before the web server starts.
Resource Stickiness	Determines how likely a resource is to stay on its current node unless necessary (e.g., node failure or resource failure).	A service has high stickiness set, so even if a better node is available, Pacemaker avoids moving the service unless the current node fails.

These examples illustrate how administrators can use Pacemaker failover policies and constraints to ensure resource availability, optimal placement, and controlled failover behavior within a cluster environment. Each policy serves specific needs related to resource management, ensuring reliability and performance in a clustered setup.

The selection of Pacemaker failover policies depends heavily on the specific requirements and characteristics of your clustered environment, as well as the criticality of the services being managed. However, here are some general recommendations for utilizing Pacemaker failover policies effectively:

- ▶ **Priority-Based Failover:** Assign appropriate priority levels to critical resources. Higher priority resources will be allocated on nodes with better capabilities or higher availability.
- ▶ **Location Constraints:** Define preferred nodes for critical services or those with specific hardware or configurations. This ensures services are allocated to nodes best suited for their requirements.
- ▶ **Collocation Constraints:** Use collocation and anti-collocation constraints to control resource placement. Group resources that benefit from running together and avoid conflicts between others.
- ▶ **Ordering Constraints:** Establish dependencies between resources to maintain proper sequencing during start-up or shut-down procedures.
- ▶ **Resource Stickiness:** Adjust stickiness levels carefully. Higher stickiness prevents unnecessary resource movements but may hinder load balancing if set too high.

Remember, there is no one-size-fits-all approach. It's crucial to tailor these failover policies to your specific environment, considering the nature of applications, available resources, and business continuity requirements. Regular monitoring and adjustments based on real-time performance are also essential for maintaining an effective and reliable high-availability setup using Pacemaker.



Planning and Installing Your Pacemaker Environment

IBM Power Systems are renowned for their robust architecture and exceptional performance. When combined with Pacemaker, they offer a compelling solution for achieving high availability. This chapter serves as a foundational guide to understanding the intricate process of installing and configuring Pacemaker on IBM Power Systems, you will have a grasp of the key concepts, prerequisites, and initial steps required to embark on your journey towards creating a highly available environment.

This chapter contains the following topics:

- ▶ 3.1, “Prerequisites and requirements” on page 36.
- ▶ 3.2, “Installing and Configuring Pacemaker environment” on page 38.

Note: While Pacemaker is compatible with several supported Linux distributions on IBM Power System and IBM Power Virtual Server, this chapter primarily concentrates on Red Hat Enterprise Linux (RHEL) High Availability Add-On.

3.1 Prerequisites and requirements

To deploy Pacemaker on Linux on IBM Power, you'll need to consider both hardware and software prerequisites:

► **Hardware Prerequisites:**

- IBM Power Server and IBM Power Virtual Server

Ensure you have an IBM Power System server with sufficient computing resources (CPU, RAM, and storage) to support your intended workload. The specific model and configuration will depend on your workload requirements.

- Network Infrastructure

A well-configured network infrastructure is essential for cluster communication. You need at least two network interfaces for redundancy. IBM Power servers typically have built-in network interfaces, but you can add additional network adapters if needed. For more information, refer to support policies for RHEL High Availability Clusters - Cluster Interconnect Network Interfaces using the following link:

<https://access.redhat.com/articles/3068841>

- Shared Storage

For high availability clustering, shared storage is crucial. You can use technologies like IBM Storage Solutions (e.g., IBM Spectrum Scale or IBM FlashSystem) or SAN (Storage Area Network) devices to provide shared storage for your cluster. For more information about support policies, requirements, and limitations for RHEL Resilient Storage components running on High Availability clusters refer to:

<https://access.redhat.com/articles/3130101>

► **Software Prerequisites:**

- Linux Distribution:

Choose a Linux distribution that is certified and supported for use with Pacemaker on IBM Power Systems. Red Hat Enterprise Linux (RHEL), SUSE Linux Enterprise Server (SLES), and Ubuntu Server are commonly used options:

- Supported RHEL High Availability releases on IBM Power architectures:

Red Hat provides support for the following combinations of RHEL High Availability running on IBM Power architectures. Table 3-1 provides an overview of Red Hat's prerequisites for utilizing RHEL High Availability components on IBM Power platform.

Table 3-1 Supported RHEL HA Add-On on Linux on Power

IBM Power Server	Supported RHEL HA Add-On version
IBM Power8®	Supported with RHEL 8. Supported with RHEL 7.4 or later - using corosync-2.4.0-9.el7 or later and Pacemaker-1.1.16-8.el7 or later.
IBM Power9®	Supported with RHEL 8.
IBM Power10	Supported with RHEL 8.4 or later.

- Support for RHEL High Availability on IBM Power bare-metal servers

Table 3-2 describes the support of RHEL High Availability clusters running on IBM Power bare-metal running OPAL firmware.

Table 3-2 Supported RHEL HA Add-On on IBM Power bare-metal

RHEL version	Supported on IBM Power bare-metal
RHEL 8	Supported by Red Hat.
RHEL 7	Supported by Red Hat with RHEL 7.5 or later - with corosync-2.4.3-2.el7 or later and Pacemaker-1.1.18-11.el7 or later.

Note: RHEL High Availability is not officially supported on IBM PowerKVM, and Red Hat does not offer assistance for RHEL High Availability in PowerKVM virtualized environments. Additionally, running KVM virtual machines (cluster nodes) on an IBM Power host (hypervisor) is not within the scope of supported configurations.

- IBM Power Virtual Server (PowerVS) LPARs as RHEL High Availability members
Red Hat only supports IBM Power Virtual Server (PowerVS) LPARs as RHEL High Availability members on the following versions shown in Table 3-3.

Table 3-3 Supported on PowerVS

RHEL version	Supported on IBM PowerVS
RHEL 8.4.z	Errata RHBA-2023:0503 with the following package(s): fence-agents-ibm-powervs-4.2.1-65.el8_4.11. or later for RHEL 8.4.z.
RHEL 8.5.z	Not supported.
RHEL 8.6	RHEL 8.6: Errata RHBA-2023:0428 with the following package(s): fence-agents-ibm-powervs-4.2.1-89.el8_6.8 or later.
RHEL 8.7	Errata RHBA-2023:0843 with the following packages(s): fence-agents-ibm-powervs-4.2.1-103.el8_7.1 or later for RHEL 8.7.z.

- Pacemaker

Install the Pacemaker software on all cluster nodes. You can typically find Pacemaker in the distribution's package repositories. Use the package manager (e.g., yum for RHEL/CentOS or zypper for SLES) to install it.

- Corosync

Pacemaker relies on Corosync for cluster communication. Ensure Corosync is installed and configured correctly on all cluster nodes.

- Resource Agents

Depending on your applications and services, you may need resource agents specific to your workload. These agents allow Pacemaker to manage resources like IP addresses, services, and file systems. Install and configure the necessary resource agents.

- Cluster Configuration

Create a cluster configuration file (typically `/etc/corosync/corosync.conf`) to define cluster properties such as nodes, rings, and quorum settings. Also, configure Pacemaker's cluster configuration (usually in `/etc/cluster/cluster.conf`).

- Cluster Authentication

Implement secure authentication between cluster nodes. This often involves setting up shared keys or certificates to ensure secure communication.

– Fencing Devices

Configure fencing devices to ensure the reliability of your cluster. IBM Power Systems may support various fencing methods. The following are requirements around fencing, fence-devices, and STONITH in a RHEL High Availability cluster:

- STONITH/fencing must be enabled.
- Every node must be managed by a fence device.
- Clusters with shared block storage or DLM require power or storage-based devices.
- Clusters with no shared block storage or DLM may use alternative agents and manual fencing.

Note: Red Hat's support for environments that employ fence agents not officially supplied or endorsed by Red Hat is restricted. In cases where fencing actions are required, Red Hat Support may refrain from providing assistance or engaging in investigations for cluster deployments utilizing non-Red Hat supported fence agents.

– Monitoring and Management Tools

Consider using monitoring and management tools like Hawk or pcs (Pacemaker Configuration System) for easier cluster management and monitoring.

3.2 Installing and Configuring Pacemaker environment

The following are required steps to install and configure Pacemaker on IBM Power Server and IBM Power Virtual Server

1. Install RHEL HA Add-On software

- a. Check that the RHEL High Availability repository is enabled. On both nodes, use the command shown in Example 3-1.

Example 3-1 Check RHEL repositories

```
$ dnf repolist
```

- b. If the HA repository is missing, use the commands in Example 3-2 to enable it.

Example 3-2 Enable RHEL HA repositories

```
$ subscription-manager repos --enable="rhel-8-for-ppc64le-highavailability-e4s-rpms"
$ dnf clean all
$ dnf repolist
```

- c. On both nodes install the RHEL HA Add-On software packages.

- i. For on-premises IBM Power systems run the command shown in Example 3-3.

Example 3-3 Install RHEL HA Add-On on IBM Power System

```
$ dnf install -y pcs pacemaker fence-agents-all
```

- ii. If you are running on IBM PowerVS instances, use the commands in Example 3-4 on page 39.

Example 3-4 Install RHEL HA Add-On on IBM Power Virtual Server

```
$ dnf install -y pcs pacemaker fence-agents-ibm-powervs
```

d. Make sure that you install the minimal version of the fence-agents-ibm-powervs package dependent on your Red Hat Enterprise Linux release:

- RHEL 8.4 fence-agents-ibm-powervs-4.2.1-65.el8_4.11.noarch.rpm (<https://access.redhat.com/errata/RHBA-2023:0503>)
- RHEL 8.6 fence-agents-ibm-powervs-4.2.1-89.el8_6.8.noarch.rpm (<https://access.redhat.com/errata/RHBA-2023:0428>)

2. Configure Cluster Nodes

Edit the Corosync configuration file `/etc/corosync/corosync.conf` to define the cluster properties, such as the ring settings and node information. See Example 3-5.

Example 3-5 Edit Corosync config file

```
$ sudo nano /etc/corosync/corosync.conf
```

Example 3-6 is an example Corosync configuration.

Example 3-6 Example of Corosync config file

```
$ sudo nano /etc/corosync/corosync.conf
totem {
version: 2
secauth: off
cluster_name: mycluster
transport: udpu }
nodelist {
node {
ring0_addr: <Node1_IP>
name: <Node1_Name>
}
node
{ ring0_addr: <Node2_IP>
name: <Node2_Name>
}
}
quorum {
provider: corosync_votequorum
}
```

3. Configure Pacemaker Cluster

Create a Pacemaker cluster configuration. Edit `/etc/cluster/cluster.conf` as shown in Example 3-7.

Example 3-7 Edit cluster.conf config file

```
$ sudo nano /etc/cluster/cluster.conf
```

In Example 3-8 we provide an example of Pacemaker cluster configuration.

Example 3-8 Example of cluster.conf file

```
<?xml version="1.0"?>
<cluster>
<cman expected_votes="2" transport="udpu">
```

```

<interface name="eth0" ringnumber="0" mcastaddr="226.94.1.1" mcastport="5405"/> </cman>
<clusternodes>
<clusternode name="<Node1_Name>" nodeid="1">
<fence>
<method name="fence_xvm">
<device name="<Fence_Device_Name>" />
</method>
</fence>
</clusternode>
<clusternode name="<Node2_Name>" nodeid="2"> <fence> <method name="fence_xvm"> <device
name="<Fence_Device_Name>" />
</method>
</fence>
</clusternode>
</clusternodes>
<fencedevices>
<fencedevice name="<Fence_Device_Name>" agent="<Fence_Agent_Name>" /> </fencedevices>
<cfs/>
</cluster>

```

4. Start and Enable Services

Start and enable the Pacemaker and Corosync services on both cluster nodes with the commands shown in Example 3-9.

Example 3-9 Start and enable services

```

$ sudo systemctl start corosync
$ sudo systemctl enable corosync
$ sudo systemctl start pacemaker
$ sudo systemctl enable pacemaker

```

5. Configure Resource Agents

Depending on your workload, configure Pacemaker resource agents for services like IP addresses, services, and file systems. Example 3-10 shows an example of configuring a simple IP address resource.

Example 3-10 Configure Resource Agents

```

$ sudo pcs resource create my_ip ocf:heartbeat:IPaddr2 ip=<Virtual_IP>
cidr_netmask=<Subnet_Mask> op_monitor interval=30s

```

6. Create and Start the Cluster

Create and start the cluster using the **pcs** command as shown in Example 3-11.

Example 3-11 Create and start cluster

```

$ sudo pcs cluster setup --name mycluster <Node1_Name> <Node2_Name>
$ sudo pcs cluster start --all

```

7. Verify the Cluster Status

Check the cluster status to ensure it's running correctly. The command is shown in Example 3-12.

Example 3-12 Verify cluster status

```

$ sudo pcs cluster status

```

8. Resource Management

You can manage resources with the `pcs` command. For example, to start the IP address resource use the command shown in Example 3-13.

Example 3-13 Start IP address resource

```
$ sudo pcs resource enable my_ip
```



4

High Availability for Your IBM Power Virtual Server Environment

This chapter provides an introduction to the IBM Power Virtual Server offering. IBM Power Virtual server is a cloud solution based on IBM Power server which is deployed within IBM data centers across the world – collocated with other IBM Cloud components.

The IBM Power Virtual Server offering is capable of hosting IBM AIX, IBM i and Linux on IBM Power (Red Hat, SUSE and Ubuntu) workloads on IBM Power systems. The IBM Power systems hosting these workloads uses their own dedicated networking infrastructure within the IBM Cloud location. The IBM Power VS offering uses Cisco ACI networking within IBM Cloud facilities to connect with the rest of the IBM Cloud services internally and to connect externally to customer on-premises infrastructure components. The storage connected to IBM Power Virtual Server offering is based on IBM FlashSystem storage subsystems that are directly connected using Fibre channel protocol via dual Brocade SAN fabrics.

This chapter contains following list of topics and discuss in detail to provide a complete understanding of IBM Power Virtual Server offering and its building blocks that provide Power Virtual Server as a Service within IBM Cloud including high availability options within the IBM Power Virtual Server offering:

- ▶ 4.1, “Introduction to Power Virtual Server offering” on page 44
- ▶ 4.2, “IBM Power Virtual Server High Availability” on page 46
- ▶ 4.3, “IBM Power Virtual Server Networking Options” on page 48
- ▶ 4.4, “IBM Power Virtual Server storage options” on page 58

4.1 Introduction to Power Virtual Server offering

IBM Power Virtual Server offering is a hosted infrastructure that is fully virtualized and fully integrated with the IBM Cloud Console which provides the public access to consume all services offered in IBM Cloud. The offering provides IBM Power compute based virtual instances, associated IBM storage and virtual network capabilities that is able to host and run all types of IBM and other Operating Systems such as AIX, IBM i and Linux (Red Hat, SUSE and Ubuntu) on fully virtualized Logical **PAR**titions (LPARs). The IBM Power Virtual Server offering is a fully managed Infrastructure as a Service (IaaS) offering based on IBM Power Systems and associated IBM Storage.

The IBM Power Virtual Server offering uses dedicated Cisco ACI networking coupled with virtual networking provided by IBM PowerVM hypervisor within the colocation. The IBM Power Virtual Server networking extends to the rest of the IBM Cloud allowing full and complete set of networking capabilities to connect with other IBM Cloud Services and customer services internally as well as externally. The network expands to provide secure networking for all types of connectivity to customers – with various configuration options available – and is capable of expanding to connect to all IBM Cloud capabilities in an integrated fashion.

The storage used in IBM Power Virtual Server offering is based on IBM FlashSystems and there are two types of virtualized storage capabilities available to use for workloads. The Tier 1 storage is capable of up to 10 IOPS/GB and Tier 3 storage is capable of up to 3 IOPS/GB. Customers can choose their tier based on their workload performance requirements.

Note: Tier 3 storage may not be suitable for production workloads due to mission critical applications performance requirements.

The systems used within IBM Power Virtual Server offering are IBM Power Systems based on either Power9 or Power10 architecture. The capabilities within these IBM Power systems are no different in this offering from the IBM Power systems that are run by customers on-premise today. They have the same proven IBM advanced Reliability, Availability and Serviceability (RAS) features built in.

The compute used within IBM Power Virtual Server offering can be scaled up or down based on customer needs. The virtualization layer is fully managed by IBM Cloud Team and maintained to the highest possible availability based on the deployment model for customer workloads. The location of specific workloads can be designated using affinity or anti-affinity rules to meet production requirements. In addition it is possible to pin workloads to a specific server within if required for performance or compliance requirements.

The infrastructure is continuously patched and kept up-to-date all the time by IBM Cloud operations team, IBM is take responsibility for managing the hardware to remove any complexities. IBM is responsible for managing the hardware and virtualization layer and the customer is responsible for managing the operating system and above. IBM is also responsible for maintaining all IBM storage, SAN and networking hardware including firmware which is kept patched and up-to-date by the IBM Cloud operations team.

One of the main advantages of using IBM Power Virtual Server offering is that any enterprise software and workloads that are currently operating on customer premises are certified to run on IBM Power Virtual Server with a simple lift and shift of those workloads. In this method, the workload will continue to run similar to how it ran previously on-premises. Unlike in on-premises environments – where the customer is expected to maintain currency on hardware, virtualization on IBM Power systems and Storage – the IBM Cloud operations team will maintain and provide timely updates and upgrades with no disruption to running workloads allowing customers to

focus more on actual business workloads rather than the infrastructure from Operating Systems and above on each LPAR.

Figure 4-1 provides a responsibility matrix showing the services management roles and responsibilities by either the IBM Cloud team or Client responsibilities for IBM Power VS offering components.

On-Premises	Power Systems Virtual Servers	Platform as a Service	Software as a Service
Applications	Applications	Applications	Applications
Data	Data	Data	Data
Runtime	Runtime	Runtime	Runtime
Middleware	Middleware	Middleware	Middleware
Operating system	Operating system	Operating system	Operating system
Virtualization	Virtualization	Virtualization	Virtualization
Servers	Servers	Servers	Servers
Storage	Storage	Storage	Storage
Networking	Networking	Networking	Networking

Client manages
 IBM manages

Figure 4-1 IBM Power Systems Virtual Server responsibility assignment matrix¹

4.1.1 IBM Power Virtual Server Compute Hardware Building Blocks

IBM Power Virtual Server deployments are built using selected hardware models based on following IBM Power Hardware systems:

- IBM Power9 Enterprise E980 Systems (9080-M9S)
- IBM Power9 Scaleout S922 Systems (9009-22A or 9009-22G)
- IBM Power10 Enterprise E1080 Systems (9080-HEX)
- IBM Power10 Scaleout S1022 Systems (9105-22A)

The latest list of available IBM Power systems within IBM Power Virtual Server offering can be found for each IBM Cloud location in this [IBM Cloud document](#).

Note: The availability of IBM Power systems may vary from one IBM Cloud location to another. Additional IBM Power systems may get deployed over time to each location.

4.1.2 IBM Power Virtual Server Network Hardware Building Blocks

IBM Power Virtual Server deployments are using the following Network infrastructure devices and capabilities as building blocks to provide all network services within IBM Power Virtual Server Offering.

- Cisco Nexus 9000 (25Gb²)
- Cisco Nexus 9000 (10G)

¹ The latest service matrix can be found in link <https://cloud.ibm.com/docs/power-iaas?topic=power-iaas-about-virtual-server>

² Latest Deployments on S1022 and E1080 systems will be using 100Gb capable SR-IOV adapters.

- Cisco Nexus 9000 (1G)
- Cisco ASR1001-HX Router
- Cisco UCS - APIC controller
- Avocent ACS 8048, ACS 8016, ACS 8032 DAC-400

The latest list of network devices used within IBM Power Virtual Server offering can be found for each IBM Cloud location in this [IBM Cloud document](#).

Note: The deployment of network devices may vary from one IBM Cloud location to another. Additional devices and technologies may get deployed over time.

4.1.3 IBM Power Virtual Server Storage Hardware Building Blocks

IBM Power Virtual Server deployments are using 32GB Brocade SAN infrastructure with following types of devices as building blocks as storage providers.

- IBM FlashSystems Storage based on IBM FS9x00 series devices
- IBM Storwize V7000 with SSD drives
- Brocade 32Gb SAN

The latest list of storage devices used within IBM Power Virtual Server offering can be found for each IBM Cloud location in this [IBM Cloud document](#).

Note: The deployment of IBM Storage devices may vary from one IBM Cloud location to another. Additional devices and technologies may get deployed over time.

4.2 IBM Power Virtual Server High Availability

This section provides an overview of IBM Power Systems Virtual Server infrastructure related High Availability capabilities and Disaster Recovery options available within IBM Cloud deployments.

The IBM Power Virtual Server offering is designed with the concept of eliminating any single point of failure (SPOF) within the underlying infrastructure deployed to host the IBM AIX, IBM i or Linux on Power (Red Hat, SUSE and Ubuntu) workloads. This is achieved through the use of redundant physical hardware components and virtualization technologies built using IBM Power Systems availability features.

The use of dual hardware adapters, each connected to different storage connections or network switches, provides the basic high availability eliminating any effect of a failure of any single hardware component. The IBM Power Virtual Server infrastructure also uses dual SAN Fabrics for storage access, providing availability in case one fabric fails as the secondary fabrics will continue to service the workloads without any outages. Similarly, there are multiple network switches used to eliminate any SPOF for network gear failures utilized by the IBM Power Virtual Server offering.

The Network and Storage IO high availability within each system is achieved using PowerVM based dual Virtual IO Servers that work in pairs within each IBM Power system building block. This allows each system to service workloads with no outages during service maintenance. If any VIO Server fails, all services will continue to operate and be available through the remaining VIO Server. In the unlikely case of dual failures or entire IBM Power system failure, the workloads are able to be remotely restarted on another server allowing the workloads to resume services automatically using the remote restart capability built in to the IBM Power

Virtual Server offering. This level of workload remote restart capability provides basic disaster recovery capabilities to continue your business critical applications and ensures that services are automatically relocated and restarted within IBM Power Systems Virtual Server offering.

For any critical workloads you may want to deploy on IBM Power Virtual Server with more advanced High Availability solutions, locally within a site across multiple availability zones or across IBM Clouds with Geographic or Metro mirror distances. IBM Power Virtual Servers supports the deployment of IBM PowerHA SystemMirror providing automatic application failovers for workloads in your IBM Power Virtual Server environment running on AIX or IBM i. Similarly, for Linux workloads, Pacemaker based cluster solutions can be deployed with advanced application and middleware configurations to automatically monitor and failover workloads based on your application HA-DR requirements.

When applications or databases are deployed to run on IBM Power Virtual Server, there are many options available to maintain high availability and business continuity for such workloads. Figure 4-2 describes different business continuity plans that can be used with the IBM Power Virtual Server Offering.

	PowerVS Remote Restart	HA Clusters	Global Replication Service	Backup and Restore
Description	PowerVS remotely restarts VMs from a failed host to another host.	OS Clustering, one OS in the cluster has access to the data, multiple active OS instances on all nodes in the cluster. Application is restarted on a secondary node upon outage event.	Global Mirror Change Volumes provides the technology that provides asynchronous replication and advanced network configuration for fast data transfer.	No clustering technologies required. Restore the VM in case a critical failure occurs.
Outage Types	Hardware planned and unplanned.	Hardware or software, planned or unplanned.	Hardware planned and unplanned.	Hardware unplanned.
Responsibility	PowerVS	Customer	Customer	Customer
OS Integration	OS agnostic	Inside of OS	OS agnostic	OS agnostic
Resources/Licenses	N + 0	N + N	N + N	N + 0

Figure 4-2 IBM Power Virtual Server HA and DR Business Continuity Plan

The High Availability solutions can be created in the area of HA Clusters in above chart for many applications and services running on Linux operating systems on IBM Power VS solution.

One of the major limitation when deploying HA clusters with Pacemaker or PowerHA within IBM Power VS offering is that there is no way to stop an active node like you would in an on-premise environment where you have access to the Hardware Management Consoles, Novalink or VIO Servers within your infrastructure. You need to be able to stop a node (fence) to control the cluster resources and prevent a split-brain scenario where nodes are isolated and no mechanism to kill one or the other node by the quorum of active nodes. IBM PowerVS supports a fence agent – *fence_ibm_powervs* – which is the only supported agent for a STONITH device on Power Virtual Server clusters. The fence agent connects to the Power Cloud API using parameters to identify an IBM Power Virtual Server instance to stop the resource as required.

4.3 IBM Power Virtual Server Networking Options

IBM Power Virtual Server networking is built using multiple virtualization and physical technologies to securely communicate and operate with multiple tenants in a secure fashion maintaining their segregation and isolation as a priority at all levels of networking.

When an IBM Power Virtual Server instance is created, it can be placed on two types of networks at an OS level.

► **Public Network:**

Easy and fast method to connect over the internet and protected by a firewall allowing only few connections through to the OS instance:

- SSH
- 5250 Terminal emulation with SSL
- HTTPS
- Ping

This network connection is used for initial configuration over a public network to a Operating System Instance where there is no private connectivity available. It is common to remove this interface after the initial deployment of a Virtual Server instance and then only use private network interfaces.

► **Private Network:**

The private network allows to have an interface separate from the public interface which is allowed to communicate with rest of the OS instances within same private VLAN once enabled.

This network can be used to communicate with various IBM Cloud services and application services as needed by the customer. Further, this can be extended to communicate with external networks in a secure fashion as required.

Note: The public network interface can be removed from the IBM Cloud Console at any time and added if needed.

4.3.1 IBM Power Systems Virtualization

IBM Power Systems are built with a rich set of virtualization features that are capable of running enterprise class workloads in a secure way utilizing the IBM PowerVM Hypervisor that enables these virtualized networking features and capabilities on IBM Power hardware.

Some of the relevant key features that are used within PowerVM virtualization technology are discussed in this section to help you understand the features available in the IBM Power Virtual Server offering.

► **Network Virtualization:**

IBM PowerVM is implemented with IEEE VLAN compatible virtual switches in the hypervisor that are capable of providing up to 4,094 VLANs per virtual switch. This provides an abstraction layer between physical network hardware and logical partitions and allows the sharing of one or more physical resources across multiple LPARs while maintaining the segregation and separation between each VLAN as per industry standards.

The VLAN concept allows two LPARs to communicate with each other using the channels between them over the hypervisor vSwitch without leaving the hypervisor – this is the

default mode of operation. The virtual switches are able to fundamentally operate in two different modes:

– **Virtual Ethernet Bridge (VEB) Mode:**

In this mode, all LPARs within a chassis on a virtual switch are allowed to communicate with each other freely within its respective VLAN using IP addresses at layer 3. However, the VLAN boundary is not allowed to be crossed, maintaining the VLAN separation at Layer 2. This is the default mode of operation for vSwitches when created in the Power Hypervisor. Figure 4-3 illustrates VEB mode operations with a vSwitch.

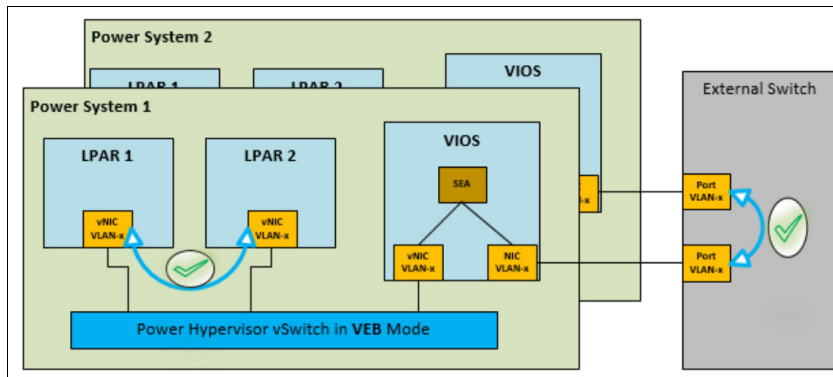


Figure 4-3 IBM Power Hypervisor vSwitch in VEB mode

– **Virtual Ethernet Port Aggregator (VEPA) Mode:**

In this mode, none of the LPARs within a chassis on a virtual switch are able to communicate with each other, even if they are placed on the same VLAN ID with an IP address. The only way to communicate with another LPAR within the same chassis on same vSwitch is to send the network packets to an external switch and to use reflective relay capability on external switches to send the packet back into the same vSwitch destined for another LPAR. The reflective relay allows the network traffic to perform a hair pin turn around and the packets return back to the vSwitch where it originated and then delivered to its intended destination at layer 3. Without using this reflective relay capability, two LPARs placed on a VEPA mode virtual switch will not be able to communicate successfully. Figure 4-4 illustrates VEPA mode restrictions imposed by the vSwitch.

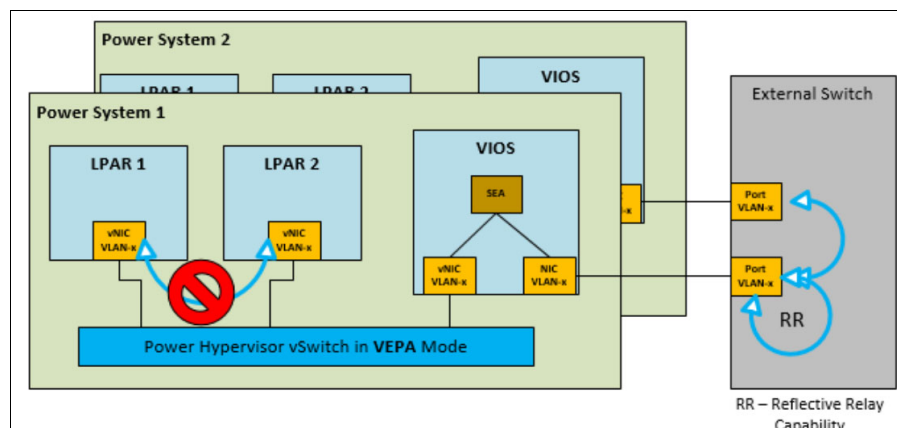


Figure 4-4 IBM Power Hypervisor vSwitch in VEPA mode

The network virtualization of physical hardware adapters is achieved by the use of Virtual IO Servers (VIOS). There are many different options available to implement and connect

using VIOS to allow Ethernet bridging or trunking to external networks. Single Root IO Virtualization (SR-IOV) is often used for client LPARs for additional adapter sharing together with the Shared Ethernet Adapter (SEA) within IBM Power Virtual Server deployments in IBM Cloud.

The IBM Power Virtual Server implementation uses two slightly different networking deployments dependent on the models of the IBM Power System being utilized.

– IBM S922 Systems:

As S922 Systems are limited in the number of available adapter slots, we use three 2-port adapters that are shared using SR-IOV Virtual Function (VF) between the two VIO Servers that service management, public and private networking for an S922 System.

The SEA configuration uses multiple SR-IOV based VFs to aggregate 2 ports coming from different adapters to eliminate any SPOF on each VIO Servers. In addition, this aggregation provides a higher bandwidth to the clients on a system

Refer to Figure 4-5 to understand how they are networking connections are deployed an how they are used for different purposes in the S922 systems.

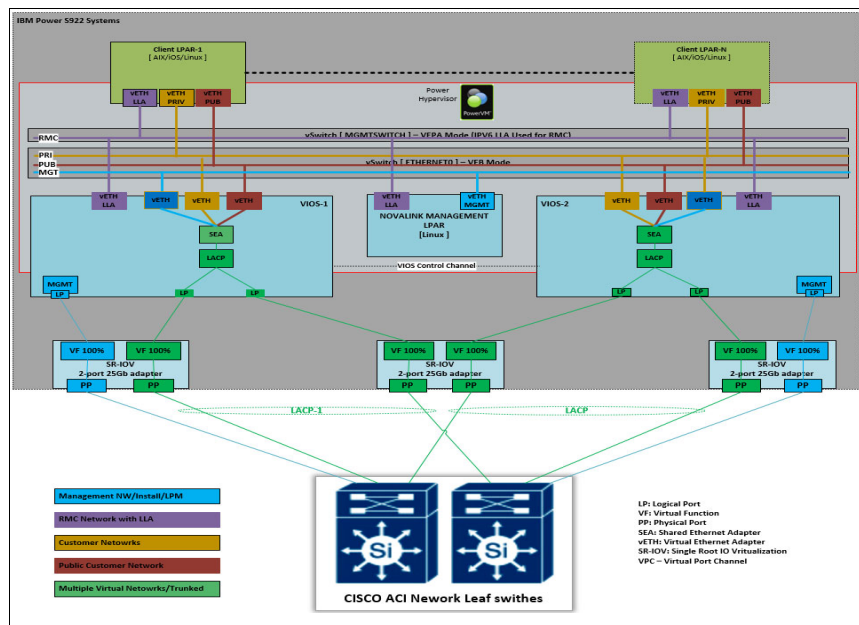


Figure 4-5 IBM Power VS S922 Network Virtualization

– IBM S1022 Systems:

The new deployments within IBM Power Virtual Server offering with IBM Power S1022 systems are using 100Gb SR-IOV adapters and has a slightly different architecture from IBM Power S922 systems. The SR-IOV capabilities are not used and the Physical ports from SR-IOV adapters are Link Aggregated instead. The same level of high availability is achieved by failing over to secondary VIO Server in case of an adapter failure. This architecture is illustrated in Figure 4-6 on page 51. Power 10 based S1022 systems are currently available in MAD02, MAD04, WDC07 and DAL10. Additional S1022 systems will be rolled out to existing and new data centers over time, based on client demand.

