

High Availability Clustering with Red Hat® Enterprise Linux 8

Solution Assurance

Bodo Brand, Daniel Kaiser, David Stark

Document version: 03-2022

Notices and disclaimers

© 2021 International Business Machines Corporation. No part of this document may be reproduced or transmitted in any form without written permission from IBM.

U.S. Government Users Restricted Rights – use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM.

Information in these presentations (including information relating to products that have not yet been announced by IBM) has been reviewed for accuracy as of the date of initial publication and could include unintentional technical or typographical errors. IBM shall have no responsibility to update this information. **This document is distributed "as is" without any warranty, either express or implied. In no event, shall IBM be liable for any damage arising from the use of this information, including but not limited to, loss of data, business interruption, loss of profit or loss of opportunity.** IBM products and services are warranted per the terms and conditions of the agreements under which they are provided.

IBM products are manufactured from new parts or new and used parts. In some cases, a product may not be new and may have been previously installed. Regardless, our warranty terms apply.

Any statements regarding IBM's future direction, intent or product plans are subject to change or withdrawal without notice.

Performance data contained herein was generally obtained in a controlled, isolated environments. Customer examples are presented as illustrations of how those customers have used IBM products and the results they may have achieved. Actual performance, cost, savings or other results in other operating environments may vary.

References in this document to IBM products, programs, or services does not imply that IBM intends to make such products, programs or services available in all countries in which IBM operates or does business.

Workshops, sessions and associated materials may have been prepared by independent session speakers, and do not necessarily reflect the views of IBM. All materials and discussions are provided for informational purposes only, and are neither intended to, nor shall constitute legal or other guidance or advice to any individual participant or their specific situation.

It is the customer's responsibility to ensure its own compliance with legal requirements and to obtain advice of competent legal counsel as to the identification and interpretation of any relevant laws and regulatory requirements that may affect the customer's business and any actions the customer may need to take to comply with such laws. IBM does not provide legal advice or represent or warrant that its services or products will ensure that the customer follows any law.

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 about this publication 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. IBM does not warrant the quality of any third-party products, or the ability of any such third-party products to interoperate with IBM's products. **IBM expressly disclaims all warranties, expressed or implied, including but not limited to, the implied warranties of merchantability and fitness for a purpose.**

The provision of the information contained herein is not intended to, and does not, grant any right or license under any IBM patents, copyrights, trademarks or other intellectual property right.

IBM, the IBM logo, ibm.com and IBM Z, IBM z15, RACF, and z/VM are 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: www.ibm.com/legal/copytrade.shtml

Agenda

❖ Introduction

- ❖ References
- ❖ Cluster Types – Single site
- ❖ Cluster Types – Multi site
- ❖ IBM Z run levels
- ❖ Components

❖ Concepts

- ❖ Resources
- ❖ Quorum
- ❖ Planned/Unplanned Outage
- ❖ Fencing/STONITH
- ❖ Advanced Concepts

❖ Use Case:

- ❖ LPAR HA Cluster with KVM as resource

❖ Appendix

Introduction

Introduction - References

Documentation:

Official Red Hat® **documentation:**

- RHEL7: [LINK](#)
- RHEL8: [LINK](#)

Official Red Hat® **support statements:**

- [z/VM specific](#)
- [further support statements](#)

Official Red Hat® **version changes:**

- RHEL7: [LINK](#)
- RHEL8 Release Notes: [LINK](#)

[Redbooks® publication - HA on Linux](#)

- HA services or applications **uptime** approaches 100%
- **HA withstands** failures that are caused by **planned or unplanned outages**

[Official Pacemaker documentation](#)

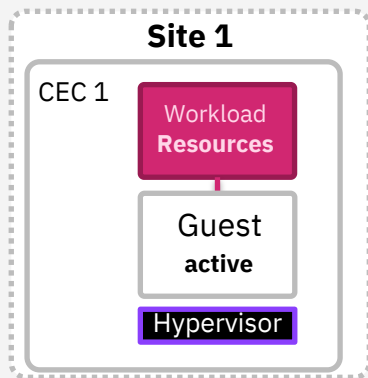
Site A



Site B



Introduction – Cluster Types – Single site

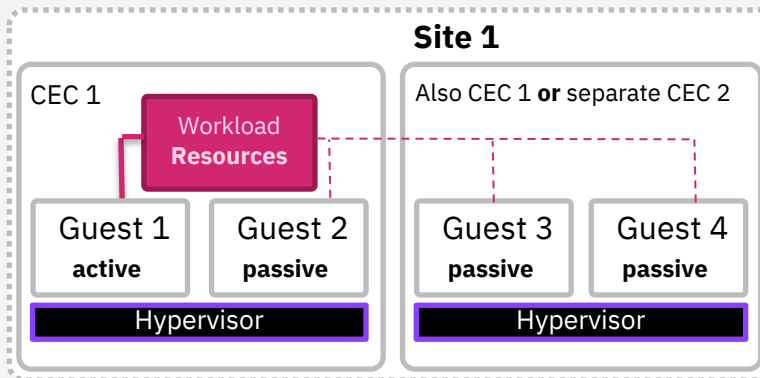


No Red Hat High Availability

Hypervisor can be:

- ❖ LPAR
- ❖ z/VM
- ❖ KVM

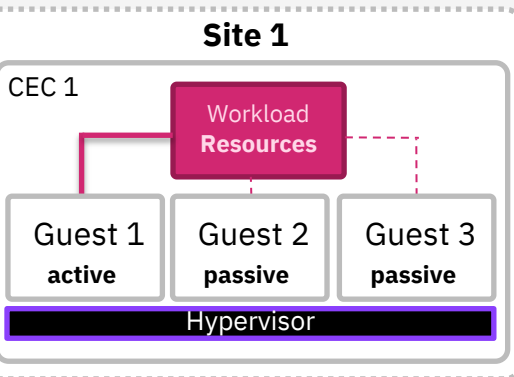
Workload Resources can be anything e.g. an apache webserver



Cluster spanning two Hypervisors&CECs

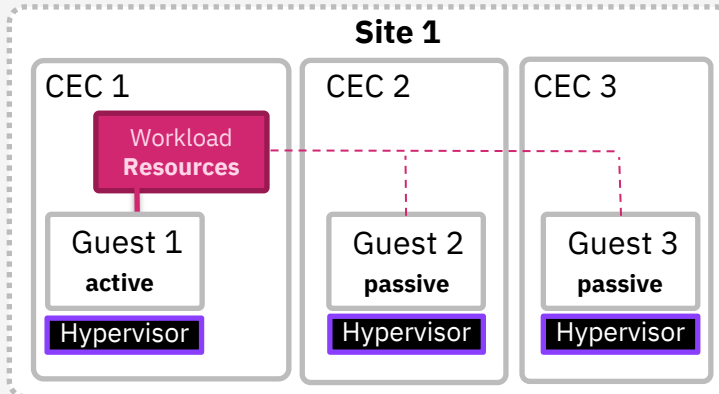
Notes:

- ❖ A quorum server running separately might be needed



Active / Passive Failover

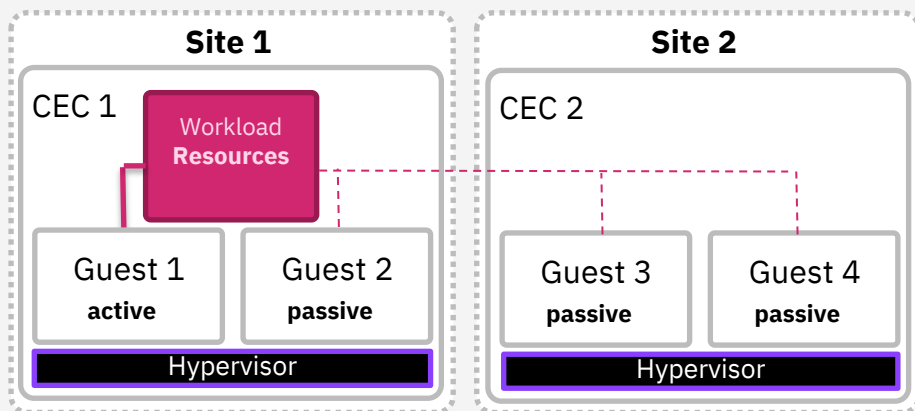
- ❖ Minimal setup to get started with HA.
- ❖ Does not sustain CEC or Hypervisor failures



Cluster spanning over more than two Hypervisors/CECs

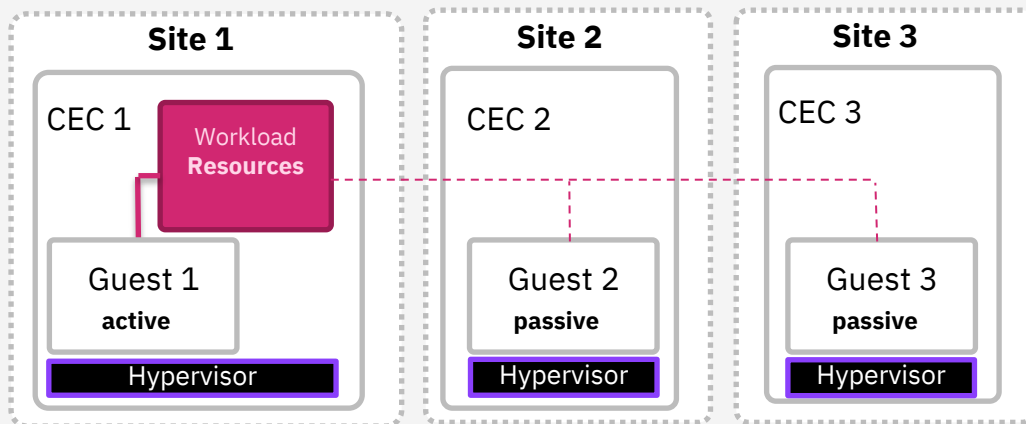
- ❖ Quorum server might not be needed!

Introduction – Cluster Types – Multi site



Multi Site Cluster with two sites

- ❖ You may have to use Replicated Storage instead of Shared Storage
- ❖ The HA settings have to keep the additional latency in mind which might be introduced
- ❖ Reachability gets more complicated



Multi Site Cluster with more than two sites

- ❖ With that many CECs a separate quorum server might not be needed anymore
- ❖ With storage replication you might want to have synchron replication between two sites and asynchron replication with a third site

Introduction – IBM Z run levels

LPAR HA-Cluster

LPAR1

LPAR2

LPAR3

Cluster spanning over three LPARs

- PR/SM Hypervisor with optionally DPM

KVM Guest

LPAR HA-Cluster

LPAR10
(KVM Hyper.)

LPAR11
(KVM Hyper.)

LPAR12
(KVM Hyper.)

Cluster spanning over three LPARs with KVM Guests

- **KVM Guest** is seen as any other workload in the cluster

z/VM HA-Cluster

z/VM Guest 1

z/VM Guest 2

z/VM Guest 3

Cluster spanning over three z/VM Guests

- **KVM Guest** is seen as any other workload in the cluster

LPAR4
(z/VM Hypervisor)

KVM HA-Cluster

KVM Guest 1

KVM Guest 2

KVM Guest 3

Cluster spanning over three KVM Guests

- **KVM Guest** is seen as any other workload in the cluster

LPAR5
(KVM Hypervisor)

Introduction – Components (Minimal)

HA-Cluster

LPARs 1 / KVM Guest 1 / z/VM Guest 1

Admin applications

pcs “cmdline tool to
interact/configure the cluster”
pcsd web ui

Daemons

pcsd
pacemaker
corosync

Optional **Daemons**
SBD

Optional kernel modules

diag288
watchdog



Scripts

Resource Agents
“manage Cluster
Services (start/stop/...)”
Fencing Agents
“isolate cluster
members”

Libraries

libQB
“for logging, tracing,
IPC and polling”

LPARs 2 / KVM Guest 2 / z/VM Guest 2

3

...

