# How to handle special HA cases in RHEV

## PREPARED FOR - Volkswagen IT Group Cloud

Ben Haubeck

Version 1.1

**Table of Contents**

## 1. History and Revisions

| Version | Date | Authors | Changes |
|---------|------|---------|---------|
| 0.1 | 19/11/2015 | Ben Haubeck bhaubeck@redhat.com | Initial Draft |
| 0.2 | 20/11/2015 | Ben Haubeck bhaubeck@redhat.com | added details for RHEV version and environment |
| 0.9 | 20/11/2015 | Ben Haubeck bhaubeck@redhat.com | added feedback from Adrian Bradshaw and Martin Tessun |
| 0.9.1 | 20/11/2015 | Ben Haubeck bhaubeck@redhat.com | added feedback Martin Tessun |
| 1.0 | 20/11/2015 | Ben Haubeck bhaubeck@redhat.com | first version, ready to hand over to the customer |
| 1.1 | 20/11/2015 | Ben Haubeck bhaubeck@redhat.com | added feedback from Eike Holtz, clarified scenario 3.2 |

## 2. Preface

### 2.1. Confidentiality, Copyright, and Disclaimer

This is a Customer-facing document between Red Hat, Inc. and Volkswagen ("Client"). Copyright © 2015 Red Hat, Inc. All Rights Reserved. No part of the work covered by the copyright herein may be reproduced or used in any form or by any means – graphic, electronic, or mechanical, including photocopying, recording, taping, or information storage and retrieval systems – without permission in writing from Red Hat except as is required to share this information as provided with the aforementioned confidential parties.

### 2.2. About This Document

Double (or even more) failures at the same time or a disaster can lead to VMs and / or hosts in a state in that RHEV cannot decide on its own, how the situation can be definitely solved without causing the risk of data corruption. The purpose of the document is to describe the procedures that might be necessary to solve such a situation in which RHEV could not decide on its own, which action(s) has to be taken to prevent data corruption. Single failures are covered by automatic procedures. Single failures are the loss of 1 connectivity or the failure of 1 component (i. e. 1 power supply unit).
This document is focusing on the RHEV installation at Volkswagen with RHEV in version 3.5.4 and RHEL 7.1 as hypvervisors and with its configuration that is described in detail in the document "rhevm-setup-and-configuration.pdf".

### 2.3. Audience

The document is intended for those team members on site at Volkswagen, who will be responsible for RHEV support and administration.

### 2.4. Additional Background and Related Documents

Numerous other documents have also been provided by Red Hat Consulting, explaining tasks such as installation and configuration of RHEV, backup etc.

### 2.5. Terminology

Some of the acronyms using in this document are included in the table below

*Table 1. Terminology Table*

| Term | Definition |
|------|-----------|
| RHEV | Red Hat Enterprise Virtualisation |
| RHEV-M | Red Hat Enterprise Virtualisation Manager |
| RHEL-H | Red Hat Enterprise Linux Hypervisor |
| iRMC | integrated Remote Management Board (Fujtsu Server) |

## 3. Description of procedures to solve HA related issues after two or more concurrent failures

The following scenarios are related to multiple simultaneous failures. IE loss of multiple/all SAN paths, loss of all power, loss of both pairs in a bond etc. In these cases, some manual intervention is required to bring back affected VMs.

Red Hat's absolute priority is to guarantee the consistency of the data: Everything else can be relatively easily corrected but a data corruption is generally impossible to correct and a recovery can only rely on a possibly outdated backup. This is the reason why, when in doubt, manual intervention will be required if the state of the VMs can't be guaranteed, either through status query or through fencing (i.e. killing the server to make sure it can't concurrently write data and hence corrupt it).

### 3.1. Complete power loss for server and remote management board

**Scenario**

All four power supply sources lost at the same time, server & iRMC without any power, so both ssh soft fencing and iRMC fencing are both not available

*Symptoms*

The host is marked as "non-responsive" in the UI and RHEV-M is trying to fence the host after a grace period. In this scenario this is not possible anymore as neither a command via ssh is possible nor iRMC can be used to power off the host or to determine its state.