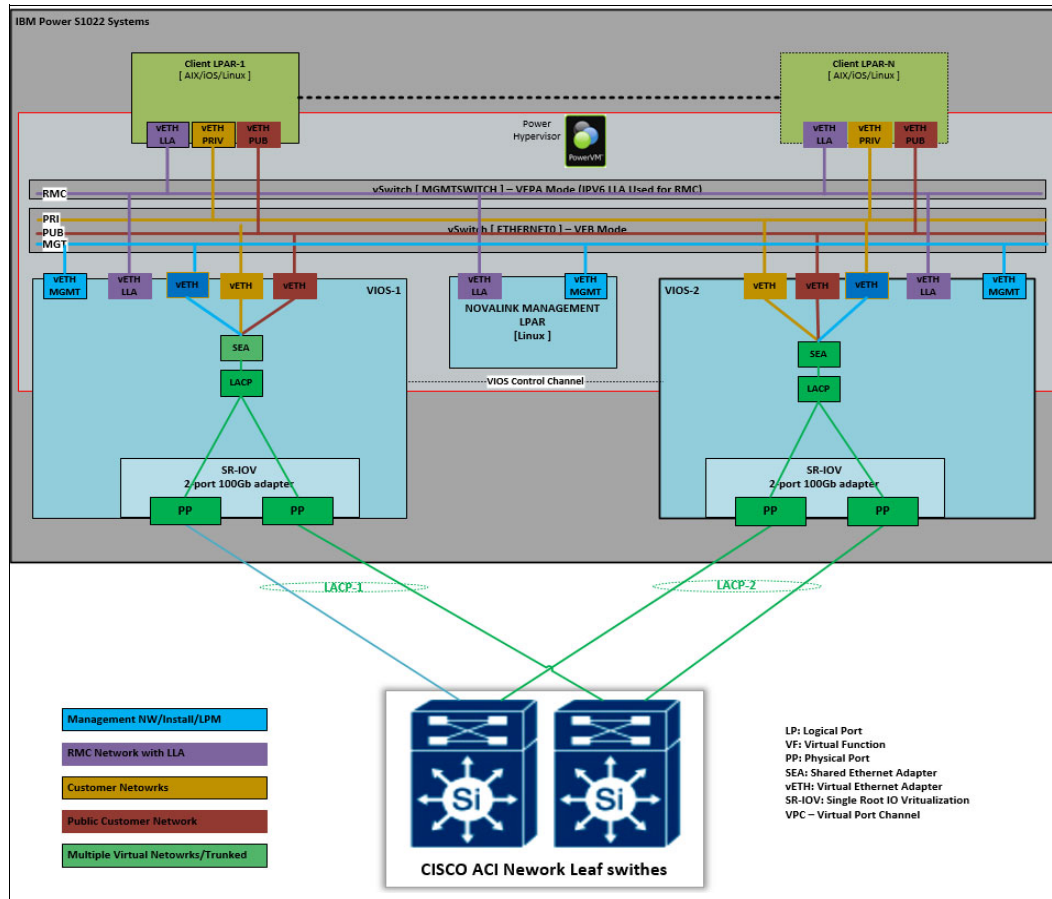


Figure 4-6 IBM Power Virtual Server S1022 Network Virtualization

– IBM E980 Systems:

The E980 systems follow a similar pattern, but with more physical adapters as it capable of having more adapter slots. There are 10 x 2-port SR-IOV adapters in use for each of these systems and it uses each port as a SR-IOV VF to create an SEA that services the management, public and private networks in each frame using dual VIO Servers.

The use of additional adapters/ports increases the network availability and bandwidth per VIO Server when connecting to the external networks. Multiple adapters are aggregated using IEEE LACP using multiple VFs provided by SR-IOV adapter ports to the external Cisco ACI network.

Refer to Figure 4-7 on page 52 to understand how the networking connections are setup and how they are used for different networking solutions with each E980 system.

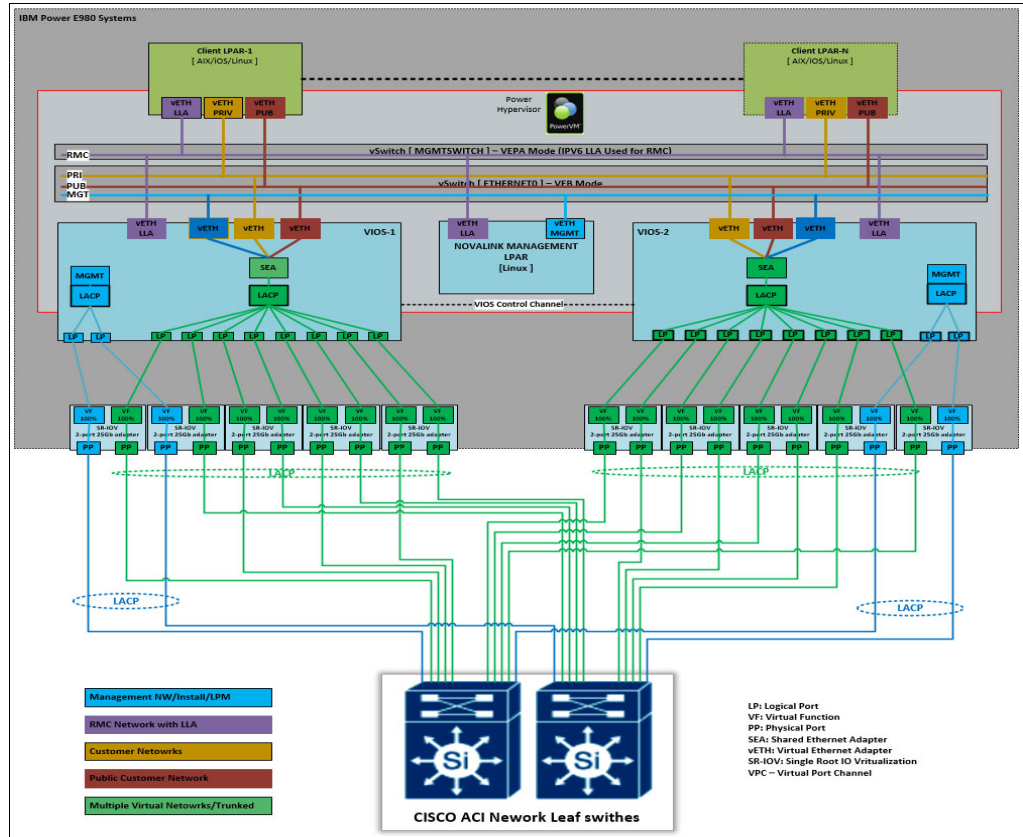


Figure 4-7 IBM Power VS E980 Network Virtualization

– IBM E1080 Systems:

The new deployments within IBM Power Virtual Server offering for E1080 systems are to be used with multiple 100Gb SR-IOV adapters. The architecture is being currently tested at the time of writing this publication. The most likely configuration will be based on the use of multiple 100Gb physical ports using SR-IOV adapters on an LACP aggregation. Deployment of E1080 systems will be based on customer demand.

► Storage Virtualization:

There are three main methods that a client LPAR could map external or internal storage used in IBM Power Systems using VIO Servers.

- Shared Storage Pools [SSP]
- Virtual SCSI [vSCSI]
- N-Port Id Virtualization [NPIV]

Each of above virtualization methods uses a different way to provide virtualized storage to LPARs hosted on a IBM power System. In the IBM Power VS offering deployments, SSPs are not in use for any purpose. Hence, SSPs are not discussed.

Virtual SCSI [vSCSI]

The vSCSI option is used between VIO Server and Client LPARs to map disks that are visible on VIO Servers. Each VIO Server will map the same SAN disk device provisioned from external SAN storage. The use of vSCSI is limited for to the Novalink LPAR and not used with any client related LPARs. The boot disk for Novalink LPAR is mapped from each VIO Server using vSCSI. The disk is a protected highly available volume provided by the IBM Storage subsystem eliminating any requirement to setup mirroring on the Novalink partition.

N-Port ID Virtualization [NPIV]

NPIV is an industry standard technology that provides the capability to assign a physical Fibre Channel adapter to multiple unique world wide port names (WWPN). To access physical storage from a SAN, the physical storage is mapped to logical units (LUNs) and the LUNs are mapped to the ports of physical Fibre Channel adapters. Then the PowerVM Virtual I/O Servers uses the maps to connect the LUNs to the virtual Fibre Channel adapter of the virtual I/O client as a pass through virtual adapter using PowerVM virtual slots.

The IBM PowerVS offering uses NPIV for all client LPARs to map storage from external storage providers. With the use of dual PowerVM VIO Servers, each mapped LUN is mapped providing highly available volumes to the client LPARs through Power VM virtualization layer maintaining a minimum of two paths from VIO Servers to maintain availability and performance with the use of multipath device drivers at the client Operating system level. The use of multipath device drivers will depend on the type of operating systems.

Refer to Figure 4-8 for further details on how vSCSI and NPIV based connectivity is deployed using physical adapters on IBM Power S922/S1022 Systems

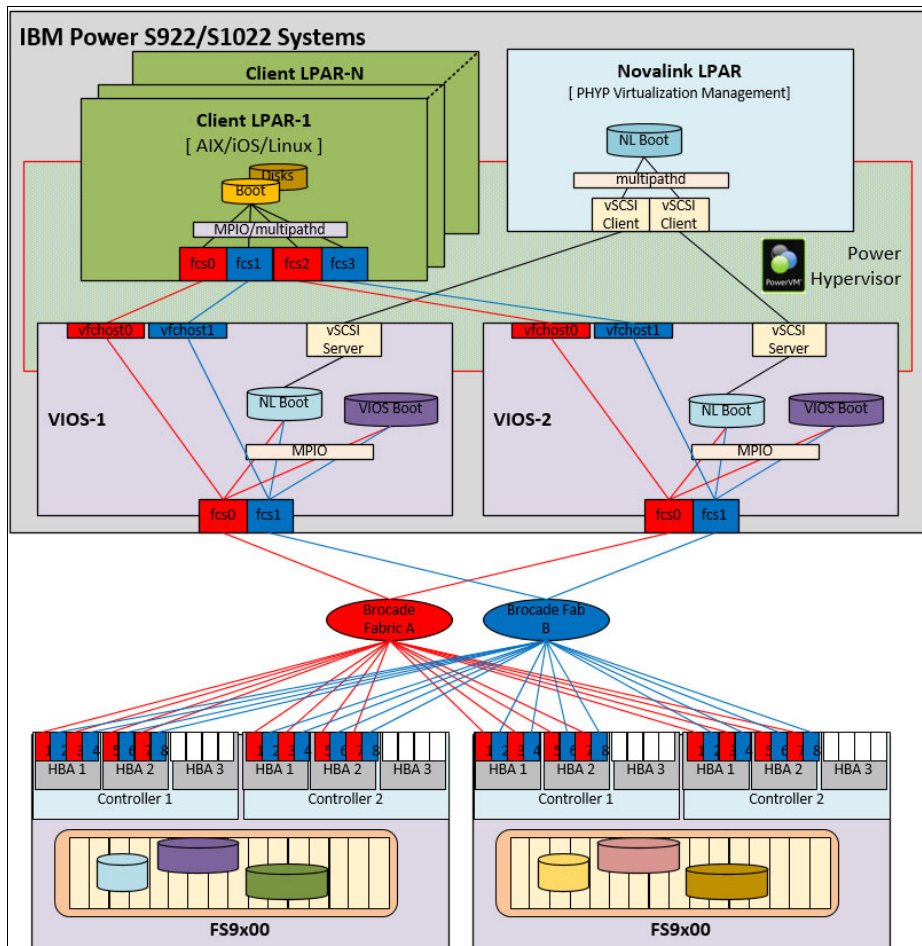


Figure 4-8 IBM Power Virtual Server vSCSI and NPIV FC Connectivity on IBM Power S922/S1022 Systems

Refer to Figure 4-9 on page 54 for further details on how vSCSI and NPIV based connectivity is deployed using physical adapters on IBM Power S922/S1022 Systems

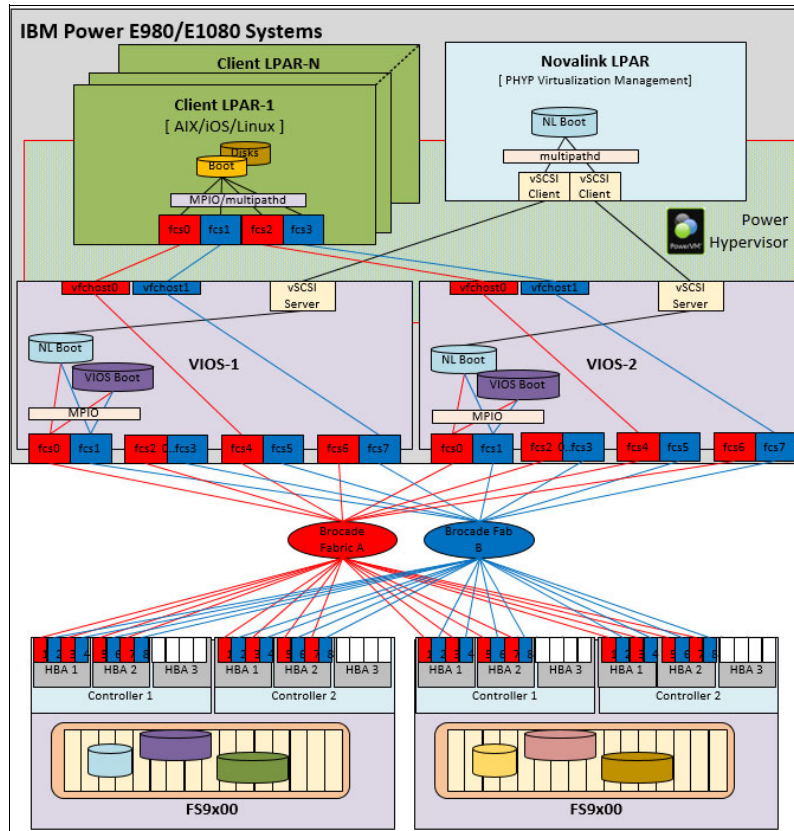


Figure 4-9 IBM Power Virtual Server vSCSI and NPIV FC Connectivity on IBM Power E980/E1080 Systems

► IBM PowerVM Novalink:

NovaLink is a software interface that is used for virtualization management. The NovaLink is installed on a PowerVM server within a Linux LPAR. PowerVM NovaLink enables highly scalable modern cloud management and deployment of critical enterprise workloads. The PowerVM NovaLink is used to provision large numbers of virtual machines on PowerVM servers quickly and in a robust way using OpenStack services. The provisioning and PowerVM management is enabled using You can manage the server through a representational state transfer application programming interface (REST API) or through a command-line interface (CLI). Novalink seamlessly integrate with IBM and non IBM Cloud managers such as PowerVC or IBM Cloud Console using REST API and Openstack.

The NovaLink runs on a Linux logical partition on a POWER8, POWER9, or Power10 processor-based systems that is virtualized by PowerVM. Each frame must have this Novalink LPAR deployed for Virtualization management using Novalink capabilities.

IBM PowerVM Novalink provides a way of distributing the virtualization management to each frame rather than from a single point like when the Hardware Management Consoles (HMC) are used. Using Novalink enables a larger number of LPARs and IBM Power systems to be managed compared to utilizing HMC management. This feature allows the maximum scalability to provide more virtual machines managed with better distribution across the managed servers.

PowerVM Novalink is a single instance that is running on each IBM Power System and may fail /restart without impacting to any active and running workloads on such frames. Hence, it is not seen as a single point of failure (SPOF). If the Novalink Partition is not operational, there may be limits to the ability make changes to any LPAR on the frame with dynamic resource changes, console access, live partition mobility actions and any power

control on an effected Power system. The full details of the Novalink capabilities are documented in IBM Knowledge Center in the following reference link:

<https://www.ibm.com/docs/en/power10/9105-42A?topic=environment-powervm-novalink>

4.3.2 Cisco ACI Networking

IBM Power Virtual Server offering uses Cisco's premier SDN solution which is Application Centric Infrastructure (ACI) networking which enables multi cloud environments with agility and resiliency. The IBM Power VS offering uses the Cisco ACI Spine and Leaf switch configuration to provide ACI capability.

As per Cisco definition, "Cisco ACI, the industry's most secure, open, and comprehensive software-defined networking (SDN) solution, enables automation that accelerates infrastructure deployment and governance, simplifies management to easily move workloads across a multifabric and multicloud frameworks, and proactively secures against risk arising from anywhere. It radically simplifies, optimizes, and expedites the application deployment lifecycle"³. Figure 4-10 shows a typical ACI networking construct as referenced from Cisco documentation.⁴

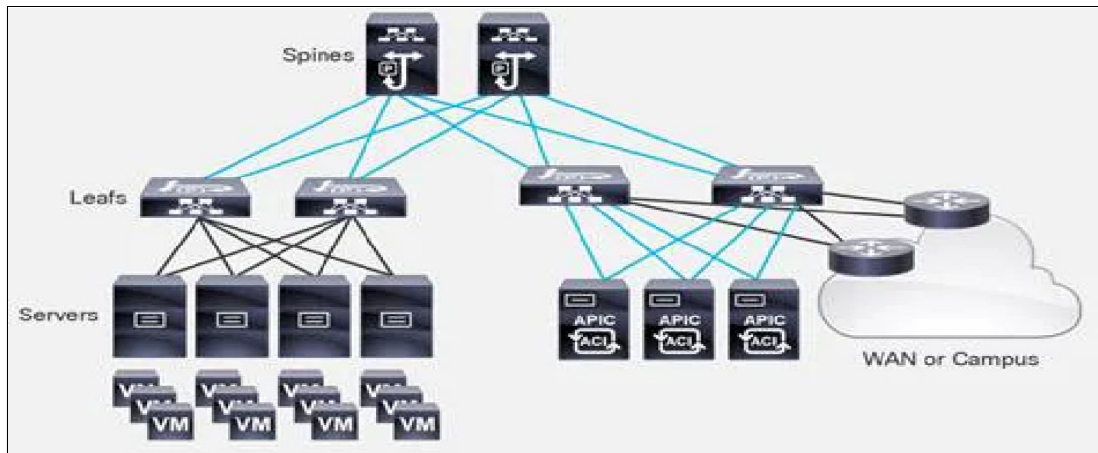


Figure 4-10 Cisco ACI Fabric⁴

4.3.3 External Networking

This section describes the available external network connectivity for IBM Power VS offering and it is not an exhaustive list of connections methods that are available today. The deployment of Power VS offering makes use of either Public or Private network interface within an LPAR.

► Internet:

Access the internet through resources that are hosted in any of the following infrastructure environments using layer 3 IPv4 traffic:

- IBM Cloud Classic
- IBM Virtual Private Cloud (VPC)
- IBM PowerVS

³ <https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/solution-overview-c22-741487.html>

⁴ https://www.cisco.com/c/dam/en/us/td/docs/dcn/whitepapers/cisco-application-centric-infrastructure-design-guide.docx/_jcr_content/renditions/cisco-application-centric-infrastructure-design-guide_1.jpg

► **Remote:**

Connect remote networks to your IBM Cloud networks using the following services which allow customers to connect to a remote network in IBM Cloud:

– **Internet VPN:**

Access over the public internet to connect remote networks through to IBM Cloud networks using a VPN. The VPN is terminated on gateway devices or a service within IBM Cloud.

– **Megaport or Direct Link 2.0 Connectivity:**

Direct Link is a suite of offerings that enable the creation of direct, private connections between your remote, on-premises network and IBM Cloud, without traversing the public internet.

Figure 4-11 illustrates Megaport or Direct Link 2.0 connectivity from IBM Power Virtual Server offering to on-premise networks.

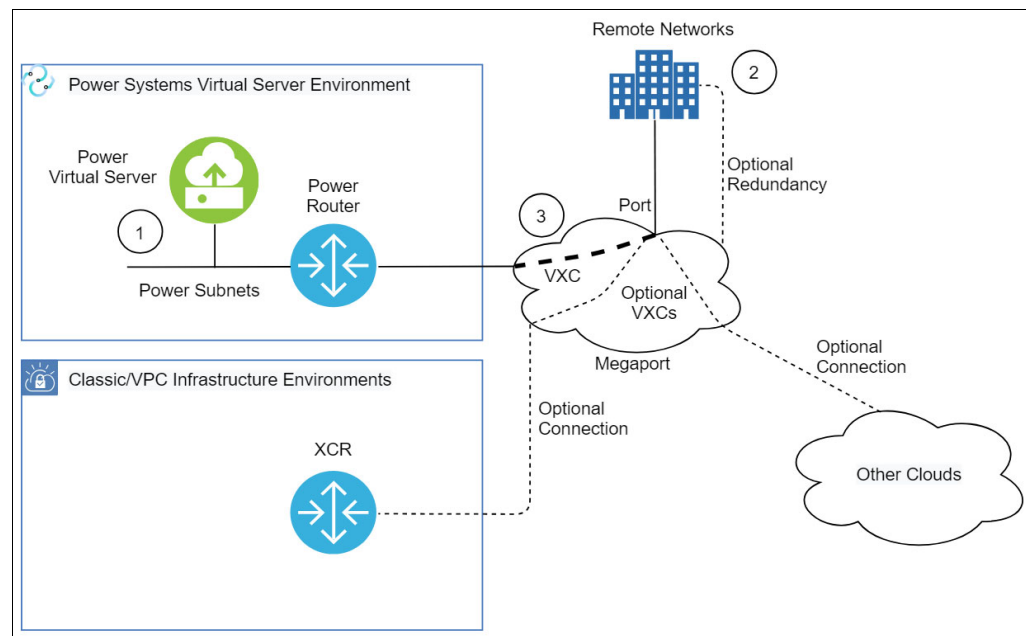


Figure 4-11 Connect to on-premise networks

The use of Direct Link 2.0 or Megaport seamlessly connects your on-premises resources to your cloud resources. The speed and reliability of Direct Link extends your organization's data center network and offers more consistent, higher-throughput connectivity, keeping traffic within the IBM Cloud network.

The Direct Links can connect to either a local or remote IBM Cloud Transit Gateway, which allows the on-premises network to access all networks that are connected to the IBM Cloud Transit Gateway.

Your Direct Link 2.0 or Megaport connections are location-specific. By default, IBM Cloud Direct Link is not a redundant service. You must order a separate Direct Link Connect instance for redundancy.

Transit Gateway enables you to connect your otherwise disconnected private networks, such as classic, VPC, and Direct Link. This allows to establish connection between multiple Power Systems Virtual Server workspaces across different data centers. Figure 4-12 on page 57 illustrates connectivity between multiple Power Systems Virtual Server locations with high availability (HA) and disaster recovery (DR) solutions.

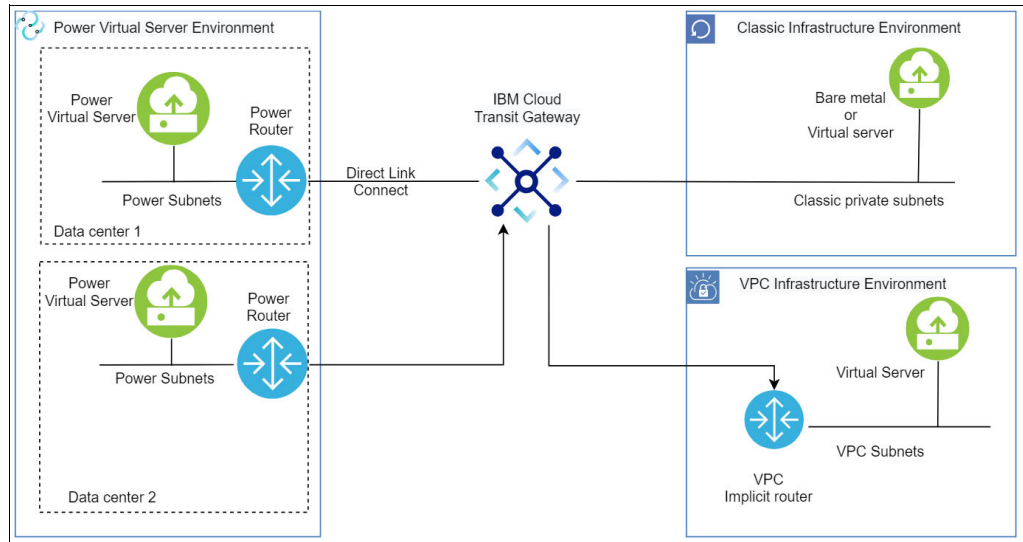


Figure 4-12 Transit Gateway deployment scenario with DL2.0 or Megaport

More details are found in this [IBM Cloud document](#).

► **Power Edge Router (PER):**

A Power Edge Router (PER) is a high-performance router that provides advanced routing capabilities for IBM Power Virtual Server offering. One of the benefits of the PER solution is that it makes it easier for a Power Systems Virtual Server instance to access other IBM Cloud services, such as IBM Cloud DNS, NTP, and Cloud Object Storage. This allows connectivity to these services without having to use proxies or virtual routers, as the PER solution includes a Network Address Translation (NAT) device that simplifies the access process.

Figure 4-13 illustrates the use of PER in IBM Cloud.

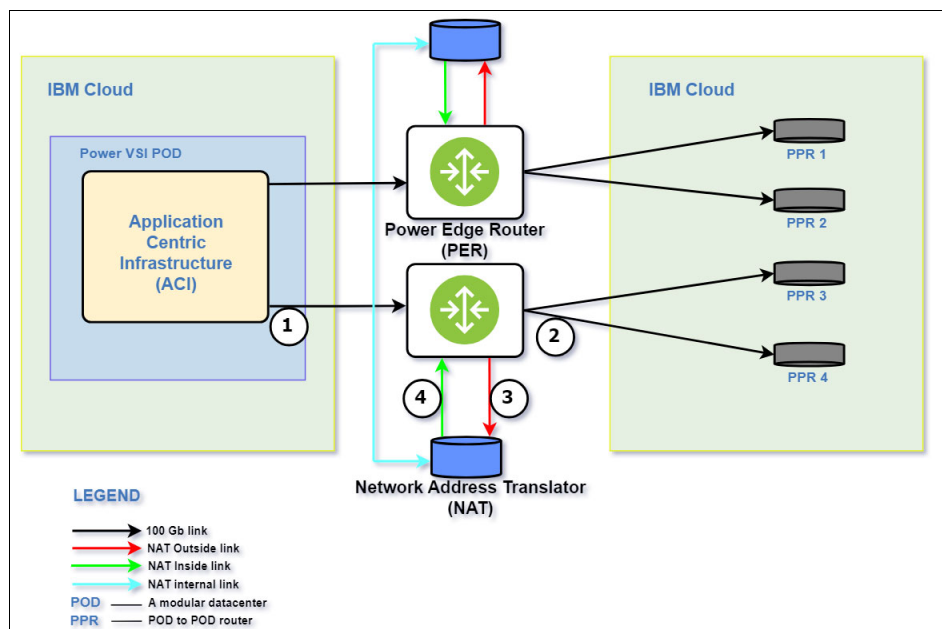


Figure 4-13 Power Edge Router Integration to IBM Cloud

PER improves network communication across different parts of IBM network. The PER solution creates a direct connection to the IBM Cloud MPLS (Multi Protocol Label Switching) backbone, making it easy for different parts of the IBM network to communicate with each other.

Note: At the time this document was written the PER solution is currently only available in DAL10. PER will be deployed in other data centers over time.

For a complete and up-to-date list of available external network connectivity details to IBM Power VS offering, refer to this [IBM Cloud document](#).

Important: Not all IBM Cloud locations provide same external connectivity and may depend on what is available depending on the available capabilities of local operators.

4.4 IBM Power Virtual Server storage options

This section explains various storage options for achieving high availability for Pacemaker clusters. Storage itself has a technology to achieve High availability through replication or IBM FlashCopy® features. High availability with Pacemaker cluster can be achieved using various technology in cloud as well on prem.

Some of the storage options are GRS (Global Replication Services), GFS2 (Global file system 2), Shared Storage technology and IBM Storage scale.

4.4.1 GFS2

GFS2 is a shared disk file system for Linux clusters. GFS2 allows all the nodes in cluster to have direct access to the same shared storage through the use of a locking protocol. All nodes in a GFS2 Linux cluster work as peer nodes.

Within Red Hat Enterprise Linux, GFS2 is provided by the Red Hat Enterprise Linux Resilient Storage Add-On which requires the RHEL High Availability Add-On to provide the Pacemaker clustering option.

Using GFS2 in the cluster requires the correct hardware to allow access to shared storage and a locking mechanism to control access to the storage. The shared storage is usually provided by a Storage Area Network using Fibre Channel, NVMe over Fabrics, or iSCSI. The locking manager is built into your Linux distribution and requires the use of a cluster manager such as Pacemaker.

More details about GFS2 can be found in this [GFS2 overview](#).

4.4.2 IBM Storage Scale

IBM Storage scale is a high performance, highly available, clustered file system and associated management software, available on different platforms. IBM storage scale can scale in several dimensions, including performance, capacity and number of nodes that can mount or access the file system. IBM storage scale addresses the needs of applications whose performance to capacity ratio demands cannot be achieved by traditional scale up storage systems, and IBM storage scale therefore deployed for many I/O demanding enterprise applications that requires high performance or scale. IBM Storage scale provide

various configuration options, access methods and many other features such as snapshots, compression and encryption.

For more details about IBM Storage Scale refer to this [IBM Storage Scale overview](#).

4.4.3 Global Replication Services (GRS)

Data replication is key to business resiliency because simply put, data drives decision-making. Data informs and feeds into mission-critical processes, analytics, systems and, ultimately, business insights. Organizations must guarantee that it is constantly available and accessible to users in near real-time. As enterprises develop across geographies and platforms, replication enables them to scale in tandem with their expanding data requirements while maintaining performance. GRS reflects IBM's commitment to enabling business continuity planning, operational excellence and optimization in cost with IBM Power Systems.

Global Replication Service (GRS) is based on well-known, industry-standard IBM Storwize Global Mirror Change Volume Asynchronous replication technology. Global Replication Service on IBM Power Virtual Server exposes the application programming interface (API)/command line interface (CLI) to create and manage replication enabled volumes.

4.4.4 Configuration of High Availability Storage with Pacemaker

High availability Storage with Pacemaker is a solution for ensuring that critical data and applications remain available and accessible even in case of hardware failure or other disruptions. A Pacemaker cluster manager allows you to configure and manage clusters of computers that work together to provide high availability for critical data and applications respectively.

The High availability storage with Pacemaker solution typically involves cluster of multiple nodes, of which each nodes have access to shared storage. This shared disk from the storage, of which each nodes have access to shared storage. This shared disk from the storage are used for storing critical application data, which can be accessed by any nodes of cluster. In case of any failure, the Pacemaker cluster automatically fails to another node in the cluster, ensuring that the data and application remain available.

This section explains how to configure High Availability storage with Pacemaker on RHEL.

As a prerequisite, all the nodes of cluster should have appropriate shared disk from same backend storage.

1. Install Pacemaker on all nodes in the cluster using the command shown in Example 4-1.

Example 4-1 Install Pacemaker

```
$ sudo dnf install pacemaker pcs fence-agents-all
```

This command will install the Pacemaker software, the `pcs` command line and collection of fencing agents.

2. Configure the firewall to allow the access to Pacemaker cluster on all nodes using the commands shown in Example 4-2.

Example 4-2 Configure firewall

```
$ sudo firewall -cmd - -permanent - - add- service=high-availability
$ sudo firewall -cmd - -reload
```

3. Create a Pacemaker cluster by running the command shown in Example 4-3 on one of the nodes.

Example 4-3 Create Pacemaker cluster

```
$ sudo pcs cluster setup - -name <cluster_name> node1 node2 node3
```

This command will configure the Pacemaker configuration files and generate the cluster configuration file that can be used to create the cluster on all nodes.

4. Start the Pacemaker cluster on all the nodes using the command shown in Example 4-4.

Example 4-4 Start the Pacemaker cluster

```
$ sudo pcs cluster start - -all
```

This command will start the cluster services on all the nodes.