Concepts

(Managed) Resources

A cluster contains one or multiple resources. Each resource has following properties:

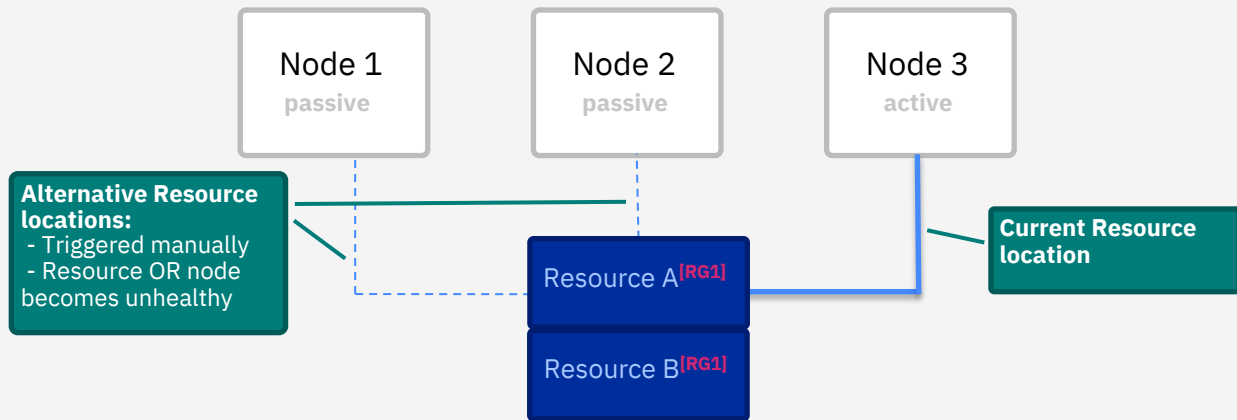
- ❖ type (e.g., apache) and resource identifier (e.g., Website)

Resource:

- ❖ Implemented as Resource Agent (RA)
- ❖ executable/service conforming to a standard (usually **ocf** or **systemd**)
- ❖ handles all **operations**: (**start**, **stop**, **monitor**)
- ❖ **attributes** for configuration (e.g., configfile=a.conf)
- ❖ **constraints** (**location**, **order**, **colocation**)

Resource Groups [RG1]:

- ❖ Resources in a resource group **start and stop in order**
- ❖ When one of the resources moves in the group, the other resources in that group move with it



Predefined Resources:

- ❖ List all predefined Resources
`# pcs resource list`
- ❖ Look into the Resource docu
`# pcs resource describe ...`
- ❖ Add Resources to cluster
`# pcs resource create ...`

Define your own Resource

- ❖ Article from Red Hat ®: [LINK](#)
- ❖ OCF compliant RA: [LINK](#)
 - ❖ XML definition file
 - ❖ Operations implemented in any programming language
 - ❖ Exit codes standardized
- ❖ Add to pacemaker search location: `/usr/lib/ocf`

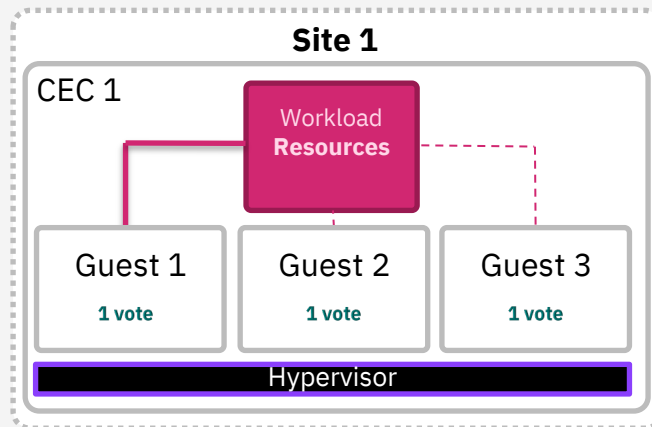
Quorum

Quorum

Corosync Votequorum

- ❖ Quorum decides how many guests can fail before the cluster becomes non-operational
- ❖ Quantity of **votes** are assigned to the **systems**
- ❖ Only when a **majority of votes are present** the cluster **operations are allowed** to proceed.
- ❖ With 1 CEC an uneven amount of nodes ensures quorum when 1 node fails.
- ❖ With 2 CECs you either need:
 - ❖ A third cluster member on a neutral/third side
 - ❖ Or a quorum server on a neutral/third side which only votes but does not participate otherwise
- ❖ (See next slides for reasons on why a third side is needed)

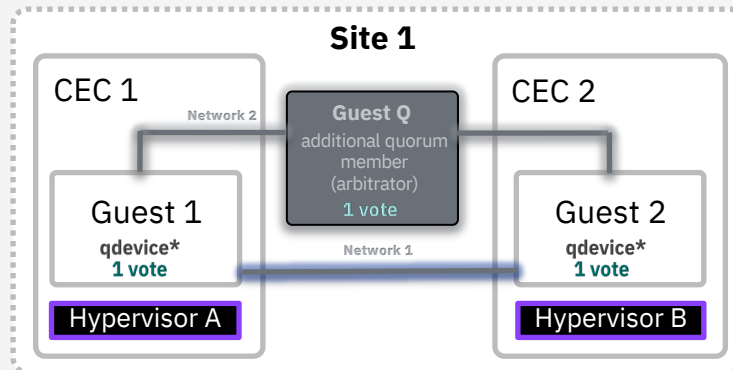
Uneven number of cluster members



Notes:

- Usually 1 vote each
- cluster is quorate / functional with 2 votes.

Even number of cluster members



Notes:

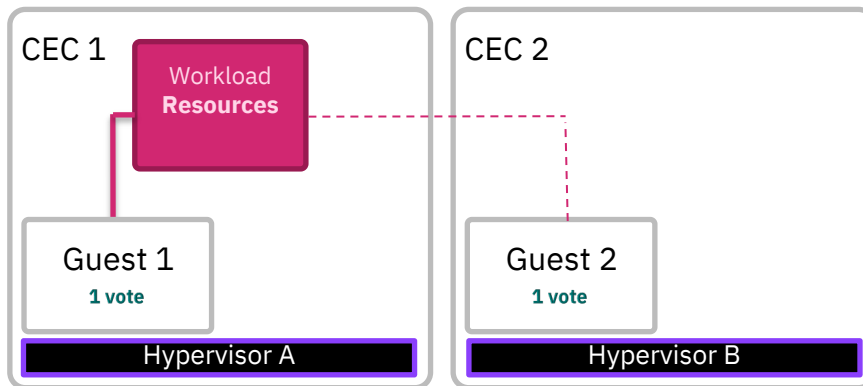
- Guest Q should be reachable through different network connection
- can be used by multiple independent clusters

Quorum – 2/3-Nodes Challenges

2-node cluster

split brain challenge

Operative with 1 vote



Does not work because in case of a network failure both CECs would be quorate or inquorate which results in running the workload twice or not at all. (**split brain**)

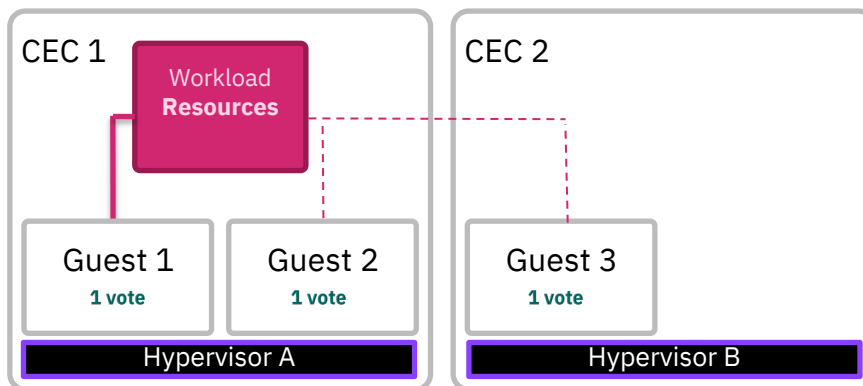
Note:

- When a CEC/Hypervisor dies, it affects the network and guest at the same time. A Tie-Breaker where the “lower node id” (e.g. guest 1) wins is not be able to handle this situation well.

3-node cluster

quorum challenge

Operative with 2 vote



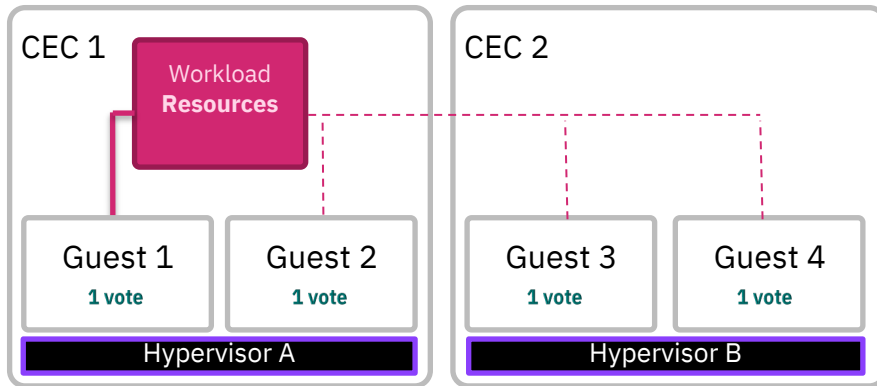
Does not work because in case of CEC 1 failure CEC 2 alone is not quorate!

Quorum – 4 Nodes Challenges and Solution

4-nodes cluster

split brain +
quorum
challenge

Operative with
3 vote

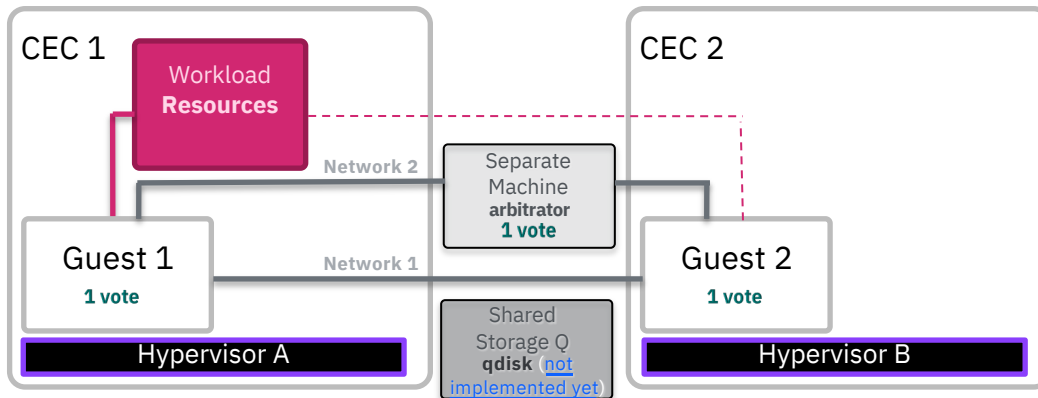


Does not work because in case of a network failure in between the CECs both CECs are not quorate.

Solution:

2 / 4 nodes
+ additional
quorum
member

Operative with
2 vote



A **separate machine** which participates as quorum member **prevents split brain** and other quorum challenges.