*Reaction by RHEV*

Tries to fence the host and determine its state. No VMs are getting restarted, that were formerly running on the host in the unknown state.

*Reason for this behaviour*

As the state of the host cannot be determined RHEV cannot ensure, that the VMs are down and are not writing to their disk.

*Solution*

Click on "Confirm host has been rebooted" in the RHEV UI.
VMs, that marked as HA will be started on other hosts automatically. The other VMs are changing their state to "down" and can be started manually. (It is also possible to configure RHEV, that all VMs that were running, were started automatically, but Volkswagen has decided to configure the cluster that way, that only VMs, that marked as "HA" are automatically restarted.)

Additionally it would be possible to automatically confirm "Host has been rebooted" after some grace period. But in that case data loss can occur and is within the responsibility of Volkswagen. If needed / wanted this can be implemented by Red Hat. We need the confirmation of Volkswagen for accepting the additional risk of data corruption before implementing this.

### 3.2. Loss of 50% or more of all hosts in one cluster or splitt-brain

**Scenario**

Catastrophic loss of entire data center or catastrophic loss of connection to multiple hosts - resulting in 50% of all hosts unreachable and therefore all fencing options unavailable.

*Symptoms*

After a grace period at least 50% of the hosts of a cluster are marked as "down" or "non-responsive".

*Reaction by RHEV*

No fencing actions and no automatic migrations are triggered.

*Reasons for this behaviour*

RHEV cannot exclude a split brain scenario and in that case restarting VMs can lead to severe data corruption.

*Solution*

If hosts in state "non-responsive" see scenario 3.1.
If hosts in state "down":
If the customer decides to work on further with half the capacity while the outage is still ongoing: Temporarily change the cluster policy for fencing by enabling "Skip fencing on cluster connectivity issues - Threshold: 50%". (Remember to revert this once the DC has returned to full operation.) The policies can be changed in the configuration of the cluster settings ("Edit Cluster") and there in the tab "Fencing Policy".
To prevent this kind of problem in general: Distribute the hosts across at least three data centers.

## 3.3. Complete loss of frontend ("Nutznetz") network connectivity

**Scenario**

Loss of all network connections in the network bond that contains the "Nutznetz".

*Symptoms*

RHEV is recognizing the bonded interface as "down", VMs remain marked as "up". No automatic migrations are triggered.

*Reasons for this behaviour*

Simultaneous outage of two cables (double failure) are not covered by automatic procedures.

*Solution*

This event should cause an alarm event in the monitoring solution. The VMs can easily migrated to hosts with working frontend network. To solve this kind of scenario the monitoring solution can be configured to trigger the migration of VMs that are unreachable via the RHEV-M API.

## 3.4. Complete loss of SAN connectivity

**Scenario**

Loss of all **8** SAN paths simultaneously

*Symptoms*

All VMs on affected hypervisor are changing in state "paused". No fencing and no automatic migrations are triggered.

*Reasons for this behaviour*

Restart of VMs can lead to data corruption, re-establish of SAN connectivity will reactivate the VMs automatically.

*Solution*

It is still being discussed at Volkswagen if the VMs should be restarted elsewhere and the DB recovered or if it is better to wait until the SAN connectivity is re-established. Hence RHEV offers two solutions:

- Use the power management in the RHEV UI to restart the host. VMs, that marked as HA will be started on other hosts

automatically. The other VMs are changing their state to "down" and can be started manually. (It is also possible to configure RHEV, that all VMs that were running, were started automatically, but Volkswagen has decided to configure the cluster that way, that only VMs, that marked as "HA" are automatically restarted.)
This can lead to data corruption, the DB needs recovering after restart of the VM on another host.

or

- Waiting until SAN connectivity is re-established and VMs resume automatically.