5. Add the shared storage devices to the cluster by running the command in Example 4-5.

Example 4-5 Add the shared storage devices to the cluster

```
$ sudo pcs cluster cib stonith_cfg
```

The above command will add the XML snippet to the Pacemaker configuration file as shown in Example 4-6 to define shared storage devices as the storage device to be used for the cluster fencing mechanism.

Example 4-6 Modify Pacemaker config file to add storage devices

```
<configuration>
  <fencing>
    <method name="sbd">
      <device class="stonith" name="sbd">
        <instance_attributes stonith-timeout="60"/>
      </device>
    </method>
  </fencing>
</configuration>
```

Save the file and close it.

6. Enable the sbd (Stonith Block Device) service on all nodes using the command:

Example 4-7 Enable the sbd service

```
$ sudo systemctl enable sbd
```

7. Configure the Pacemaker cluster to use the shared storage device for high availability storage by executing the command shown in Example 4-8.

Example 4-8 Configure the Pacemaker cluster to use shared storage device for HA

```
$ sudo pcs resource create shared_storage ocf:heartbeat:Filesystem device="/dev/sdb1" \
directory="/mnt/shared_storage" fstype="xfs" options="nofail" op monitor interval="30s" \
on-fail="standby" clone
```

The command will create a Pacemaker resource agent that monitors the shared disk and automatically fails over to another node in the Cluster if a problem is detected. Replace /dev/sdb1 with the path to required shared disk path and also change the mount as per existing mount in your setup.

8. Verify that the shared disk is all available to the nodes of cluster by running the command:

Example 4-9 Verify that the shared disk is available to the nodes of the cluster

```
$ sudo pcs status resources
```

This command will display the status of all resources in the Pacemaker cluster, including the shared disks.

9. Test the high availability storage configuration by simulating a node failure using the command shown in Example 4-10.

Example 4-10 Test the HA storage configuration

```
$ sudo pcs cluster simulate node1
```

This command will fail node1 and it will move the access to other node.

For more informations refer to the following link:

<https://unixcop.com/high-availability-storage-with-pacemaker/>



High Availability for Database Servers

Pacemaker plays a pivotal role in ensuring continuous availability and reliability for critical services, including IBM Db2 databases. Leveraging Pacemaker for IBM Db2 deployments involves implementing a well-thought-out strategy that maximizes uptime, minimizes downtime, and guarantees data integrity.

Pacemaker orchestrates the deployment and management of IBM Db2 in a clustered environment – allowing seamless failover and resource allocation across multiple nodes. Its suite of failover policies, constraints, and monitoring capabilities empowers administrators to design and maintain a highly available database system that can withstand node failures, network interruptions, or planned maintenance events.

The integration of Pacemaker with IBM Db2 introduces mechanisms for handling failovers, managing resource dependencies, & ensuring the efficient allocation of computing resources.

This chapter explores the intricacies of deploying and configuring Pacemaker to orchestrate high availability for IBM Db2 databases. It delves into failover policies, resource constraints, and best practices for setting up a fault-tolerant environment tailored specifically for IBM Db2.

In this chapter we discuss the following:

- ▶ 5.1, “Db2 database introduction” on page 64.
- ▶ 5.2, “Db2 High Level deployment Architecture with Pacemaker” on page 66.
- ▶ 5.3, “Using Pacemaker to create highly available Db2 databases” on page 67.
- ▶ 5.4, “Prerequisites for an integrated solution using Pacemaker” on page 70.
- ▶ 5.5, “Overview of Installing Db2 on Red Hat Enterprise Linux and configuration of Pacemaker” on page 71
- ▶ 5.6, “Using the Db2 cluster manager utility for high availability” on page 72.

5.1 Db2 database introduction

Db2 is a family of database management system (RDBMS) products from IBM that serves several different operating system (OS) platforms. Used by organizations of all sizes, Db2 provides a data platform for both transactional and analytical operations, as well as continuous availability of data to keep transactional workflows and analytics operating efficiently. In addition to being a relational DBMS, Db2 also offers integrated support for several NoSQL capabilities, including XML, graph store & JavaScript Object Notation, or JSON.

5.1.1 Db2 database concepts

A Db2 database is a group of data treated collectively as a unit. A database is a large, structured set of persistent data, and its purpose is to store, retrieve and manipulate related information. The data is stored in tables. The Db2 DBMS operates as the server to manage data in databases across a multiuser environment, enabling many concurrent users to access the same data simultaneously. The Db2 DBMS also prevents unauthorized access, provides utilities for backing up and recovering data and offers performance tools and data management capabilities.

Db2 databases have logical structures and physical structures which the DBMS manages separately. The physical storage of data can be managed without affecting the access to logical storage structures. Db2 databases are created using Data Definition Language, or DDL, commands and are composed of tablespaces, tables with rows and columns, views, indexes, stored procedures, and other supporting structures. DBA or developer can use a Db2 database and its underlying structures to create, read, update, and delete data to support an organization's business requirements. SQL Language is used by applications to interact with the database to perform create or change data.

5.1.2 Db2 high availability and disaster recovery

High availability disaster recovery (HADR) provides a high availability solution to protect against both partial and complete site failures. HADR protects against data loss by replicating data changes from a source database, called the *primary database*, to the target databases, called the *standby databases*. HADR supports up to three remote standby servers.

A partial site failure can be caused by a hardware, network, or software failure. Without HADR, a partial site failure requires restarting the database management system (DBMS) server that contains the database. The length of time that it takes to restart the database and the server where it is located is unpredictable. It can take several minutes before the database is brought back to a consistent state and made available. With HADR, a standby database can take over in seconds. Further, you can redirect the clients that used the original primary database to the new primary database by using automatic client reroute or retry logic in the application.

A complete site failure can occur when a disaster, such as a fire, causes the entire site to be destroyed. However, because HADR uses TCP/IP for communication between the primary and standby databases, they can be situated in different locations. For example, the primary database might be at your head office in one city, and a standby database might be at your sales office in another city. If a disaster occurs at the primary site, data availability is maintained by having the remote standby database take over as the primary database with full Db2 functionality. After a takeover operation occurs, you can bring the original primary database back up and return it to its primary database status; this procedure is known as *failback*. You can initiate a failback if you can make the old primary database consistent with the new primary database. After you reintegrate the old primary database into the HADR setup as a standby database, you can switch the roles of

the databases. This operation would enable the original primary database to be the primary database again.

Figure 5-1 shows a Db2 database utilizing HADR. The example shows HADR deployed with one Primary and two Standby instances. One database instance in Data Center 1 is running in a Primary role, a Standby database instance is running in Data Center 1 and a second Standby database instance is available for disaster recovery in Data Center 2.

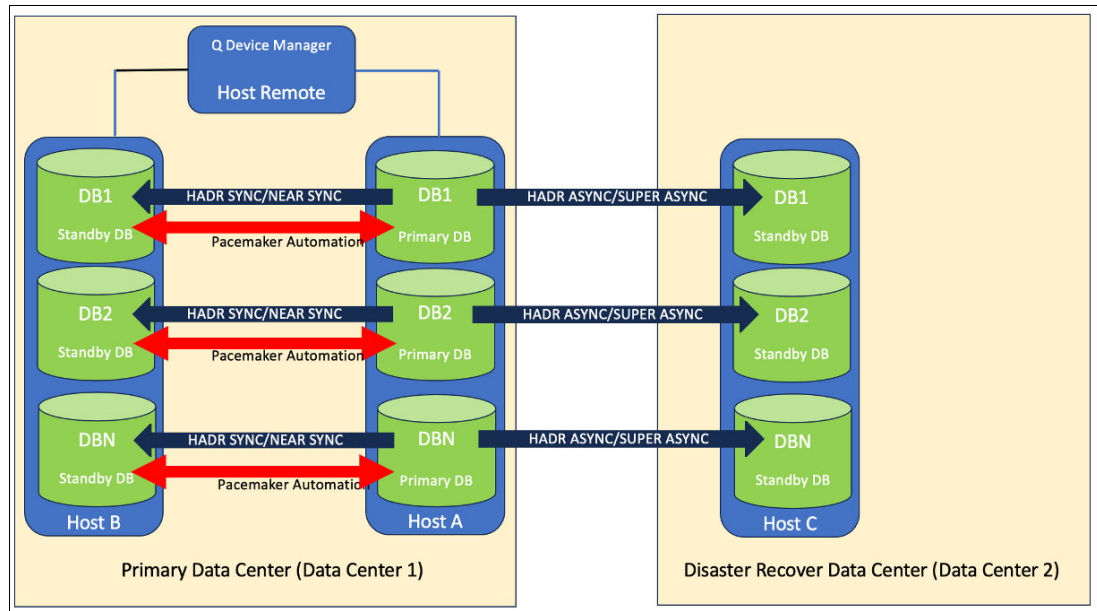


Figure 5-1 Db2 HADR Architecture with Pacemaker across Data Centers

5.1.3 Db2 databases in High Availability deployment with HADR

When Db2 database servers are deployed for high availability, the database will be running on multiple Db2 server hosts. Data is only allowed to be changed on one of the DB2 instances, the data is replicated to the other instances which are running in standby mode. The HADR feature of Db2 is responsible for managing the data replication. The normal implementation of Db2 HADR is uses One Primary Db2 instance and One Standby DB2 instance in a production site and optionally another Standby database will be run in the DR site if DR protection is desired.

For Db2 running on Linux, Pacemaker, the Linux clustering component, will be deployed to automate the database service availability. Additional requirements for Pacemaker are listed in section 5.4, “Prerequisites for an integrated solution using Pacemaker” on page 70. A cluster domain will be created in the primary data center since it cannot span across the distance to the DR data center. Pacemaker will maintain and manage a Virtual IP address (VIP) which will be used to access whatever the current primary database instance is. The application will use this VIP to connect to the database. The Pacemaker cluster will manage the Db2 server’s HADR Role. The Db2 HADR feature will manage the replication of data between the instances.

Figure 5-2 shows the components involved in running Db2 HADR database. Db2 with the HADR feature manages the data replication. Pacemaker ensures that if the Db2 primary instance fails, that the secondary instance is promoted to primary and continues handling data requests. Note that this cluster uses the Qdevice quorum setup where

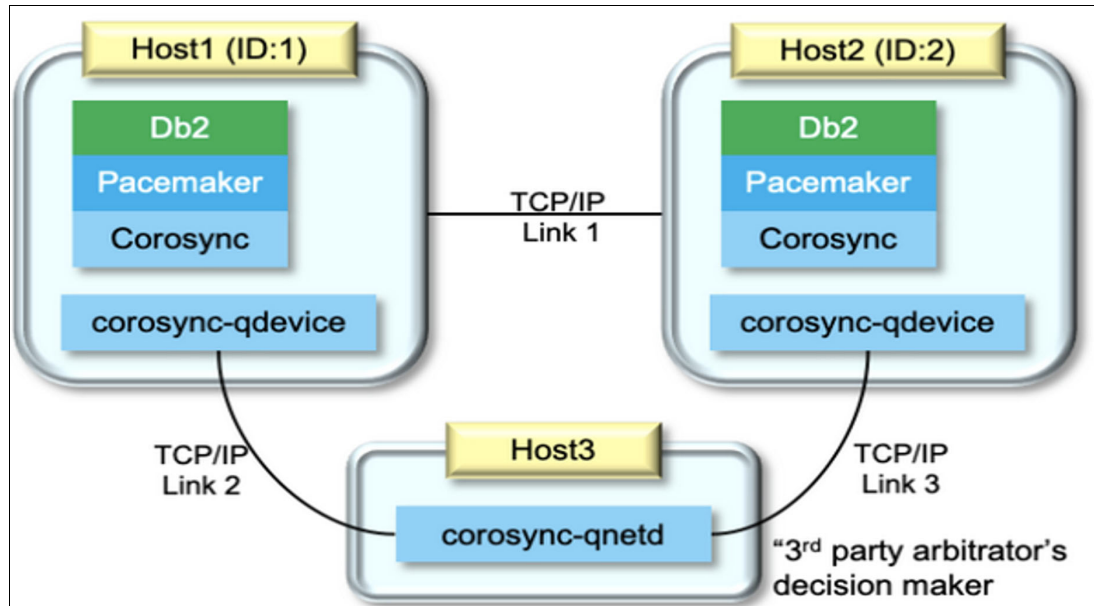


Figure 5-2 DB2 Pacemaker typical High Level Deployment Architecture

The DB2 database server could be typically deployed in following example pattern:

- ▶ Primary data center will have primary database and Principal standby database hosted on two database server hosts.
- ▶ Time synchronization needs to be maintained between the servers in different sites.
- ▶ Auxiliary standby database can be deployed in Standby data center.
- ▶ Service failover in primary data center will be automatically managed by Pacemaker. The server will failover from primary and standby database and managed by Pacemaker.
- ▶ In case of planned or unplanned downtime the database service will have to be manually started in DR Site (Data Centre) in case of failure.
- ▶ Allocate two NIC adapters on each Db2 Server of production environment to separate workload for client interactions and workload from database data replication (HADR).
 - NIC 1 for corporate application client connectivity.
 - NIC 2 for DB2 HADR replication network and connectivity.

5.2 Db2 High Level deployment Architecture with Pacemaker

A cluster is a group of connected machines that work together as a single system. When one machine in a cluster fails, cluster management software transfers the workload of the failed machine onto other machines.

Pacemaker is an open source high-availability cluster resource manager software that runs on a set of nodes. Together with Corosync – an open-source group communication system that provides ordered communication delivery, cluster membership, quorum enforcement, and other features among the nodes – it helps detect component failures and orchestrate the

necessary failover procedures to minimize interruptions to applications. Pacemaker provides a framework to manage the availability of resources. Resources are services on a host that needs to be kept highly available.

DB2 Pacemaker utilizes Corosync to setup and manage the failover between the machines managed by the quorum arbitrator explained in section 5.3.2, “The Quorum device support on Pacemaker” on page 69.

5.3 Using Pacemaker to create highly available Db2 databases

Pacemaker is seamlessly integrated with Db2 on Linux platforms including Linux for Power. The unified solution extends its support across multiple databases and multiple instances, all orchestrated via the Db2 cluster manager utility (*db2cm*).

Starting from Db2 version 11.5.5 and beyond, Pacemaker enables the automation of HADR (High Availability Disaster Recovery) configurations as the cluster manager. This advanced capability encompasses the automation of multiple HADR databases per instance and facilitates the automation of multiple HADR instances, all seamlessly coordinated within a shared Pacemaker/Corosync cluster environment.

Moreover, in Db2 versions 11.5.8 and later, Pacemaker introduces support for Mutual Failover (MF) automation. This integration empowers Pacemaker as the cluster manager to facilitate Mutual Failover, enhancing the disaster recovery capabilities of the system.

Note:

- In Db2 version 11.5.4, Pacemaker is provided solely as a technology preview, suitable for development, testing, and proof-of-concept purposes.
- Beginning with Db2 11.5.5, Pacemaker becomes part of the offering and is made accessible for deployment in production environments.
- In Db2 11.5.6 and subsequent versions, the Pacemaker cluster manager is bundled and installed along with Db2 to facilitate automated fail-over to HADR standby databases.
- From Db2 version 11.5.8 and onward, the utilization of Pacemaker as the integrated cluster manager allows for Mutual Failover high availability.

5.3.1 Pacemaker core: foundational components & key elements

In the comprehensive high availability (HA) setup integrating Pacemaker, the cluster software stack comprises multiple essential components, all vital for the effective operation of Pacemaker.

Resources

These represent specific entities within Db2 that require monitoring, initiation, or termination of states. This encompasses various elements (resources), such as the Db2 member process, HADR-capable databases for High Availability Disaster Recovery (HADR), mount points necessary for Mutual Failover functionality, Ethernet network adapters, and virtual IP addresses.

Constraints

Constraints serve as rules configured during cluster establishment to enhance the behavior of processes:

- ▶ **Location Constraint:** Determines where specific resources are permitted to operate within the cluster.
- ▶ **Ordering Constraint:** Specifies the sequence in which certain resource actions need to take place.
- ▶ **Co-location Constraint:** Defines the interdependency of one resource's location with another resource's placement within the cluster environment.

Resource Set

A collection of resources that fall under the influence of a specific constraint.

Resource Model

Within Pacemaker, the resource model for Db2 encompasses the predefined relationships and constraints of all resources. This model is established during cluster setup using the *db2cm* utility with the *-create option*. Any modifications or deviations from this model without approval from Db2 will render the model unsupported.

Resource Agents

In Pacemaker, resource agents are Db2 user exits, consisting of a set of shell scripts developed and maintained by Db2. These agents execute actions on the resources defined within the resource model. Table 5-1 identifies available resource agents.

Table 5-1 *Resource agents*

Resource agent	Definition	Level
db2ethmon	The resource agent to monitor the defined Ethernet network adapter.	Host level.
db2inst (HADR only)	The resource agent to monitor, start, and stop the Db2 member process.	Db2 instance level.
db2hadr (HADR only)	The resource agent to monitor, start, and stop individual HADR-enabled databases.	Db2 database level.
db2partition (Mutual Failover)	The resource agent to monitor, start, and stop the Db2 partition process.	Db2 instance level.
db2fs (Mutual Failover)	The resource agent to monitor, start, and stop a file system.	N/A.

Cluster Topology and Communication Layer

Effective high-availability cluster management requires consistent cluster topology among nodes. Pacemaker employs the Corosync Cluster Engine, an open-source group communication system, to ensure uniform cluster topology across nodes. Corosync guarantees a reliable messaging infrastructure, ensuring that events occur in the same sequence on each node, and enforces quorum constraints.

Cluster Domain Leader

Within the cluster, one node is elected as the Domain Leader (also known as the Designated Controller or DC in Pacemaker terminology). The Pacemaker controller daemon on the Domain Leader assumes responsibility for making critical cluster decisions. If the current Domain Leader node fails, a new Domain Leader is elected to assume this pivotal role.

Networks

To set up components within your Pacemaker cluster domain associated with networking, use the Db2 cluster manager (*db2cm*) utility. This utility allows the addition of a physical network to your cluster domain. This physical network typically consists of network interface cards (NICs), IP addresses, and subnet masks.

QDevice

The QDevice (Quorum Device) refers to a specialized quorum configuration that utilizes an external device or mechanism to determine the cluster's availability and make decisions regarding its operation. QDevice quorum provides an additional layer of reliability by introducing an external resource (such as a storage device or network device) to participate in the cluster's quorum. Key aspects of QDevice quorum include:

- ▶ QDevice quorum involves integrating an external device or mechanism into the cluster's quorum calculations. This device acts as an additional voter in determining the cluster's state and making decisions regarding failover or resource management.
- ▶ By adding an external quorum device, the cluster's quorum decisions are less reliant solely on the number of cluster nodes. Instead, the device's availability or state contributes to the overall quorum calculation, enhancing reliability and reducing the risk of split-brain scenarios or incorrect failover decisions.
- ▶ Quorum devices can vary and might include shared disks, dedicated quorum servers, or network-based devices specifically designed to participate in quorum calculations. These devices are configured to communicate with the cluster nodes and provide input to determine the cluster's operational state.
- ▶ The QDevice's availability contributes to the quorum vote count, influencing the cluster's ability to continue operations or trigger failover actions. If the external device is accessible and operational, it adds its quorum vote to the overall decision-making process.
- ▶ In scenarios where the majority of cluster nodes become unavailable or network partitions occur, the presence and operational status of the QDevice can be decisive in making proper failover decisions to ensure cluster integrity and data consistency.

Note: QDevice quorum is the recommended quorum mechanism for a production system.

5.3.2 The Quorum device support on Pacemaker

A quorum device helps a cluster manager make cluster management decisions when the cluster manager's normal decision process does not produce a clear choice. To select an action to take, a cluster manager counts the number of cluster domain nodes that support each of the potential actions. The cluster manager then selects the action that is supported by most of the cluster domain nodes. If the same number of cluster domain nodes supports more than one choice, then the cluster manager refers to a quorum device to make the choice. The Pacemaker cluster stack supports the following quorum devices:

- ▶ Two-node quorum
- ▶ Majority quorum
- ▶ QDevice quorum

Two node Quorum and Majority Quorum

The two-node quorum is the default mechanism. Because no tie-breaker mechanism exists, the two-node quorum is prone to the split-brain scenario. It is not intended for production environments.

Majority quorum is intended to use with cloud providers such as AWS, GCP, or AZURE.

QDevice quorum

This is recommended quorum management method for Db2 HADR. QDevice requires an external resource which must be accessible by all hosts in the Pacemaker cluster. QDevice is much more reliable because the quorum decision logic is more robust than a simple TCP/IP ping to the external IP address.

QDevice requires that the resource is placed on a separate host, like a majority quorum requirement. However, setup is simplified, as the host with the QDevice does not need to be part of the Pacemaker quorum. The QDevice is an excellent choice as a quorum solution as it combines reliability and simplicity.

5.4 Prerequisites for an integrated solution using Pacemaker

The following are the prerequisites for an integrated solution using Pacemaker for Db2 on IBM Power Linux. For more information refer to this [IBM Document](#).

Supported IBM Power Linux distributions

Table 5-2 identifies the supported Power Linux distributions for running a Pacemaker managed integrated high availability solution.

Table 5-2 Supported Power Linux for IBM Db2

Db2 version	Supported RHEL version	Supported SUSE Version
Db2 version 11.5.8	Version 8.4 and higher.	Version 15 SP3 and above.
Db2 version 11.5.7	Version 8.2 and higher.	Version 15 SP3 and above.
Db2 before 11.5.7	Not supported.	Not supported.

User Configuration

Ensure that the necessary users (such as instance owner, fenced user), associated groups, and other users essential for the setup are created.

SSH Configuration

Set up passwordless SSH between HADR nodes for both root and instance users to enable seamless communication.

Storage Requirements

Allocate adequate local storage on each node for cluster-related software:

- Approximately 50 MB for cluster storage RPMs and extracted files.
- Around 200 MB for a full installation of cluster-related software.
- Minimum 1GB space in /var to store cluster software log files.
- At least 150 MB in /usr for RHEL and 300 MB for SUSE distributions.

RPM Requirements

Please enable Yum repository for installation.

Approximately

- /usr/share/pacemaker
- /usr/share/doc/packages
- /usr/share/licenses
- /usr/share/man/man7
- /usr/share/man/man8

- /usr/lib/pacemaker
- /usr/lib/ocf/resource.d/pacemaker
- /usr/lib/systemd/system
- /usr/lib/debug/dwz
- /usr/lib64
- /usr/lib64/pkgconfig
- /usr/sbin

Firewall Configuration

If a firewall is in place, open specific ports (3121, 5403, 5404 - 5405) on each host or in the network to facilitate communication.

Required Packages

Ensure the presence of KornShell (*ksh*) and *python3-dnf-plugin-versionlock* packages for Pacemaker. The latter is used by the Db2 installer to lock all relevant Pacemaker and Corosync RPMs.

Virtual IP Setup

In Db2 HADR, set up a Virtual IP for each HADR-enabled database to enable automatic client rerouting during failover.

QDevice Quorum

To establish QDevice quorum, install corosync-qnetd software on a third host acting as the arbitrator. This host doesn't need the Db2 server and serves as a mediator.

db2prereqcheck command

Executing this command allows you to check if your system meets the prerequisites prior to commencing the installation. Remember, you'll need to download the installation media and run [db2prereqcheck](#) for Db2 to proceed with the installation.

5.5 Overview of Installing Db2 on Red Hat Enterprise Linux and configuration of Pacemaker

The following are the required steps to install and configure IBM Db2 with Pacemaker on RHEL:

- ▶ Download the latest version of IBM Db2 from [IBM Passport Advantage® Online for Customers web site](#) or [IBM PartnerPlus Software Access Catalog](#) page.
- ▶ Ensure that you have the supported version of RHEL as described on Table 5-2 on page 70.
- ▶ Extract the archive and unpack each of the parts into the working directory that you created on your two LPARs where you will install the Primary and Standby IBM Db2 databases. The archive will contain the following directory tree:
 - Db2/
 - Db2agents/
 - RPMS/
 - RPMS/ppcle
 - RPMS/noarch
 - SRPMS/

- ▶ To check whether the system meets the prerequisites for Db2 version 11.5.8, issue the command shown in Example 5-1.

Example 5-1 db2prereqcheck command

```
$ db2prereqcheck -v 11.5.8.0 -s
```

- ▶ Install Pacemaker using the `db2_install` or `db2setup`¹ command: you can run the commands without any parameters specific to Pacemaker. The default behavior is to install Pacemaker. Setup Db2 according to the [IBM documentation](#).

5.6 Using the Db2 cluster manager utility for high availability

The Db2 cluster manager (*db2cm*) tool facilitates the setup and management of robust Db2 databases within a Pacemaker cluster to ensure high availability (HA). Here are two key configurations available:

- ▶ Use the HADR high availability feature within a Db2 instance operating on a Pacemaker-controlled Linux cluster.
- ▶ Use the Mutual Failover high availability capability within a Db2 instance running on a Pacemaker-controlled Linux cluster.

5.6.1 The HADR option for Pacemaker-managed Db2 instances

High availability disaster recovery (HADR) offers a solution for maintaining operational continuity during both partial and complete site outages. It safeguards against data loss by replicating data modifications from a source database (the primary) to designated databases (the standbys). HADR facilitates up to three remote standby servers to enhance resilience.

Prior to setting up HADR, confirm the presence of the subsequent system prerequisites:

- ▶ All hosts within the cluster must have the Pacemaker cluster software stack installed.
- ▶ Configure and ensure the online status of Db2 instances and the HADR-enabled database.
- ▶ In this scenario, 192.168.1.100 and 192.168.1.200 represent the host names of the cluster nodes, `yslhadomain` stands as the domain name, and `eth0` serves as the network interface (device) name allocated to each host.
- ▶ Create the Pacemaker cluster and the public network resources by running the command shown in Example 5-2.

Example 5-2 Create Pacemaker cluster command

```
$ home/db2inst1/sqllib/bin/db2cm -create -cluster -domain yslhadomain -host
ip-192.168.1.100 -publicEthernet eth0 -host ip-192.168.1.200 -publicEthernet eth0
Created db2_ip-192.168.1.100_eth0 resource.
Created db2_ip-192.168.1.200_eth0 resource.
Cluster created successfully.
```

¹ `db2setup` utility, which provides silent install with a response file and a graphical user interface, supports installing Pacemaker in Db2 11.5.8 and later

- Create the instance cluster domain by running the commands shown in Example 5-3.

Example 5-3 Create instance cluster domain command

```

$/home/db2inst1/sqllib/bin/db2cm -create -instance db2inst1 -host ip-192.168.1.100
Created db2_ip-192.168.1.100_db2inst1_0 resource.
Instance resource for db2inst1 on ip-192.168.1.100 created successfully.

```

```

$ /home/db2inst1/sqllib/bin/db2cm -create -instance db2inst1 -host ip-192.168.1.200
Created db2_ip-192.168.1.200_db2inst1_0 resource.
Instance resource for db2inst1 on ip-192.168.1.200 created successfully.

```

- Verify the cluster by running the *crm* status command as shown in Example 5-4.

Example 5-4 Verify cluster status command

```

$crm status
Stack: corosync
Current DC: ip-192.168.1.200 (version 2.0.2-1.e18-744a30d655) - partition with quorum
Last updated: Tue Dec 05 21:49:57 2023
Last change: Tue Dec 05 21:39:45 2023 by root via cibadmin on ip-192.168.1.100
2 nodes configured
4 resources configured
Online: [ ip-192.168.1.200 ip-192.168.1.100 ]
Full list of resources:
db2_ip-192.168.1.100_eth0 (ocf::heartbeat:db2ethmon): Started ip-192.168.1.100
db2_ip-192.168.1.200_eth0 (ocf::heartbeat:db2ethmon): Started ip-192.168.1.200
db2_ip-192.168.1.100_db2inst1_0 (ocf::heartbeat:db2inst): Started ip-192.168.1.100
db2_ip-192.168.1.200_db2inst1_0 (ocf::heartbeat:db2inst): Started ip-192.168.1.200

```

- Create a new database and configure HADR on the new database.
- Create the HADR database resources as shown in Example 5-5.

Example 5-5 Create HADR database resources

```

$/home/db2inst1/sqllib/bin/db2cm -create -db YSLDB -instance db2inst1
Database resource for YSLDB created successfully

```

- Create the VIP resources for the newly created database as shown in Example 5-6.

Example 5-6 Create VIP resources

```

$/home/db2inst1/sqllib/bin/db2cm -create -primaryVIP 192.168.1.201 -db YSLDB -instance
db2inst1

```

- Verify the cluster by running the *crm* status command after creating Db2 databases, HADR resources and VIP resources. This is shown in Example 5-7.

Example 5-7 Verify cluster status

```

$ crm status
Stack: corosync
Current DC: ip-192.168.1.200 (version 2.0.2-1.e18-744a30d655) - partition with quorum
Last updated: Tue Dec 05 22:20:52 2023
Last change: Tue Dec 05 22:14:29 2023 by root via cibadmin on ip-192.168.1.200
2 nodes configured
7 resources configured
Online: [ ip-192.168.1.200 ip-192.168.1.100 ]
Full list of resources:
db2_ip-192.168.1.100_eth0 (ocf::heartbeat:db2ethmon): Started ip-192.168.1.100

```

```
db2_ip-192.168.1.200_eth0 (ocf::heartbeat:db2ethmon): Started ip-192.168.1.200
db2_ip-192.168.1.100_db2inst1_0 (ocf::heartbeat:db2inst): Started ip-192.168.1.100
db2_ip-192.168.1.200_db2inst1_0 (ocf::heartbeat:db2inst): Started ip-192.168.1.200 Clone
Set: db2_db2inst1_db2inst2_YSLDB-clone [db2_db2inst1_db2inst2_YSLDB] (promotable)
Masters: [ ip-192.168.1.200 ]
Slaves: [ ip-192.168.1.100 ]
db2_db2inst1_db2inst1_YSLDB-primary-VIP (ocf::heartbeat:IPAddr2): Started ip-192.168.1.200
```

► Install and configure a QDevice quorum

- On the primary and standby hosts, install the *corosync-qdevice* packages as shown in Example 5-8.

Example 5-8 Install the corosync-qdevice packages

```
$dnf install
<Db2_image>/db2/<platform>/pcmk/Linux/<OS_distribution>/<architecture>/corosync-qdevice*
```

- On the QDevice host, install the Corosync QNet software. This is shown in Example 5-9.

Example 5-9 Install Corosync QNet software

```
$dnf install
<Db2_image>/db2/<platform>/pcmk/Linux/<OS_distribution>/<architecture>/corosync-qnetd*
```

- As the root user, run the `db2cm` command to setup the QDevice from one of the cluster nodes as shown in Example 5-10.

Example 5-10 Setup QDevice

```
$ <Db2_image>/bin/db2cm -create -qdevice qserver.powerm.ma
```

- On the primary and standby hosts, verify that the quorum was setup correctly. See Example 5-11.

Example 5-11 Verify Quorum Setup

```
$corosync-qdevice-tool -s
Qdevice information -----
Model:          Net
Node ID:        1
Configured node list:
0   Node ID = 1
1   Node ID = 2
Membership node list:  1, 2
Qdevice-net information -----
Cluster name:    yslhadomain
QNetd host:     qserver.powerm.ma:5403
Algorithm:      LMS
Tie-breaker:    Node with lowest node ID
State:          Connected
```

- On the QDevice host, verify that the quorum device is running correctly as shown in Example 5-12.

Example 5-12 Verify Quorum

```
$ corosync-qnetd-tool -l
```

- ▶ Initiate the takeover of a primary HADR host from the standby host on a Pacemaker-managed Db2 instance.
 - Verify Cluster Status as shown in Example 5-13.

Example 5-13 Verify Cluster

```
$ crm status
Stack: corosync
Current DC: dbsecondary (version 2.0.2-1.e18-744a30d655) - partition with quorum
Last updated: Tue Dec 05 22:22:32 2023
Last change: Tue Dec 05 22:15:43 2023 by root via cibadmin on dbprimary
2 nodes configured
7 resources configured
Online: [ dbprimary dbsecondary ]
Full list of resources:
db2_dbprimary_eth1 (ocf::heartbeat:db2ethmon): Started dbprimary db2_dbsecondary_eth1
(ocf::heartbeat:db2ethmon): Started dbsecondary db2_dbprimary_db2inst1_0
(ocf::heartbeat:db2inst): Started dbprimary db2_dbsecondary_db2inst1_0
(ocf::heartbeat:db2inst): Started dbsecondary
Clone Set: db2_db2inst1_db2inst1_YSLDB-clone [db2_db2inst1_db2inst1_YSLDB] (promotable)
Masters: [ dbprimary ]
Slaves: [ dbsecondary ]
db2_db2inst1_db2inst1_YSLDB-primary-VIP (ocf::heartbeat:IPaddr2): Started dbprimary
db2_db2inst1_db2inst1_YSLDB-standby-VIP (ocf::heartbeat:IPaddr2): Started
dbsecondary
```

- As the Db2 instance owner, run a user-initiated takeover from the standby host as shown in Example 5-14.