Planned / Unplanned Outage

Planned / Unplanned Outage

Planned Outage

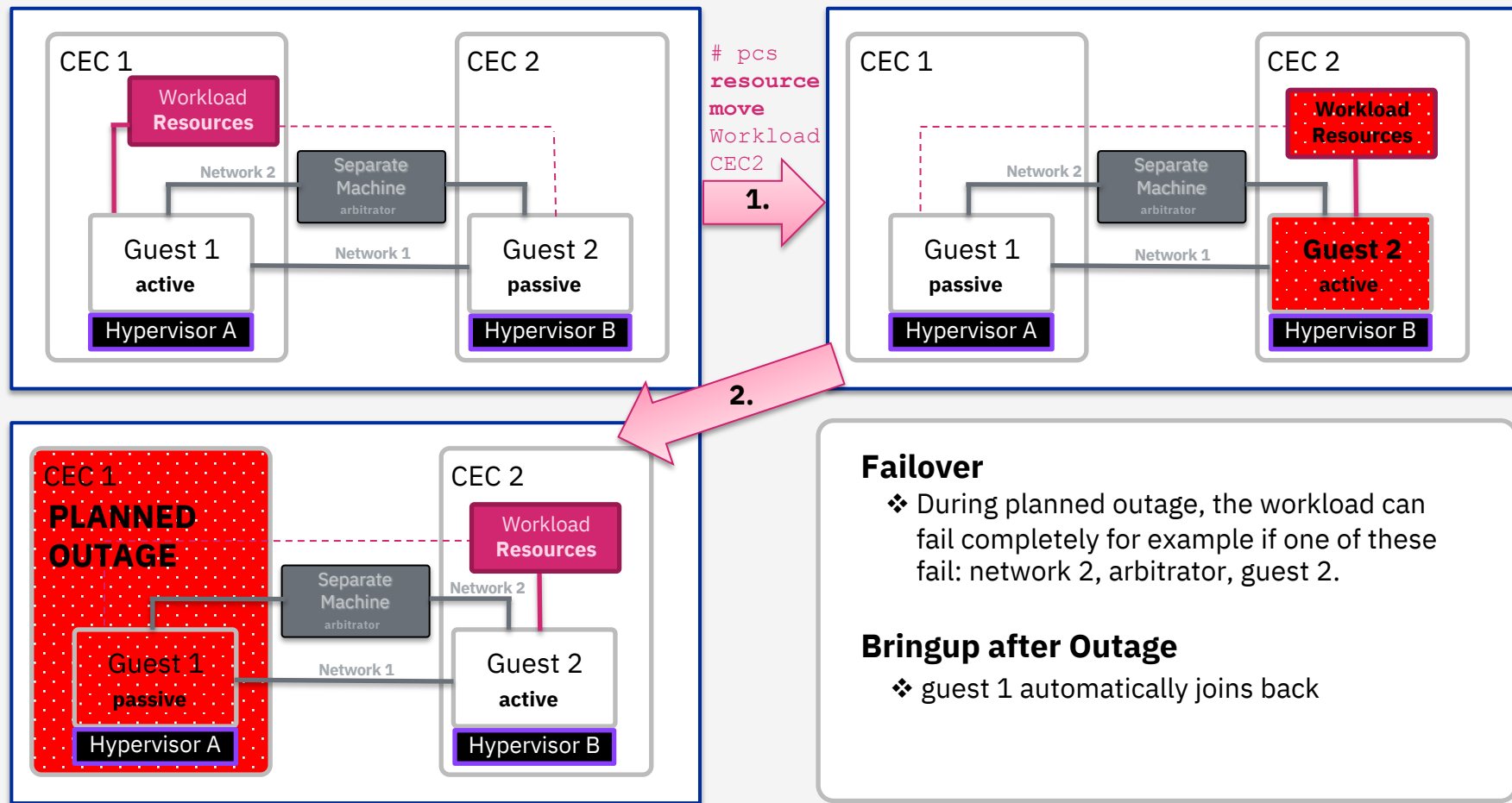
- ❖ Red Hat HA allows you to manually trigger the movement of workload to another cluster member
- ❖ Live Guest Relocation (LGR) (**z/VM**) / Live Migration (**KVM**)
 - ❖ **z/VM Guests Cluster:** For Live Guest Relocation to work you need a **SSI cluster**.
 - ❖ Only LGR of passive guests might be supported ([LINK](#))
 - ❖ **KVM Guests Cluster:** For Live Migration to work you usually need Shared Storage (or Replicated Storage) which is mounted read and write between the Hypervisors.
 - ❖ **LPARs Cluster with KVM Guest Workload:** Live Migration of the Resource is automatically tried when all requirements are met.

Unplanned Outage

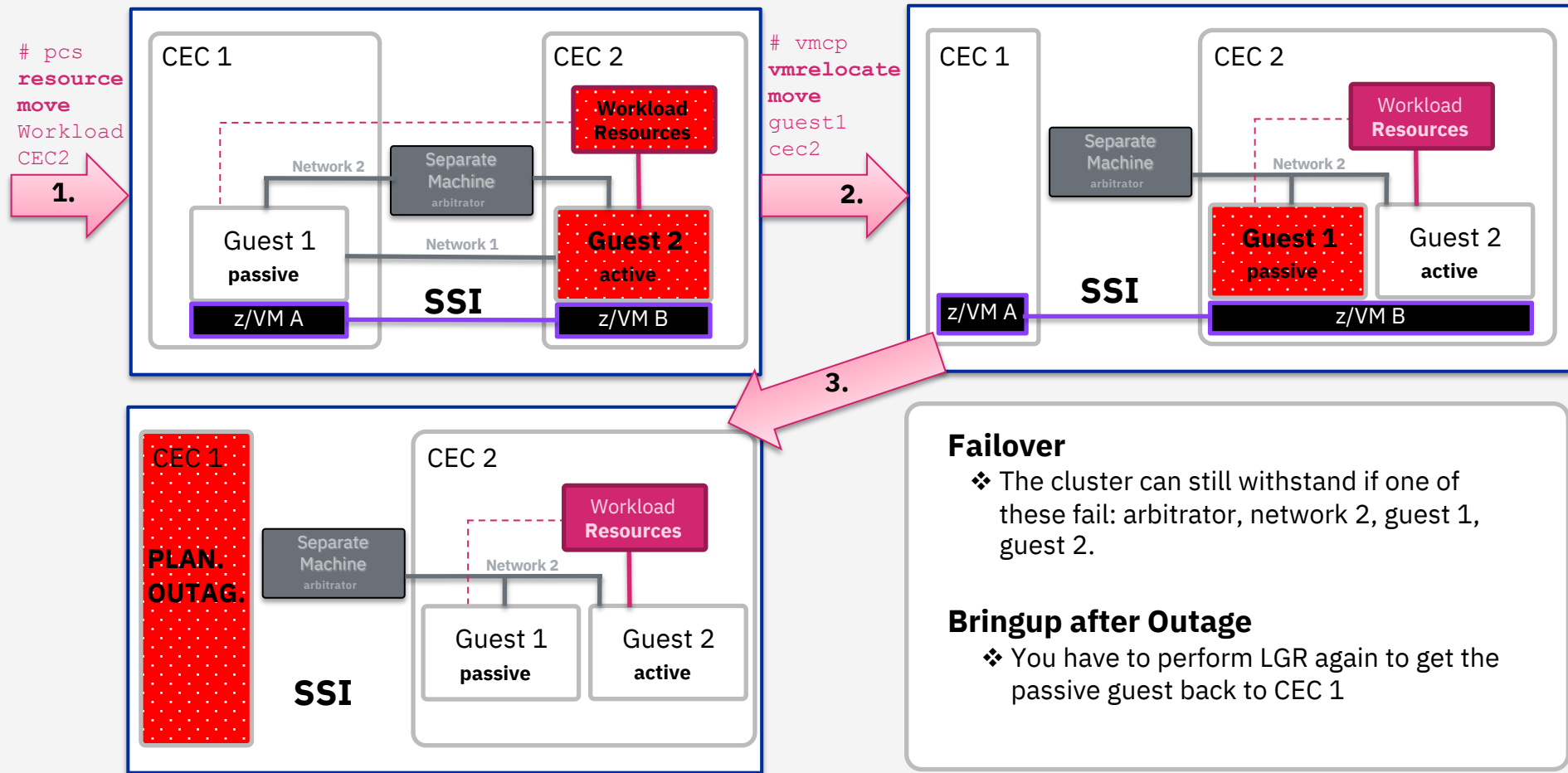
- ❖ Red Hat HA automatically fails over in case of failure as soon as the node released all Resources (see Fencing/STONITH concept).
- ❖ Live Guest Relocation (**z/VM**) / Live Migration (**KVM**)
 - ❖ Cannot be used as the Guest as you would move a corrupted/broken guest in this case.

➤ See following slides for graphic illustrations.

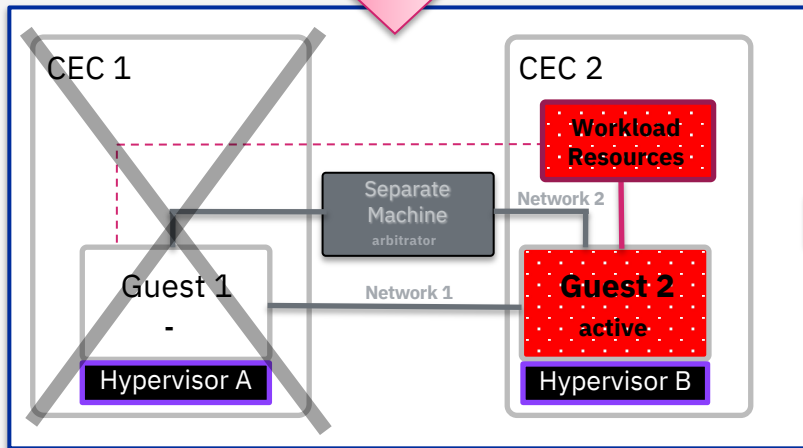
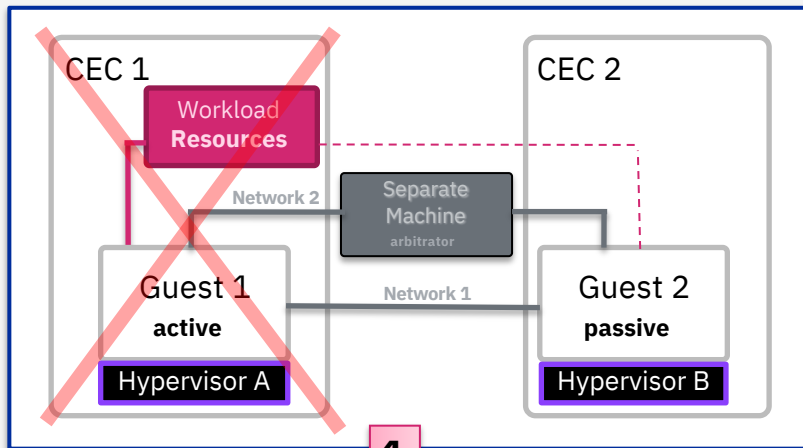
Planned Outage – 2 CECs Example



Planned Outage – 2 CECs Example with SSI (z/VM)



UN-Planned Outage – 2 CECs Example

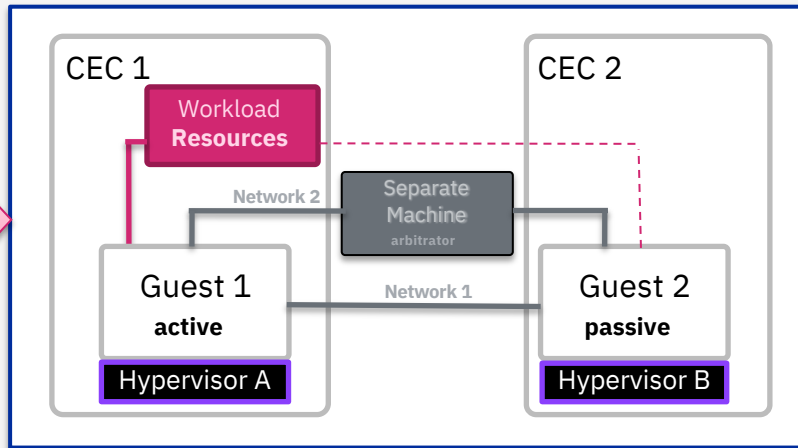


Failover

- ❖ Workload is automatically moved to guest 2 on CEC 2
- ❖ Fencing of guest 1 fail
- ❖ **No big differences to non-SSI Cluster**

Bringup after Outage:

- ❖ Guest 1 automatically joins back
- ❖ Workload resources move back to Guest 1 by default (configurable)



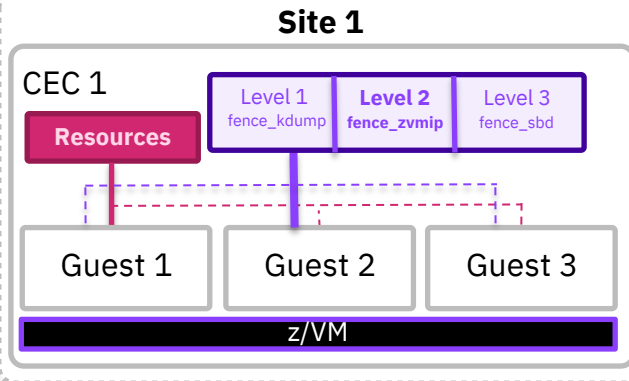
Fencing / STONITH

Fencing / STONITH

Concept of Fencing/STONITH

- ❖ Ensures that it is not possible for a guest to run resources if the guest is not intended to do so
- ❖ Depending on your HA Setup you might want to use a combination of the following available fencing agents:
 - ❖ **fence_zvmip (via z/VM [SMAPI](#))**: In a 2 CEC setup you have to use 2 fence-agents specifying the other side. This fence agent is not SSI aware which means you would have to change both fence-agents every time you do LGR. ([Instructions](#))
 - ❖ **fence_sbd (via [SBD](#))**: SBD watches the cluster health locally and triggers self fencing if needed. Additionally, SBD watches a shared disk where the fence-agent can write a poison pill to which also triggers fencing.
 - ❖ **fence_ibmz (via [HMC API](#))**: Performs a deactivate, activate and load operation on a LPAR. Both Classic (PR/SM) and DPM is included. (available with RHEL8.6+ and RHEL9+)
 - ❖ **fence_kdump**: This fence agent just detects if the failing guest is currently taking a kdump. If yes fencing is considered complete.
 - ❖ **fence_virsh (via [KVM - virsh](#))**: Simply ssh's to the Hypervisor and fences the guest through virsh commands.

