



Metro/Stretched Storage Clustering Best Practices

with Virtual Volumes (vVols)
on VMware vSphere 8

November 20, 2024

Table of Contents

Disclaimer.....	4
Overview.....	5
Intended Audience.....	5
Introduction.....	5
Technical Requirements and Constraints.....	6
Virtual Volumes, VASA, and Stretched Storage Clustering.....	6
Active/Active vs Active/Passive Access to vVols.....	8
Uniform Versus Non-Uniform Configurations.....	8
Site Affinity of Virtual Machines and Datastores.....	10
Infrastructure Architecture.....	10
vSphere Configuration.....	12
vSphere High Availability.....	12
Admission Control.....	12
vSphere HA Heartbeats.....	14
Permanent Device Loss and All Paths Down Scenarios.....	17
Restart Ordering and Priority.....	20
Proactive HA.....	23
Distributed Resource Scheduler (DRS).....	24
Site Affinity.....	24
Advanced Settings.....	26
Correcting Affinity Rule Violations.....	27
Storage DRS.....	27
Known Issues and Considerations.....	28
Failure Scenarios.....	28
Single Host Failure in Frimley Data Center.....	29
Single Host Isolation in Frimley Data Center.....	30
Storage Partition.....	31
Data Center Partition.....	32
Disk Shelf Failure in Frimley Data Center.....	34
Full Storage Failure in Frimley Data Center.....	35
Permanent Device Loss.....	36
Full Compute Failure in Frimley Data Center.....	38
Loss of Frimley Data Center.....	39

Loss of VASA Providers in Frimley Data Center	40
Loss of VASA Providers in Both Data Centers	41
Summary & Conclusion.....	42

Disclaimer

This document is intended to provide general guidance for organizations that are considering Broadcom solutions. The information contained in this document is for educational and informational purposes only. This document is not intended to provide advice and is provided “AS IS.” Broadcom makes no claims, promises, or guarantees about the accuracy, completeness, or adequacy of the information contained herein. Organizations should engage appropriate legal, business, technical, and audit expertise within their specific organization for review of requirements and effectiveness of implementations.

Overview

VMware Metro Storage Cluster (vMSC) is a specific storage configuration that is commonly referred to as Stretched Storage Cluster (SSC) or Metro Storage Cluster (MSC). This applies to both Virtual Volume (vVol) and non-vVol implementations. These configurations are typically implemented in environments where disaster and downtime avoidance is a key requirement.

This recommended best practice document provides insight and information for operating an SSC infrastructure on vVol. We use the term SSC for vVol to distinguish it from both Virtual Machine File System (VMFS) MSC and another vVol feature that uses the Metro Storage term to achieve efficient (replication-assisted) vMotions. For vVol stretched storage, we will use the term vVol Stretched Storage Clustering (vVol SSC) to avoid confusion with other features.

This paper explains how VMware vSphere handles specific failure scenarios and discusses various design considerations and operational procedures that apply to vVol, which differs slightly from VMFS MSC operations. However, the solutions for vVol are largely similar to VMFS MSC operations, which is intentional as it helps prevent new pitfalls and allows VMFS MSC users to feel comfortable with vVol SSC once they understand the vVol-specific nuances.

Unlike VMFS MSC, the vVol solution is supported by both VMware and the partner through a certification program. When customers use a combination of VMware products (vCenter, ESXi, High Availability/VM Component Protection, Distributed Resource Scheduler, and other components) with an array certified for vVol SSC on the given VMware release, VMware support will accept all incidents. VMware will work to resolve issues either internally (if identified as VMware-specific) or in collaboration with the storage partner (if determined to be storage partner-specific or requiring fixes from both parties). This differs from VMFS MSC architecture, which is entirely partner-verified and supported.

Intended Audience

This document is intended for technical professionals who design, deploy, or manage a vVol Stretched Storage Cluster infrastructure. The intended audience includes, but is not limited to, technical consultants, infrastructure architects, IT managers, implementation engineers, partner engineers, sales engineers, and operations staff.

Readers should be familiar with vSphere, vCenter Server, vSphere HA, vSphere DRS, vSphere Storage DRS, and replication and storage clustering technologies and terminology.