Example 5-14 Initiate takeover

```
$ db2 takeover hadr on db YSLDB
DB20000I The TAKEOVER HADR ON DATABASE command completed
successfully.
```

- Validate that the primary role fails over in the crm output. Also, verify that the reads on the standby VIP and primary VIP have switched hosts as expected. This can be seen in Example 5-15.

Example 5-15 Verify Cluster

```
$ crm status
Stack: corosync
Current DC: dbsecondary (version 2.0.2-1.e18-744a30d655) - partition with quorum Last
updated: Tue Dec 05 22:26:42 2023
Last change: Tue Dec 05 22:17:35 2023 by root via cibadmin on dbprimary
2 nodes configured
7 resources configured
Online: [ dbprimary dbsecondary ]
Full list of resources:
db2_dbprimary_eth1 (ocf::heartbeat:db2ethmon): Started dbprimary
db2_dbsecondary_eth1 (ocf::heartbeat:db2ethmon): Started dbsecondary
```

```

db2_dbprimary_db2inst1_0      (ocf::heartbeat:db2inst):      Started dbprimary
db2_dbsecondary_db2inst1_0  (ocf::heartbeat:db2inst):      Started dbsecondary Clone Set:
db2_db2inst1_db2inst1_YSLDB-clone [db2_db2inst1_db2inst1_YSLDB] (promotable)
Masters: [ dbsecondary ]
Slaves: [ dbprimary ]
db2_db2inst1_db2inst1_YSLDB-primary-VIP (ocf::heartbeat:IPAddr2): Started dbsecondary
db2_db2inst1_db2inst1_YSLDB-standby-VIP (ocf::heartbeat:IPAddr2): Started
dbprimary

```

To maintain a Pacemaker-managed cluster domain for an HADR Db2 instance refer to [this document](#).

5.6.2 The Mutual Failover option for Pacemaker-managed Db2 instances

Db2 versions 11.5.8 and beyond offer an alternative high availability feature known as Mutual Failover when utilized within a Pacemaker-controlled Linux cluster. The setup of this cluster is established using the Db2 cluster manager (*db2cm*) utility.

Db2 Mutual Failover represents a high availability configuration where two hosts possess the same Db2 instance installation and a shared mount, allowing only one host to be active at any given moment

In this scenario, *192.168.1.100* and *192.168.1.200* represent the host names of the cluster nodes, *yslhadomain* stands as the domain name, and *eth0* serves as the network interface (device) name allocated to each host.

- ▶ On the first host *192.168.1.100*:
 - As the root user, create a directory on *192.168.1.100* for the file system that is to be shared by the two Db2 Mutual Failover hosts (*192.168.1.100* and *192.168.1.200*). The command is shown in Example 5-16.

Example 5-16 Create directory for the shared file system

```
$mkdir /db2ha
```

- Format a partition on the target device with the appropriate file system for the cluster as shown in Example 5-17.

Example 5-17 Formation partition

```
$mkfs.ext3 /dev/sda
```

- Get the Universal Unique Identifier (UUID) of the mount device that you created as seen in Example 5-18.

Example 5-18 Formation partition

```

$blkid /dev/sda1
...
UUID="31a28d09-381c-4249-8b88-1d5db75c3b81" BLOCK_SIZE="4096" TYPE="ext3"

```

- Create a new entry in the */etc/fstab* file on both hosts. This is shown in Example 5-19.

Example 5-19 Create fstab entry

```
/dev/sda1:UUID=31a28d09-381c-4249-8b88-1d5db75c3b81/db2ha ext3 acl,user_xattr,noauto 0 0
```

- Mount the new file system on this host as shown in Example 5-20.

Example 5-20 Mount file system

```
$mount /db2ha
```

- Create administration groups, an instance user, and a fenced user on both hosts in your cluster. This is shown in Example 5-21.

Example 5-21 Create group and users

```
groupadd -g 990 db2iadm1
groupadd -g 989 db2fadm1
useradd -u 1002 -g db2iadm1 -m -d /home/db2inst1 -s /usr/bin/ksh db2inst1
useradd -u 1003 -g db2fadm1 -m -d /home/db2fenc1 -s /usr/bin/ksh db2fenc1
```

- Install Db2 on both hosts in your cluster using the command shown in Example 5-22.

Example 5-22 Install Db2

```
./db2_install -b /opt/ibm/db2/V115M8A -p SERVER -y
```

- Create an empty directory for the db2inst1 user and copy the existing .profile and .kshrc files over to the folder. See Example 5-23.

Example 5-23 copy .profile and .kshrc

```
mkdir /db2ha/db2home
chown -R db2inst1.db2iadm1 /db2ha
usermod -d /db2ha/db2home db2inst1
cp -f /home/db2inst1/.profile /db2ha/db2home
cp -f /home/db2inst1/.kshrc /db2ha/db2home
```

- Create the Db2 instance using the command shown in Example 5-24.

Example 5-24 Create Db2 instance

```
./opt/ibm/db2/V115M8A/instance/db2icrt -u db2fenc1 db2inst1
```

- Configure passwordless SSH by first adding the Db2 instance user's SSH key to the authorized_keysfile on 192.168.1.100 as shown in Example 5-25.

Example 5-25 Configure passwordless

```
$cat /db2ha/db2home>/.ssh/id_rsa.pub >> /db2ha/db2home /.ssh/authorized_keys
```

- ▶ On the second host 192.168.1.200.
 - Umount the shared file system from the first host using the command shown in Example 5-26.

Example 5-26 umount file system

```
$umount /db2ha
```

- Create file system directory on the second host and then mount the shared file system to the directory. See Example 5-27.

Example 5-27 Create directory and mount file system

```
$mkdir /db2ha
$ mount /db2ha
```

- Confirm that the users are the same on both hosts, mount the shared file system and add the Db2 instance user's SSH key to the authorized_keys file, and configure the home directory on the second host and then remove all files from the sqllib directory. These commands are shown in Example 5-28.

Example 5-28 Configure home directory and remove all files from sqllib directory

```
chown -R db2inst1.db2iadm1 /db2ha
usermod -d /db2ha/db2home db2inst1
rm -rf /db2ha/db2home/sqllib
```

- Create the Db2 instance on the second host using the command in Example 5-29.

Example 5-29 Create Db2 instance

```
$/opt/ibm/db2/V115M8A/instance/db2icrt -u db2fenc1 db2inst1
```

- ▶ Configure Mutual Failover with *db2cm* command:
 - Run *db2cm* to create your Db2 cluster as seen in Example 5-30.

Example 5-30 Create Db2 cluster

```
$ home/db2inst1/sqllib/bin/db2cm -create -cluster -domain yslhadomain -host
ip-192.168.1.100 -publicEthernet eth0 -host ip-192.168.1.200 -publicEthernet eth0
Created db2_ip-192.168.1.100_eth0 resource.
Created db2_ip-192.168.1.200_eth0 resource.
Cluster created successfully.
```

- Create a partition resource on the first host as shown in Example 5-31.

Example 5-31 Configure partition resource

```
$/db2cm -create -partition 0 -instance db2inst1
```

- Create a VIP resource on the first host as shown in Example 5-32.

Example 5-32 Create VIP resource

```
$/db2cm -create -primaryVIP 192.168.1.201 -partition 0 -instance db2inst1
```

- Create a Quorum Device on the first host as shown in Example 5-33.

Example 5-33 Create Quorum Device

```
$/db2cm -create -qdevice qserver
```

- Verify the configuration of your Mutual Failover, Pacemaker-managed Db2 cluster using the `db2cm` command as shown in Example 5-34.

Example 5-34 Verify configuration

```
./db2cm -list
```

- ▶ Initiate a failover and move a mutual failover partition from one host to the other.
 - As root, find the location of the partition that you want to move. See the command in Example 5-35.

Example 5-35 Find partition location

```
./db2cm -list
```

- As root, from the install path run `db2cm -move` to move your partition. This is shown in Example 5-36.

Example 5-36 Move partition

```
$/opt/ibm/db2/V11.5/bin/db2cm -move -instance db2inst1 -partition 0 -host dbsecondary  
Partition resource db2_db2inst1_0 moved to dbsecondary successfully.
```

For additional information on maintaining a Pacemaker-managed cluster domain for a Mutual Failover Db2 instance refer to this [IBM Document](#).



Using Pacemaker for Highly Available SAP Landscapes

This chapter describes the use of Pacemaker on Linux on Power systems either on IBM Power Systems Virtual Server Offering or directly on IBM Power Systems to provide highly available SAP landscapes.

There are two main application areas of SAP is used with IBM Power systems either on IBM cloud or on-premise together with Pacemaker.

- Highly Available deployments of SAP HANA Databases
- Highly Available SAP NetWeaver and S/4HANA Applications

Following topics will discuss and details out options on how to deploy SAP HANA Databases and SAP NetWeaver or S/4HANA Applications on power systems.

- ▶ 6.1, “Introduction to SAP systems” on page 82
- ▶ 6.2, “Using Pacemaker for S/4HANA or NetWeaver application resources” on page 85
- ▶ 6.3, “Using Pacemaker to create highly available SAP HANA database instances” on page 109

6.1 Introduction to SAP systems

SAP SE is a German multinational software company that develops enterprise software to manage business operations and customer relations. The company is a leading enterprise resource planning (ERP) software vendor. Founded in 1972, the company was initially called System Analysis Program Development (Systemanalyse Programmentwicklung), later abbreviated to SAP. Since then, it has grown from a small, five-person endeavor to a multinational enterprise headquartered in Waldorf, Germany, with more than 105,000 employees worldwide.¹

SAP creates solutions for all types of enterprises, from small to large corporations. SAP's main goal is centralizing all business processes to provide real-time information, enabling enterprises to map and design various business processes inside their companies. SAP delivers software solutions for over 25 industries helping customers to work with data and allowing businesses to achieve their goals and improve their data management, decision-making and planning.

SAP ECC

SAP ECC (also known as SAP ERP or NetWeaver) is SAP's legacy suite of enterprise applications and was designed to run on a third-party database such as Oracle or IBM DB2. ECC is the core business product inside SAP Business Suite, and it provides an integrated and updated overview of an organization's business processes – from financial and warehousing to human resources and logistics. To better understand just how much of a far reaching impact ECC has on a business, consider the breadth of use cases that are solved using some of the commonly used modules which include:

- Financial Accounting and Controlling (often referred to as FICO)
- Sales and Distribution / Customer Service
- Materials Management
- Production Planning
- Quality Management
- Warehousing
- Logistics
- Plant Maintenance
- Human Resources Management

SAP Business Suite added more capabilities with advanced modules that integrated with that core, such as Customer Relationship Management, Supplier Relationship Management, Supply Chain Management and Product Lifestyle Management.

Important: in 2020, SAP announced that the end of maintenance support for its ECC product would occur in 2027 (this was an extension of the original 2025 date). Any new SAP implementations should be done on SAP S/4HANA and customers running SAP ECC should be planning migration to SAP S/4HANA to maintain support.

SAP BW

In 1998 SAP released the first version of SAP BW, providing a model-driven approach to Enterprise Data Warehouse that made data warehousing easier and more efficient, particularly for SAP R/3 data. Since then, SAP BW has evolved to become a key component for thousands of companies. It can transform and consolidate business information from virtually any source system and it ran on industry standard relational databases like Oracle

¹ <https://www.sap.com/about/what-is-sap.html>

and IBM DB2 until version 7.3 at which point it began to transition to use SAP HANA - SAP's in-memory database manager with the release of version 7.4.

SAP HANA

In 2010, SAP launched SAP HANA (High-performance ANalytic Appliance). SAP HANA is a multi-model database that stores data in memory instead of keeping it on a disk. The column-oriented in-memory database design provides the ability to run advanced analytics alongside high-speed transactions – in a single system. This allows companies to process massive amounts of data with near-zero latency, query data quickly, and become truly data-driven.

By storing data in column-based tables in main memory and bringing online analytical processing (OLAP) and online transactional processing (OLTP) together, SAP HANA is unique – and significantly faster than other database management systems (DBMS) on the market today. Initially, the legacy SAP ECC product was enabled to connect to the SAP HANA database as well as the previously supported databases.

SAP S/4HANA

In 2015, SAP announced its next generation of SAP ERP software, SAP S/4HANA. SAP S/4HANA is designed to provide real-time analytics, a streamlined user interface, and a simplified data model to help businesses run more efficiently and effectively. S/4HANA also offers a personalized user experience with SAP Fiori and delivers the faster response times and performance.

S/4HANA is built on the advanced SAP HANA in-memory database platform, which allows it to process large amounts of data quickly and efficiently. This enables businesses to make more informed decisions based on real-time insights—it is also one of the most important differences when compared to SAP ECC, which, as mentioned earlier, runs on a third-party database.

S/4HANA not only provides the end-to-end integration of key business components – like financial management, procurement and supply chain management, sales and distribution, production planning and control, and project management – it also serves as a critical backbone for business operations that want to tap into the power of next generation technologies like robotic process automation, machine learning, big data management, IoT, AI, and much more.

Overall, S/4HANA is the digital core enabling SAP's intelligent enterprise model—helping businesses improve productivity, reduce costs, and increase their agility in a rapidly changing market.

SAP BW/4HANA

In 2016 SAP introduced its enhanced enterprise data warehouse solution – SAP BW/4HANA. Like the SAP S/4HANA offering, the SAP BW/4HANA is written to take advantage of the SAP HANA database capabilities and as such only runs on SAP HANA.

6.1.1 High availability concepts for SAP solutions

SAP ERP solutions comprise multiple components and services that interact with each other to form a comprehensive network of functions, including but not limited to:

- ABAP Central Services (ASCS)
- Java Central Services (SCS)
- Enqueue Replication Server (ERS)
- Primary Application Server (PAS)

In addition to these services, the SAP landscape relies on a database element, such as SAP HANA, which also needs to be considered for HA protection. Your SAP environment may also include NetWeaver, and S4/HANA.

Figure 6-1 shows how these components work together and how your HA solution must support both the application components and the database components.

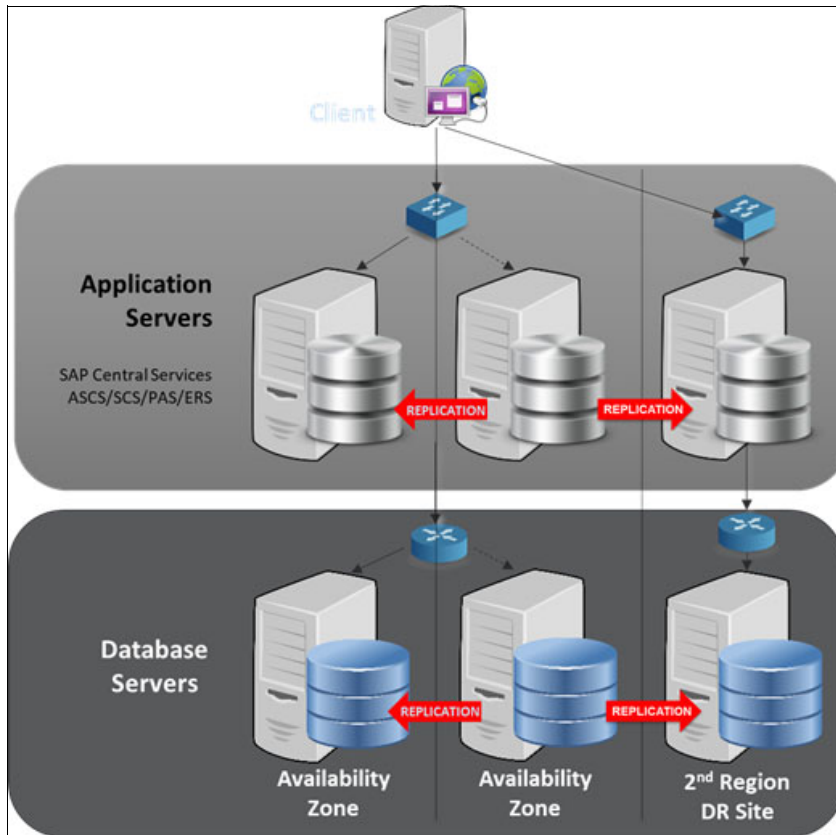


Figure 6-1 HA and DR solutions for SAP environment

To be effective, an HA solution must monitor the health of all these services, and in the event that a service does not perform as expected – or even fails completely – the high-availability solution will take the appropriate action to recover that service as quickly as possible. Further, the solution will bring the services online in the secondary node(s) and in the correct order to ensure SAP functions resume quickly and that data is not corrupted.

The service restart may occur on the same node or on an alternative node within the cluster. In the latter case, clustering software is used to configure secondary nodes to be kept on standby as ready-to-start-services that will continue operating if there is an issue with the primary node.

As discussed earlier, to meet stringent recovery point objectives (RPOs), this failover cluster should also be configured so that either the primary and secondary nodes share the same storage (e.g. SAN), or so that the storage for the primary node and secondary nodes are mirrored with efficient replication. Synchronized (mirrored) storage has several key advantages. It allows you to create a failover cluster in cloud environments where shared storage is not an option. It also allows you to use local solid-state drives (SSD) for cost-efficient, high-performance storage.

Failover clustering protection ensures minimal downtime of that service, and thus minimal impact to end-users. In parallel, the database is protected by replicating the data to be available on the secondary node within the same HA cluster, eliminating the SAN single point of failure. This data replication protection can be applied regardless of the type of database, including MaxDB, DB2, and SAP HANA. Also consider the need to geographically separate the nodes in the cluster for protection from site-wide or regional disasters. Use a clustering solution that allows failover between cluster nodes on different cloud subnets or in different geographically separated locations.

6.2 Using Pacemaker for S/4HANA or NetWeaver application resources

There are two deployment options to consider for running your SAP systems on IBM Power:

- ▶ The first use case involves using on-premise IBM Power systems hardware.
- ▶ The second use case uses IBM Power servers in the IBM Cloud using the IBM Power Systems Virtual Server offering.

Running SAP NetWeaver or S/4HANA on IBM Power Systems, whether using on-premises hardware or IBM Cloud PowerVS, offers a consistent robust platform for SAP applications providing best-in-class performance, resilience for critical workloads, and a flexible infrastructure. SAP has certified both of these options and they are certified for use with either SUSE Linux or Red Hat Linux distributions.

The implementation of Pacemaker (HA Add-on Cluster) is fully certified for SAP systems running on IBM Power hardware. The implementation is logically the same across the on-premises option and the IBM PowerVS options but will have minor differences in the commands used to implement them. There are also some minor differences between the implementation when using SUSE Linux or Red Hat Linux.

Note: There are some minor differences when using Red Hat or SUSE Linux and relevant OS vendor documentation must be referenced. The SAP notes will provide specific details relating to the implementation differences between the two operating systems.

Refer to following SAP Notes for further details on how to plan, deploy and configure these options on IBM Power Systems.

- *SAP Note - 2378874 - Install SAP Solutions on Linux on IBM Power Systems (little endian)*
- *SAP Note - 2855850 - SAP Applications on IBM Power Virtual Servers*
- *SAP Note - 2923984 - SAP on IBM Power Virtual Servers: Support prerequisites*
- *SAP Note - 2772999 - Red Hat Enterprise Linux 8.x: Installation and Configuration*

Note: A valid SAP “S-User” id is required to display/access above SAP Notes.

The examples we provide will be using Red Hat Enterprise Linux. However, it is important to note that there is a parallel and mostly identical deployment using SUSE Linux if that is your preference.

IBM PowerVS configuration options

Table 6-1 on page 86 shows the current options for running SAP NetWeaver on PowerVS instances.

Table 6-1 Current configuration options for SAP NetWeaver on PowerVS

Power System	SAPS per CPU Core	SAPS per CPU Thread (SMT-8)
S922	5,570	696.25
E980	6,000	750

As additional configuration options might be available in the future, it would be a good idea to check the latest information in the [IBM Cloud documentation](#).

Custom sizes for the IBM PowerVS instances are supported for use with SAP NetWeaver or S/4HANA in accordance with existing SAP recommendations. Sizing for your SAP NetWeaver or S/4HANA instances should follow the standard SAP sizing guidelines provided by using the SAPS benchmark calculations. For more information refer to *SAP Note 2855850*.

All SAP NetWeaver Application Server ABAP-based products and SAP NetWeaver Application Server Java-based products are supported on IBM Power Virtual Servers. As some of the older versions are no longer supported, refer to respective SAP Notes to validate and deploy supported combinations of these components on these SAP landscapes.

To setup RHEL 8.x Linux operating systems suitable for SAP NetWeaver or S/4HANA application landscapes, refer to *SAP Note 772999 - Red Hat Enterprise Linux 8.x: Installation and Configuration*

6.2.1 HA Solution for NetWeaver or S/4 Based on ABAP Platform 1709 or older

The Standalone Enqueue Server (ENSA1) is an older ASCS implementation. When using ENSA1 the Enqueue Client is required to fail over to the cluster node where the active ERS instance is running limiting your configuration flexibility. It must be ensured that the ERS instance is shut down and moved to the node where the ASCS instance was running when it was online because the ASCS instance has to access the shared memory where the ERS instance maintains a copy of the enqueue lock table to continue to keep track of the enqueue locks of active transactions.

It is important to note that when using stand alone enqueue server (ENSA1), the ASCS instance must always move to the node where the ERS instance is running and can only be deployed as a two node Pacemaker cluster. The detailed configuration guides for Red Hat are found in [this link](#). The cluster configuration is illustrated in Figure 6-2.

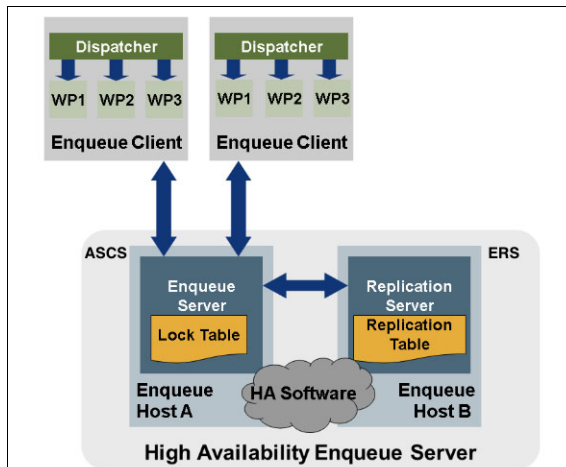


Figure 6-2 Standalone Enqueue Server 1 (ENSA1)²

It is recommended that this deployment architecture be upgraded to ENSA2 which is the current standard default deployment configuration and an improvement over the older setups. Refer to section 6.2.2, “HA Solution for S/4HANA based on ABAP Platform 1809 or newer” on page 87 for further details.

6.2.2 HA Solution for S/4HANA based on ABAP Platform 1809 or newer

The Standalone Enqueue Server has evolved into generation 2 with ABAP Platform 1809 and later, known as Standalone Enqueue Server 2, or ENSA2. In ENSA2, if the ASCS fails, it can start on a separate node in a cluster and the new instance will copy the lock entries from the enqueue replicator 2.

Figure 6-3 provides an illustration of the Standalone Enqueue Server which is a component of Application Server ABAP. ENSA2 provides a mechanism to ensure the high availability of the lock table and its entries.

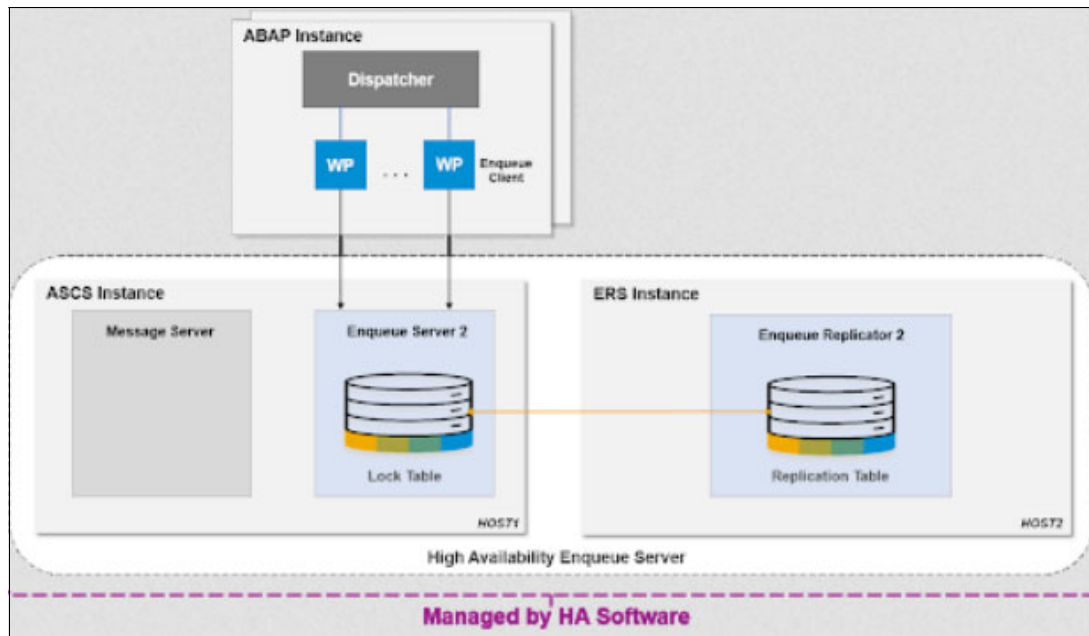


Figure 6-3 ENSA2 Enqueue Replicator³

Starting with ABAP Application 1809, ENSA2 based deployment has become the default installation configuration for new installations.

There are two main supported scenarios as per Table 6-2.

Table 6-2 Cluster configuration options

Supported Scenario	Description
Multi-Node Cluster	Multi-node clusters are supported since there is no need for ENSA2 ASCS to “follow” ERS
Two-node cluster	The existing or upgrades from ENSA1 to ENSA2 deployments are easily adjustable and usable with two-node clusters when switching from ENSA1 to ENSA2

² Red Hat HA Solutions for SAP

³ Red Hat documentation on HA solutions.

Brand new installations can take advantage of the new features when designing the architecture, now having the option of choosing between a two-node cluster or a multi-node cluster. Figure 6-4 shows a reference architecture diagram of a typical 3-node cluster.

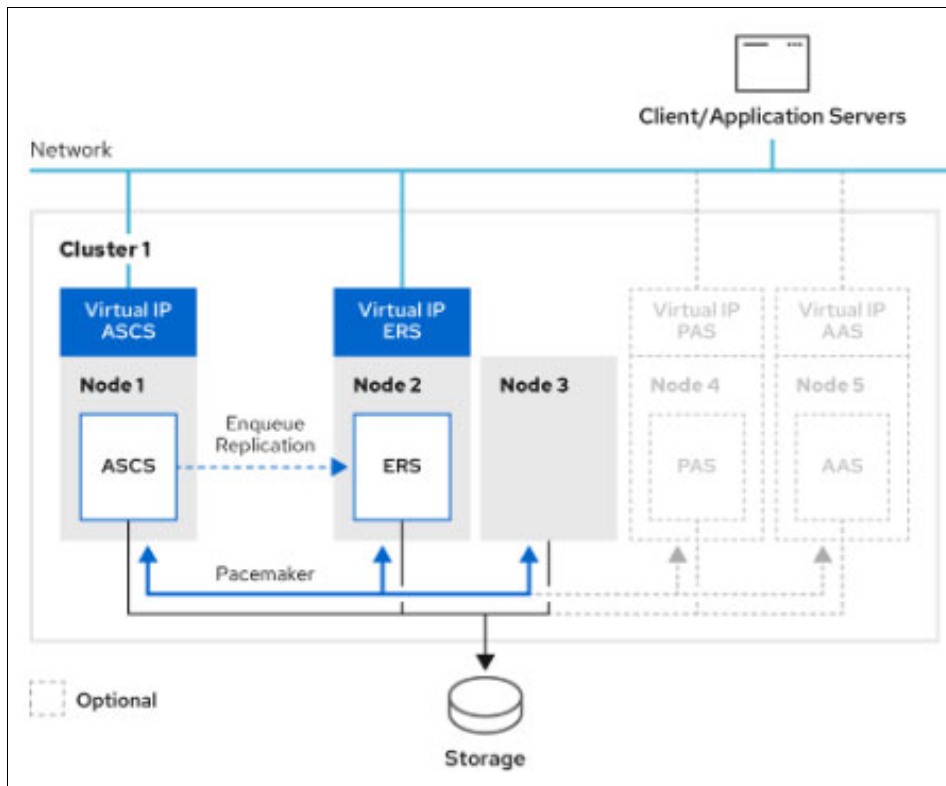


Figure 6-4 ENSA2 Deployment Reference Architecture⁴

Note: More nodes can be added to this design, depending upon the specific datacenter requirements.

6.2.3 Configuring high availability for SAP S/4HANA (ASCS and ERS) in a RHEL HA Add-On cluster

The Red Hat documentation [Red Hat HA Solutions for SAP HANA, S/4HANA, and NetWeaver based SAP Applications](#) discusses various configuration aspects for SAP solutions in a RHEL High Availability Add-On cluster utilizing Pacemaker. The following information describes the configuration of ABAP SAP Central Services (ASCS) and Enqueue Replication Service (ERS) instances in an RHEL HA Add-On cluster on IBM Power Systems Virtual Server. The focus of these example configuration is on the second generation of the Standalone Enqueue Server, or ENSA2.

As mentioned before, starting with the release of SAP S/4HANA 1809, ENSA2 is installed by default and can be configured in either a two-node cluster or a multi-node cluster. This example uses the ENSA2 setup for a two-node RHEL HA Add-On cluster.

⁴ [Red Hat HA Solutions for SAP](#)

If the ASCS service fails in a two-node cluster, it restarts on the node where ERS is running. The lock entries for the SAP application are protected when they rebuild from the copy of the lock table that's available in the ERS. When the failed cluster node reactivates, the ERS instance moves to the other node (anti-collocation) to protect the lock table copy.

Note: This cluster is designed to protect the ASCS and ERS services. It is recommended that you install the SAP database instance and other SAP application servers on virtual server instances in a separate cluster.

Installing on PowerVS

The following information outlines the configuration of ASCS and ERS in a RHEL HA Add-On cluster with Red Hat Enterprise Linux 8.4 (RHEL) by using IBM Power Systems Virtual Server instances as cluster nodes.

For more information, see the following Red Hat knowledge base articles:

- [Configuring SAP S/4HANA ASCS/ERS with Standalone Enqueue Server 2 \(ENSA2\) in Pacemaker](#)
- [Support Policies for RHEL High Availability Clusters - Management of SAP S/4HANA in a cluster](#)

Prerequisites

Before starting, please ensure that you understand and have satisfied the following prerequisites:

- ▶ A Red Hat Customer Portal account.
- ▶ An IBM Cloud account.
- ▶ A valid RHEL for SAP Applications or RHEL for SAP Solutions subscription is required to enable the repositories that you need to install SAP HANA and the resource agents for HA configurations.
- ▶ Virtual server instances need to fulfill the hardware and resource requirements for the SAP instances in scope.
- ▶ This information describes a setup that uses shareable storage volumes accessible on both cluster nodes. Certain file systems are created on shareable storage volumes so that they can be mounted on both cluster nodes. This setup applies to both instance directories:
 - `/usr/sap/<SID>/ASCS<Inst#>` of the ASCS instance.
 - `/usr/sap/<SID>/ERS<Inst#>` of the ERS instance.
- ▶ Make sure that the storage volumes that were created for those file systems are attached to both virtual server instances. During SAP instance installation and RHEL HA Add-On cluster configuration, each instance directory must be mounted on its appropriate node. HA-LVM makes sure that each of the two instance directories is mounted on only one node at a time.

Important: Different storage setups for the instance directories, for example, NFS mounts, are possible. Choosing a different storage configuration will require changes in the setup steps and creation of cluster file system resources as appropriate.

- ▶ The virtual hostname for ASCS instance and ERS instance must meet the requirements as documented in *Hostname of SAP ABAP Platform servers*. Make sure that the virtual IP addresses for the SAP instances are assigned to a network adapter and that they can communicate in the network.
- ▶ Make sure that you subscribed to either RHEL for SAP Applications or RHEL for SAP Solutions and that you enabled the Red Hat Enterprise Linux for SAP Solutions for Power LE -

Update Services for SAP Solutions 8.4 ppc64le (rhel-8-for-ppc64le-sap-solutions-e4s-rpms) repository. You need to set the release to 8.4 through the subscription manager.