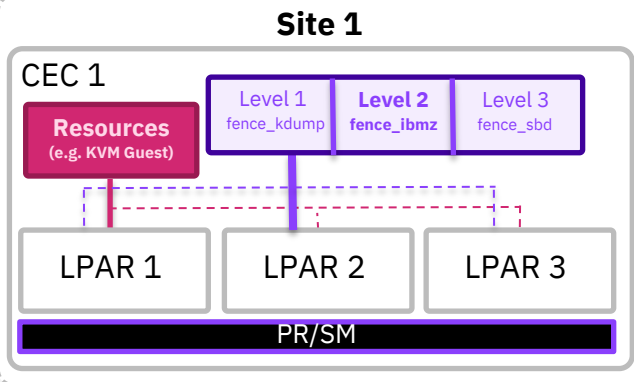
z/VM Guest Cluster – fencing example



Specifics:

- ❖ Level 3 (fence_sbd): You can specify the z/VM command executed on SBD fencing

LPARs Cluster – fencing example



Specifics:

- ❖ Level 2 (fence_ibmz): Check support for Redhat version. Should be able to handle both:
 - ❖ classic (PR/SM)
 - ❖ DPM

Advanced Concepts

Advanced concepts for reference

Remark

- ❖ Not covered technically in this presentation

Cluster notifications & Error conditions

- ❖ When errors happen in the cluster the cluster might not proceed without manual intervention
- ❖ Getting notified of cluster problems in time can be crucial for High Availability

Cluster User Permissions

- ❖ To make sure a specific role (e.g. a cluster operator) can only perform actions specific to his job role you can configure ACLs (Access Control Lists)

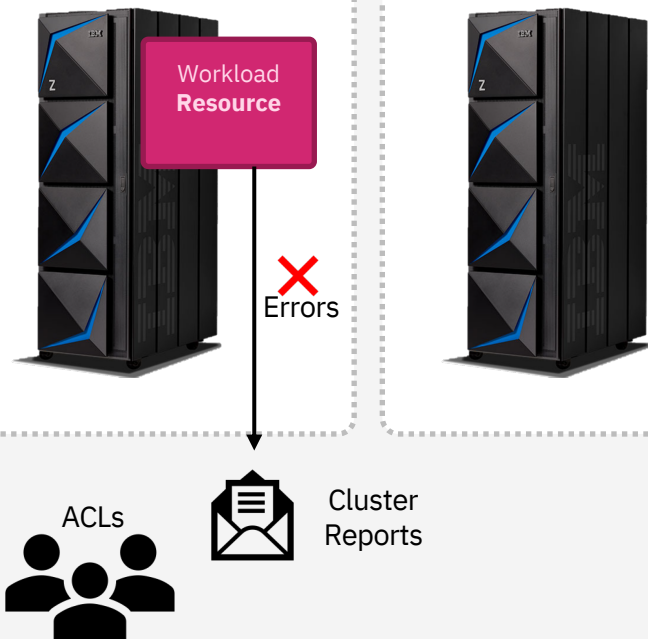
Deal with Multi Site Clusters

- ❖ To prevent "split-brain" in multi-site clusters the booth ticket manager spans an overlay cluster over existing clusters on different sites

Disaster Recovery

- ❖ A secondary cluster can be specified as recovery site

Site A



LPAR HA Cluster with KVM as resource

Agenda

❖ Architecture and Requirements

❖ Steps for setup creation

1. First Steps – Cluster Setup
2. Quorum
3. Fencing/STONITH
 4. fence_ibmz
 5. fence_sbd
 6. fence_kdump
 7. Fencing levels
8. GFS2 (Shared Storage)
9. VirtualDomain (KVM Guest)
10. Cluster Testing

Guidance Notes

- Some of the operations must be run on all nodes and some only one node.
- The "**Run on**" graphic on the right indicate on which of the nodes you must run the command.

Run on:

LPAR 1

LPAR 2

...

- "#" at the beginning of the line indicate a privileged bash command.

Run on:

LPAR 1

LPAR 2

...

- The graphic on the right is used for illustration purposes.

Architecture and Requirements

Architecture and Requirements

Our Test Setup:

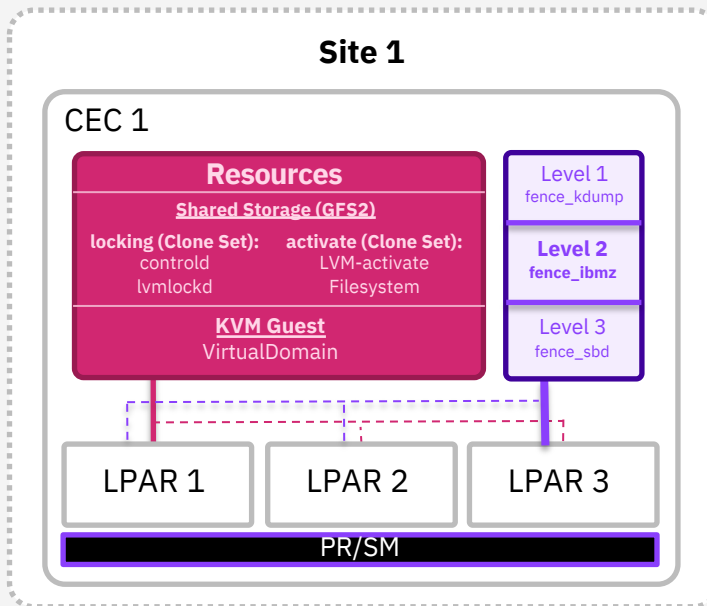
- ❖ z15™ and DS8000®
- ❖ 3 LPARs (on one CEC)
 - ❖ 3 ECKD DASD (LPAR guest OS)
- ❖ Shared Storage (for SBD and KVM)
 - ❖ 2+ ECKD DASD

LPARs:

- ❖ PR/SM mode **or** DPM mode
- ❖ Distro: RHEL 8.5

Legend:

--- Failover Paths



- ❖ **Clone Set:**
Is a Resource Group which is cloned to all other cluster members.

First Steps – Cluster Setup

Installation and Firewall

Step 1.0 – Installation

```
rhel8# subscription-manager register --auto-attach
```

```
rhel8# dnf config-manager --set-enabled \  
rhel-8-for-s390x-highavailability-rpms
```

```
rhel8# yum update -y
```

```
rhel8# yum install -y pcs
```

```
rhel8# yum install -y pacemaker \  
fence-agents-all
```

Note:

- fence_ibmz will only be installed with later Red Hat Enterprise Linux versions.
Manual installation from upstream will be described later in this Use Case.

Run on:

LPAR 1

LPAR 2

LPAR 3

Step 1.1 – Firewall Configuration

```
rhel8# firewall-cmd --permanent \  
--add-service=high-availability
```

```
rhel8# firewall-cmd --reload
```

Run on:

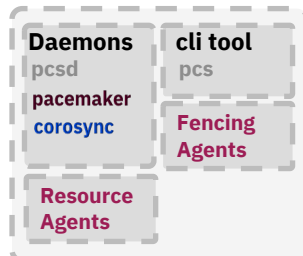
LPAR 1

LPAR 2

LPAR 3

Site 1

CEC 1



LPAR 1

LPAR 2

LPAR 3

PR/SM

Cluster Setup

Step 1.2 – Prepare Cluster

- ❖ Create a linux® user used by the cluster

```
rhel18# passwd hacluster
```

Note:

- Same password for each node is recommended.

- ❖ Enable the cluster controlling and configuration daemon

```
rhel18# systemctl enable pcsd.service --now
```

Run on:

LPAR 1

LPAR 2

LPAR 3

Step 1.3 – Auth Nodes

```
rhel18# pcs host auth LPAR1 LPAR2 LPAR3
```

Username: hacluster

Password: ...

Run on:

LPAR 1

Step 1.4 – Setup Cluster

```
rhel18# pcs cluster setup \  
my_cluster LPAR1 LPAR2 LPAR3
```

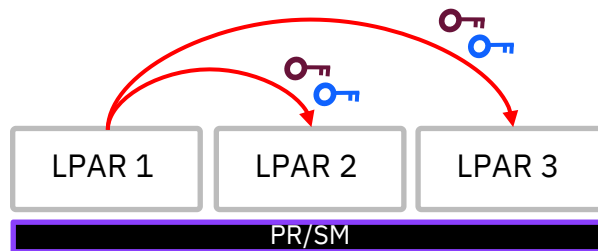
Run on:

LPAR 1

Site 1

CEC 1

/etc/pacemaker/authkey
/etc/corosync/authkey



Startup of the Cluster and Status

Step 1.5 – Start and Enable Services

Run on:

LPAR 1

```
rhel18# pcs cluster start --all
rhel18# pcs cluster enable --all
```

Step 1.6 – Status and CLI Tools

Run on:

LPAR 1

❖ The pcs CLI tool allows you to configure the cluster (pacemaker + corosync) and view the status

❖ Full cluster status:

```
rhel18# pcs status --full
```

❖ Pacemaker configuration

```
rhel18# pcs cluster cib
```

❖ Pacemaker and corosync have their own cli tools:

❖ Pacemaker configuration

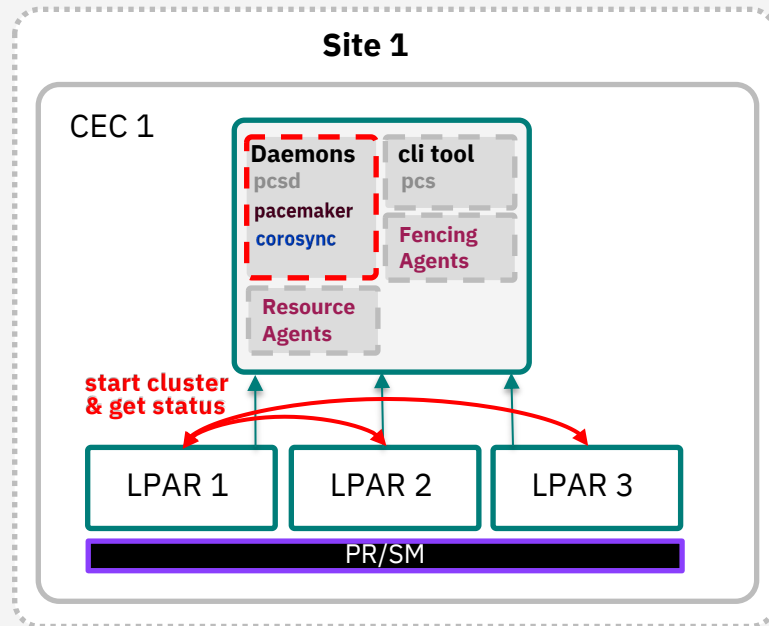
```
rhel18# cibadmin -Q
```

❖ Show corosync object database

```
rhel18# corosync-cmapctl
```

❖ Dump live corosync flight data

```
rhel18# corosync-blackbox
```



Quorum

Quorum

Step 2.0 – Considerations

Run on:

LPAR 1

- ❖ With `wait_for_all` enabled the whole cluster only becomes quorate/functional for the first time when all cluster members are available

```
rhel8# pcs quorum update wait_for_all=1
```

Note:

- For example, when starting three LPARs consecutively. Without enabling `wait_for_all` the last LPAR might be fenced from the two already available LPARs.