Introduction

A VMware Metro/Stretched Storage Cluster configuration relies on array-based Metro/Stretched Storage Clustering (SSC) to stretch datastores across arrays. SSC solutions are typically deployed in environments where the distance between data centers is limited, such as metropolitan or campus environments—hence the term "metro" in the name.

VMware MSC infrastructures extend the benefits of vSphere HA clusters from a local site to a geographically dispersed model with two data centers in different locations. A VMware SSC infrastructure is essentially a stretched cluster. The architecture extends the definition of "local" for network, storage, and compute resources, enabling these subsystems to span geographical locations while presenting a single, common infrastructure to the vSphere cluster at both sites.

The primary benefit of a stretched cluster model is enabling fully active and workload-balanced data centers that can be used to their full potential while allowing for extremely fast recovery during host or site failures. The ability to provide active resource balancing should be the primary design and implementation goal. Although often associated with disaster recovery, stretched cluster infrastructures are not recommended as primary solutions for pure disaster recovery.

This document does not address the distinctions between disaster recovery and downtime or disaster avoidance solutions.

Stretched cluster solutions offer the following benefits:

- Workload mobility
- Cross-site automated load balancing
- Enhanced downtime avoidance

- Disaster avoidance
- Fast recovery

Technical Requirements and Constraints

Before implementing a stretched cluster, the following technical requirements must be met:

- Storage connectivity using SCSI Fibre Channel and iSCSI is supported with qualified vendors. Currently, the solution is "SCSI only," though support for additional protocols and transports is planned.
- The maximum supported network latency between sites for ESXi management networks is 11ms round-trip time (RTT).
- The maximum supported latency for synchronous storage replication links is 11ms RTT. However, storage vendor documentation should be consulted, as vendor-specific maximum tolerated latency may be lower. Users should implement the lower value between the vSphere limit (11ms) and the storage vendor's specified limit.
- The vMotion network requires 250 Mbps of dedicated bandwidth per concurrent vMotion session.
- Fault Tolerance (FT) is not supported with vVol Stretched Storage Cluster.
- vCenter HA is not yet supported with Virtual Volumes, limiting deployment to a single vCenter instance. The vCenter appliance should not be deployed on vVol due to circular dependencies, but it can be placed on VMFS MSC-based storage.
- Storage DRS and Storage IO Control features are not supported on a vVol datastore.
- While replication to and from stretchable containers should technically function, no current vendor supports this capability, and it remains untested.

The above constraints generally apply to both VMFS and vVol stretched storage implementations. The minimum write latency for a stretched storage object is always at least one RTT, though some arrays may require more. For reads, some arrays allow I/O operations from either array, while others require one array to handle all I/O operations at a given time, resulting in similar latency for reads and writes as RTT typically dominates I/O latency on modern storage arrays.

The 11ms RTT represents the minimum additional latency for any cross-site I/O or VASA operation. Multiple RTTs may be required for certain operations. For example, when a host on site A writes to a controller on the array at site B, synchronization back to the array on site A requires a minimum of two RTTs on top of existing latency.

Other I/O operations (such as Compare And Write) may require additional round-trips between arrays, though the expected range is typically 1-5 round-trips for any I/O or VASA (control plane) operation for vVol. This creates an upper bound of approximately 60ms for network latency alone. With additional processing time, most operations remain minimally affected since ESXi generally operates with acceptable latencies measured in seconds, even for the shortest I/O operations like VMFS Heartbeats.

Virtual Volumes, VASA, and Stretched Storage Clustering

We use the term vVol Stretched Storage Cluster (SSC) in this document when discussing the vVol MSC solution to avoid confusion with both VMFS MSC and the "virtual Metro Cluster" feature for vVol (which covers replication-assisted vMotion). Note that vVol SSC support begins with vSphere 8 Update 3 on the host side.

The storage requirements for vVol are more complex. To explain these details, we first need to provide a high-level overview of vVol operation.

A storage array supporting vVol requires a VASA Provider (VP) that vCenter and ESXi hosts can connect to for various control path operations (through the VASA API, which exists in different versions as features are added). The VASA Provider is essentially a web server that serves special protocols for vVol, using a client/server scheme to process requests from

vCenter/ESXi hosts. These requests handle vVol object creation/destruction, vVol Container (Datastore) management, and other functions to provide an abstract datastore of vVol objects.