- ▶ SAP application server instances require a common shared file system SAPMNT */sapmnt/<SID>* with read and write access, and other shared file systems such as SAPTRANS */usr/sap/trans*. These file systems are typically provided by an external NFS server. The NFS server must be installed on virtual servers that are not part of the ENSA2 cluster.
- ▶ Configuring an active-passive NFS server in a Red Hat High Availability cluster describes the implementation of an active-passive NFS server in a RHEL HA Add-On cluster with Red Hat Enterprise Linux 8 by using virtual server instances in Power Systems Virtual Server.
- ▶ Make sure that all SAP installation media is available.

Preparing nodes for SAP installation

This section describes how to prepare the nodes for an SAP installation.

Preparing environment variables

To simplify the setup, prepare the following environment variables for user root on both cluster nodes. These environment variables are used in subsequent commands in the remainder of the instructions.

1. On both nodes, create a file with the environment variables shown in Example 6-1 making sure to adapt them to your configuration.

Example 6-1 File “sap_env.sh” to be created on both nodes

```
export SID=<SID>                # SAP System ID (uppercase)
export sid=<sid>                # SAP System ID (lowercase)

export ASCS_nr=<INSTNO>        # Instance Number for ASCS
export ERS_nr=<INSTNO>        # Instance Number for ERS

# virtual hostnames
export ASCS_vh=<virtual hostname> # virtual hostname for ASCS
export ERS_vh=<virtual hostname> # virtual hostname for ERS
export ASCS_ip=<IP address>      # virtual IP address for ASCS
export ERS_ip=<IP address>      # virtual IP address for ERS

# LVM storage for instance file systems
export ASCS_vg=<vg name>        # volume group name for ASCS
export ERS_vg=<vg name>        # volume group name for ERS

export ASCS_lv=<lv name>        # logical volume name for ASCS
export ERS_lv=<lv name>        # logical volume name for ERS
```

Important: It is recommended to use meaningful names for the volume groups and logical volumes that designate their content. For example, include the SID and *ascs* or *ers* in the name. Don't use hyphens in the volume group or logical volume names. For example:

- *s01ascsvg* and *s01ascslv*
- *s01ersvg* and *s01erslv*

2. Next apply the variables in both nodes by running the command shown in Example 6-2.

Example 6-2 Source “sap_env.sh” and set environment variables on both nodes.

```
source sap_envs.sh
```

Assigning virtual IP addresses

3. Check whether the virtual IP address for the SAP instance is present. Otherwise, you need to identify the correct network adapter to assign the IP address. The command is shown in Example 6-3.

Example 6-3 On both nodes, check the list of currently active IP addresses

```
ip -o -f inet address show | '/scope global/ {print $2, $4}'
```

Sample output from the previous command is shown in Example 6-4.

Example 6-4 Sample output of IP addresses

```
# ip -o -f inet address show | awk '/scope global/ {print $2, $4}'
env2 10.51.0.66/24
env3 10.52.0.41/24
env4 10.111.1.28/24
```

The device name of the network adapter appears in the first column. The second column lists the active IP addresses and the number of bits that are reserved for the netmask - which are separated by a slash.

4. If the virtual IP address for the SAP instance is not present, make sure that it isn't erroneously set on another virtual server instance by running the command in Example 6-5.

Example 6-5 Ping to ensure VIP not assigned

```
ping -c 3 ${ASCS_vh}
```

The results of the ping should show that the address is unreachable as shown in Example 6-6.

Example 6-6 Ping Sample Output

```
# ping -c 2 cl-sap-scs
PING cl-sap-scs (10.111.1.248) 56(84) bytes of data.
From cl-sap-1.tst.ibm.com (10.111.1.28) icmp_seq=1 Destination Host Unreachable
From cl-sap-1.tst.ibm.com (10.111.1.28) icmp_seq=2 Destination Host Unreachable

--- cl-sap-ers ping statistics ---
2 packets transmitted, 0 received, +2 errors, 100% packet loss, time 2112ms
pipe 3
```

Note: If the ping output shows Destination Host Unreachable, the IP address is available, and you can assign the IP alias to the virtual server instance. Use the correct device name of the network adapter that matches the subnet of the IP address.

5. Add the IP address definition on both of the nodes using the command in Example 6-7.

Example 6-7 Add IP addresses on Nodes

```
On Node1
# ip addr add ${ASCS_ip} dev env4
On Node2
# ip addr add ${ERS_ip} dev env4
```

Note: Adjust the device name for the network adapter based on your configuration.

Preparing volume groups, logical volumes, and shared file systems

Shared storage is an important resource in an ENSA2 cluster. ASCS and ERS must be able to run on both nodes, and their runtime environment is stored in the shared storage volumes. All cluster nodes need to access the shared storage volumes, but only one node has exclusive read and write access to a volume.

Preparing High Availability Logical Volume Manager settings

1. Edit file `/etc/lvm/lvm.conf` to include the system ID in the volume group as shown in Example 6-8.

Example 6-8 On both nodes, edit the lvm.conf file

```
On Both Nodes
# vi /etc/lvm/lvm.conf

Search for parameter system_id_source and change its value to "uname".
# system_id_source = "uname"
```

Identifying World Wide Names of shared storage volumes

2. Identify the World Wide Name (WWN) for each storage volume that is in the shared volume groups.
 - Log in to IBM Cloud to the Storage volumes view of Power Systems Virtual Server.
 - Select your workspace.
 - Filter for the volume prefix in the Storage volumes list, and identify all the World Wide Names of the volumes that are in scope for ASCS and ERS (the World Wide Name is a 32-digit hexadecimal number).

Note: Make sure that the attribute Shareable is On for those volumes.

3. In the Virtual server instances view, go to both virtual server instances of the cluster. Verify that all volumes that are in scope for ASCS and ERS are attached to both virtual server instances.

If you need to attach a new storage volume to a virtual server instance, make sure that you rescan the SCSI bus to detect the new volume. Afterward, update the multi-path configuration of the virtual server instance. The command for this is shown in Example 6-9 which should be run on both nodes.

Example 6-9 Command to rescan for new devices and update multi-path

```
# rescan-scsi-bus.sh && sleep 10 && multipathd reconfigure
```

4. Log in to both cluster nodes, and add the WWN to the environment variables of user root as shown in Example 6-10.

Example 6-10 Add WWN to environment variables

```
On NODE1, run the following command.
# export ASCS_pvid=<WWN> # WWN of shared storage volume for ASCS

On NODE2, run the following command.
# export ERS_pvid=<WWN> # WWN of shared storage volume for ERS
```

Note: The matching WWN values can also be found with the **pvs --all** command.

Make sure that you set the environment variable by using the hexadecimal number with lowercase letters.

Creating physical volumes

5. Create physical volumes on the nodes for the attached storage volumes as shown in Example 6-11.

Example 6-11 Create on both Nodes

```
On NODE1
# pvcreate /dev/mapper/${ASCS_pvid}

On NODE2
# pvcreate /dev/mapper/${ERS_pvid}
```

Sample output from this command is shown in Example 6-12.

Example 6-12 pvcreate output on each node

```
# On NODE1
pvcreate /dev/mapper/${ERS_pvid}
Physical volume "/dev/mapper/36005076810810335700000000002e31" successfully created.

# On NODE2
pvcreate /dev/mapper/${ASCS_pvid}
Physical volume "/dev/mapper/36005076810810335700000000002ddc" successfully created.
```

Creating volume groups

6. Create volume groups for ASCS and ERS using the commands shown in Example 6-13.

Example 6-13 Create on both Nodes

```
On NODE1
# vgcreate ${ASCS_vg} /dev/mapper/${ASCS_pvid}
Verify that the System ID is set.
# vgs -o+systemid

On NODE2
# vgcreate ${ERS_vg} /dev/mapper/${ERS_pvid}
Verify that the System ID is set.
# vgs -o+systemid
```

Sample output is shown in Example 6-14.

Example 6-14 vgs output on each node

```
# On NODE1
vgs -o+systemid
VG          #PV #LV #SN Attr   VSize  VFree  System ID
s01ascsvg   1  0  0 wz--n- <50.00g <50.00g cl-sap-1

# On NODE2
vgs -o+systemid
VG          #PV #LV #SN Attr   VSize  VFree  System ID
s01ersvg    1  0  0 wz--n- <50.00g <50.00g cl-sap-2
```

Creating logical volumes and file systems

7. Create the logical volume for the ASCS and ERS, then format it as XFS file system as shown in Example 6-15.

Example 6-15 Create Filesystems on both Nodes

```
On NODE1
# lvcreate -l 100%FREE -n ${ASCS_lv} ${ASCS_vg}
# mkfs.xfs /dev/${ASCS_vg}/${ASCS_lv}

On NODE2
# lvcreate -l 100%FREE -n ${ERS_lv} ${ERS_vg}
# mkfs.xfs /dev/${ERS_vg}/${ERS_lv}
```

Making sure that a volume group is not activated on multiple cluster nodes

8. Volume groups that are managed by the cluster must not activate automatically on startup. On both nodes, edit file `/etc/lvm/lvm.conf` and modify the `auto_activation_volume_list` entry to limit auto activation to specific volume groups. This is shown in Example 6-16.

Example 6-16 Edit /etc/lvm/lvm.conf on both Nodes, recreate boot image and reboot

```
On Both Nodes
# vi /etc/lvm/lvm.conf

Search for parameter auto_activation_volume_list and add the
volume groups, other than the volume group that you defined for
the NFS cluster, as entries in that list.
# auto_activation_volume_list = [ "rhel_root" ]

On both NODES, Rebuild the initramfs boot image to make sure that
the boot image does not activate a volume group that is controlled
by the cluster.
# dracut -H -f /boot/initramfs-$(uname -r).img $(uname -r)

Reboot both nodes.
# shutdown -t now -r
```

Note: For RHEL 8.5 and later, you can disable auto activation for a volume group when you create the volume group by specifying the `--setautoactivation n` flag for the `vgcreate` command.

Mounting the file systems for SAP installation

9. Activate the volume groups and mount the SAP instance file systems. Run the commands shown in Example 6-17,

Example 6-17 Run on Node1(ASCS) and Node 2(ERS)

```
On NODE1 (ASCS)
# vgchange -a y ${ASCS_vg}
# mkdir -p /usr/sap/${SID}/ASCS${ASCS_nr}
# mount /dev/${ASCS_vg}/${ASCS_lv} /usr/sap/${SID}/ASCS${ASCS_nr}

On NODE2 (ERS)
# vgchange -a y ${ERS_vg}
# mkdir -p /usr/sap/${SID}/ERS${ERS_nr}
# mount /dev/${ERS_vg}/${ERS_lv} /usr/sap/${SID}/ERS${ERS_nr}
```

Mounting the required NFS file systems

10. On both nodes, make sure that the NFS file systems `/sapmnt` and `/usr/sap/trans` are mounted by running the commands shown in Example 6-18.

Example 6-18 Check NFS mounts

```
On NODE1 (ASCS)
# mount | grep nfs
On NODE2 (ERS)
# mount | grep nfs
```

Installing SAP instances

Use the SAP Software Provisioning Manager (SWPM) to install all instances.

Next it is time to install SAP instances on the cluster nodes. This is done using the SAP installer SWPM.

- Install ASCS instance on NODE1 by using the virtual hostname `${ASCS_vh}` that is associated with the virtual IP address for ASCS:


```
<swpm>/sapinst SAPINST_USE_HOSTNAME=${ASCS_vh}
```
- Install ERS instance on NODE2 by using the virtual hostname `${ERS_vh}` that is associated with the virtual IP address for ERS:


```
<swpm>/sapinst SAPINST_USE_HOSTNAME=${ERS_vh}
```
- Install instances outside the cluster:
 - DB instance
 - PAS instance
 - AAS instances

Installing and setting up the RHEL HA Add-On cluster

Now that your SAP landscape is installed and ready to go it is time to enable the Pacemaker function to manage the high availability functions. For additional details on this process see [Implementing a RHEL HA Add-On cluster on IBM Power Virtual Server](#).

Preparing ASCS and ERS instances for the cluster integration

Use the following steps to prepare the SAP instances for the cluster integration.

1. Removing the ASCS and ERS entries from the SAP services files

Adjust the SAP services file `/usr/sap/sapservices` to prevent automatic start of the `sapstartsrv` instance agent for both ASCS and ERS instance after a reboot as shown in Example 6-19.

Example 6-19 On both nodes comment the sapstartsrv entries for ASCS and ERS

```
On NODE1 (ASCS)
# vi /usr/sap/sapservices

On NODE2 (ERS)
# vi /usr/sap/sapservices
```

Example 6-20 shows an example of what line to comment out.

Example 6-20 Example Commented sapstartsrv entries

```
On NODE1 and NODE2
```

```
#LD_LIBRARY_PATH=/usr/sap/S01/ASCS01/exe:$LD_LIBRARY_PATH;export
LD_LIBRARY_PATH;/usr/sap/S01/ASCS01/exe/sapstartsrv
pf=/usr/sap/S01/SYS/profile/S01_ASCS01_c1-sap-scs -D -u s01adm
```

2. Creating mount points for the instance file systems on the takeover node

Create the mount points for the instance file systems and adjust their ownership using the commands shown in Example 6-21.

Example 6-21 Create mount points for ASCS and ERS

```
On NODE1 (ASCS)
# mkdir /usr/sap/${SID}/ASCS${ASCS_nr}
# chown ${sid}adm:sapsys /usr/sap/${SID}/ASCS${ASCS_nr}

On NODE2 (ERS)
# mkdir /usr/sap/${SID}/ERS${ERS_nr}
# chown ${sid}adm:sapsys /usr/sap/${SID}/ERS${ERS_nr}
```

3. Installing permanent SAP license keys

When the SAP ASCS instance is installed on a Power Systems Virtual Server instance, the SAP license mechanism relies on the partition UUID. For more information, see *SAP note 2879336 - Hardware key based on unique ID*.

On both nodes, run the commands shown in Example 6-22 as user `<sid>adm` to identify the **HARDWARE KEY** of the node.

Example 6-22 On both Nodes Get SAP License hardware Key

```
On Both NODEs
# sudo -i -u ${sid}adm -- sh -c 'saplikey -get'
```

Sample output is provided in Example 6-23.

Example 6-23 Sample Hardware Key Output

```
$ sudo -i -u ${sid}adm -- sh -c 'saplikey -get'

saplikey: HARDWARE KEY = H1428224519
```

Note the **HARDWARE KEY** of each node as you need both hardware keys to request two different SAP license keys. Check the following SAP notes for more information about requesting SAP license keys:

- *SAP Note - 2879336 - Hardware key based on unique ID*
- *SAP Note - 2662880 - How to request SAP license keys for failover systems*

4. Installing SAP resource agents

Install the required software packages. The **resource-agents-sap** includes the SAP Instance cluster resource agent for managing the SAP instances.

Unless **sap_cluster_connector** is configured for the SAP instance, the RHEL HA Add-On cluster considers any state change of the instance as an issue. If other SAP tools such as **sapcontrol** are used to manage the instance, then **sap_cluster_connector** grants permission to control SAP instances that are running inside the cluster. If the SAP instances are managed by only cluster tools, the implementation of **sap_cluster_connector** is not necessary.

Install the packages for the resource agent and the SAP Cluster Connector library. For more information, see How to enable the SAP HA Interface for SAP ABAP application server instances managed by the RHEL HA Add-On

On both nodes, run the command shown in Example 6-24. If needed, use **subscription-manager** to enable the SAP NetWeaver repository.

Example 6-24 Enable Red Hat subscriptions

```
On Both NODEs
# subscription-manager repos --enable="rhel-8-for-ppc64le-sap-netweaver-e4s-rpms"
```

5. Configuring SAP Cluster Connector

Add user `${sid}adm` to the `haclient` group using the command in Example 6-25.

Example 6-25 Add user to haclient group

```
On Both NODEs
# usermod -a -G haclient ${sid}adm
```

6. Adapting the SAP instance profiles

Modify the start profiles of all SAP instances that are managed by SAP tools outside the cluster. Both ASCS and ERS instances can be controlled by the RHEL HA Add-On cluster and its resource agents. Adjust the SAP instance profiles to prevent an automatic restart of instance processes using the commands shown in Example 6-26.

Example 6-26 Update SAP Profiles for ASCS and ERS

```
On NODE1
Change all occurrences of Restart_Program to Start_Program in the instance profile of both ASCS and ERS.
```

```
# cd /sapmnt/${SID}/profile
# sed -i -e 's/Restart_Program_\([0-9][0-9]\)/Start_Program_\1/'
${SID}_ASCS${ASCS_nr}_${ASCS_vh}
# sed -i -e 's/Restart_Program_\([0-9][0-9]\)/Start_Program_\1/'
${SID}_ERS${ERS_nr}_${ERS_vh}
```

Add the following two lines at the end of the SAP instance profile to configure `sap_cluster_connector` for the ASCS and ERS instance, and activate the integration.

```
service/haLib = $(DIR_EXECUTABLE)/saphascriptco.so
service/haLib_cluster_connector = /usr/bin/sap_cluster_connector
```

Configuring ASCS and ERS cluster resources

Up to this point, the following are assumed:

- ▶ A RHEL HA Add-On cluster is running on both virtual server instances and fencing of the nodes was tested.
- ▶ The SAP System is running.
 - SAP ASCS is installed and active on node 1 of the cluster.
 - SAP ERS is installed and active on node 2 of the cluster.
- ▶ All steps in Prepare ASCS and ERS instances for the cluster integration are complete.

1. Configuring resource for sapmnt share

Create a cloned Filesystem cluster resource to mount the SAPMNT share from an external NFS server to all cluster nodes. Make sure that the environment variable `${NFS_vh}` is set to the virtual hostname of your NFS server `${NFS_vh}`, and

`{NFS_options}` according to your mount options. Use the commands shown in Example 6-27.

Example 6-27 Set Mount Options

```
On Both NODEs
# export NFS_options="rw,sec=sys"

On NODE1
# pcs resource create fs_sapmnt Filesystem \
  device="{NFS_vh}:/${SID}" \
  directory="/sapmnt/${SID}" \
  fstype='nfs' \
  "options={NFS_options}" \
  clone interleave=true
```

Configuring the ASCS resource group

2. Create a resource for the virtual IP address of the ASCS as shown in Example 6-28.

Example 6-28 Create ASCS Resource Group

```
On NODE1
# pcs resource create ${sid}_vip_asc${ASCS_nr} IPAddr2 \
  ip={ASCS_ip} \
  --group ${sid}_asc${ASCS_nr}_group
```

In Example 6-29 we show creating resources for an HA-LVM file system on a shared storage volume, we create resources for LVM-activate and for the instance file system of the ASCS.

Example 6-29 Create Resource with LVM Filesystem

```
On NODE1
# pcs resource create ${sid}_fs_asc${ASCS_nr}_lvm LVM-activate \
  vgname="{ASCS_vg}" \
  vg_access_mode=system_id \
  --group ${sid}_asc${ASCS_nr}_group

# pcs resource create ${sid}_fs_asc${ASCS_nr} Filesystem \
  device="/dev/mapper/{ASCS_vg}-{ASCS_lv}" \
  directory="/usr/sap/${SID}/ASC${ASCS_nr}" \
  fstype=xfs \
  --group ${sid}_asc${ASCS_nr}_group
```

In the alternative example shown in Example 6-30, the instance file system of the ASCS is provided by an HA NFS server, only the file system resource is required. Make sure that you defined the environment variable `{NFS_vh}` according to your NFS server, and that you created a directory `{SID}/ASC` under the NFS root directory during the SAP installation of the ASCS instance.

Example 6-30 Create Resource with NFS Filesystem

```
On NODE1
# pcs resource create ${sid}_fs_asc${ASCS_nr} Filesystem \
  device="{NFS_vh}:/${SID}/ASC" \
  directory="/usr/sap/${SID}/ASC${ASCS_nr}" \
  fstype=nfs \
  force_unmount=safe \
```

```

op start interval=0 timeout=60 \
op stop interval=0 timeout=120 \
--group ${sid}_ascs${ASCS_nr}_group

```

3. Create a resource for managing the ASCS instance. This is shown in Example 6-31.

Example 6-31 Create Resource for Managing ASCS Instance

```

On NODE1
# pcs resource create ${sid}_ascs${ASCS_nr} SAPInstance \
  InstanceName="${SID}_ASCS${ASCS_nr}_${ASCS_vh}" \
  START_PROFILE=/sapmnt/${SID}/profile/${SID}_ASCS${ASCS_nr}_${ASCS_vh} \
  AUTOMATIC_RECOVER=false \
  meta resource-stickiness=5000 \
  migration-threshold=1 failure-timeout=60 \
  op monitor interval=20 on-fail=restart timeout=60 \
  op start interval=0 timeout=600 \
  op stop interval=0 timeout=600 \
  --group ${sid}_ascs${ASCS_nr}_group

```

The meta **resource-stickiness=5000** option is used to balance the failover constraint with ERS so that the resource stays on the node where it started and doesn't migrate uncontrollably in the cluster.

4. Add a resource stickiness to the group to make sure that the ASCS remains on the node as shown in Example 6-32.

Example 6-32 Add a resource with stickiness to remain on node

```

On NODE1
# pcs resource meta ${sid}_ascs${ASCS_nr}_group resource-stickiness=3000

```

Configuring ERS resource group

5. Create a resource for the virtual IP address of the ERS as shown in Example 6-33.

Example 6-33 Create ERS Resource Group

```

On NODE1
# pcs resource create ${sid}_vip_ers${ERS_nr} IPAddr2 \
  ip=${ERS_ip} \
  --group ${sid}_ers${ERS_nr}_group

```

Example 6-34 shows creating resources for an HA-LVM file system on a shared storage volume, you create resources for LVM-activate and for the instance file system of the ERS.

Example 6-34 Create Resource with LVM Filesystem

```

On NODE1
# pcs resource create ${sid}_fs_ers${ERS_nr}_lvm LVM-activate \
  vgname="${ERS_vg}" \
  vg_access_mode=system_id \
  --group ${sid}_ers${ERS_nr}_group
# pcs resource create ${sid}_fs_ers${ERS_nr} Filesystem \
  device="/dev/mapper/${ERS_vg}-${ERS_lv}" \
  directory=/usr/sap/${SID}/ERS${ERS_nr} \
  fstype=xfs \
  --group ${sid}_ers${ERS_nr}_group

```

In the alternative example shown in Example 6-35, the instance file system of the ERS is provided by an HA NFS server, only the file system resource is required. Make sure that you defined the environment variable `#{NFS_vh}` according to your NFS server, and that you created a directory `#{SID}/ERS` under the NFS root directory during the SAP installation of the ERS instance.

Example 6-35 Create Resource with NFS Filesystem

```
On NODE1

# pcs resource create ${sid}_fs_ers${ERS_nr} Filesystem \
  device="#{NFS_vh}:${SID}/ERS" \
  directory=/usr/sap/${SID}/ERS${ERS_nr} \
  fstype=nfs \
  force_unmount=safe \
  op start interval=0 timeout=60 \
  op stop interval=0 timeout=120 \
  --group ${sid}_ers${ERS_nr}_group
```

6. Create a resource for managing the ERS instance which is shown in Example 6-36.

Example 6-36 Create Resource for Managing ERS Instance

```
On NODE1

# pcs resource create ${sid}_ers${ERS_nr} SAPInstance \
  InstanceName="#{SID}_ERS${ERS_nr}_#{ERS_vh}" \
  START_PROFILE=/sapmnt/${SID}/profile/${SID}_ERS${ERS_nr}_#{ERS_vh} \
  AUTOMATIC_RECOVER=false \
  IS_ERS=true \
  op monitor interval=20 on-fail=restart timeout=60 \
  op start interval=0 timeout=600 \
  op stop interval=0 timeout=600 \
  --group ${sid}_ers${ERS_nr}_group
```

Configuring constraints

7. A collocation constraint prevents resource groups `#{sid}_ascs${ASCS_nr}_group` and `#{sid}_ers${ERS_nr}_group` from being active on the same node whenever possible. The stickiness score of **-5000** makes sure that they run on the same node if only a single node is available. See Example 6-37 for the details.

Example 6-37 Add collocation constraints and Order constraints

```
Add constraint
# pcs constraint collocation add \
  ${sid}_ers${ERS_nr}_group with ${sid}_ascs${ASCS_nr}_group -5000

Add the following Order constraint to control resource group ${sid}_ascs${ASCS_nr}_group
starts before ${sid}_ers${ERS_nr}_group.
# pcs constraint order start \
  ${sid}_ascs${ASCS_nr}_group then stop ${sid}_ers${ERS_nr}_group \
  symmetrical=false \
  kind=Optional

Add following two order constraints make sure that the file system SAPMNT mounts before
resource groups ${sid}_ascs${ASCS_nr}_group and ${sid}_ers${ERS_nr}_group start.
# pcs constraint order fs_sapmnt-clone then ${sid}_ascs${ASCS_nr}_group
# pcs constraint order fs_sapmnt-clone then ${sid}_ers${ERS_nr}_group
```

At this point cluster setup is complete.

Testing SAP ENSA2 clusters

It is vital to thoroughly test the cluster configuration to make sure that the cluster is working correctly. The following information provides a few sample failover test scenarios, but is not a complete list of test scenarios.

For example, the description of each test case includes the following information.

- Component under test
- Description of the test
- Prerequisites and the initial state before failover test
- Test procedure
- Expected behavior and results
- Recovery procedure

Test1 - Testing failure of the ASCS instance

This test is testing the failover of the ASCS instance when the ASCS service fails.

► **Test1 - Description**

Simulate a failure of the SAP ASCS instance that is running on NODE1.

► **Test1 - Prerequisites**

- A functional two-node RHEL HA Add-On cluster for SAP ENSA2.
- Both cluster nodes are active.
- Cluster is started on NODE1 and NODE2.
 - Resource group `${sid}_ascs${ASCS_nr}_group` is active on NODE1.
 - Resources `${sid}_vip_ascs${ASCS_nr}`, `${sid}_fs_ascs${ASCS_nr}_lvm`, `${sid}_fs_ascs${ASCS_nr}` and `${sid}_ascs${ASCS_nr}` are **Started** on NODE1.
 - Resource group `${sid}_ers${ERS_nr}_group` is active on NODE2.
 - Resources `${sid}_vip_ers${ERS_nr}`, `${sid}_fs_ers${ERS_nr}_lvm`, `${sid}_fs_ers${ERS_nr}` and `${sid}_ers${ERS_nr}` are **Started** on NODE2.

To validate the status of the SAP instances and show where each instance is running, run the command shown in Example 6-38. The expected starting point is:

- ASCS instance is running on NODE1.
- ERS instance is running on NODE2.

Example 6-38 Cluster Status

```
# pcs status
```

A sample output is shown in Example 6-39.

Example 6-39 Sample Cluster Status

```
# pcs status
Cluster name: SAP_ASCS
Cluster Summary:
 * Stack: corosync
 * Current DC: cl-sap-1 (version 2.0.5-9.e18_4.5-ba59be7122) - partition with quorum
 * Last updated: Tue Feb 14 07:59:16 2023
 * Last change: Tue Feb 14 05:02:22 2023 by hacluster via crmd on cl-sap-1
 * 2 nodes configured
 * 11 resource instances configured
```

Node List:

```
 * Online: [ cl-sap-1 cl-sap-2 ]
```

Full List of Resources:

```
* res_fence_ibm_powervs(stonith:fence_ibm_powervs): Started cl-sap-2
* Resource Group: s01_ascs01_group:
  * s01_vip_ascs01(ocf::heartbeat:IPaddr2): Started cl-sap-1
  * s01_fs_ascs01_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-1
  * s01_fs_ascs01(ocf::heartbeat:Filesystem): Started cl-sap-1
  * s01_ascs01(ocf::heartbeat:SAPInstance): Started cl-sap-1
* Resource Group: s01_ers02_group:
  * s01_vip_ers02(ocf::heartbeat:IPaddr2): Started cl-sap-2
  * s01_fs_ers02_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-2
  * s01_fs_ers02(ocf::heartbeat:Filesystem): Started cl-sap-2
  * s01_ers02(ocf::heartbeat:SAPInstance): Started cl-sap-2
* Clone Set: fs_sapmnt-clone [fs_sapmnt]:
  * Started: [ cl-sap-1 cl-sap-2 ]
```

Daemon Status:

```
corosync: active/disabled
Pacemaker: active/disabled
pcsd: active/enabled
```

► Test1 - Test Procedure

Forcefully terminate the SAP ASCS instance by sending a SIGKILL signal as user `#{sid}adm` to the enqueue server. On NODE1, identify the PID of the enqueue server using the command shown in Example 6-40.

Example 6-40 On Node1 identify PID for Enqueue Server

```
# pgrep -af "(en|enq).sap"
```

The output of the command is shown in Example 6-41. After discovering the PID, use the `kill` command as shown in the example.

Example 6-41 Locating PID and issuing kill command

```
# pgrep -af "(en|enq).sap"
30186 en.sapS01_ASCS01 pf=/usr/sap/S01/SYS/profile/S01_ASCS01_cl-sap-scs
# kill -9 30186
```

► Test1 - Expected behavior

The following actions should occur after you issue the `kill` command for the ASCS instance.

- SAP ASCS instance on NODE1 fails.
- The cluster detects the failure of the ASCS instance.
- The cluster stops the dependent resources on NODE1 (virtual IP address, file system `/usr/sap/#{SID}/ASCS#{ASCS_nr}`, and the LVM resources), and acquires them on NODE2.
- The cluster starts the ASCS on NODE2.
- The cluster stops the ERS instance on NODE2.
- The cluster stops the dependent resources on NODE1 (virtual IP address, file system `/usr/sap/#{SID}/ERS#{ERS_nr}`, and the LVM resources), and acquires them on NODE1.
- The cluster starts the ERS on NODE1.

To validate the outcome, after a few seconds check the status using the command shown in Example 6-42.

Example 6-42 Command to check status

```
# pcs status
```

The output should look similar to what is shown in Example 6-43.

Example 6-43 Status after failover

```
# pcs status
Cluster name: SAP_ASCS
Cluster Summary:
  * Stack: corosync
  * Current DC: cl-sap-1 (version 2.0.5-9.e18_4.5-ba59be7122) - partition with quorum
  * Last updated: Tue Feb 14 08:10:18 2023
  * Last change: Tue Feb 14 05:02:22 2023 by hacluster via crmd on cl-sap-1
  * 2 nodes configured
  * 11 resource instances configured

Node List:
  * Online: [ cl-sap-1 cl-sap-2 ]

Full List of Resources:
  * res_fence_ibm_powervs(stonith:fence_ibm_powervs): Started cl-sap-2
  * Resource Group: s01_ascs01_group:
    * s01_vip_ascs01(ocf::heartbeat:IPaddr2): Started cl-sap-2
    * s01_fs_ascs01_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-2
    * s01_fs_ascs01(ocf::heartbeat:Filesystem): Started cl-sap-2
    * s01_ascs01(ocf::heartbeat:SAPInstance): Started cl-sap-2
  * Resource Group: s01_ers02_group:
    * s01_vip_ers02(ocf::heartbeat:IPaddr2): Started cl-sap-1
    * s01_fs_ers02_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-1
    * s01_fs_ers02(ocf::heartbeat:Filesystem): Started cl-sap-1
    * s01_ers02(ocf::heartbeat:SAPInstance): Started cl-sap-1
  * Clone Set: fs_sapmnt-clone [fs_sapmnt]:
    * Started: [ cl-sap-1 cl-sap-2 ]

Daemon Status:
corosync: active/disabled
Pacemaker: active/disabled
pcsd: active/enabled
```

Test2 - Testing failure of the node that is running the ASCS instance

This test shows how Pacemaker responds to a node failure when the node is running the ASCS instance.

► **Test2 - Description**

Simulate a failure of the node where the ASCS instance is running.

► **Test2 - Prerequisites**

- A functional two-node RHEL HA Add-On cluster for SAP ENSA2.
- Both cluster nodes are active.
- Cluster is started on NODE1 and NODE2.
 - Resource group `#{sid}_ascs${ASCS_nr}_group` is active on NODE2.
 - Resources `#{sid}_vip_ascs${ASCS_nr}`, `#{sid}_fs_ascs${ASCS_nr}_lvm`, `#{sid}_fs_ascs${ASCS_nr}` and `#{sid}_ascs${ASCS_nr}` are **Started** on NODE2.