- ❖ "totem token timeout" specifies in milliseconds until a token loss is declared

```
rhel8# pcs cluster config update \  
totem token=5000
```

Note:

- For totem token limits check out the [corosync support policies](#)

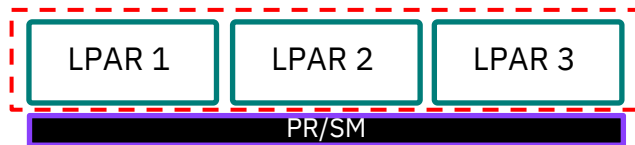
- ❖ Check totem token timeout

```
rhel8# corosync-cmapctl | \  
grep "runtime.*totem.token "
```

Site 1

CEC 1

Cluster only becomes active for the first time when all three LPARs are available.

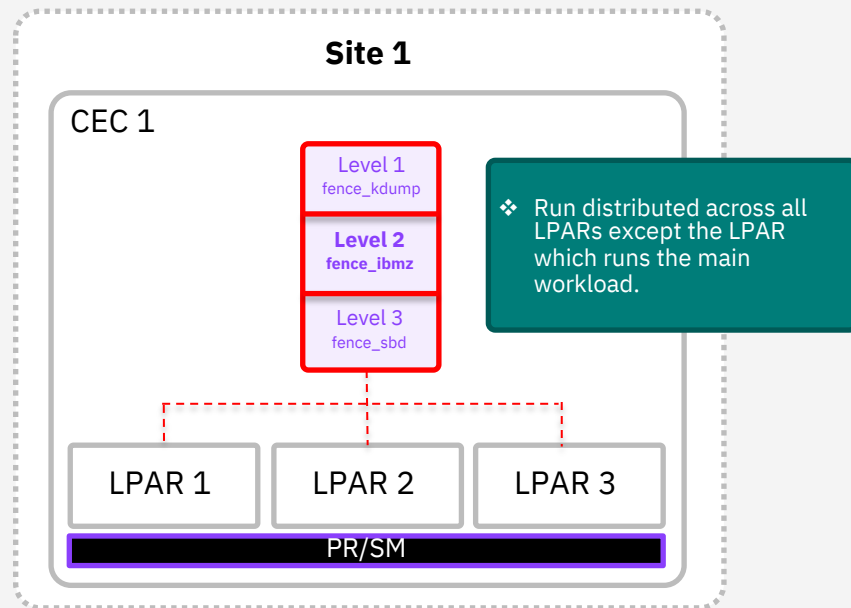


Fencing / STONITH

Fencing / STONITH

Step 3.0 – Considerations

- ❖ The main fencing method will be power fencing over the HMC:
 - ❖ **fence_ibmz (Level 2)**: Provides solid fencing because it is a power fencing method which triggers fencing externally via the HMC API.
- ❖ When the HMC is not available, SBD is used as backup fence agent:
 - ❖ **fence_sbd (Level 3)**: As last resort, self fencing is a reliable backup option which might take a bit longer but should take effect in the worst cases. The poison pill is used to speed up this fence method in some failure cases.
- ❖ For debugging purposes, we also include:
 - ❖ **fence_kdump (Level 1)**: When you take a kdump (either automatically or manually) you want to prevent other fencing methods to trigger. This way the fencing is considered successful when a LPAR kdumps.



Fencing / STONITH fence_ibmz (Level 2)

fence_ibmz – Prerequisites

Step 4.0 – Create a HMC user

Run on:

HMC - Classic mode

❖ Create minimal viable role (Set Scope and Permissions of HMC user)

❖ Minimum tasks required:
Deactivate, Activate, Load, View Activation Profiles

❖ Scope to cluster members only.

❖ Allow access to Web Services management interfaces

☒ Allow access to Web Services management interfaces

Maximum web services API sessions (0-9999):

100

Idle web services API session timeout (1-360 minutes):

15

❖ Keep timeouts configured here in mind. The default values are usually very high which should not affect fencing actions.

Session

☒ Session timeout (minutes): 300

☒ Verify timeout (minutes): 15

☒ Idle timeout (minutes): 20

Summary for user4fencing

General

Description:

Last logon:

Last mobile logon:

Email Address:

Disabled: No

Authentication

Password authentication type: Local

Password rule: Standard

Multi-factor authentication type: No MFA

Roles

Partition_Fencing

Partition_Fencing_Objects

Groups

Tasks

Activate

Deactivate

Load

View Activation Profiles

Object Types

Objects

(Defined CPC)

(LPAR Image)

(LPAR Image)

(LPAR Image)

Quick tip:

❖ You can separate the Object Scope and Task permissions into two roles.

fence_ibmz – optional HMC SCSI configuration

Step 4.1 – Use HMC activate-on-load

- ❖ Set **load during activation** allows you to:
 - ❖ skip the additional load task (saves time)
 - ❖ can be found in the activation profile of each LPAR
 - ❖ works with any supported storage type
 - ❖ currently this option is **required** for SCSI usage

☒ Load during activation

Load type:

☐ Standard load
☒ SCSI load
☐ SCSI dump
☐ NVMe load
☐ NVMe dump

☐ Enable Secure Boot for Linux

Load address:

Load parameter:

Worldwide port name:

Logical unit number:

Boot program selector:

Boot record logical block address:

Operating system specific load parameters:

☐ Use dynamically changed address
☐ Use dynamically changed parameter

fence_ibmz – TLS CA Certificates

Step 4.2 – Install and Trust HMC TLS Certificate

- ❖ Get Root Certificate and all Intermediate CA Certificates used in the chain of the server certificate (in PEM format) and then trust them by executing:

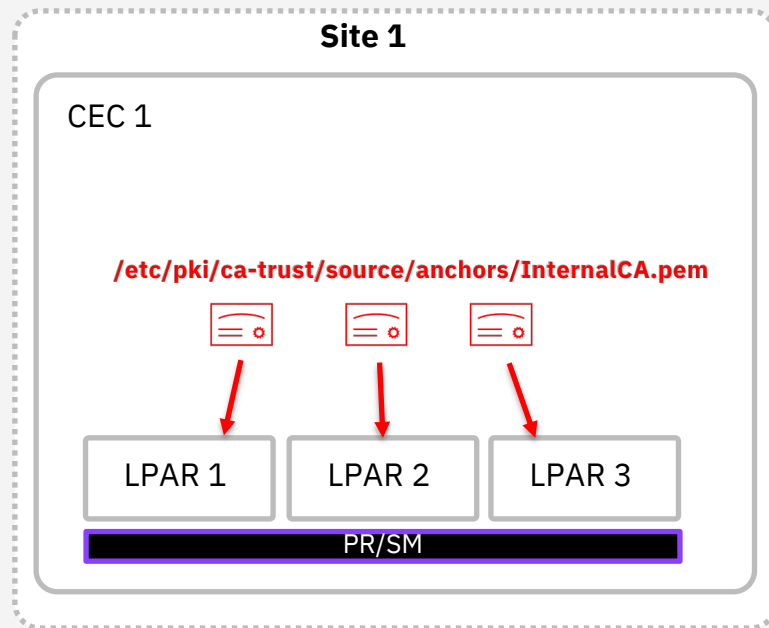
```
# cp CA_CERT.pem /etc/pki/ca-trust/source/anchors/  
# update-ca-trust
```

- ❖ Verify that the Certificates are in the trust store:

```
# trust list | less
```

- ❖ Verify that the whole certificate chain to the HMC is trusted:

```
# openssl s_client -showcerts -connect ${HMC_URL}:443 \  
-verify_return_error < /dev/null
```



fence_ibmz – Installation

Remark

- ❖ In newer RHEL versions fence_ibmz might be already installed with the fence-agents-all package

Step 4.3 – Setup fence_ibmz

❖ Installation via Upstream

```
# dnf install -y wget
# wget https://raw.githubusercontent.com/ClusterLabs/ \
fence-agents/master/agents/ibmz/fence_ibmz.py
# sed -i 's+@PYTHON@+/usr/libexec/platform-python+' \
fence_ibmz.py
# sed -i 's+@FENCEAGENTS_LIBDIR@+/usr/share/fence+' \
fence_ibmz.py
# cp fence_ibmz.py /usr/sbin/fence_ibmz
# chmod +x /usr/sbin/fence_ibmz
```

❖ Verify Installation

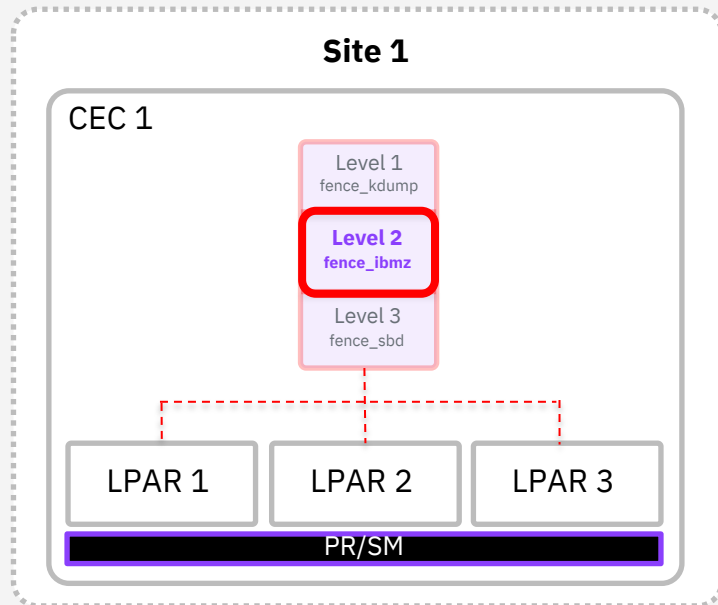
```
# pcs stonith list|grep fence_ibmz
```

Run on:

LPAR 1

LPAR 2

LPAR 3



fence_ibmz – Add to Cluster

Step 4.4 – Add fence_ibmz to cluster

Run on:

LPAR 1

❖ Add fence agent to cluster

```
# pcs stonith create fence_ibmz fence_ibmz \  
    ip=${HMC_URL} \  
    username="${HMC_USER}" \  
    password="${HMC_USER_PASSWORD}" \  
    ssl_secure=true \  
    pcmk_host_map="lpar1:CEC1/LPAR1;lpar2:CEC1/LPAR2;lpar3:CEC1/LPAR3"
```

Note:

- pcmk_host_map: the second value is case sensitive!
- It is possible to hide the password by providing a password_script: [LINK](#).

❖ Add debug log output for verification

```
# pcs stonith update fence_ibmz verbose=1 debug_file=/tmp/fence_ibmz.log
```

❖ Trigger fencing manually to verify the fence agent

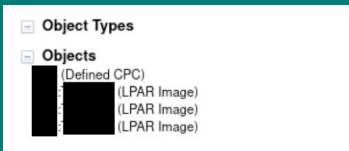
```
# pcs stonith fence lpar2  
# cat /tmp/fence_ibmz.log | less
```

Note:

- stonith-timeout and stonith-action options might be ignored when triggering manual fencing: [LINK](#).