Earlier vSphere releases managed VASA Provider high availability through vCenter. Since vCenter is also a single point of failure, the HA scheme changed in VASA 6 specification (where vVol SSC was introduced) to resemble an active/active I/O path scheme with priorities.

The primary requirement for vVol SSC best practices is VASA provider high availability for all stretched containers and any un-stretched (legacy) containers, both appearing as datastores in ESXi. For a vVol datastore to function, an ESXi host needs both datapath (FC, iSCSI) and control path availability, making VP HA crucial for SSC.

Many VASA operations work with vVol objects, which are SCSI LUNs in SCSI-based vVol arrays (currently the only supported implementation). vVol enables virtual disks and related components to be array-known objects (a LUN in SCSI, an NVMe Namespace for NVMe, or a file for NFS) that ESXi hosts can manipulate through the VASA API, leveraging array capabilities for operations like snapshots or clones.

vVol object types include:

- CONFIG vVols: VM home folders (VMFS-formatted LUNs containing VM configuration files and virtual disk descriptors)
- DATA vVols: Actual virtual disks
- SWAP vVols: ESXi VM swap space
- MEM vVols: Memory snapshot components
- OTHER vVols: Additional uses

In vVol, the array/VP represents a datastore as an abstract 'container' holding vVol objects. ESXi hosts use VASA APIs to query containers for vVol objects based on specific constraints, building datastores with CONFIG vVol mount points as needed.

VASA 5 specification (released with vSphere 8 Update 1) added support for multiple containers per array and container isolation to specific vCenter instances and their hosts. vVol SSC in VASA 6 (vSphere 8 Update 3) builds on this foundation.

For SCSI-based vVol, each vVol is a LUN. Due to potentially large numbers of vVol objects, they remain hidden until bound to a Protocol Endpoint (PE) for host I/O use. The PE appears as a visible logical storage device, acting as a proxy in the ESXi storage stack. At the driver layer, I/O redirects to the actual vVol object.

vVol SSC requires that containers/datastores span array pairs with continuous mirroring of all vVol objects. PEs indicate I/O direction, and if mirroring fails, only one array can continue exposing the container, PEs, and associated vVol objects.

This design requires simultaneous read/write datastore access from hosts at both sites, with transparent continuity during problems. While active/passive storage can work in some configurations, active/active storage is recommended (including arrays with optimized and non-optimized paths, excluding paths requiring activation).

Traditional synchronous replication solutions are incompatible due to their primary-secondary relationships. SSC requires uninterrupted simultaneous access for live VM migration between sites.

The SSC stretched datastore storage subsystem must maintain synchronous consistency between locations (sync mirrored or instant mirrored). This ensures data consistency regardless of read location and enables immediate workload continuation from the surviving site during failures.

This architecture requires significant bandwidth and low latency between cluster sites. Increased distances or latencies can degrade performance and may prevent successful vMotion migration between cluster nodes in different locations.

Active/Active vs Active/Passive Access to vVols

Similar to storage arrays with dual controllers—where a LUN is actively accessed through one controller and remains passive/standby on another—this concept applies across arrays mirroring a LUN (or vVol object, which is a LUN in SCSI-based vVol implementations).

In stretched storage configurations, one array may have one or more active controllers for a LUN while the other array exposes only standby controllers for the same LUN. This active/passive LUN access model applies per LUN backing a VMFS volume in VMFS MSC implementations.

vVol implementations differ because vVols are not directly accessible or visible in the storage stack, but are accessed through Protocol Endpoints in SCSI-based vVol. Individual vVols are bound to a Protocol Endpoint, and the PE and all bound vVols must share the same LUN accessibility state (path state). For vVol SSC, the PE is stretched (exposed by both arrays), and vVol path states are determined by the PE to which a vVol is bound.

VMware recommends active/active configurations because SAN cross-link failures in active/passive setups can cause PEs to alternate between active and passive states on each array, severely degrading performance.