- Resource group `#{sid}_ers#{ERS_nr}_group` is active on NODE1.
- Resources `#{sid}_vip_ers#{ERS_nr}`, `#{sid}_fs_ers#{ERS_nr}_lvm`, `#{sid}_fs_ers#{ERS_nr}` and `#{sid}_ers#{ERS_nr}` are **Started** on NODE1.
- Check SAP instance processes:
 - ASCS instance is running on NODE2.
 - ERS instance is running on NODE1.

► **Test2 - Test procedure**

Forcefully terminate NODE2 by sending a **fast-restart** system request as shown in Example 6-44.

Example 6-44 Fast-restart on NODE2

```
sync; echo b > /proc/sysrq-trigger
```

► **Test2 - Expected behavior**

- NODE2 restarts.
- The cluster detects the failed node and sets its state to offline (UNCLEAN).
- The cluster acquires the ASCS resources (virtual IP address, file system `/usr/sap/#{SID}/ASCS#{ASCS_nr}`, and the LVM items) on NODE1.
- The cluster starts the ASCS on NODE1.

Note: The second node is offline and both resource groups are running on the first node.

After a while, check the status with the command in Example 6-45.

Example 6-45 Check Cluster Status.

```
# pcs status
```

Example 6-46 shows the expected output showing all services running on NODE1.

Example 6-46 Output showing cluster status after Node2 failure

```
# pcs status
Cluster name: SAP_ASCS
Cluster Summary:
 * Stack: corosync
 * Current DC: cl-sap-1 (version 2.0.5-9.e18_4.5-ba59be7122) - partition with quorum
 * Last updated: Tue Feb 14 08:34:16 2023
 * Last change: Tue Feb 14 08:34:04 2023 by hacluster via crmd on cl-sap-1
 * 2 nodes configured
 * 11 resource instances configured

Node List:
 * Online: [ cl-sap-1 ]
 * OFFLINE: [ cl-sap-2 ]

Full List of Resources:
 * res_fence_ibm_powervs(stonith:fence_ibm_powervs): Started cl-sap-1
 * Resource Group: s01_ascs01_group:
   * s01_vip_ascs01(ocf::heartbeat:IPaddr2): Started cl-sap-1
   * s01_fs_ascs01_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-1
   * s01_fs_ascs01(ocf::heartbeat:Filesystem): Started cl-sap-1
   * s01_ascs01(ocf::heartbeat:SAPInstance): Started cl-sap-1
 * Resource Group: s01_ers02_group:
   * s01_vip_ers02(ocf::heartbeat:IPaddr2): Started cl-sap-1
   * s01_fs_ers02_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-1
```



```

* s01_fs_ers02(ocf::heartbeat:Filesystem): Started cl-sap-1
* s01_ers02(ocf::heartbeat:SAPInstance): Started cl-sap-1
* Clone Set: fs_sapmnt-clone [fs_sapmnt]:
* Started: [ cl-sap-1 ]
* Stopped: [ cl-sap-2 ]

```

```

Daemon Status:
corosync: active/disabled
Pacemaker: active/disabled
pcsd: active/enabled

```

► Test2 - Recovery procedure

Wait until NODE2 restarts, then restart the cluster framework using the command shown in Example 6-47.

Example 6-47 On Node1 Start Cluster.

```
pcs cluster start
```

The results of the cluster start are:

- The cluster starts on NODE2 and acquires the ERS resources (virtual IP address, file system `/usr/sap/${SID}/ERS${ERS_nr}`, and the LVM resources) on NODE2.
- The cluster starts the ERS instance on NODE2.

Wait a moment and check the status with the command as shown in Example 6-48 which shows that the ERS resource group moved to the second node.

Example 6-48 Cluster status after recovery

```

# pcs status
Cluster name: SAP_ASCS
Cluster Summary:
* Stack: corosync
* Current DC: cl-sap-1 (version 2.0.5-9.e18_4.5-ba59be7122) - partition with quorum
* Last updated: Tue Feb 14 08:41:23 2023
* Last change: Tue Feb 14 08:34:04 2023 by hacluster via crmd on cl-sap-1
* 2 nodes configured
* 11 resource instances configured

```

```

Node List:
* Online: [ cl-sap-1 cl-sap-2 ]

```

```

Full List of Resources:
* res_fence_ibm_powervs(stonith:fence_ibm_powervs): Started cl-sap-1
* Resource Group: s01_ascs01_group:
* s01_vip_ascs01(ocf::heartbeat:IPaddr2): Started cl-sap-1
* s01_fs_ascs01_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-1
* s01_fs_ascs01(ocf::heartbeat:Filesystem): Started cl-sap-1
* s01_ascs01(ocf::heartbeat:SAPInstance): Started cl-sap-1
* Resource Group: s01_ers02_group:
* s01_vip_ers02(ocf::heartbeat:IPaddr2): Started cl-sap-2
* s01_fs_ers02_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-2
* s01_fs_ers02(ocf::heartbeat:Filesystem): Started cl-sap-2
* s01_ers02(ocf::heartbeat:SAPInstance): Started cl-sap-2
* Clone Set: fs_sapmnt-clone [fs_sapmnt]:
* Started: [ cl-sap-1 cl-sap-2 ]

```

```

Daemon Status:
corosync: active/disabled
Pacemaker: active/disabled

```

pcsd: active/enabled

Test3 - Testing failure of the ERS instance

This test will show the results of the ERS service on NODE2.

► **Test3 - Description**

Simulate a failure of the ERS instance while it is running.

► **Test3 - Prerequisites**

- A functional two-node RHEL HA Add-On cluster for SAP ENSA2.
- Both cluster nodes are active.
- Cluster starts on NODE1 and NODE2.
 - Resource group `${sid}_ascs${ASCS_nr}_group` is active on NODE1.
 - Resources `${sid}_vip_ascs${ASCS_nr}`, `${sid}_fs_ascs${ASCS_nr}_lvm`, `${sid}_fs_ascs${ASCS_nr}` and `${sid}_ascs${ASCS_nr}` are **Started** on NODE1.
 - Resource group `${sid}_ers${ERS_nr}_group` is active on NODE2.
 - Resources `${sid}_vip_ers${ERS_nr}`, `${sid}_fs_ers${ERS_nr}_lvm`, `${sid}_fs_ers${ERS_nr}` and `${sid}_ers${ERS_nr}` are **Started** on NODE2.
- Check SAP instance processes:
 - ASCS instance is running on NODE1.
 - ERS instance is running on NODE2.

► **Test3 - Test Procedure**

Forcefully terminate the SAP ERS instance by sending a SIGKILL signal as shown in Example 6-49.

Example 6-49 Identify the PID of the enqueue replication server and send kill command

```
pgrep -af "(er|enqr).sap"
2527198 er.sapS01_ERS02 pf=/usr/sap/S01/ERS02/profile/S01_ERS02_cl-sap-ers NR=01
kill -9 2527198
```

► **Test3 - Expected behavior**

The following actions should occur as a result of shutting down the ERS service.

- SAP Enqueue Replication Server on NODE2 fails immediately.
- The cluster detects the stopped ERS and marks the resource as failed.
- The cluster restarts the ERS on NODE2.

Example 6-50 shows running the `pcs status` command and the expected results. The `${sid}_ers${ERS_nr}` ERS resource restarted on the second node. If you run the `pcs status` command too soon, you might see the ERS resource briefly in status **FAILED**. See Example 6-50 for sample results.

Example 6-50 Cluster status after ERS service failure

```
# pcs status
Cluster name: SAP_ASCS
Cluster Summary:
 * Stack: corosync
 * Current DC: cl-sap-1 (version 2.0.5-9.e18_4.5-ba59be7122) - partition with quorum
 * Last updated: Tue Feb 14 08:50:53 2023
 * Last change: Tue Feb 14 08:50:50 2023 by hacluster via crmd on cl-sap-2
 * 2 nodes configured
 * 11 resource instances configured

Node List:
 * Online: [ cl-sap-1 cl-sap-2 ]
```

Full List of Resources:

```
* res_fence_ibm_powervs(stonith:fence_ibm_powervs): Started cl-sap-1
* Resource Group: s01_ascs01_group:
  * s01_vip_ascs01(ocf::heartbeat:IPaddr2): Started cl-sap-1
  * s01_fs_ascs01_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-1
  * s01_fs_ascs01(ocf::heartbeat:Filesystem): Started cl-sap-1
  * s01_ascs01(ocf::heartbeat:SAPInstance): Started cl-sap-1
* Resource Group: s01_ers02_group:
  * s01_vip_ers02(ocf::heartbeat:IPaddr2): Started cl-sap-2
  * s01_fs_ers02_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-2
  * s01_fs_ers02(ocf::heartbeat:Filesystem): Started cl-sap-2
  * s01_ers02(ocf::heartbeat:SAPInstance): Started cl-sap-2
* Clone Set: fs_sapmnt-clone [fs_sapmnt]:
  * Started: [ cl-sap-1 cl-sap-2 ]
```

Daemon Status:

```
corosync: active/disabled
Pacemaker: active/disabled
pcsd: active/enabled
```

Test3 - Recovery Procedure

Run the commands in Example 6-51 to refresh the cluster.

Example 6-51 On NODE2 Run following commands

```
pcs resource refresh
pcs status --full
```

Test4 - Testing the manual move of the ASCS instance

This test moves the ASCS instance to NODE2 using SAP control interfaces as in a maintenance action.

► Test4 - Description

Use SAP Control commands to move the ASCS instance to the other node for maintenance purposes.

► Test4 - Prerequisites

- A functional two-node RHEL HA Add-On cluster for SAP ENSA2.
- The **sap_cluster_connector** is installed and configured.
- Both cluster nodes are active.
- Cluster is started on NODE1 and NODE2.
 - Resource group `#{sid}_asc#{ASCS_nr}_group` is active on NODE1.
 - Resources `#{sid}_vip_asc#{ASCS_nr}`, `#{sid}_fs_asc#{ASCS_nr}_lvm`, `#{sid}_fs_asc#{ASCS_nr}` and `#{sid}_asc#{ASCS_nr}` are **Started** on NODE1.
 - Resource group `#{sid}_ers#{ERS_nr}_group` is active on NODE2.
 - Resources `#{sid}_vip_ers#{ERS_nr}`, `#{sid}_fs_ers#{ERS_nr}_lvm`, `#{sid}_fs_ers#{ERS_nr}` and `#{sid}_ers#{ERS_nr}` are **Started** on NODE2.
- Check SAP instance processes:
 - ASCS instance is running on NODE1.
 - ERS instance is running on NODE2.

Test4 - Test Procedure

Log in to NODE1 and run **sapcontrol** to move the ASCS instance to the other node. The command is shown in Example 6-52.

Example 6-52 On NODE1 Run following commands

```
sudo -i -u ${sid}adm -- sh -c "sapcontrol -nr ${ASCS_nr} -function HAFailoverToNode"
```

► Test4 - Expected behavior

The **sapcontrol** function interacts with the cluster through the **sap-cluster-connector**. The cluster creates location constraints to move the resource. You can check the status as shown in Example 6-53.

Example 6-53 Status after move of ASCS service

```
# pcs status
Cluster name: SAP_ASCS
Cluster Summary:
 * Stack: corosync
 * Current DC: cl-sap-1 (version 2.0.5-9.e18_4.5-ba59be7122) - partition with quorum
 * Last updated: Tue Feb 14 09:03:19 2023
 * Last change: Tue Feb 14 09:01:40 2023 by s01adm via crm_resource on cl-sap-1
 * 2 nodes configured
 * 11 resource instances configured

Node List:
 * Online: [ cl-sap-1 cl-sap-2 ]

Full List of Resources:
 * res_fence_ibm_powervs(stonith:fence_ibm_powervs): Started cl-sap-1
 * Resource Group: s01_ascs01_group:
   * s01_vip_ascs01(ocf::heartbeat:IPaddr2): Started cl-sap-2
   * s01_fs_ascs01_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-2
   * s01_fs_ascs01(ocf::heartbeat:Filesystem): Started cl-sap-2
   * s01_ascs01(ocf::heartbeat:SAPInstance): Started cl-sap-2
 * Resource Group: s01_ers02_group:
   * s01_vip_ers02(ocf::heartbeat:IPaddr2): Started cl-sap-1
   * s01_fs_ers02_lvm(ocf::heartbeat:LVM-activate): Started cl-sap-1
   * s01_fs_ers02(ocf::heartbeat:Filesystem): Started cl-sap-1
   * s01_ers02(ocf::heartbeat:SAPInstance): Started cl-sap-1
 * Clone Set: fs_sapmnt-clone [fs_sapmnt]:
   * Started: [ cl-sap-1 cl-sap-2 ]

Daemon Status:
 corosync: active/disabled
 Pacemaker: active/disabled
 pcsd: active/enabled
```

Note: Keep in mind that the ASCS resource group moved to the second node. If you run the **pcs status** command too soon, you might see some resources **stopping** and **starting**.

► Test4 - Recovery Procedure

Wait until the ASCS instance is active on NODE2. After five minutes, the cluster removes the created location constraints automatically. Example 6-54 shows the constraints during the failover.

Example 6-54 On NODE2 run the pcs constraint command to show constraints

```
# pcs constraint
Location Constraints:
  Resource: s01_asc01_group
    Constraint: cli-ban-s01_asc01_group-on-cl-sap-1
    Rule: boolean-op=and score=-INFINITY
    Expression: #uname eq string cl-sap-1
    Expression: date lt 2023-02-08 09:33:50 -05:00
Ordering Constraints:
  start s01_asc01_group then stop s01_ers02_group (kind:Optional) (non-symmetrical)
  start fs_sapmnt-clone then start s01_asc01_group (kind:Mandatory)
  start fs_sapmnt-clone then start s01_ers02_group (kind:Mandatory)
Colocation Constraints:
  s01_ers02_group with s01_asc01_group (score:-5000)
Ticket Constraints:
```

Example 6-55 shows how to remove the constraints manually by running the **pcs resource clear** command.

Example 6-55 On NODE2, Clear location Constraints

```
pcs resource clear ${sid}_asc0${ASCS_nr}_group
```

Example 6-56 shows the results of clearing the constraints.

Example 6-56 On NODE2, run the following command.

```
# pcs constraint
Location Constraints:
Ordering Constraints:
  start s01_asc01_group then stop s01_ers02_group (kind:Optional) (non-symmetrical)
  start fs_sapmnt-clone then start s01_asc01_group (kind:Mandatory)
  start fs_sapmnt-clone then start s01_ers02_group (kind:Mandatory)
Colocation Constraints:
  s01_ers02_group with s01_asc01_group (score:-5000)
Ticket Constraints:
```

6.3 Using Pacemaker to create highly available SAP HANA database instances

SAP HANA production databases are generally business critical and it is essential to ensure that the highest possible uptime is maintained for these databases by deploying them in a highly available cluster configuration. Red Hat HA Add-On Solutions for SAP HANA uses Pacemaker to enable highly available cluster configurations in conjunction with HANA Server Replication (HSR).

SAP HANA System Replication (HSR) is a built-in high availability and disaster recovery feature to support business continuity. HANA System Replication is designed to copy and continuously synchronize an SAP HANA database to one or more locations. Data is constantly preloaded on the secondary system to minimize recovery time objective (RTO). While HSR provides the ability

to have an up to date copy of the SAP HANA database, it does not have the capability to automatically detect failures and take actions to failover between primary and secondary database systems.

The Red Hat High Availability Add-On capability complements HANA HSR by allowing automated failovers within such clusters so that the HANA Databases are able to operate within a highly available clusters on top of the industry leading Red Hat Enterprise Linux High Availability cluster framework. The Red Hat Enterprise Linux High Availability Add-On provides all the necessary packages that are required to configure a Pacemaker-based cluster that provides reliability, scalability, and availability to critical production level HANA Database services. Pacemaker clustering capability works in conjunction with SAP HANA HSR to reduce planned and unplanned downtime for your SAP HANA environment.

HSR has several options for setting up the replication between the database copies. These options can support various distances between the database copies, ranging from Metro distances and beyond depending on the network latencies and bandwidth. This enables the ability to have local recovery, metro distance recovery, or geographic level disaster recovery for the HANA databases.

Figure 6-5 is an illustration of a typical HANA Database deployment with HSR capability.

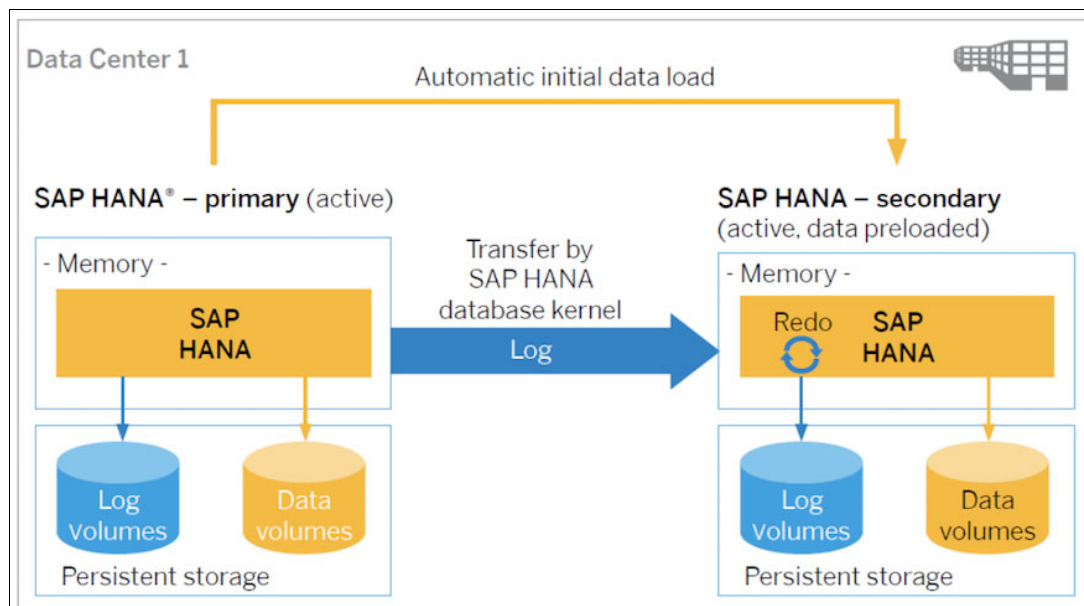


Figure 6-5 SAP HANA System Replication (HSR)⁵

Table 6-3 summarizes a list of all supported high availability scenarios for Automated SAP HANA Scale-Up System Replication options

Table 6-3 Automated SAP HANA Scale-Up System Replication options

Supported Scenario	Notes
6.3.1, "Performance Optimized" on page 111	Secondary site is not active for client/application servers
6.3.2, "Cost Optimized" on page 111	Support a QA/Test instance running on the secondary site (Cost-Optimized); QA/Test instance will be shutdown first during the failover of Prod

⁵ [Red Hat documentation on HA implementation for SAP Hana](#)

Supported Scenario	Notes
6.3.3, “Active/Active (Read Enabled)” on page 112	The secondary HANA instance can take read-only inquiries
6.3.4, “Multi-tier System Replication” on page 113	Multi tier System Replication is possible, but the tertiary site cannot be managed by the cluster
6.3.5, “Multi-target System Replication” on page 113	In addition to the standard HANA System Replication the data is replicated to additional secondary HANA instances that are not managed by the cluster

6.3.1 Performance Optimized

The Performance Optimized scenario, the secondary HANA database is configured to preload the tables into memory, thus the takeover time is normally very fast. However, as the secondary HANA database is dedicated to System Replication and doesn't accept client inquiries, this setup is expensive in terms of hardware cost.

Figure 6-6 illustrates the Performance Optimized deployment with Pacemaker cluster.

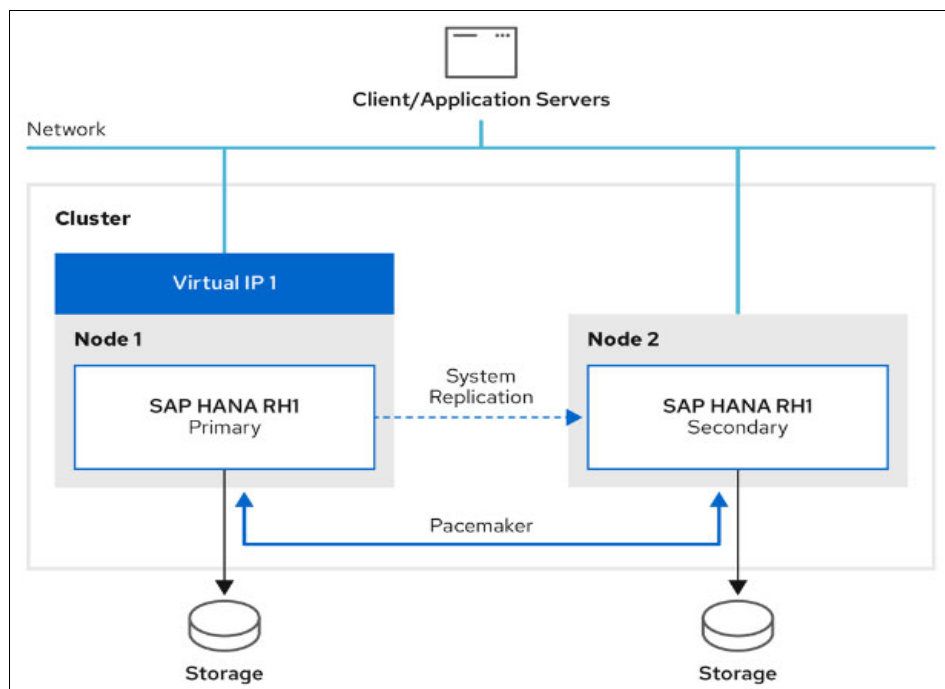


Figure 6-6 Performance Optimized - HSR Pacemaker cluster⁶

The configuration guide for Performance Optimized deployment can be found for IBM Power Virtual Server offering in this [IBM Cloud document](#).

6.3.2 Cost Optimized

The Cost Optimized scenario provides support for additional TEST/QA HANA database on the secondary site, serving client inquiries. As the hardware resources are used and utilized by the TEST/QA instances, the Production HANA database could not be preloaded. The take

⁶ [Red Hat HA Solutions for SAP](#)

over from secondary site must shutdown TEST/QA instances to allow assignment of hardware resources to the production instance so that the secondary HANA instance could be promoted as primary instance. The takeover time is effected and longer in this case for Cost Optimized deployments.

Figure 6-7 illustrates the Cost Optimized deployment with Pacemaker clusters

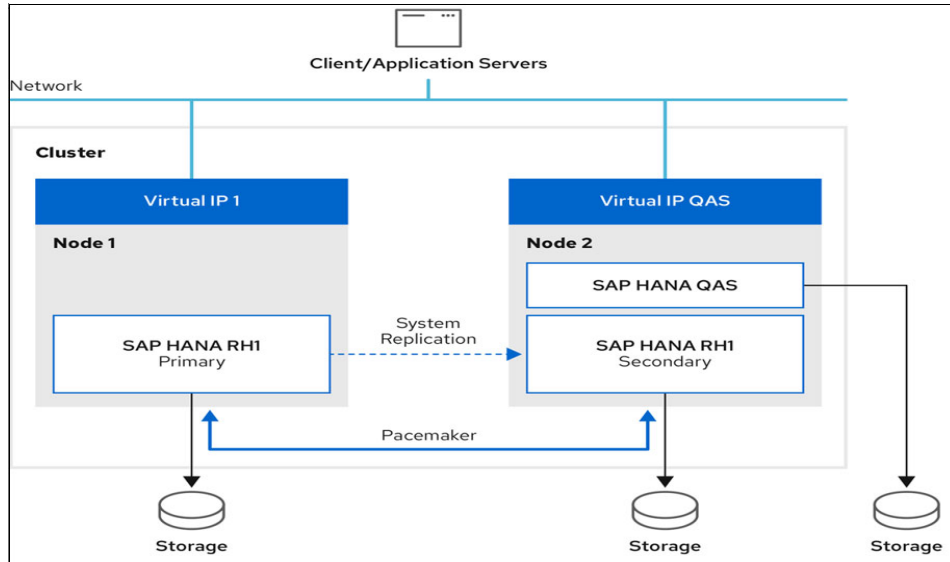


Figure 6-7 Cost Optimized - HSR Pacemaker cluster

The configuration guide for Cost Optimized deployment with Pacemaker can be found for IBM Power Virtual Server offering in this [IBM Cloud Document](#).

6.3.3 Active/Active (Read Enabled)

In this configuration, both primary site as well as secondary site are loaded and run in an active/active cluster where primary site instance operating in read/write mode where the secondary site instance is operating in read-enabled mode where secondary instance is able to service read only queries over a second Virtual IP address. Figure 6-8 illustrates the Active/Active (Read Enabled) deployment with Pacemaker clusters

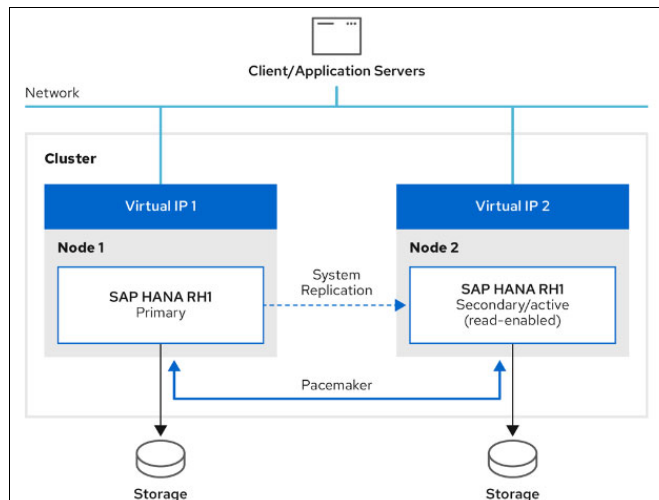


Figure 6-8 Active/Active (read-enabled) Pacemaker cluster⁷

The configuration guide for Active/Active (read-enabled) deployment with Pacemaker can be found for IBM Power Virtual Server offering in this [IBM Cloud document](#).

6.3.4 Multi-tier System Replication

Multi-tier System Replication provides a HANA database to be replicated using HSR from primary site to Secondary and then onto a 3rd site as a chain replication. In this option, any failovers or recovery activities are totally manual for the 3rd site. The Primary and Secondary sites with HANA databases are setup to automatically fail-over. However, the Pacemaker cluster will not have any configuration or the knowledge of the 3rd site. The fail-back of the cluster need to be manually performed with the cluster stopped as there is no knowledge of the 3rd site in the Pacemaker cluster.

Figure 6-9 illustrates the Multi-tier System Replication deployment with Pacemaker clusters.

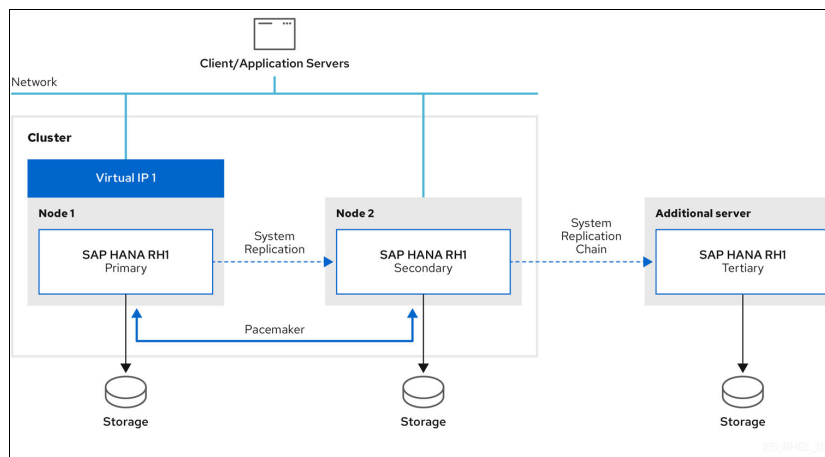


Figure 6-9 Multi-tier System Replication Pacemaker Cluster⁸

The Configuration guide for Multi-tier System Replication deployment with Pacemaker can be found for IBM Power Virtual Server offering in this [IBM Cloud Document](#).

6.3.5 Multi-target System Replication

The Multi-target System Replication HA cluster setup in Scale up deployment of HANA 2.0 databases, the primary HANA instance is replicated to a secondary HANA instance managed by the HA cluster and also, primary instance is asynchronously to an additional secondary HANA instances that is not managed by the cluster to meet additional HA-DR capabilities.

Figure 6-10 on page 114 illustrates the Multi-target System Replication deployment with Pacemaker clusters

⁷ Red Hat HA Solutions for SAP

⁸ Red Hat HA Solutions for SAP

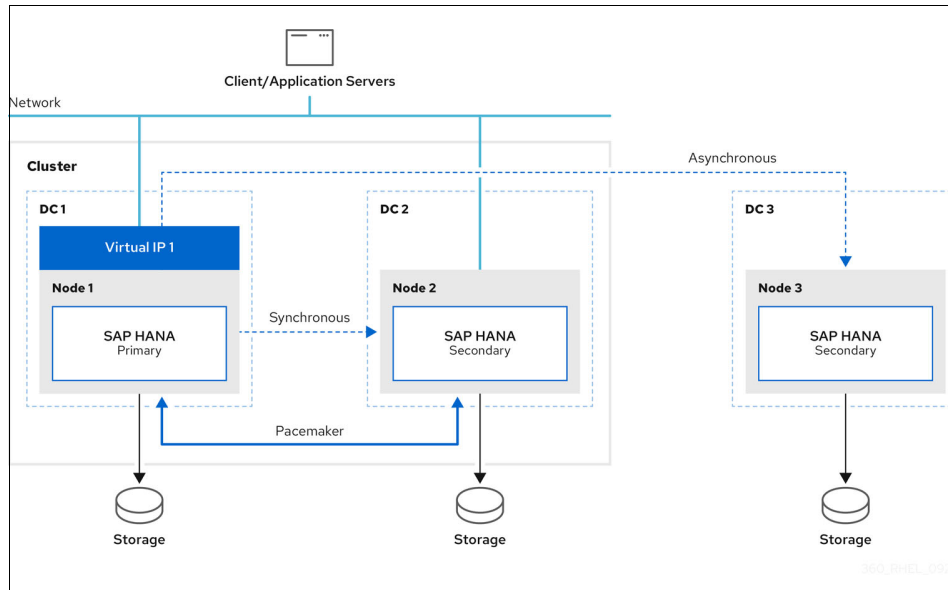


Figure 6-10 Multi-target System Replication Pacemaker cluster⁹

Configuration guide for Multi-target System Replication deployment with Pacemaker can be found for on-premise deployments in this [Red Hat Document](#).

Note: This section is not meant to replace any SAP or Red Hat Documentation. You should always refer to respective SAP or Red Hat Documentation for latest and current up to date information.

6.3.6 Configuring SAP HANA Scale-Up System Replication in a RHEL HA Add-On cluster

Overview

This section describes the configuration of a Red Hat Enterprise Linux 8 (RHEL) HA Add-On cluster for managing SAP HANA Scale-Up System Replication. The cluster uses virtual server instances within IBM Power Virtual Server offering in IBM Cloud as cluster nodes.

Prerequisites

1. A Red Hat High Availability cluster deployed on two virtual server instances in Power Systems Virtual Server.
 - a. Install and set up the RHEL HA Add-On cluster
 - b. Configure and verify fencing as described in the preceding sections
2. The virtual server instances used from IBM Power Virtual Server offering need to fulfill hardware and resource requirements for the SAP HANA systems within its scope.
3. The hostnames used for the virtual server instances must meet the SAP HANA requirement.
4. SAP HANA installed on both virtual server instances and SAP HANA System Replication is configured following standard SAP HANA installation procedures for IBM Power Systems. The following references describe this in detail:
 - a. Install SAP HANA on RHEL:
 - SAP Note 2772999 - Red Hat Enterprise Linux 8.x: Installation and Configuration

⁹ [Red Hat HA Solutions for SAP](#)

- SAP Note 2777782 - SAP HANA DB: Recommended OS Settings for RHEL 8
 - b. SAP HANA Server Installation:
 - SAP HANA Server Installation and Update Guide
 - c. SAP HANA System Replication:
 - SAP HANA System Replication Guide
 - d. Configure SAP HANA System Replication:
 - Automating SAP HANA Scale-Up System Replication by using the RHEL HA Add-On
 - SAP HANA System Replication
5. A valid RHEL for SAP Applications or RHEL for SAP Solutions subscription is required to enable the repositories that you need to install SAP HANA and the resource agents for HA configurations.

For more details refer to [Automating SAP HANA Scale-Up System Replication by using the RHEL HA Add-On](#).

Environment variables

To simplify the setup and configuration, on both nodes, as root, set the environment variables shown in Example 6-57 as they are used throughout the following example.