Quick tip:

- ❖ Each pair consists of:
HOSTNAME:CECNAME/LPARNAME
- ❖ See HMC Objects:



fence_ibmz – Considerations

Step 4.5 – Further Considerations

Run on:

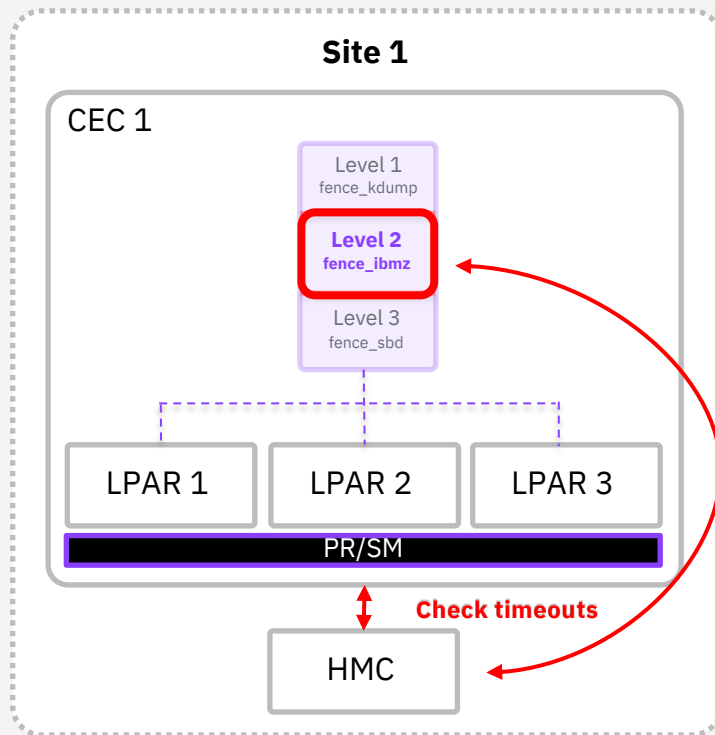
LPAR 1

- ❖ The stonith-timeout defines how long to wait for STONITH action (e.g. on, off) to complete. (default 60s)
 - ❖ Can be overwritten by pcmk_xxx_timeout on a fence agent basis
- ❖ Considering the deactivate operation on the HMC can take up to 900 seconds (by default), you can overwrite the STONITH action timeout for the fence agent:

```
# pcs stonith update fence_ibmz \  
    pcmk_reboot_timeout=1810 \  
    pcmk_off_timeout=905 \  
    pcmk_on_timeout=905
```

Note:

- pcmk_reboot_timeout should not be relevant here as the fence operation maps the reboot action to off and on internally
- When the system is at it limits the HMC task can take significantly longer



Fencing / STONITH fence_sbd (Level 3)

SBD Fencing – Watchdog

Step 5.0 – Setup Watchdog

- ❖ Enable watchdog kernel module

```
# modprobe diag288 wdt
```

- ❖ Show watchdog

```
# wdctl
```

- ❖ Make watchdog loading persistent

```
# echo "diag288_wdt" > /etc/modules-load.d/watchdog.conf
```

Note:

- Watchdog timeout cannot be lower than 15s! ([link](#))

Run on:

LPAR 1

LPAR 2

LPAR 3

Step 5.1 – Setup SBD Shared Storage

- ❖ Shared Storage is already assumed to be setup.

- ❖ Create SBD header on disk partition

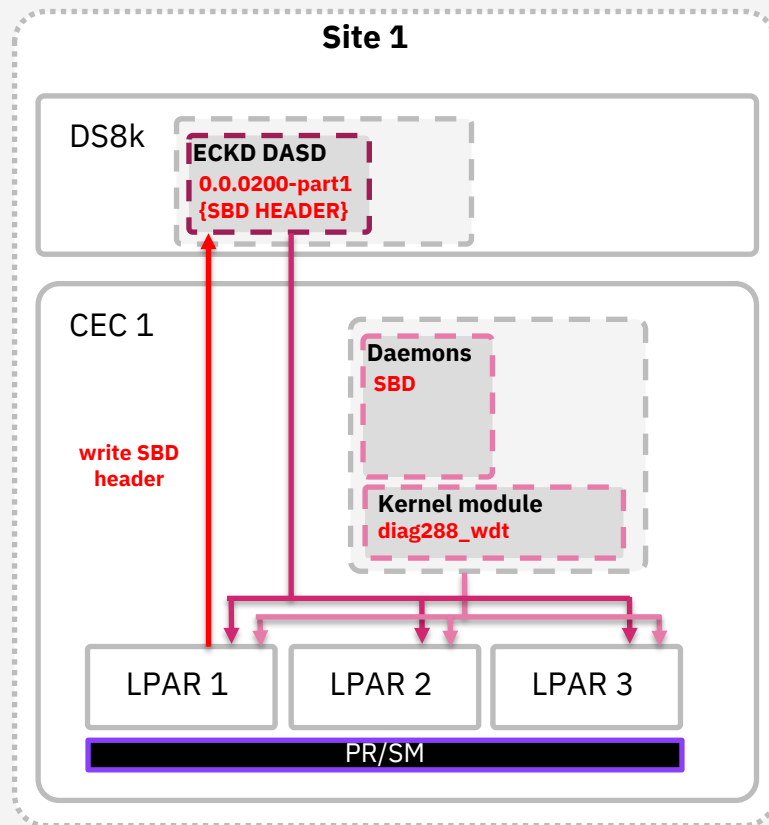
```
rhel18# pcs stonith sbd device setup \  
--device=/dev/disk/by-path/ccw-0.0.0200-part1 \  
watchdog-timeout=15 \  
msgwait-timeout=30
```

- ❖ Show SBD header

```
rhel18# pcs stonith sbd status --full
```

Run on:

LPAR 1



SBD Fencing – Daemon

Step 5.2 – Setup SBD

Run on:

LPAR 1

❖ Enable SBD systemd daemon in cluster

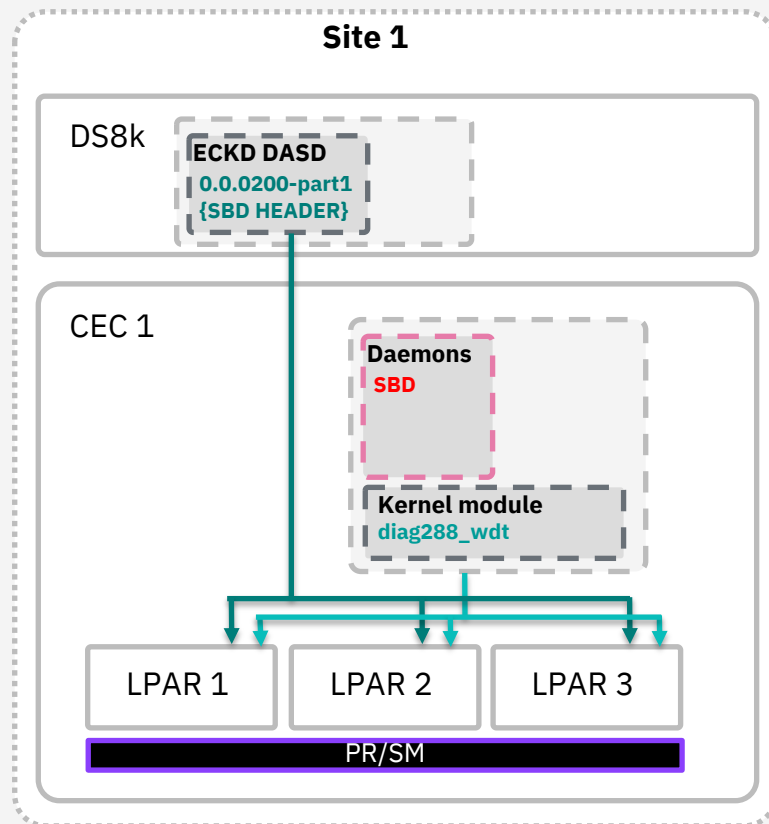
```
rhel18# pcs stonith sbd enable \  
watchdog=/dev/watchdog \  
device=/dev/disk/by-path/ccw-0.0.0200-part1 \  
SBD_DELAY_START=60 SBD_WATCHDOG_TIMEOUT=15
```

Note:

- SBD_* are environment variables for the SBD systemd service.
- SBD_WATCHDOG_TIMEOUT **only applies** when SBD runs in diskless mode.
 - > when disks are defined the watchdog timer written to the disk header is used.
- **The diag288 watchdog minimum timeout is 15 seconds. ([LINK](#))**
- SBD_DELAY_START postpones the start of the pacemaker systemd daemon
- SBD_DELAY_START should be longer then: corosync token timeout (5) + consensus timeout (6) + pcmk_delay_max (0) + msgwait (30) = 41 seconds. Otherwise, you might run into the issue that pacemaker starts with exit code 100.

❖ Restart cluster

```
rhel18# pcs cluster stop --all  
rhel18# pcs cluster start --all
```



SBD Fencing – Fence Agent

Step 5.3 – Setup SBD Fence Agent

Run on:

LPAR 1

- ❖ Show default power_timeout which indicates how long the fencing process waits before pacemaker must be up again after fencing

```
rhel18# fence_sbd -o metadata|grep -A 2 power_timeout
```

- ❖ Create SBD fence agent and set power_timeout

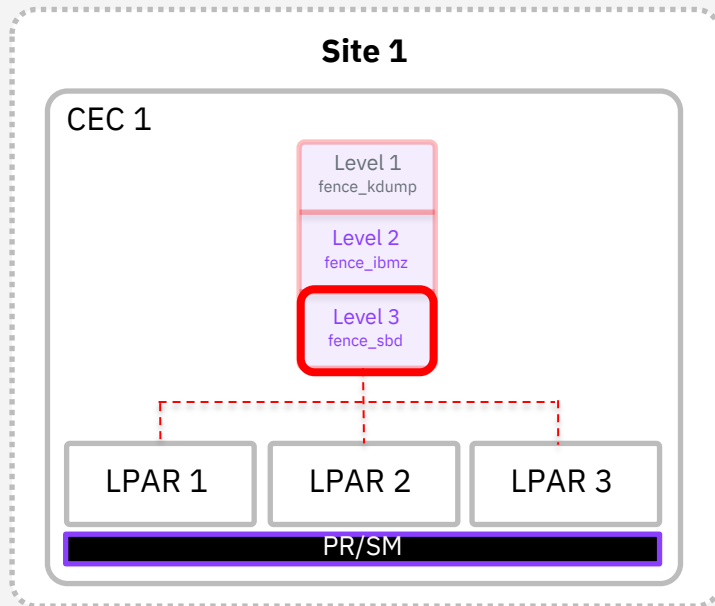
```
rhel18# pcs stonith create fence_sbd fence_sbd \  
    devices="/dev/disk/by-path/ccw-0.0.0200-part1" \  
    power_timeout=45
```

Note:

- [power_timeout should be bigger than msgwait timeout](#)

- ❖ Show settings of SBD fence agent

```
rhel18# pcs stonith config
```



SBD Fencing – Testing

Step 5.4 – Test SBD fencing

Run on:

LPAR 2

❖ Test sending of messages through A

```
rhel18# sbd -d /dev/disk/by-path/ccw-0.0.0200-part1 \  
message lpar1 test
```

❖ Look at the Log of the SBD systemd service

Run on:

LPAR 1

```
rhel18# journalctl -u sbd -f
```

❖ Send poison pill from lpar2 to lpar1

Run on:

LPAR 2

```
rhel18# pcs stonith disable fence_ibmz  
rhel18# time pcs stonith fence lpar1  
rhel18# pcs stonith enable fence_ibmz
```

Note:

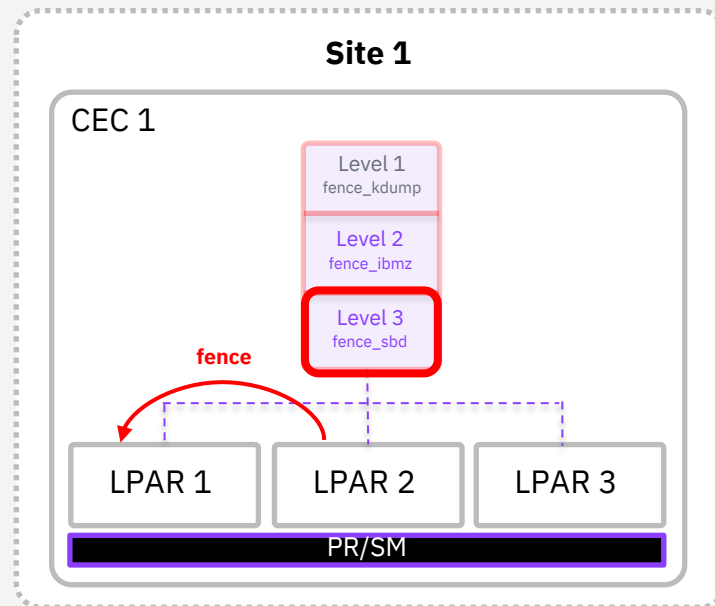
- other system should reboot and the pacemaker systemd service should be delayed by the SBD systemd service by msgwait-timeout seconds.