Uniform Versus Non-Uniform Configurations

SSC solutions are classified into two distinct types based on how hosts access storage. Understanding these types is crucial as they influence design considerations:

- **Uniform host access configuration:** vSphere hosts from both sites connect to storage nodes in the storage cluster across all sites. Paths presented to vSphere hosts come from both local and remote arrays/sites.
- **Non-uniform host access configuration:** vSphere hosts at each site connect only to storage nodes at the same site. Paths presented to vSphere hosts from storage nodes are limited to the local array/site.

In uniform host access configuration, hosts in data center A have access to the storage systems in data center B. This effectively stretches the Storage Area Network (SAN) between sites, allowing all hosts to access either array's storage when no faults exist.

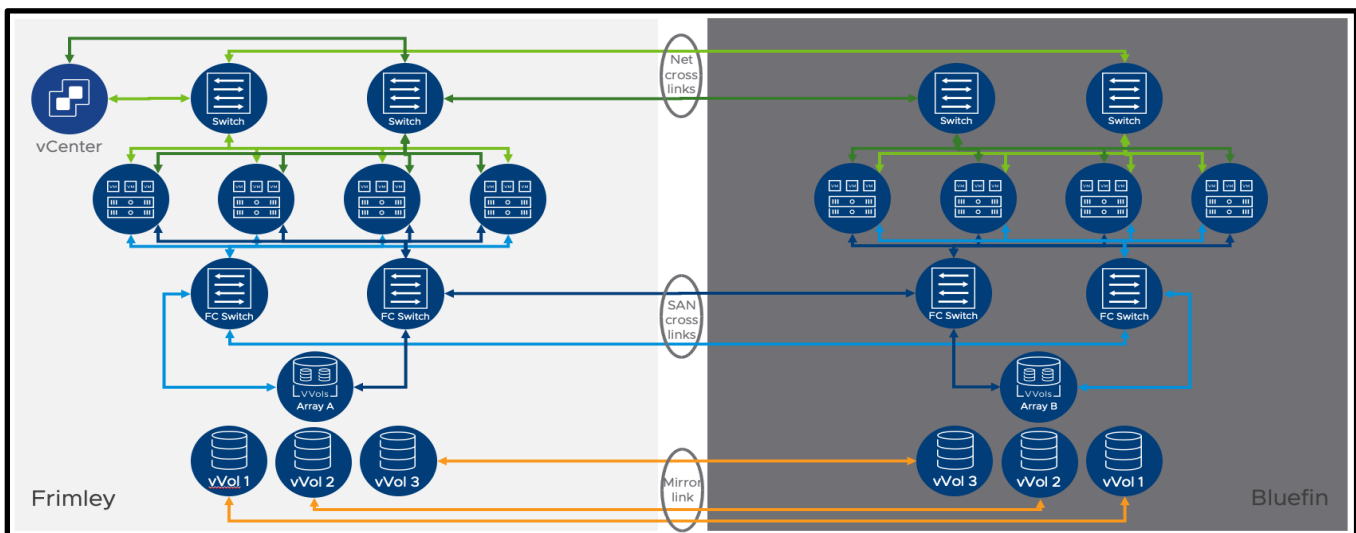


Figure 1 - Uniform Configuration

In non-uniform host access configuration, hosts in data center A access only the array within their local data center. Each array, along with its peer in the opposite data center, is responsible for providing datastore access in one data center. The arrays maintain sync mirroring for all writes to the datastore, preserving write capture on both arrays.

Non-uniform configurations differ from uniform configurations in their error tolerance scenarios. If mirror links between arrays break (or mirroring fails for any reason), non-uniform configurations lose storage container (datastore) access from one site. Uniform configurations can continue serving storage to both sites through SAN cross-links. However, if SAN cross-links also fail, the uniform setup degrades to non-uniform behavior, and both configurations exhibit similar failure characteristics.