Example 6-57 Environment Variables

```
export SID=<SID>           # SAP HANA System ID (uppercase)
export sid=<sid>           # SAP HANA System ID (lowercase)
export INSTNO=<INSTNO>     # SAP HANA Instance Number

export DC1=<Site1>         # HANA System Replication Site Name 1
export DC2=<Site2>         # HANA System Replication Site Name 2

export NODE1=<Hostname 1> # Hostname of virtual server instance 1
export NODE2=<Hostname 2> # Hostname of virtual server instance 2
```

Installing SAP HANA Resource Agents

Run the command shown in Example 6-58 to install the RHEL HA Add-On resource agents for SAP HANA System Replication.

Example 6-58 Install SAP HANA Resource Agents

```
dnf install -y resource-agents-sap-hana
```

Starting SAP HANA System

Start SAP HANA and verify that HANA System Replication is active. On both nodes, run the commands shown in Example 6-59. For more information, refer to [“Checking SAP HANA System Replication state”](#).

Example 6-59 Start HANA and Check Status

```
sudo -i -u ${sid}adm -- HDB start

sudo -i -u ${sid}adm -- <<EOT
    hdbnsutil -sr_state
    HDBSettings.sh systemReplicationStatus.py
EOT
```

Enabling the SAP HANA srConnectionChanged() hook

Recent versions of SAP HANA provide hooks so SAP HANA can send out notifications for certain events. For more information, refer to [“Implementing a HA/DR Provider”](#)

The srConnectionChanged() hook improves the ability of the cluster to detect a status change of HANA System Replication that requires an action from the cluster. The goal is to prevent data loss and corruption by preventing accidental takeovers.

Activating the srConnectionChanged() hook on all SAP HANA instances

First, stop the cluster as shown in Example 6-60.

Example 6-60 Stop the Cluster

```
# On NODE1, run the following command.
pcs cluster stop --allEOT
```

Now, install the hook script. Run the commands shown in Example 6-61 on each HANA node in the cluster.

Example 6-61 Install the hook script for each HANA instance, and set the required ownership.

```
# On Both Nodes, run the following commands.

mkdir -p /hana/shared/myHooks
cp /usr/share/SAPHanaSR/srHook/SAPHanaSR.py /hana/shared/myHooks
chown -R ${sid}adm:sapsys /hana/shared/myHooks
```

Update the global.ini file on each HANA node as seen in Example 6-62. This enables the hook script.

*Example 6-62 Update the **global.ini** file on each HANA node to enable the hook script.*

```
# On Both Nodes, run the following command.

sudo -i -u ${sid}adm -- <<EOT
  python \${DIR_INSTANCE}/exe/python_support/setParameter.py \
    -set SYSTEM/global.ini/ha_dr_provider_SAPHanaSR/provider=SAPHanaSR \
    -set SYSTEM/global.ini/ha_dr_provider_SAPHanaSR/path=/hana/shared/myHooks \
    -set SYSTEM/global.ini/ha_dr_provider_SAPHanaSR/execution_order=1 \
    -set SYSTEM/global.ini/trace/ha_dr_saphanasr=info
EOT
```

Validate the changes to the global.ini file using the command in Example 6-63.

*Example 6-63 Verify the changed **global.ini** file.*

```
# On Both Nodes, run the following commands.
cat /hana/shared/${SID}/global/hdb/custom/config/global.ini
```

Create sudo settings for SAP HANA OS user which is required to allow the right permissions for the user during configuration changes. See Example 6-64.

Example 6-64 Create sudo settings for SAP HANA OS user

```
# You need the following sudo settings to allow the ${sid}adm user
# Script can update the node attributes when the srConnectionChanged()
# hook runs.
```

```

#
# Create a file with the required sudo aliases and user specifications.
# On Both Nodes, run the following commands.

cat >> /etc/sudoers.d/20-saphana << EOT
Cmdnd_Alias DC1_SOK = /usr/sbin/crm_attribute -n hana_${sid}_site_srHook_${DC1} -v SOK -t
crm_config -s SAPHanaSR
Cmdnd_Alias DC1_SFMAIL = /usr/sbin/crm_attribute -n hana_${sid}_site_srHook_${DC1} -v SFMAIL
-t crm_config -s SAPHanaSR
Cmdnd_Alias DC2_SOK = /usr/sbin/crm_attribute -n hana_${sid}_site_srHook_${DC2} -v SOK -t
crm_config -s SAPHanaSR
Cmdnd_Alias DC2_SFMAIL = /usr/sbin/crm_attribute -n hana_${sid}_site_srHook_${DC2} -v SFMAIL
-t crm_config -s SAPHanaSR
${sid}adm ALL=(ALL) NOPASSWD: DC1_SOK, DC1_SFMAIL, DC2_SOK, DC2_SFMAIL
Defaults!DC1_SOK, DC1_SFMAIL, DC2_SOK, DC2_SFMAIL !requiretty
EOT

# Adjust the permissions and check for syntax errors.

chown root:root /etc/sudoers.d/20-saphana
chmod 0440 /etc/sudoers.d/20-saphana
cat /etc/sudoers.d/20-saphana
visudo -c

```

Important: Any problems that are reported by the `visudo -c` command must be corrected.

Verify that the hook functions.

To verify the configured hook works,

- Restart both HANA instances and verify that the hook script works as expected.
- Perform an action to trigger the hook, such as stopping a HANA instance.
- Check whether the hook logged anything in the trace files.

These actions are shown in Example 6-65.

Example 6-65 Restart and then stop the HANA instance.

```

# On Both Nodes, run the following commands.

sudo -i -u ${sid}adm -- HDB restart
sudo -i -u ${sid}adm -- HDB stop

# Check messages in trace files.

sudo -i -u ${sid}adm -- sh -c 'grep "ha_dr_SAPHanaSR.*crm_attribute"
$DIR_INSTANCE/$VTHOSTNAME/trace/nameserver_* | cut -d" " -f2,3,5,17'

# Start the HANA instance.

sudo -i -u ${sid}adm -- HDB start

```

Note: After you verify that the hooks function, you can restart the HA cluster.

Start the cluster using the commands shown in Example 6-66.

Example 6-66 Start cluster.

```
# On NODE1, run the following commands.
pcs cluster start --all

# Check the status of the cluster.
pcs status --full
```

Configuring general cluster properties

To avoid failovers of the resources during initial testing and post production, set the following default values (shown in Example 6-67) for the resource-stickiness and migration-threshold parameters. Keep in mind that defaults don't apply to resources that override them with their own defined values.

Example 6-67 Start cluster.

```
# On NODE1, run the following commands.
pcs resource defaults update resource-stickiness=1000
pcs resource defaults update migration-threshold=5000
```

Creating a cloned SAPHanaTopology resource

The SAPHanaTopology resource gathers status and configuration of SAP HANA System Replication on each node. It also starts and monitors the local SAP Host Agent, which is required for starting, stopping, and monitoring SAP HANA instances. Create this resource with the commands shown in Example 6-68.

Example 6-68 Create the SAPHanaTopology resource.

```
# On NODE1, run the following commands.

pcs resource create SAPHanaTopology_${SID}_${INSTNO} SAPHanaTopology \
  SID=${SID} InstanceNumber=${INSTNO} \
  op start timeout=600 \
  op stop timeout=300 \
  op monitor interval=10 timeout=600 \
  clone clone-max=2 clone-node-max=1 interleave=true

# Check the configuration and the cluster status by running the
# following commands.

pcs resource config SAPHanaTopology_${SID}_${INSTNO}
pcs resource config SAPHanaTopology_${SID}_${INSTNO}-clone
pcs status --full
```

Creating master and slave SAPHana resources

The SAPHana resource manages two SAP HANA instances that are configured as HANA System Replication nodes. Create these resources as shown in Example 6-69.

Example 6-69 Create the SAPHana resource

```
# On NODE1, run the following commands.
pcs resource create SAPHana_${SID}_${INSTNO} SAPHana \
  SID=${SID} InstanceNumber=${INSTNO} \
  PREFER_SITE_TAKEOVER=true \
```

```

DUPLICATE_PRIMARY_TIMEOUT=7200 \
AUTOMATED_REGISTER=false \
op start timeout=3600 \
op stop timeout=3600 \
op monitor interval=61 role="Slave" timeout=700 \
op monitor interval=59 role="Master" timeout=700 \
op promote timeout=3600 \
op demote timeout=3600 \
promotable notify=true clone-max=2 clone-node-max=1 interleave=true

# Check the configuration and the cluster status.
pcs resource config SAPHana_${SID}_${INSTNO}
pcs status --full

```

Creating a virtual IP address resource

Review the information in Reserving virtual IP addresses and reserve a virtual IP address for the SAP HANA System Replication cluster.

Use the reserved IP address to create a virtual IP address resource. This virtual IP address is used to reach the System Replication primary instance.

Assign the reserved IP address to a VIP environment variable and create the virtual IP address cluster resource as shown in Example 6-70.

Example 6-70 Create VIP

```

# On NODE1, run the following commands.
export VIP=<reserved IP address>
echo $VIP
pcs resource create vip_${SID}_${INSTNO} IPAddr2 ip=$VIP

# Check the configured virtual IP address and the cluster status.
pcs resource config vip_${SID}_${INSTNO}
ip addr show

pcs status --full

```

Creating constraints

Make sure that “**SAPHanaTopology**” resources are started before you start the “**SAPHana**” resources.

The virtual IP address must be present on the node where the primary resource of “SAPHana” is running. This constraint mandates the start order of these resources. Example 6-71 shows the order constraint being set.

Example 6-71 Create the SAPHanaTopology order constraint.

```

# On NODE1, run the following commands.
pcs constraint order SAPHanaTopology_${SID}_${INSTNO}-clone \
  then SAPHana_${SID}_${INSTNO}-clone symmetrical=false

# Check the configuration.
pcs constraintip addr show

```

Example 6-72 shows setting the collocation constraint.

Example 6-72 Create constraint to collocate the virtual IP address with primary.

```
# On NODE1, run the following commands.
pcs constraint collocation add vip_${SID}_${INSTNO} \
  with master SAPHana_${SID}_${INSTNO}-clone 2000

# Check the configuration and the cluster status.
pcs constraint
pcs status --full
```

Note: This constraint collocates the virtual IP address resource with the SAPHana resource that was promoted as primary.

Enabling automated registration of secondary instance

You need to set the parameter `AUTOMATED_REGISTER` according to your operational requirements. If you want to keep the ability to revert to the state of the previous primary SAP HANA instance, then `AUTOMATED_REGISTER=false` avoids an automatic registration of the previous primary as a new secondary.

If you experience an issue with the data after a takeover that was triggered by the cluster, you can manually revert if `AUTOMATED_REGISTER` is set to `false`.

If `AUTOMATED_REGISTER` is set to `true`, the previous primary SAP HANA instance automatically registers as secondary, and cannot be activated on its previous history. The advantage of `AUTOMATED_REGISTER=true` is that high-availability capability is automatically reestablished after the failed node reappears in the cluster.

For now, it is recommended to keep `AUTOMATED_REGISTER` set to the default value – `false` – until the cluster is fully tested and that you verify that the failover scenarios work as expected.

Tip: The `pcs resource update` command is used to modify resource attributes and `pcs resource update SAPHana_${SID}_${INSTNO} AUTOMATED_REGISTER=true` sets the attribute to `true`.

Testing the SAP HANA System Replication cluster

It is vital to thoroughly test the cluster configuration to make sure that the cluster is working correctly. The following information provides a few sample failover test scenarios, but is not a complete list of test scenarios.

For example, the description of each test case includes the following information.

- Component that is being tested
- Description of the test
- Prerequisites and the cluster state before you start the failover test
- Test procedure
- Expected behavior and results
- Recovery procedure

Test1 - Testing failure of the primary database instance

Use the proceeding information to test the failure of the primary database instance.

► Test1 - Description

This test will simulate a failure of the primary HANA database instance that is running on NODE1.

► **Test1 - Prerequisites**

- A functional two-node RHEL HA Add-On cluster for HANA system replication.
- Both cluster nodes are active.
- Cluster that is started on NODE1 and NODE2.
- Cluster Resource SAPHana_\${SID}_\${INSTNO} that is configured with AUTOMATED_REGISTER=false.
- Check SAP HANA System Replication status:
 - Primary SAP HANA database is running on NODE1
 - Secondary SAP HANA database is running on NODE2
 - HANA System Replication is activated and in sync

► **Test1 - Test procedure**

Forcefully terminate the SAP HANA primary instance by sending a SIGKILL signal as user \${sid}adm. This is shown in Example 6-73.

Example 6-73 Kill the primary HANA database

```
# On NODE1, run the following commands.
sudo -i -u ${sid}adm -- HDB kill-9
```

► **Test1 - Expected behavior**

- SAP HANA primary instance on NODE1 fails.
- The cluster detects the stopped primary HANA database and marks the resource as failed.
- The cluster promotes the secondary HANA database on NODE2 to take over as new primary.
- The cluster releases the virtual IP address on NODE1, and acquires it on the new primary on NODE2.
- If an application, such as SAP NetWeaver, is connected to a tenant database of SAP HANA, the application automatically reconnects to the new primary.

► **Test1 - Recovery procedure**

As cluster resource SAPHana_\${SID}_\${INSTNO} is configured with the setting AUTOMATED_REGISTER=false, the cluster doesn't restart the failed HANA database, and doesn't register it against the new primary. Which means that the status on the new primary (NODE2) also shows the secondary in status 'CONNECTION TIMEOUT'. To recover register the previous primary database as a secondary as shown in Example 6-74.

Example 6-74 Re-register the previous primary as new secondary

```
# On NODE1, run the following commands.
```

```
sudo -i -u ${sid}adm -- <<EOT
  hdbnsutil -sr_register \
  --name=${DC1} \
  --remoteHost=${NODE2} \
  --remoteInstance=00 \
  --replicationMode=sync \
  --operationMode=logreplay
EOT
#
# Verify the system replication status:
#
```

```

sudo -i -u ${sid}adm -- <<EOT
  hdbnsutil -sr_state
  HDBSettings.sh systemReplicationStatus.py
EOT
#
# On NODE1, run the following commands.
#
pcs resource refresh SAPHana_${SID}_${INSTNO}
pcs status --full

```

Note: After the manual register and resource refreshes, the new secondary instance restarts and shows up in status synced (SOK).

Test2 - Testing failure of the node that is running the primary database

Use the following information to test the failure of the node that is running the primary database.

► Test2 - Description

This test will simulate a failure of the node that is running the primary HANA database.

► Test2 - Preparation

Make sure that Cluster Resource SAPHana_\${SID}_\${INSTNO} is configured with AUTOMATED_REGISTER=true. This is shown in Example 6-75.

Example 6-75 On NODE1, run the following command.

```

pcs resource update SAPHana_${SID}_${INSTNO} AUTOMATED_REGISTER=true
pcs resource config SAPHana_${SID}_${INSTNO}

```

► Test2 - Prerequisites

- A functional two-node RHEL HA Add-On cluster for HANA system replication.
- Both nodes active.
- Cluster is started on NODE1 and NODE2.
- Check SAP HANA System Replication status.
- Primary SAP HANA database is running on NODE2
- Secondary SAP HANA database is running on NODE1
- HANA System Replication is activated and in sync

► Test2 - Test procedure

Forcefully terminate the primary on NODE2 by sending a shutoff system request as shown in Example 6-76.

Example 6-76 On NODE2, run the following command.

```

sync; echo o > /proc/sysrq-trigger

```

► Test2 - Expected behavior

- NODE2 shuts down.
- The cluster detects the failed node and sets its state to OFFLINE.
- The cluster promotes the secondary HANA database on NODE1 to take over as new primary.
- The cluster acquires the virtual IP address on NODE1 on the new primary.

- If an application, such as SAP NetWeaver, is connected to a tenant database of SAP HANA, the application automatically reconnects to the new primary.

► **Test2 - Recovery procedure**

Log in to the IBM Cloud Console and start the NODE2 instance. Wait until NODE2 is available again, then restart the cluster framework. This is shown in Example 6-77.

Example 6-77 On NODE2, run the following commands.

```
pcs cluster start
pcs status --full
```

Note: As cluster resource SAPHana_\${SID}_\${INSTNO} is configured with AUTOMATED_REGISTER=true, SAP HANA restarts when NODE2 rejoins the cluster and the former primary re-registers as a secondary.

Test3 - Testing the failure of the secondary database instance

Use the following information to test the failure of the secondary database instance.

► **Test3 - Description**

Simulate a failure of the secondary HANA database.

► **Test3 - Prerequisites**

- A functional two-node RHEL HA Add-On cluster for HANA system replication.
- Both nodes are active.
- Cluster is started on NODE1 and NODE2.
- Cluster Resource SAPHana_\${SID}_\${INSTNO} is configured with AUTOMATED_REGISTER=true.
- Check SAP HANA System Replication status:
 - Primary SAP HANA database is running on NODE1
 - Secondary SAP HANA database is running on NODE2
 - HANA System Replication is activated and in sync

► **Test3 - Test Procedure**

Forcefully terminate the SAP HANA secondary instance by sending a SIGKILL signal as user \${sid}adm. On NODE2, run the command shown in Example 6-78.

Example 6-78 Kill secondary HANA database

```
sudo -i -u ${sid}adm -- HDB kill-9
```

► **Test3 - Expected behavior**

- SAP HANA secondary on NODE2 stops.
- The cluster detects the stopped secondary HANA database and marks the resource as failed.
- The cluster restarts the secondary HANA database.
- The cluster detects that the system replication is in sync again.

► **Test3 - Recovery procedure**

Wait until the secondary HANA instance starts and syncs again (SOK), then cleanup the failed resource actions as shown in `pcs status`. This is shown in Example 6-79.

Example 6-79 On NODE2, run the following command.

```
pcs resource refresh SAPHana_${SID}_${INSTNO}
pcs status --full
```

Test4 - Testing the manual move of a SAPHana resource to another node

This test will manually move a SAPHana resource to another node.

► Test4 - Description

Use cluster commands to move the primary instance to the other node for maintenance purposes.

► Test4 - Prerequisites

- A functional two-node RHEL HA Add-On cluster for HANA system replication.
- Both nodes are active.
- Cluster is started on NODE1 and NODE2.
- Cluster Resource SAPHana_\${SID}_\${INSTNO} is configured with AUTOMATED_REGISTER=true.
- Check SAP HANA System Replication status:
 - Primary SAP HANA database is running on NODE1
 - Secondary SAP HANA database is running on NODE2
 - HANA System Replication is activated and in sync

Test4 - Test procedure

Move SAP HANA primary to other node by using the `pcs resource move` command as shown in Example 6-80.

Example 6-80 Move SAP primary

```
pcs resource move SAPHana_${SID}_${INSTNO}-clone
```

► Test4 - Expected behavior

- The cluster creates location constraints to move the resource.
- The cluster triggers a takeover to the secondary HANA database.
- If an application, such as SAP NetWeaver, is connected to a tenant database of SAP HANA, the application automatically reconnects to the new primary.

► Test4 - Recovery procedure

The automatically created location constraints must be removed to allow automatic failover in the future. Wait until the primary HANA instance is active and remove the constraints. When the restraints are removed, the cluster registers and starts the HANA database as a new secondary instance. This is shown in Example 6-81.

Example 6-81 Remove constraints

```
pcs constraint
pcs resource clear SAPHana_${SID}_${INSTNO}-clone
pcs constraint
pcs status --full
```



Additional Use Cases for Pacemaker High Availability

Pacemaker stands as a cornerstone in various scenarios ensuring system resilience and continuity. This chapter delves into additional use cases where Pacemaker plays a pivotal role in maintaining high availability (HA) within diverse environments. It's important to note that while these use cases have undergone testing and scrutiny within this publication, caution and meticulous attention are advised when implementing these solutions in production environments. Each scenario detailed here showcases the potential of Pacemaker but warrants careful consideration and validation to ensure seamless deployment and operation in live, production-grade settings.

This chapter describes the following topics:

- ▶ 7.1, “IBM MQ Replicated Data Queue Manager” on page 126
- ▶ 7.2, “Using VM Recovery Manager to Replicate Pacemaker cluster” on page 132
- ▶ 7.3, “NSF Server High Availability” on page 138
- ▶ 7.4, “IBM MQ active/passive architecture using Pacemaker and NFS on IBM Power Virtual Server” on page 139
- ▶ 7.5, “IBM PowerVC three-node architecture using Pacemaker” on page 140

7.1 IBM MQ Replicated Data Queue Manager

An Replicated Data Queue Manager (RDQM)¹ setup consists of a trio of servers meticulously configured into a High Availability (HA) group, each hosting an instance of the queue manager. Among these instances, one serves as the active queue manager, diligently replicating its data to the other two instances in real-time. In the unfortunate event of the server hosting the active queue manager encountering an issue, another instance of the queue manager swiftly takes over, ensuring seamless operation with up-to-date data.

Furthermore, the three queue manager instances can opt to share a floating IP address, simplifying client configurations by requiring just a single IP address. Importantly, even if the HA group experiences network disruptions leading to a partition, only one instance of the queue manager can run concurrently. The server responsible for running the queue manager is referred to as the 'primary,' while the other two servers are known as 'secondaries'.

To avert the possibility of a split-brain scenario, this configuration employs three nodes. In a two-node High Availability system, the risk of split-brain arises when the connection between the two nodes is severed. In such a situation, both nodes could potentially run the queue manager simultaneously, resulting in the accumulation of disparate data sets. When connectivity is eventually restored, these conflicting data versions create a 'split-brain,' necessitating manual intervention to determine which dataset to retain and which to discard.

IMPORTANT:

Be Aware that the RDQM (Replicated Data Queue Manager) component is supported solely on Linux for x86-64 (64 bit) on RHEL 7.3 or later versions.

However, for the purpose of demonstrating Pacemaker functionality, we are configuring MQ RDQM on Red Hat Linux running on IBM Power. It's crucial to emphasize that this setup is intended solely for demonstration purposes and does not imply any support for the configuration.

For detailed and updated information, we strongly recommend referring to the latest IBM MQ documentation.

Distributed Replicated Block Device

A Distributed Replicated Block Device, often abbreviated as DRBD², is a sophisticated storage technology designed to ensure data redundancy, high availability, and fault tolerance within a distributed computing environment.

At its core, DRBD operates by mirroring data blocks across multiple storage devices or nodes in a network. Each node typically hosts a DRBD resource, and these resources are kept in sync through real-time replication. This replication can occur synchronously or asynchronously, depending on the specific configuration and requirements of the system:

- ▶ DRBD ensures that data is duplicated across multiple nodes, providing a safety net in case of hardware failures or data corruption. This redundancy minimizes the risk of data loss.
- ▶ DRBD contributes to high availability solutions by enabling failover capabilities. In the event of a node failure, another node can seamlessly take over, ensuring uninterrupted access to the stored data.

¹ <https://www.ibm.com/docs/en/ibm-mq/9.3?topic=configurations-rdqm-high-availability>

² https://linbit.com/drbd-user-guide/drbd-guide-9_0-en/

- ▶ DRBD is designed to provide good performance by optimizing data transfer and replication mechanisms. It strives to minimize latency and maximize throughput, making it suitable for various demanding applications.
- ▶ DRBD can be configured in various ways to suit specific use cases. Administrators can choose between synchronous and asynchronous replication modes, adjust replication policies, and configure storage layouts to meet their performance and availability requirements.
- ▶ It is commonly used with Linux-based systems and integrates well with existing Linux storage management tools, making it a popular choice for organizations that rely on Linux servers.
- ▶ DRBD can be used in clusters with multiple nodes, allowing organizations to scale their storage infrastructure to accommodate growing data needs.
- ▶ DRBD includes features to ensure data consistency and integrity, such as checksums and error handling mechanisms.

Note: DRBD 9 undergoes testing for compatibility on various CPU architectures, including:

- amd64
- arm64
- ppc64le
- s390x

Please note, recent iterations of DRBD 9 are exclusively verified for construction on 64-bit CPU architectures.

Note that DRBD functionality is integrated into the MQ RDQM architecture. Figure 7-1 depicts a standard configuration where there is an RDQM operating on each of the three nodes within the HA group. In the normal mode, one of the nodes acts as primary and the other two nodes act as secondary nodes.

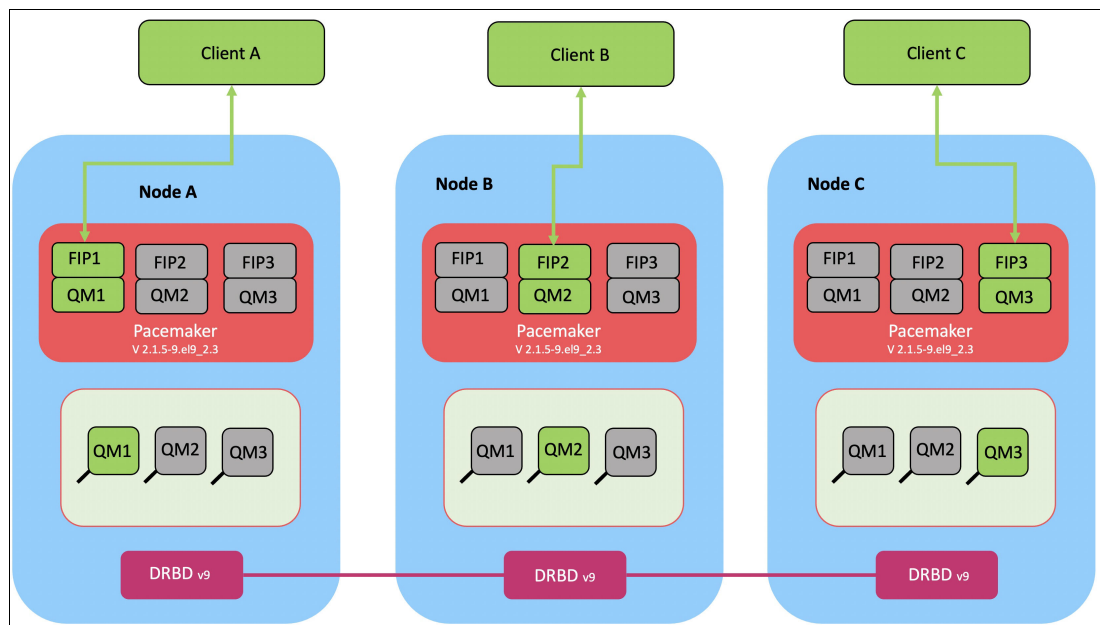


Figure 7-1 Normal RDQM configuration

In Figure 7-2 we show where Node C experiences a failure, resulting in the loss of Pacemaker links, and consequently, queue manager QM3 is now running on Node B.

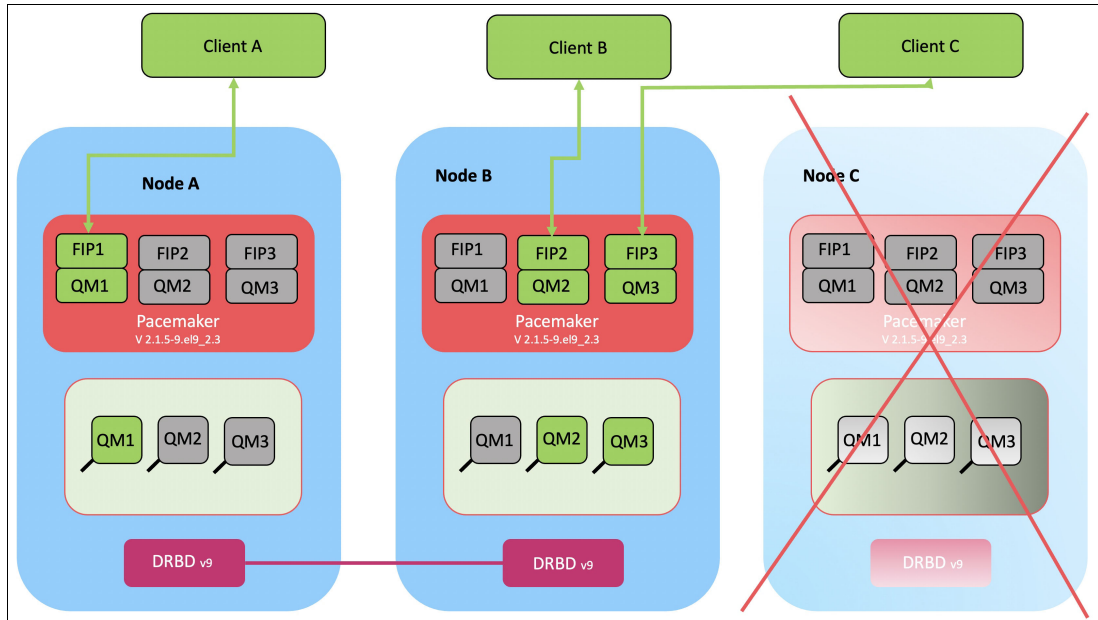


Figure 7-2 Node C failed

Configuring IBM MQ RDQM

To configure IBM MQ RDQM follow these instructions:

- Verify Red Hat Enterprise Linux version using the command shown in Example 7-1.

Example 7-1 Verifying OS version

```
$ hostnamectl
Static hostname: p1325-pvm1
Icon name: computer-vm
Chassis: vm ?
Machine ID: 500ede37aabe4bc4bf54cd41745ec3fa
Boot ID: 50dd26cdedc5471e923472a91ca8666a
Virtualization: povervm
Operating System: Red Hat Enterprise Linux 9.2 (Plow)
CPE OS Name: cpe:/o:redhat:enterprise_linux:9::baseos
Kernel: Linux 5.14.0-284.11.1.el9_2.ppc64le
Architecture: ppc64-le
```

- List available repositories using the command shown in Example 7-2.

Example 7-2 List available repositories

```
$ subscription-manager repos --list
You are attempting to run "subscription-manager" which requires administrative
privileges, but more information is needed in order to do so.
Authenticating as "root"
Password:
+-----+
Available Repositories in /etc/yum.repos.d/redhat.repo
+-----+
Repo ID:   rhel-9-for-ppc64le-highavailability-rpms
Repo Name: Red Hat Enterprise Linux 9 for Power, little endian - High Availability (RPMs)
```



```
Repo URL:
https://rhns.pbm.ihost.com/pulp/content/Default_Organization/Library/content/dist/rhel9/$re
le
           asever/ppc64le/highavailability/os
Enabled: 0
```

```
Repo ID:  codeready-builder-for-rhel-9-ppc64le-rpms
Repo Name: Red Hat CodeReady Linux Builder for RHEL 9 Power, little endian (RPMs)
Repo URL:
https://rhns.pbm.ihost.com/pulp/content/Default_Organization/Library/content/dist/rhel9/$re
le
           asever/ppc64le/codeready-builder/os
Enabled: 1
```

```
Repo ID:  rhel-9-for-ppc64le-baseos-rpms
Repo Name: Red Hat Enterprise Linux 9 for Power, little endian - BaseOS (RPMs)
Repo URL:
https://rhns.pbm.ihost.com/pulp/content/Default_Organization/Library/content/dist/rhel9/$re
le
           asever/ppc64le/baseos/os
Enabled: 1
```

```
Repo ID:  rhel-9-for-ppc64le-sap-solutions-rpms
Repo Name: Red Hat Enterprise Linux 9 for Power, little endian - SAP Solutions (RPMs)
Repo URL:
https://rhns.pbm.ihost.com/pulp/content/Default_Organization/Library/content/dist/rhel9/$re
le
           asever/ppc64le/sap-solutions/os
Enabled: 0
```

```
Repo ID:  rhel-9-for-ppc64le-appstream-rpms
Repo Name: Red Hat Enterprise Linux 9 for Power, little endian - AppStream (RPMs)
Repo URL:
https://rhns.pbm.ihost.com/pulp/content/Default_Organization/Library/content/dist/rhel9/$re
le
           asever/ppc64le/appstream/os
Enabled: 1
```

```
Repo ID:  rhel-9-for-ppc64le-supplementary-rpms
Repo Name: Red Hat Enterprise Linux 9 for Power, little endian - Supplementary (RPMs)
Repo URL:
https://rhns.pbm.ihost.com/pulp/content/Default_Organization/Library/content/dist/rhel9/$re
le
           asever/ppc64le/supplementary/os
Enabled: 1
```

-
- ▶ Enable Red Hat High Availability repository using the command shown in Example 7-3.

Example 7-3 Enable Red Hat HA Repository

```
$ sudo subscription-manager repos --enable=rhel-9-for-ppc64le-highavailability-rpms
Repository 'rhel-9-for-ppc64le-highavailability-rpms' is enabled for this system.
```

- ▶ Download and install DRDB9 Binary Packages as described in the following link:
https://linbit.com/drbd-user-guide/drbd-guide-9_0-en/#s-linbit-packages

Before proceeding with the setup of the RDQM High Availability (HA) group, specific configurations are necessary on each of the three servers designated for inclusion in the HA group. On each node, the following prerequisites must be met:

- ▶ A volume group named "drbdpool" should be created.