❖ Helpful debugging options

❖ Increase SBD verbosity: add "-v" to SBD_OPTS in
/etc/sysconfig/sbd

❖ Look at systemd startup

```
rhel18# systemd-analyze critical-chain
```



Fencing / STONITH fence_kdump (Level 1)

Fence_kdump

Step 6.0 – Configure kdump

❖ Setup kdump

```
rhel18# systemctl enable kdump --now
```

Run on:

LPAR 1

❖ Firewall rules

```
rhel18# firewall-cmd --add-port=7410/udp --permanent
rhel18# systemctl reload firewalld
rhel18# systemctl restart firewalld
```

Run on:

LPAR 1

LPAR 2

LPAR 3

❖ Edit kdump configuration

```
rhel18# vim /etc/kdump.conf
```

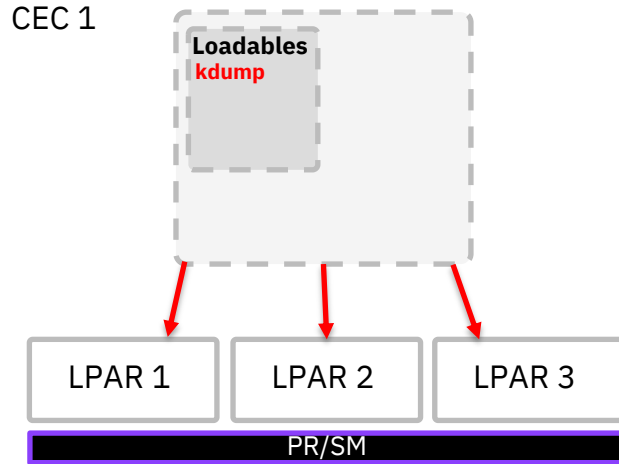
- Set Port to send the kdump message (see `man fence_kdump_send`)
Set interval to every 10 seconds (forever):
`fence_kdump_args -p 7410 -f auto -c 0 -i 10`
- all hostnames (cluster members) to send the kdump message to:
`fence_kdump_nodes lpar1 lpar2 lpar2`

❖ Restart kdump

```
rhel18# systemctl restart kdump
```

Site 1

CEC 1



Additional documentation

- ❖ Red Hat article - fence_kdump: [LINK](#)
- ❖ Solution Assurance guide - kdump: [LINK](#)

Fence_kdump – fence agent

Step 6.1 – Add kdump fence agent

Run on:

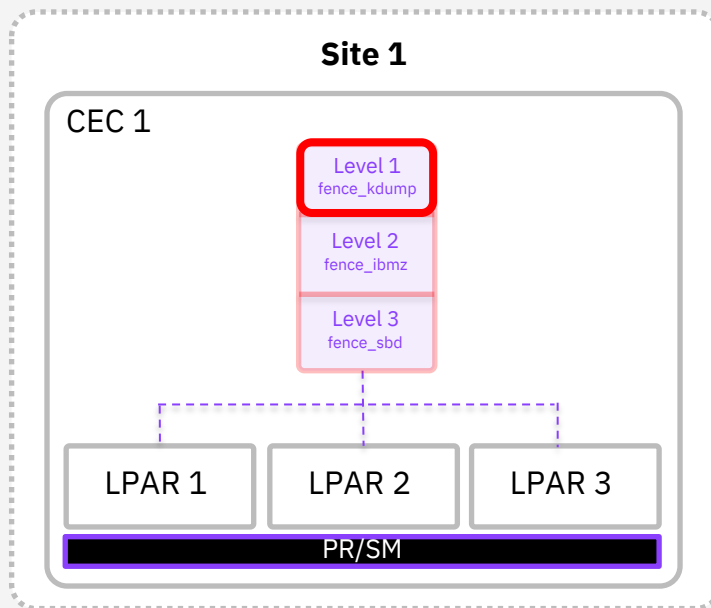
LPAR 1

❖ Create kdump fence agent

```
rhel18# pcs stonith create kdump fence_kdump \  
    pcmk_reboot_action="off" \  
    pcmk_host_list="lpar1 lpar2 lpar3" \  
    verbose=1
```

❖ Will be tested in the cluster testing chapter

❖ See following chapter to add fence levels to your cluster.



Fencing / STONITH

Add fencing levels

Fencing levels

Step 7.0 – fence-levels:

Run on:

LPAR 1

- ❖ To order the execution of fence agents, fence levels are introduced. They are executed from low to high.

```
rhel18# pcs stonith level add 1 regexp%lpar[0-9] fence_kdump
rhel18# pcs stonith level add 2 regexp%lpar[0-9] fence_ibmz
rhel18# pcs stonith level add 3 regexp%lpar[0-9] fence_sbd
```

- ❖ Show and verify fencing levels:

```
rhel18# pcs stonith level
rhel18# pcs stonith level verify
```

- ❖ When you want to remove levels:

```
rhel18# pcs stonith level remove 1
```

- ❖ Additional considerations:

- When your fencing is misconfigured, or the node has still a healthy cluster communication (e.g. when using fabric fencing) the node to be fenced is notified of its own fencing. In this case the fence-reaction property decides what happens. Panic is the safest choice and reboots the node.

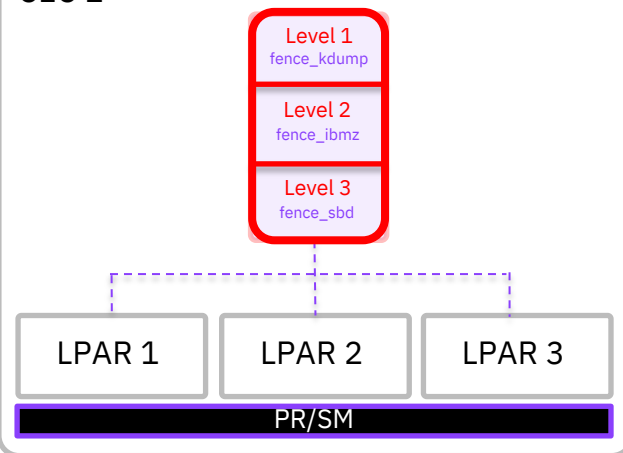
```
rhel18# pcs property set fence-reaction=panic
```

- To prevent multiple fencing operations in parallel you can disable concurrent-fencing. In a 3-node cluster that can only withstand the failure of one node we might not need concurrent-fencing:

```
rhel18# pcs property set concurrent-fencing=false
```

Site 1

CEC 1



GFS2

(Shared Storage)

GFS2 – Installation and Locking

Step 8.0 – [GFS2](#) Packages

```
rhel8# subscription-manager repos \
--enable=rhel-8-for-s390x-resilientstorage-rpms
rhel8# yum update -y
rhel8# yum install -y lvm2-lockd gfs2-utils dlm
```

Run on:

LPAR 1

LPAR 2

LPAR 3

Step 8.1 – GFS2 properties and locking

By default, the cluster stops all resources when the quorum is lost. The GFS2 resource cannot be stopped because it relies on the quorum. For this reason, the behavior must be changed to freeze all resources instead:

```
rhel8# pcs property set no-quorum-policy=freeze
```

DLM is used by **lvmlockd** for basic locking (read/write):

```
rhel8# pcs resource create dlm --group locking \
ocf:pacemaker:controld op monitor interval=30s \
on-fail=fence
rhel8# pcs resource clone locking interleave=true
```

Lvmlockd locks lvm metadata, validates caching of lvm metadata and prevents activation conflicts :

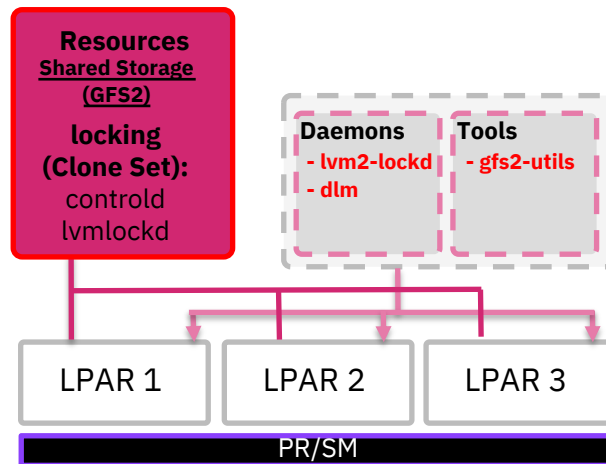
```
rhel8# pcs resource create lvmlockd --group locking \
ocf:heartbeat:lvmlockd op monitor interval=30s \
on-fail=fence
```

Run on:

LPAR 1

Site 1

CEC 1



GFS2 – Filesystem

Step 8.2 – Create GFS2 Filesystem

Run on:

LPAR 1

- ❖ Shared Storage is already assumed to be setup (e.g. with [chzdev](#) -e dasd 0.0.0201)

❖ Create VG, LV and make GFS2 filesystem

Run on:

LPAR 1

```
rhel18# vgcreate --shared shared_vg1 \  
/dev/disk/by-path/ccw-0.0.0201-part1
```

Note:

- Now the VG should be already visible on all LPARs

❖ Start the lock space on the other LPARs

Run on:

LPAR 2-3

```
rhel18# vgchange --lock-start shared_vg1
```

❖ Create shared Logical Volume

Run on:

LPAR 1

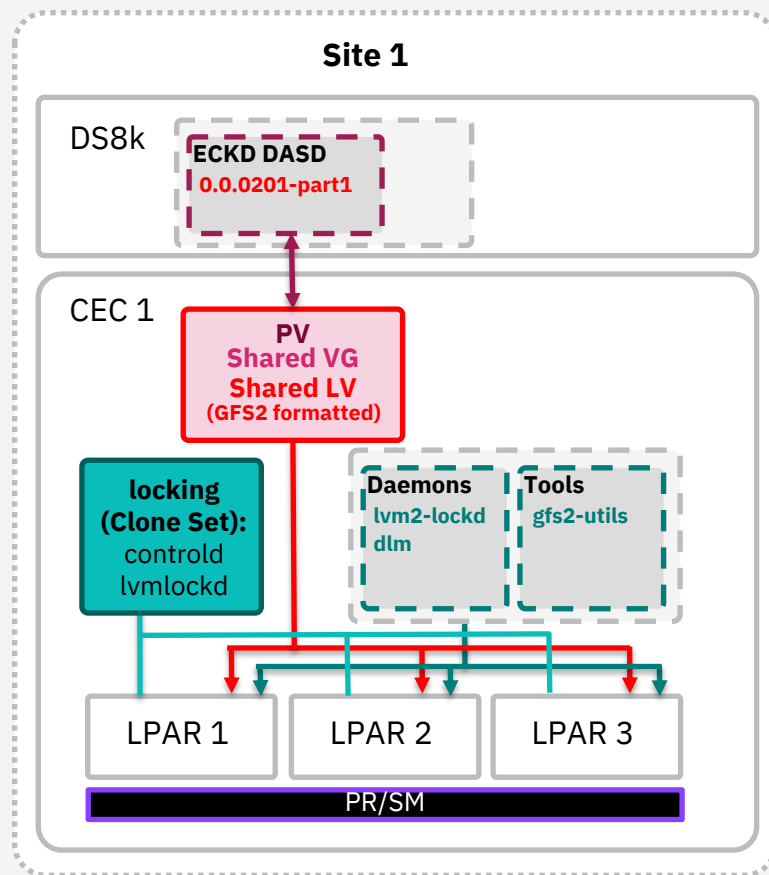
```
rhel18# lvcreate --activate sy -l 100%FREE \  
-n shared_lv1 shared_vg1
```

❖ Create GFS2 filesystem

```
rhel18# mkfs.gfs2 -j3 -p lock_dlm \  
-t my_cluster:gfs2-demo1 \  
/dev/shared_vg1/shared_lv1
```

Note:

- Make sure to create 3 journals (1 journal for each cluster member)