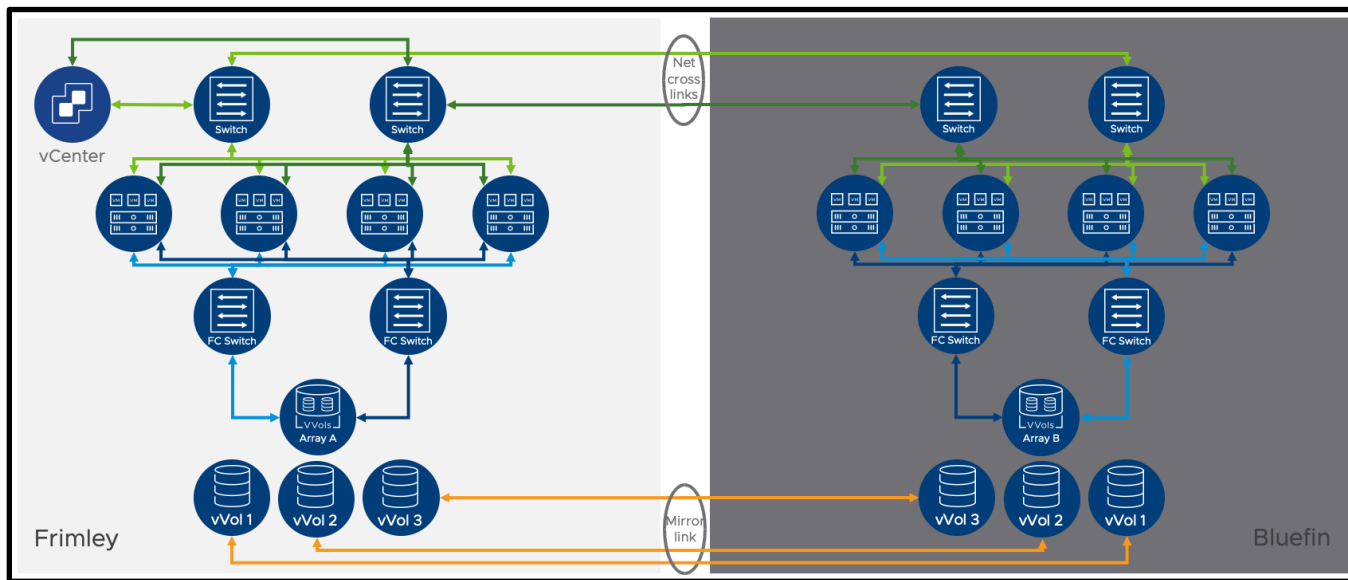


Figure 2 - Non-Uniform Configuration

The initial release of vVol SSC is in vSphere 8 Update 3, which supports only uniform configuration. While non-uniform configurations should work technically, support for this configuration will follow in a later release. The design considerations for uniform configurations apply to both types, and we note exceptions where they exist.

Note that arrays can expose both stretched and un-stretched (legacy) containers to hosts. While this is permitted, VM Component Protection (VMCP) applies only to stretched containers, not legacy containers.

In the figures throughout this document, the infrastructure layers are represented by different colors:

- Green links represent the vSphere networking infrastructure layer, with vCenter shown in the leftmost part of this layer at the Frimley site
- Blue links represent the SAN infrastructure layer, showing missing connectivity to the alternate site array in non-uniform configurations
- Orange links represent the array-to-array mirroring network layer

In some uniform setups, the mirroring network may combine with SAN cross-links. However, we show them separately as they can operate independently, allowing dedicated bandwidth for mirroring.

Regarding configuration, creating a stretched container is array/vendor specific. A stretched container may not always be stretched to another array at a given time. The container continues to appear as stretched in the vSphere environment, with its current state reported by the array. VMCP and other SSC features for vVol remain applicable, though data and VP redundancy would have reduced I/O paths and/or VP paths in this case.

These containers are sometimes called "stretch-capable" or "stretchable" as they maintain the properties enabling VMCP functionality. This makes them valuable and able to become fully stretched at any time through the array management interface, using array/vendor-specific methods.

Site Affinity of Virtual Machines and Datastores

vCenter and ESXi use the concept of "host affinity." VM site affinity can be achieved by establishing affinity between a VM and all hosts within a site. Similarly, a stretched datastore can have site affinity when the array management interface allows a preference for datastore accessibility on a specific array/site if the inter-array mirroring link fails. This site affinity is also known as site bias, datastore affinity, or datastore locality and, like using an external witness for datastore failover decisions, is vendor-specific.

With proper configuration and array support, administrators can ensure that VMs with site affinity to data center A have datastores with matching site affinity. By implementing multiple datastores with different site affinities, each containing VMs with matching site affinity, the configuration can avoid VM restarts when mirror links fail.