- ▶ Each node should be equipped with up to three interfaces, each serving distinct purposes within the RDQM setup:
 - A primary interface dedicated to Pacemaker for monitoring the HA group.
 - An alternate interface intended for Pacemaker to monitor the HA group.
 - A designated replication interface used for synchronous data replication.
- ▶ Additionally, if a firewall exists between the nodes participating in the HA group, it is imperative to ensure that the firewall allows traffic to flow seamlessly between the nodes across a specified range of ports.
- ▶ In cases where the system employs SELinux in Enforcing mode, it is essential to execute the command shown in Example 7-4 for proper configuration.

Example 7-4 semanage command

```
$ semanage permissive -a drbd_t
```

Note: Should you decide to deploy the nodes across different data centers, it's crucial to take into consideration the following constraints:

- Performance experiences a rapid deterioration as latency between data centers increases. While IBM provides support for latency of up to 5 ms, you may discover that your application's performance becomes intolerant when latency exceeds 1 to 2 ms.
- Data transmitted across the replication link does not undergo any additional encryption beyond what may already be in place through the use of IBM MQ AMS.

In the subsequent sections, we configure the RDQM HA group with the root user. However, it is strongly recommended to employ a user within the mqm group, equipped with the appropriate permissions. Detailed instructions for setting up the user mqm with the correct permissions can be found in the following link:

<https://www.ibm.com/docs/en/ibm-mq/9.3?topic=availability-requirements-rdqm-ha-solution>

- ▶ You have the option to establish SSH access without the need for passwords, allowing you to issue configuration commands on just one node within the HA group. Alternatively, if you prefer, you can manually copy the commands to each individual node as an alternative approach.
- ▶ The HA group operates as a Pacemaker cluster, and configuring this Pacemaker cluster involves several steps. Initially, you define the Pacemaker cluster by making adjustments to the `/var/mqm/rdqm.ini` file. In this particular configuration, each interface is allocated a distinct IP address, as demonstrated in Example 7-5.

Example 7-5 Interface configuration

```
Node:
Name=node-a.powerm.ma
HA_Primary=192.168.100.1
HA_Alternate=192.168.101.1
HA_Replication=192.168.102.1
Node:
Name=node-b.powerm.ma
HA_Primary=192.168.100.2
HA_Alternate=192.168.101.2
HA_Replication=192.168.102.2
Node:
Name=node-c.powerm.ma
HA_Primary=192.168.100.3
```

```
HA_Alternate=192.168.101.3
HA_Replication=192.168.102.3
```

Once the configuration is complete, follow these steps:

- ▶ Edit the `/var/mqm/rdqm.ini` file on one of the three servers to define the cluster.
- ▶ Duplicate the modified file onto the remaining two servers that will serve as nodes within the Pacemaker cluster.
- ▶ Execute the command shown in Example 7-6 as the root user on each of the three servers:

Example 7-6 rdqmadm command

```
$ rdqmadm -c
```

- ▶ Next, create a Replicated Data Queue Manager with Root Privilege by executing the command in Example 7-7 on the secondary nodes (Node B and Node C). This creates the secondary queue manager RDQM1.

Example 7-7 Create secondary queue manager on nodes Node B and Node C

```
$ ./crtmqm -sxs -fs 3072M RDQM1
Creating replicated data queue manager configuration.
IBM MQ secondary queue manager created.
```

- ▶ On the primary node (Node A), use the command shown in Example 7-8 to create the primary queue manager RDQM1.

Example 7-8 Create primary queue manager

```
$ ./crtmqm -sx -fs 3072M RDQM1
Creating replicated data queue manager configuration.
IBM MQ queue manager created.
Directory '/var/mqm/vols/rdqm1/qmgr/rqdm1' created.
The queue manager is associated with installation 'Installation1'.
Creating or replacing default objects for queue manager 'RDQM1'.
Default objects statistics : 84 created. 0 replaced. 0 failed.
Completing setup.
Setup completed.
Enabling replicated data queue manager.
Replicated data queue manager enabled.
```

- ▶ To retrieve the status of the RDQM, use the command shown in Example 7-9.

Example 7-9 Checking status of RDQM

```
$ rdqmstatus -m RDQM1
Node: node-a.powerm.ma
Queue manager status: Running
CPU: 0.07%
Memory: 267MB
Queue manager file system: 67MB used, 2.9GB allocated [2%]
HA role: Primary
HA status: Normal
HA control: Enabled
HA current location: This node
HA preferred location: This node
(...)
```

```
Node: node-b.powerm.ma  
HA status: Normal
```

```
Node: node-c.powerm.ma  
HA status: Normal
```

7.2 Using VM Recovery Manager to Replicate Pacemaker cluster

This section provides an overview on how to configuration VM Recovery Manager DR to provide recovery of a Pacemaker cluster in a second data center in order to provide a disaster recovery solution.

The production Pacemaker cluster is running in the primary data center and Pacemaker is providing fail over support for nodes within the cluster to provide high availability. However, if the primary data center is faced by a critical failure or disaster that takes the whole data center off line, then the workloads need to be recovered in a secondary data center which is not affected by the disaster. One solution to this issue can be provided by using IBM VM Recovery Manager Disaster Recovery (VMRM DR).

VMRM DR is an IBM solution that utilizes storage replication and virtual machine restart technologies to restart all of the Pacemaker cluster nodes to a known state in a controlled fashion in your backup site.

7.2.1 VM Recovery Manager DR Overview and Architecture

DR of applications and services is a key component to provide continuous business services. The VMRM DR for IBM Power solution is a DR solution that is easy to deploy and provides automated operations to recover the production site. The VMRM DR solution is based on the IBM Geographically Dispersed Parallel Sysplex (IBM GDPS®) offering concept that optimizes the usage of resources. This solution does not require you to deploy backup VMs for DR. Thus, the VMRM DR solution reduces the software license and administrative costs.

Unlike the clustered-based technology that uses redundant hardware and software components to provide fail over in a component failure, VMRM DR is developed by using VM restart-based technology, where an out-of-band monitoring and management component is used to restart the VMs when a failure of the host infrastructure occurs.

Figure 7-3 on page 133 shows an overview of VMRM DR. IBM VMRM DR was formerly known as IBM Geographically Dispersed Resiliency (IBM GDR) and provides a disaster recovery solution using virtual machine restart technology.

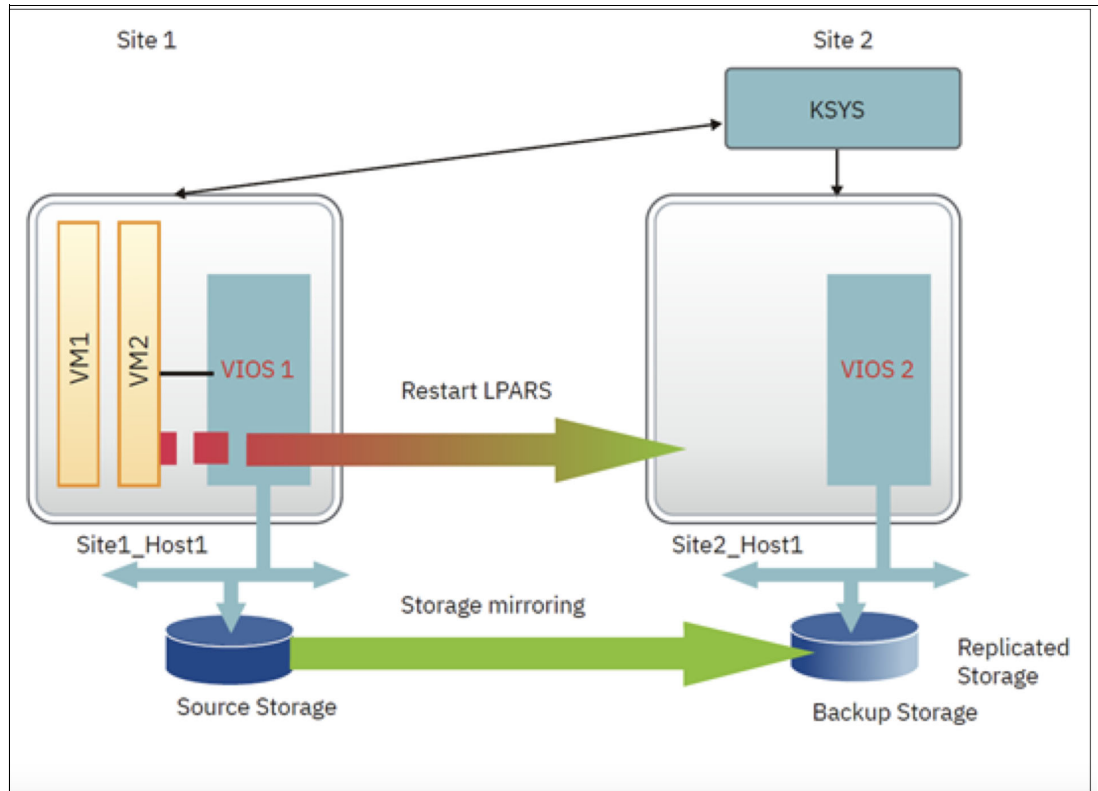


Figure 7-3 Disaster recovery Architecture for VM Recovery Manager.

In this architecture VM1 and VM2 can be assumed to be members of your Pacemaker cluster. In case all the nodes of Pacemaker cluster are failed, it can be restarted on the backup site with replicated storage. VMRM-DR is a DR solution that is easy to deploy and provides automated operations to recover the production workloads at the remote site.

VMRM DR is based on the IBM Geographically Dispersed Parallel Sysplex® (IBM GDPS) offering that optimizes the usage of resources. This solution does not require you to deploy the backup VMs for DR. Thus, the VMRM DR solution reduces the software license and administrative costs.

The VM restart-based HADR solution relies on an out-of-band monitoring and management component that restarts the VMs on other hardware when the host infrastructure fails which is shown in Figure 7-3 as the KSYS. This solution provides an easy deployment model that uses the KSYS to monitor the entire VM environment and take the appropriate recovery action based on user defined and flexible fail over policies along with management of the storage replication.

Table 7-1 on page 134 identifies the differences between a conventional cluster-based DR model and the VMRM DR solution.

Table 7-1 Clustered DR versus VM Recovery Manager DR

Parameters	Cluster-based DR model	VM restart DR model that is used by the VM Recovery Manager DR solution
Deployment method	Redundant hardware and software components are deployed at the beginning of the implementation to provide near real-time failovers when some of the components fail.	With virtualization technology, many images of the OS are deployed in a system. These VMs are deployed on physical hardware by the hypervisor that allocates and manages the CPU, memory, and I/O physical resources that are shared among the VMs.
Dependency	This solution relies on the monitoring and heartbeat capabilities within the cluster to monitor the health of the cluster and take recovery action if a failure condition is detected.	This solution relies on out-of-band monitoring software that works closely with the hypervisor to monitor the VM environment and to provide a DR mechanism for the VM environment.
Workload startup time	The workload startup time is faster because the VMs and the software stack are already available.	The VMs require more time to restart in the backup environment.
Cluster administration required	Yes.	No.
Error coverage	Comprehensive. This solution monitors the entire cluster for any errors.	Limited. This solution monitors the servers and the VMs for errors.
Deployment simplicity	This solution must be set up in each VM.	Aggregated deployment at the site level.
Protected workload type	Critical workloads can be protected by using this solution.	Critical workloads can be protected by using this solution.
Software license and administrative costs	This solution costs more because redundant software and hardware are required to deploy this solution.	This solution costs less because of optimized usage of resources.

Demonstration: A demonstration of VMRM DR under its original name of IBM GDR is available at [YouTube](#).

VMRM DR has the added feature of supporting DR testing. A copy of the storage at the remote data center is copied and zoned to the DR test VMs, which then can be started in a sandbox environment for testing.

7.2.2 Installation and Configuration requirements

To perform an installation task, root authority is required.

The KSYS LPAR must have at least one core CPU and 8 GB of memory. These requirements can be higher if you have a large environment of more than 100 LPARs in the data center.

As a best practice, you must deploy AIX rules in the VIOS. The VIOS must have enough available space in the / (root), /var, and /usr file systems. Extra CPU and memory resources are needed in each VIOS for VMRM HA management.

The VIOS must have enough available space in the / (root), /var, and /usr file systems and add at least 1-core CPU and 4 GB memory to the sizing that you are planning to deploy on your production environment, and have at least 1-core CPU and 10 GB memory in a scalable environment.

Ensure that you have enough space in the LPAR so that KSYS file sets can be installed successfully. You must have 30 MB of disk space in the /opt directory and 200 MB of disk space in the /var directory. Check whether a KSYS installation is already in progress by using the **ksysmgr q cl** command. If the KSYS software was installed previously, you must uninstall the KSYS software.

Refer to <https://www.ibm.com/docs/en/vmrmdr/1.7?topic=planning> for planning requirement and installation requirement for VM Recovery Manager DR solution implementation.

7.2.3 Disaster recovery mechanism for the VM recovery manager DR solution

After you plan the details about how the VM Recovery Manager DR solution can be integrated into your current environment, review the following flow chart that contains the high-level steps that are involved in the VMRM DR implementation.

Figure 7-4 provides a summary of the entire VMRM DR mechanism for disaster recovery:

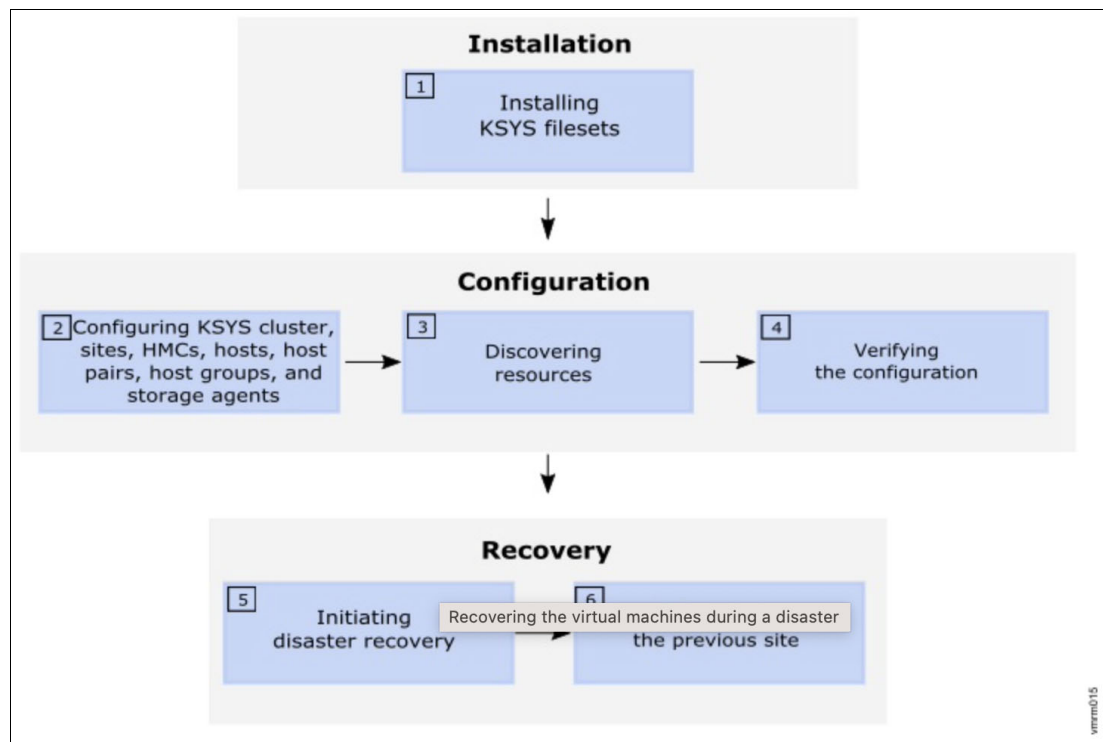


Figure 7-4 VM Recovery manager DR solution: Disaster recovery mechanism

The boxes in the flow chart are described as:

1. Installation

The controlling system (KSYS) is the fundamental component in the VM Recovery Manager DR solution. Therefore, the KSYS filesets must be installed first. The KSYS runs in an AIX 7.2.1 (or later) logical partition in the disaster recovery site. It controls the entire cloud environment for the VM Recovery Manager DR solution. To manage the servers and data replication, the KSYS must be connected to all the managed servers through the HMC and out-of-band storage subsystem connectivity to all associated primary and secondary disks.

Refer to this [VMRM DR document](#) for details on the installation of VM Recovery manager DR solution on controller node.

2. Configuration

After the KSYS is installed and a one-node KSYS cluster is set up, you must configure all the other entities by using the KSYS interface. You must complete the following procedures by using the **ksysmgr** command:

The details on the ksysmgr command including syntax can be found in this [VMRM document](#).

- Create a one-node cluster for the KSYS node.
- Create sites.
- Add HMCs to the corresponding sites.
- Add hosts to the corresponding sites.
- Identify host pairs across the sites.
- Create host groups and work groups.
- Add storage agents to the corresponding sites.
- Add contacts details for error notification.

Refer to this [KSYS configuration document](#) for configuring KSYS subsystem.

3. Discovering the resources.

After the initial configuration is complete, the KSYS discovers all the hosts from all the host groups that are managed by the HMCs in both sites and displays the status.

During discovery, the KSYS subsystem monitors the discovery of all logical partitions (LPARs) or virtual machines (VMs) in all the managed hosts in the active site. The KSYS collects the configuration information for each LPAR, and displays the status, and also logs the status in the log files at the `/var/ksys/log/directory`.

The KSYS discovers the disks of each VM and checks whether the VMs are configured currently for the storage devices mirroring. If the disks are not configured for mirroring properly, KSYS notifies you about the volumes that are not mirrored. All volumes of a VM must be mirrored. Disks can be virtualized by using N_Port ID virtualization (NPIV), virtual SCSI (vSCSI), or combination of all these modes.

The KSYS subsystem collects information about the hosts, VIOS, and logical partitions that can be managed from the HMC during the discovery phase. For example, the KSYS subsystem collects information about the system processor, system memory, hardware, and worldwide port name (WWPN) of the physical Fibre Channel adapter. The KSYS subsystem also checks for the VIOS capability for disaster recovery operations. The KSYS subsystem also collects the dynamic information about the host state, LPAR state, VIOS state, and IP addresses of the host, VIOS, and LPAR.

Command to trigger discover:

```
ksysmgr -t discover site <source_site>.
```


Discovery can be triggered at site level, host group level or workgroup level.

Note: After change in the configuration, running discovery is mandatory to reflect the changes in configuration.

Refer to this [VMRM Document](#) for more details on VM Recovery manager DR discovery operation.

4. Verification of the resource for DR readiness

In addition to the configuration validation that you initiated, the KSYS verifies and validates the environment periodically. The KSYS also verifies the configuration as part of the recovery process. In the verification phase, the KSYS fetches information from the HMC to check whether the backup site is capable to host the VMs during a disaster. The KSYS also verifies storage replication-related details and accessibility of the target disks. The verification is successful only if the storage area network (SAN) zones are configured properly on the target side.

If the verification fails as a part of the recovery process, the failure is considered as recovery failure.

Command to execute verification:

```
ksysmgr -t verify site <Source site>.
```

Note: Verification can be triggered at site level, host group level and at workgroup level.

Refer to this [VMRM document](#) for more details on verification operations.

5. Recovery of the resource in case of production failure

When any planned or unplanned outages occur, you must manually initiate the recovery by using the **ksysmgr** command that moves the virtual machines to the backup site. If you initiate a planned recovery, the storage replication direction is reversed from the current active site to the previously active site. If you initiate an unplanned recovery, the storage is failed over to the backup site and you must manually resynchronize the storage after the previously active site becomes operational.

For planned recovery trigger command:

```
ksysmgr -t move site from<source_site> to=<backup_site>
```

For unplanned disaster recovery failure triggered command:

```
ksysmgr -t move site from=<source_site> to=<backup_site> dr_type=unplanned
```

Note: The move command can be triggered at site level, host group level or workgroup level.

Refer to <https://www.ibm.com/docs/en/vmrmdr/1.7?topic=recovering-vms> for more detailed information

6. Cleanup of resource after DR movement.

After the disaster recovery phase, when the virtual machines have been moved to the target site, you must clean up the source site. The site cleanup process is different for planned recovery operation and unplanned recovery operation. In case of a planned recovery, the KSYS automatically cleans up the source site of all the disk mapping and adapter information. In the case of an unplanned recovery, the source site components such as HMC and storage subsystems might not be operational. So, you must manually

clean up the source site when the HMC and hosts in the previously active site have become operational.

Latest release of VM recovery manager is 1.7. Refer to this [product guide](#) for more information.

Conclusion

By using VM Recovery manager and registering Power system virtual machines having Pacemaker cluster, can be monitored by KSYS controller node. On disaster or critical hardware failure, Pacemaker cluster nodes on the target site can be recovery using replication disk on target site having Pacemaker cluster node details.

7.3 NSF Server High Availability

Configuring NFS Server High Availability using Pacemaker involves several steps to ensure a robust and fail-safe setup. This setup aims to maintain continuous access to NFS shares even if one server fails. Below are the general steps involved in this configuration:

- ▶ Begin by installing Pacemaker on each server that will be part of the high availability setup. Configure these tools to facilitate communication and coordination between the servers.
- ▶ Set up the NFS server on each node, ensuring that the NFS shares are correctly configured and accessible by the clients. Verify that the NFS service is working as expected before proceeding further.
- ▶ Create Pacemaker resource agents for the NFS server. These agents will manage the NFS services, monitor their status, and control their behavior in case of failures. Using Pacemaker, define the NFS server as a cluster resource. This includes specifying the NFS service, IP addresses, and any other necessary resources that need to be managed in a high availability scenario.
- ▶ Establish constraints and failover policies within Pacemaker to control how the resources should behave during various failure scenarios. Define rules for resource placement and recovery in case of node failures or network issues.
- ▶ Simulate failures, node reboots, and network interruptions to ensure that Pacemaker responds correctly and fails over the NFS service seamlessly without disrupting client access.
- ▶ Implement monitoring tools or scripts to continuously monitor the health and performance of the NFS servers. Regularly maintain and update the high availability configuration as needed.

Before configuring NFS Server High Availability using Pacemaker, certain prerequisites must be fulfilled. These include:

- ▶ Red Hat Enterprise Linux 8 and 9: Ensure that the servers intended for NFS Server High Availability are running Red Hat Enterprise Linux (RHEL) versions 8 or 9. Verify compatibility and adherence to the specific requirements for Pacemaker and associated tools on these RHEL versions.
- ▶ Enable the High Availability or Resilient Storage Add-Ons for RHEL on each server participating in the high availability cluster. These add-ons provide the necessary tools and features to establish a robust, fail-safe environment using the commands shown in Example 7-10 on page 139 and Example 7-11 on page 139.

Example 7-10 Enable Red Hat High Availability Add-On Repository

```
$ sudo subscription-manager repos --enable=rhel-9-for-ppc64le-highavailability-rpms
```

Repository 'rhel-9-for-ppc64le-highavailability-rpms' is enabled for this system.

Example 7-11 Enable Red Hat Resilient Storage Add-On Repository

```
$ sudo subscription-manager repos --enable=rhel-9-for-ppc64le-resilientstorage-rpms
```

Repository 'rhel-9-for-ppc64le-resilientstorage-rpms' is enabled for this system.

Note: The [Red Hat Enterprise Linux Resilient Storage Add-On](#) enables simultaneous shared storage access within a highly available cluster. Each server within the cluster necessitates direct access to a shared block device, achievable via a local storage area network (SAN) or provided by various public cloud providers. This Add-On encompasses:

- The Global File System 2.
- The Red Hat Enterprise Linux High Availability Add-On.
- A clustered volume manager, with RHEL 8.4 and higher versions supporting encrypted volumes. An active/active clustered Samba (SMB/CIFS) server.
- Performance Co-Pilot.

- ▶ Set up an Active/Passive configuration for shared storage among the nodes in the cluster. This configuration ensures that only one node (the active node) manages the shared storage while the others (passive nodes) remain on standby, ready to take over in case of a failure. Verify that the shared storage is accessible by all nodes and can be seamlessly switched between active and passive modes without compromising data integrity.

The subsequent links provide detailed instructions for setting up NFS Server High Availability using Pacemaker, applicable to either Linux on IBM Power or Linux on IBM Power Virtual Server.

- <https://cloud.ibm.com/docs/sap?topic=sap-ha-rhel-nfs>
- <https://access.redhat.com/solutions/3796221>
- https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux/9/html/configuring_and_managing_high_availability_clusters/index

7.4 IBM MQ active/passive architecture using Pacemaker and NFS on IBM Power Virtual Server

To configure IBM MQ High Availability in an active-passive configuration, you can use the Red Hat Enterprise Linux High Availability Add-On. This setup requires the following components:

- ▶ IBM Power Virtual Server with a Linux Operating System deployed in either a single or multiple IBM Cloud locations.
- ▶ Allocation of storage to the Logical Partitions (LPARs) within the IBM Power Virtual Server environment.
- ▶ Implementation of an IBM network linking IBM Cloud locations, especially if multiple data centers are selected. The High Availability repository correctly enabled.
- ▶ Deployment of IBM MQ version 9.2 or above

For comprehensive guidance and detailed instructions, please refer to the publication titled *IBM Power Systems High Availability and Disaster Recovery Updates: Planning for a Multicloud Environment*, REDP-5663.³

7.5 IBM PowerVC three-node architecture using Pacemaker

IBM PowerVC (Power Virtualization Center) is a software solution designed to simplify the management and deployment of virtualized resources on IBM Power Systems servers. It serves as a comprehensive tool for managing virtual environments, offering capabilities for provisioning, monitoring, and optimizing virtualized infrastructure. Key features of IBM PowerVC include:

- ▶ PowerVC enables users to create, deploy, and manage virtual machines (VMs) on IBM Power Systems servers. It provides a user-friendly interface for VM creation, resource allocation, and configuration.
- ▶ It allows administrators to pool and efficiently manage physical resources such as CPU, memory, storage, and network, optimizing their utilization across multiple virtual machines.
- ▶ PowerVC offers automation capabilities for streamlining routine tasks, such as VM provisioning and resource allocation. It also supports self-service portals, enabling users to request and deploy VMs within defined resource constraints.
- ▶ It integrates with various virtualization technologies, including PowerVM, IBM's hypervisor technology for Power Systems. PowerVC supports multiple operating systems and is designed to work seamlessly with Power Systems architecture.
- ▶ It provides monitoring tools to track the performance and health of virtualized environments. It offers reporting functionalities for resource usage, allowing administrators to optimize resource allocation.

PowerVC versions 2.0.2 and higher offer support for a three-node architecture. The multi-node setup encompasses specific elements aimed at delivering a High Availability (HA) solution, outlined in subsequent sections here. In this deployment model, the majority of services operate in an Active-Active mode, overseen by Pacemaker or Corosync for monitoring.

For a comprehensive, step-by-step deployment approach, please consult the following publication and link:

- *IBM Power Systems High Availability and Disaster Recovery Updates: Planning for a Multicloud Environment*, REDP-5663³.
- <https://www.ibm.com/docs/en/powervc/2.1.1?topic=powervc-high-availability-scal-e-architecture>

³ <https://www.redbooks.ibm.com/redpapers/pdfs/redp5663.pdf>

Related publications

The publications listed in this section are considered particularly suitable for a more detailed discussion of the topics covered in this book.

IBM Redbooks

The following IBM Redbooks publications provide additional information about the topic in this document. Note that some publications referenced in this list might be available in softcopy only.

- ▶ *IBM Power Systems High Availability and Disaster Recovery Updates: Planning for a Multicloud Environment*, REDP-5663
- ▶ *High Availability and Disaster Recovery Options for IBM Power Cloud and On-Premises*, REDP-5656
- ▶ *IBM Virtual Machine Recovery Manager for IBM Power Cookbook*, SG24-8539
- ▶ *SAP HANA on IBM Power Systems Virtual Servers: Hybrid Cloud Solution*, REDP-5693

You can search for, view, download or order these documents and other Redbooks, Redpapers, Web Docs, draft and additional materials, at the following website:

ibm.com/redbooks

Online resources

These websites are also relevant as further information sources:

- ▶ Cluster Labs Pacemaker Documentation
<https://clusterlabs.org/pacemaker/>
- ▶ Red Hat Enterprise Linux Configuring and managing high availability clusters: Using the Red Hat High Availability Add-On to create and maintain Pacemaker clusters
https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux/9/html-single/configuring_and_managing_high_availability_clusters/index#doc-wrapper
- ▶ Red Hat HA Solutions for SAP HANA, S/4HANA and NetWeaver based SAP Applications
https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux_for_sap_solutions/9/html-single/red_hat_ha_solutions_for_sap_hana_s4hana_and_netweaver_based_sap_applications/index
- ▶ Automating SAP HANA Scale-Up System Replication using the RHEL HA Add-On (RHEL8)
https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux_for_sap_solutions/8/html/automating_sap_hana_scale-up_system_replication_using_the_rhel_ha_add-on
- ▶ Automating SAP HANA Scale-Up System Replication using the RHEL HA Add-On (RHEL9)
https://access.redhat.com/documentation/en-us/red_hat_enterprise_linux_for_sap_solutions/9/html/automating_sap_hana_scale-up_system_replication_using_the_rhel_ha_add-on

- ▶ DB2 High Availability strategies
<https://www.ibm.com/docs/en/db2/11.5?topic=availability-high-strategies>
- ▶ DB2 High availability disaster recovery (HADR)
<https://www.ibm.com/docs/en/db2/11.5?topic=server-high-availability-disaster-recovery-hadr>
- ▶ High availability with Db2 server
<https://www.ibm.com/docs/en/db2/11.5?topic=availability-high-db2-server>

Help from IBM

IBM Support and downloads

[ibm.com/support](https://www.ibm.com/support)

IBM Global Services

[ibm.com/services](https://www.ibm.com/services)

To determine the spine width of a book, you divide the paper PPI into the number of pages in the book. An example is a 250 page book using Plainfield opaque 50# smooth which has a PPI of 526. Divided 250 by 526 which equals a spine width of .4752". In this case, you would use the .5" spine. Now select the Spine width for the book and hide the others: **Special>Conditional Text>Show/Hide>SpineSize(->Hide)>Set** . Move the changed Conditional text settings to all files in your book by opening the book file with the spine.fm still open and **File>Import>Formats** the Conditional Text Settings (ONLY!) to the book files.

Draft Document for Review December 20, 2023 6:09 pm

8557spine.fm 143



Redbooks

Using Pacemaker to Create Highly Available Linux Solutions

SG24-8557-00

ISBN DocISBN



(1.5" spine)

1.5" <-> 1.998"

789 <-> 1051 pages



Redbooks

Using Pacemaker to Create Highly Available Linux Solutions on IBM

SG24-8557-00

ISBN DocISBN



(1.0" spine)

0.875" <-> 1.498"

460 <-> 788 pages

Redbooks

Using Pacemaker to Create Highly Available Linux Solutions on

SG24-8557-00

ISBN DocISBN



(0.5" spine)

0.475" <-> 0.873"

250 <-> 459 pages

Redbooks

Using Pacemaker to Create Highly Available Linux Solutions on IBM Power

(0.2" spine)

0.17" <-> 0.473"

90 <-> 249 pages

(0.1" spine)
0.1" <-> 0.169"
53 <-> 89 pages

To determine the spine width of a book, you divide the paper PPI into the number of pages in the book. An example is a 250 page book using Plainfield opaque 50# smooth which has a PPI of 526. Divided 250 by 526 which equals a spine width of .4752". In this case, you would use the .5" spine. Now select the Spine width for the book and hide the others: **Special>Conditional Text>Show/Hide>SpineSize(->Hide)>Set** . Move the changed Conditional text settings to all files in your book by opening the book file with the spine:fm still open and **File>Import>Formats** the Conditional Text Settings (ONLY!) to the book files.

Draft Document for Review December 20, 2023 6:09 pm

8557spine.fm 144



Using Pacemaker to Create Highly Available

SG24-8557-00

ISBN DocISBN

(2.5" spine)
2.5" <-> mmm.n"
1315 <-> mmm pages



Using Pacemaker to Create Highly Available Linux Solutions on IBM Power

SG24-8557-00

ISBN DocISBN

(2.0" spine)
2.0" <-> 2.498"
1052 <-> 1314 pages





SG24-8557-00

ISBN DocISBN

Printed in U.S.A.

Get connected