GFS2 – Create Resources

Step 8.3 – Add LVM Resource to Cluster

❖ Create LVM Resource and required constraints

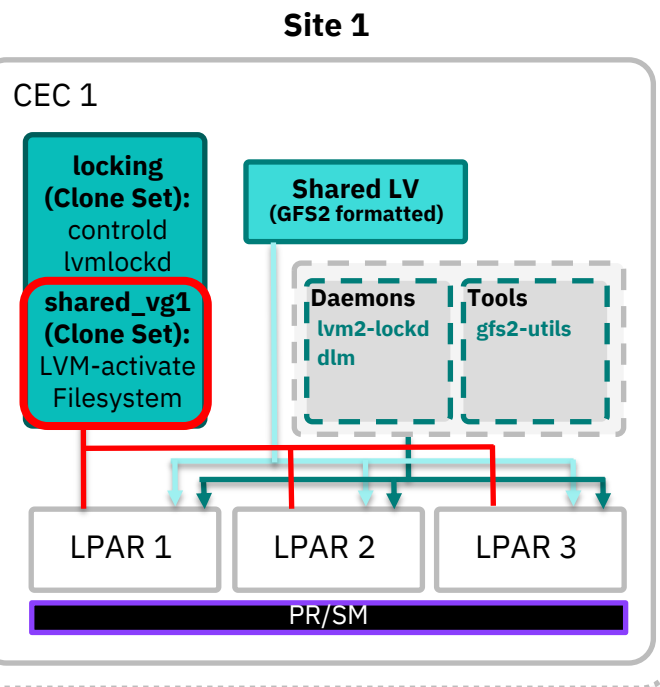
```
rhel18# pcs resource create sharedlv1 --group shared_vg1 \  
ocf:heartbeat:LVM-activate lvname=shared_lv1 \  
vgname=shared_vg1 activation_mode=shared \  
vg_access_mode=lvmlockd  
  
rhel18# pcs resource clone shared_vg1 interleave=true  
  
rhel18# pcs constraint order start locking-clone then shared_vg1-clone  
rhel18# pcs constraint colocation add shared_vg1-clone with locking-clone
```

❖ Create Filesystem Resource

```
rhel18# pcs resource create sharedfs1 --group shared_vg1 \  
ocf:heartbeat:Filesystem \  
device="/dev/shared_vg1/shared_lv1" \  
directory="/shared-fs-1" fstype="gfs2" \  
options=noatime,nodiratime,context=system_u:object_r:svirt_image_t:s0 \  
op monitor interval=10s on-fail=fence
```

Run on:

LPAR 1



VirtualDomain

VirtualDomain – Creation

Step 9.0 – Create KVM Guest

- ❖ Check the Redhat Virtualization documentation: [LINK](#).
When following commands passes you are ready to create KVM guests:

```
rhel18# virt-host-validate
```

- ❖ Create a minimal kvm guest. E.g.:

```
rhel18# virt-install \
  --name rhel8-guest1 \
  --memory 3000 \
  --vcpus 1 \
  --disk "size=4" \
  --location ${INSTALL_SERVER_URL}/s390x/RHEL8.5/DVD/ \
  --os-variant "rhel8.5" \
  --network "network=default" \
  --initrd-inject "/root/${YOUR_KICKSTART_FILE}.ks" \
  --extra-args "ks=file:///${YOUR_KICKSTART_FILE}.ks" \
  --noautoconsole
```

- ❖ Dump the guest configuration to the GFS2 shared directory

```
rhel18# virsh dumpxml rhel8-guest1 > /shared-fs-1/rhel8-guest1.xml
```

- ❖ Move the qcow image to the shared directory.

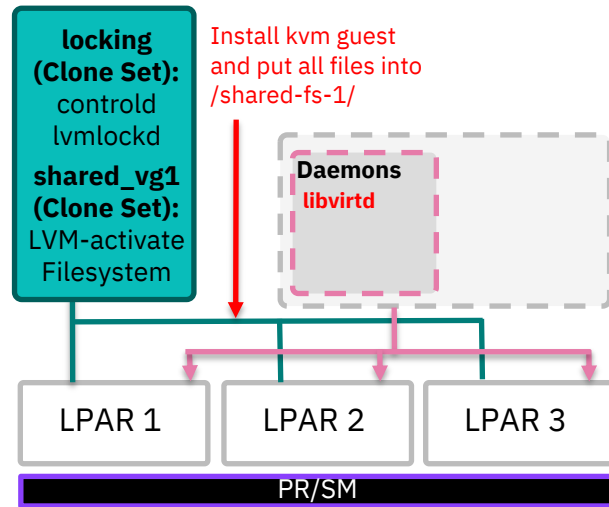
```
rhel18# mkdir /shared-fs-1/images
rhel18# virsh shutdown rhel8-guest1
rhel18# mv /var/lib/libvirt/images/* /shared-fs-1/images/
rhel18# restorecon -vr '/shared-fs-1'
```

Run on:

LPAR 1

Site 1

CEC 1



VirtualDomain – Add KVM Guest

Step 9.1 – Add VirtualDomain Resource

Run on:

LPAR 1

- ❖ Make sure the kvm hypervisor is reachable via ssh from all cluster members (without password prompting).

```
rhel8# virsh -c qemu+ssh://lpar2/system
```

- ❖ Add VirtualDomain Resource and required constraint:

```
rhel8# pcs resource create rhel8-guest1 \  
    ocf:heartbeat:VirtualDomain \  
    config="/shared-fs-1/rhel8-guest1.xml" \  
    hypervisor="qemu:///system" \  
    migration_transport="ssh" \  
    meta allow-migrate=true
```

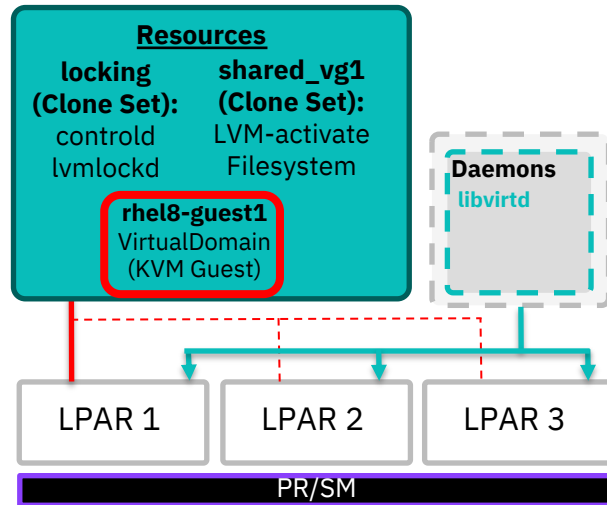
```
rhel8# pcs constraint order start shared_vg1-clone then rhel8-guest1
```

Note:

- The allow-migrate option allows live migration of the KVM guest when you move it manually.

Site 1

CEC 1



Cluster Testing

Cluster Testing

Step 10.0 – Test Workload Failover

❖ Watch the cluster status (live)

```
# watch -n1 pcs cluster status
```

Run on:

LPAR 2

❖ Trigger a kernel panic

```
# echo c > /proc/sysrq-trigger
```

Note:

- you should see a kdump written to the /var/crash directory

Run on:

LPAR 1

❖ Check systemd logs (errors, warnings...)

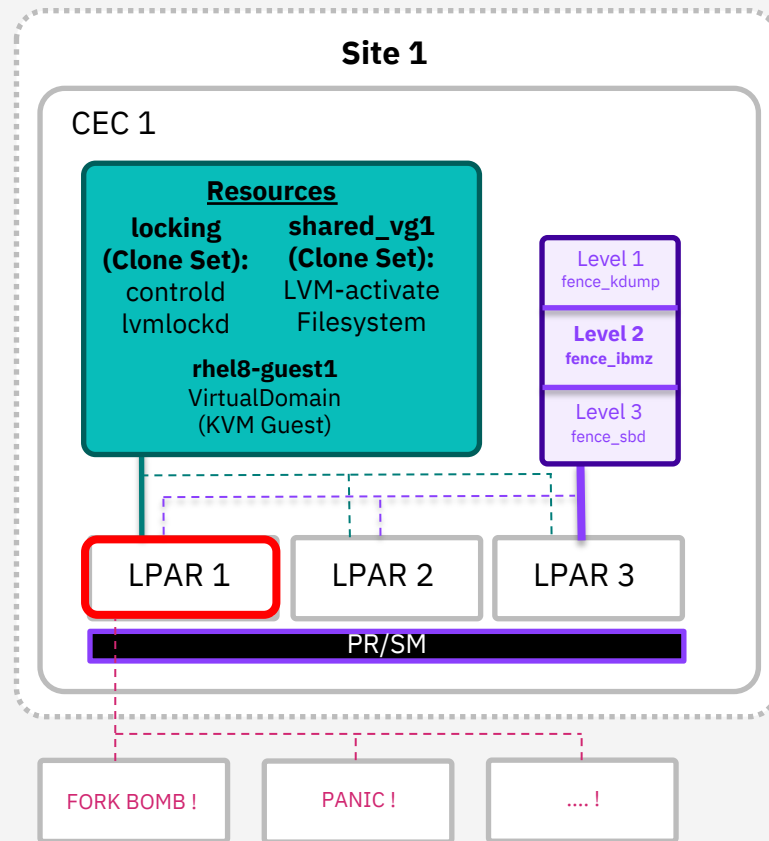
```
# journalctl -u corosync  
# journalctl -u pacemaker  
# journalctl -u sbd
```

❖ Fork Bomb:

```
# :() { :|:& }; :
```

Note:

- It might take some time for the cluster to hang
- The deactivate step in the fence_ibmz fence agent might take longer



Appendix A – Shared Storage – z/VM Shared Storage

Create Fullpack Minidisk in z/VM

This might be required to be executed multiple times when setting up shared storage in z/VM.

❖ Create new user LINSHARE

```
USER LINSHARE NOLOG  
MDISK 0200 3390 DEVNO 1111 MWV
```

Run on:

z/VM

❖ Link Minidisks

```
# vmcp 'link * 0200 0200 rw'
```

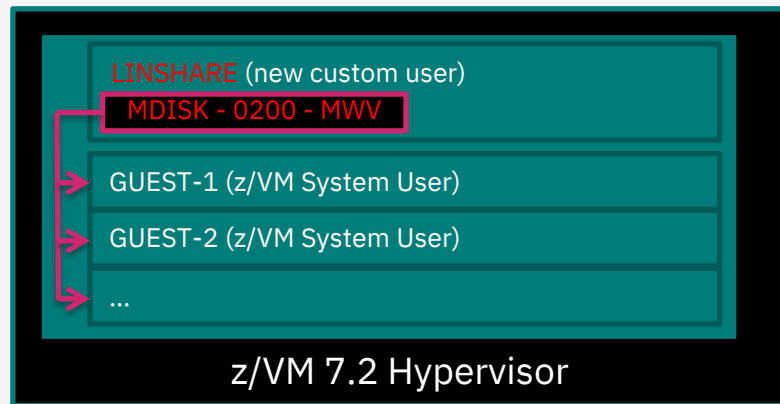
every
guest
in
cluster

❖ De-ignore, enable and make DASD persistent:

```
# chzdev -e dasd 0.0.0200
```

❖ Format dasd and create partition over whole dasd

```
# dasdfmt -b 4096 -d cdl \  
-p /dev/disk/by-path/ccw-0.0.0200  
# fdasd -a /dev/disk/by-path/ccw-0.0.0200
```



Appendix B – Introduction – High Availability Stack

	Layers	HA related examples
Workload, Automation, Orchestration	Applications / DBs & more	Oracle RAC
Operating Sys.	Linux	RHEL HA (Pacemaker + Corosync)
Virtualization	z/VM / KVM / PR/SM	z/VM SSI LGR KVM Live Migration
Networking	Networking	Multipath (path redundancy) Copy/Mirror Metro/Global Mirror IBM Hyperswap IBM FlashCopy
Storage	DS8k	
Physical	IBM LinuxONE	<ul style="list-style-type: none">- Redundant Power Supplies + Battery- Memory/Processor sparing- IBM System Recovery Boost (fixed-duration performance boost)