If the array mirror link between sites fails, only the storage system on the preferred site for a given datastore will maintain datastore access. VMs may require restart by vSphere HA in these scenarios:

- VMs running outside the location with datastore access (in non-uniform configurations)
- VMs running outside the location with datastore access (in uniform configurations where site cross-links are also lost)
- VMs running on hosts that have lost all storage access

These restarts cause brief VM outages and should be avoided when possible through proper site affinity configuration.

Infrastructure Architecture

This section describes the basic architecture referenced in this document and discusses the configuration and performance of various vSphere features. For detailed explanations of individual features, refer to the vSphere Availability Guide and the vSphere Resource Management Guide in the product documentation. We provide specific recommendations based on VMware best practices and operational guidance where applicable. The failure scenarios section explains how these best practices prevent or limit downtime.

The described infrastructure consists of a single vSphere cluster with 16 ESXi hosts, managed by a VMware vCenter instance. The first site, Frimley, and the second site, Bluefin, are connected by a stretched layer 2 network. The sites are minimally distant, typical of campus cluster scenarios.

Layer 2 network stretching is necessary because VMs must be able to run on either site, requiring the VM Network to stretch across both locations. While the vSphere management network can be routed or stretched, do not rely on solutions like HCX or a VM to route or stretch between networks, as this creates a network domain single point of failure. Layer 2 stretching remains the recommended best practice.

Each site has eight vSphere hosts. The vCenter Server instance can reside at either site or at a third location. While vCenter Server should remain highly available for management operations and could potentially use stretched storage, placing it on vVol storage is not recommended. This avoidance prevents a circular dependency when powering on vCenter during array VP issues. In our figures, vCenter Server uses vSphere VM-Host affinity with hosts in the Bluefin data center, potentially on a separate VMFS metro-clustered volume.

A vVol stretched cluster environment uses only a single vCenter Server instance, similar to vMSC setup best practices. While vCenter supports HA with two nodes, this feature is not yet supported with vVol. Since the HA model is active-passive, it provides limited benefit, as only one vCenter can be active at any time, meaning only one site maintains vCenter management during site disconnections.

Figure 3 shows three vVol Datastores/Containers, though the tested setup includes different datastores (see Table 1). Multiple stretched datastores are important, as demonstrated in the failure examples. If Container1 has site affinity to Frimley, VMs with Frimley ESXi host affinity should reside on this Container/Datastore or another stretched datastore with Frimley affinity. This configuration helps prevent automatic VM restarts during many failure scenarios.

Regarding VM datastore placement, keep each VM within a single datastore when using stretched storage. If a VM spans multiple containers/datastores (Container1 and Container2, for example) and failures cause Container1 to remain exposed only to Frimley while Container2 remains exposed only to Bluefin, vSphere cannot operate the VM as no single host can access both datastores.

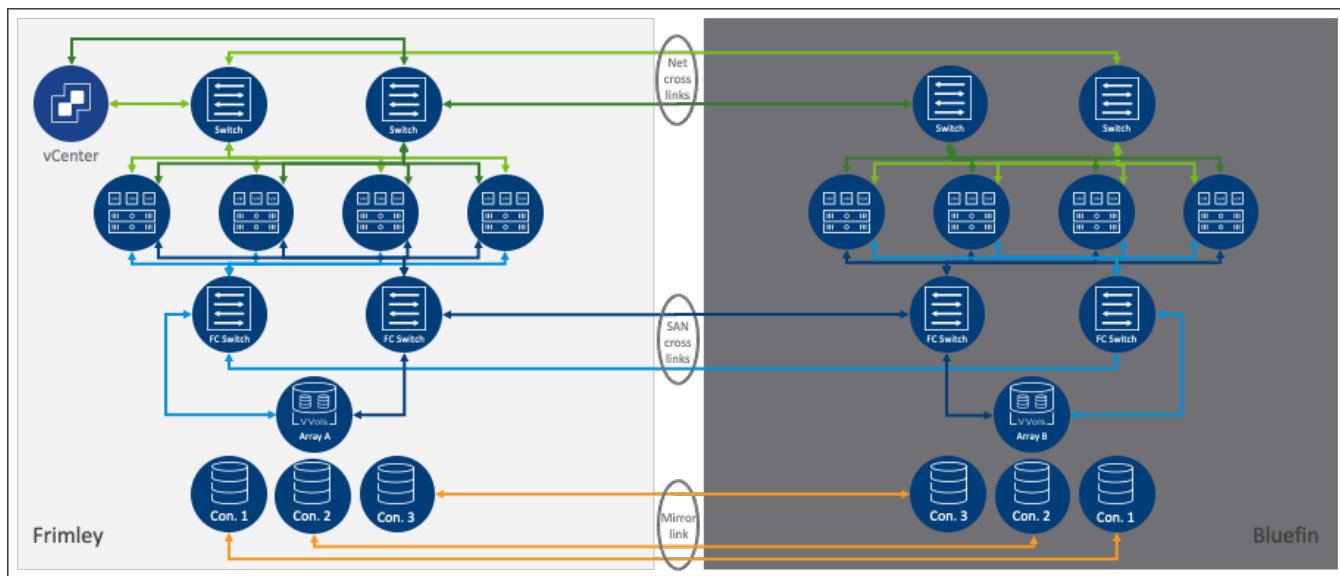


Figure 3 - Test Environment

Location	Hosts	Datastores
Frimley	199.10.107.202	StretchedDS01
	199.10.107.203	StretchedDS02
Bluefin	199.10.107.204	StretchedDS01
	199.10.107.205	StretchedDS02

Table 1 - Test Environment Details

The vSphere cluster connects to a stretched storage system in a fabric configuration with a uniform device access model.

Note: While this setup uses VASA-6 throughout, rolling upgrades may include VASA-5-connected hosts. vCenter upgrades first, and since VASA-5 support in vSphere 8 Update 1, vCenter's VASA version represents the maximum available to underlying hosts. Hosts connect using the maximum VASA version they support within vCenter's version. During rolling upgrades, some hosts may still use VASA-5 and cannot see stretched containers. Though tested, minimize time in this configuration as it represents a minority test case.

Starting with VASA-5, upgrading vSphere to use a higher VASA version with an array should not involve unregistering and re-registering VASA providers. Instead, use the 'Upgrade' button in vCenter UI for the VP to ensure continuous operation during upgrades.

vSphere Configuration

This document examines vSphere HA, vSphere Distributed Resource Scheduler (DRS), and vSphere Storage DRS in stretched cluster deployments. While storage architecture typically dominates planning discussions, workload management and provisioning require equal consideration for optimal performance.

A key motivation for implementing stretched clusters is workload balance and disaster avoidance. This raises important questions:

- How do we maintain proper environment balance without compromising availability or significantly increasing operational costs?
- How do we incorporate requirements into our provisioning process and validate ongoing compliance?

Neglecting these requirements leads to administrative confusion and reduced predictability during failure scenarios when the environment should provide maximum benefit.

Each of these vSphere features has specific configuration requirements and can enhance environment resiliency and workload availability. We provide architectural recommendations throughout this section based on our failure scenario testing results.

vSphere High Availability

Our environment consists of 16 hosts and a uniform stretched storage solution. A resilient architecture must account for various failure scenarios:

- Full site failures
- Host failures
- Network failures
- Storage failures

vSphere HA is crucial in any environment for providing availability, particularly in SSC configurations. VMware recommends enabling vSphere HA and carefully reviewing the following guidelines for optimal vSphere HA configuration in SSC-based infrastructures.

Note that vSphere HA settings apply at the cluster level (created by right-clicking a datacenter and selecting "New cluster..."). Configure these settings under "vSphere Availability" in the cluster "Services" configuration tab.

Admission Control

VMware recommends enabling vSphere HA Admission Control. As workload availability drives most stretched cluster implementations, providing sufficient capacity for full site failure is recommended. With hosts equally divided across sites, configure the admission control policy to 50 percent for both memory and CPU to ensure all workloads can restart using vSphere HA on a single site.

The percentage-based policy is recommended for its flexibility and reduced operational overhead. This approach eliminates the need to adjust percentages when adding new hosts and prevents skewed consolidation ratios that might occur with VM-level reservations. For detailed information about admission control policies and their algorithms, consult the vSphere Availability Guide.

You can specify acceptable Performance Degradation levels for workloads. While the default setting is 100%, adjust this based on your business service level agreement (SLA). Consider an environment with:

- 75GB of memory available in a three-node cluster (25GB per host)
- One host failure tolerance specified
- 60GB of memory actively used by VMs
- 0% resource reduction tolerance set in the UI

vSphere HA accounts for a single host failure in this cluster. After losing one host's memory (25GB), available memory decreases from 75GB to 50GB. With 60GB of memory in use and 0% resource reduction tolerance, the full 60GB remains required. Since only 50GB is available post-failure, vSphere issues a warning. Note that this warning does not prevent new VM provisioning or VM power-on operations—that function belongs to Admission Control.

Figure 4 demonstrates a vSphere HA cluster with Admission Control enabled. Configure specific admission control settings after cluster creation.

New Cluster

1 Basics
2 Image
3 Review

Basics [X]

Name	Cluster
Location	vcqaDC
vSphere DRS	<input checked="" type="checkbox"/>
vSphere HA	<input checked="" type="checkbox"/>
vSAN	<input type="checkbox"/>
	<input type="checkbox"/> Enable vSAN ESA i

Manage all hosts in the cluster with a single image [i](#)

Choose how to set up the cluster's image

- Compose a new image
- Import image from an existing host in the vCenter inventory
- Import image from a new host

Manage configuration at a cluster level [i](#)

CANCEL NEXT

Figure 4 - vSphere HA Configuration

VMware recommends configuring admission control to reserve 50% of CPU and memory resources for High Availability (HA). This configuration ensures your virtual machines (VMs) can restart if an entire site fails. In the vSphere Client, set the number of host failures to tolerate, which the system converts to a percentage of reserved resources, as shown in Figure 5. Select the Cluster Resource Percentage admission control policy, as it offers the greatest flexibility when managing cluster resources. For additional configuration guidance, consult the vSphere Availability Guide.

The screenshot shows the 'Edit Cluster Settings' dialog for a cluster, specifically the 'Admission Control' section. The dialog has a title bar with 'Edit Cluster Settings | Cluster' and a close button (X). Below the title bar is a descriptive text: 'Admission control is a policy used by vSphere HA to ensure failover capacity within a cluster. Raising the number of potential host failures will increase the availability constraints and capacity reserved.' The main configuration area includes: 'Host failures cluster tolerates' set to '1' with a note 'Maximum is one less than number of hosts in cluster.'; 'Define host failover capacity by' set to 'Cluster resource Percentage'; a toggle for 'Override calculated failover capacity.' which is turned off; 'Reserved failover CPU capacity' set to '50 % CPU'; 'Reserved failover Memory capacity' set to '50 % Memory'; a toggle for 'Reserve Persistent Memory failover capacity' which is turned off and has an information icon; a toggle for 'Override calculated Persistent Memory failover capacity' which is turned off; 'Reserve' set to '0 % of Persistent Memory capacity'; and 'Performance degradation VMs tolerate' set to '100 %' with a note: 'Percentage of performance degradation the VMs in the cluster are allowed to tolerate during a failure. 0% - Raises a warning if there is insufficient failover capacity to guarantee the same performance after VMs restart. 100% - Warning is disabled.' At the bottom right are 'CANCEL' and 'OK' buttons.

Figure 5 - vSphere HA Admission Control

vSphere HA Heartbeats

VMware vSphere High Availability (HA) validates host status through two heartbeat mechanisms: network heartbeating and datastore heartbeating. Network heartbeating serves as the primary method, while datastore heartbeating acts as the secondary validation method when network heartbeating fails.

When a host stops receiving heartbeats, it initiates a fail-safe process to determine if it is isolated from its primary node or experiencing complete network isolation. The host pings its default gateway as an initial check. For increased reliability, VMware recommends configuring at least two additional isolation addresses, with each address located in a different physical site.

In this deployment, one isolation address is placed in the Frimley data centre, and another in the Bluefin data centre. This configuration allows vSphere HA to detect complete network isolation, even during inter-site connection failures. Figure 6 demonstrates the configuration of multiple isolation addresses using the `das.isolationaddress` advanced setting. For configuration instructions, refer to the Broadcom Knowledge Base article "[Setting Multiple Isolation Response Addresses for vSphere High Availability](#)."

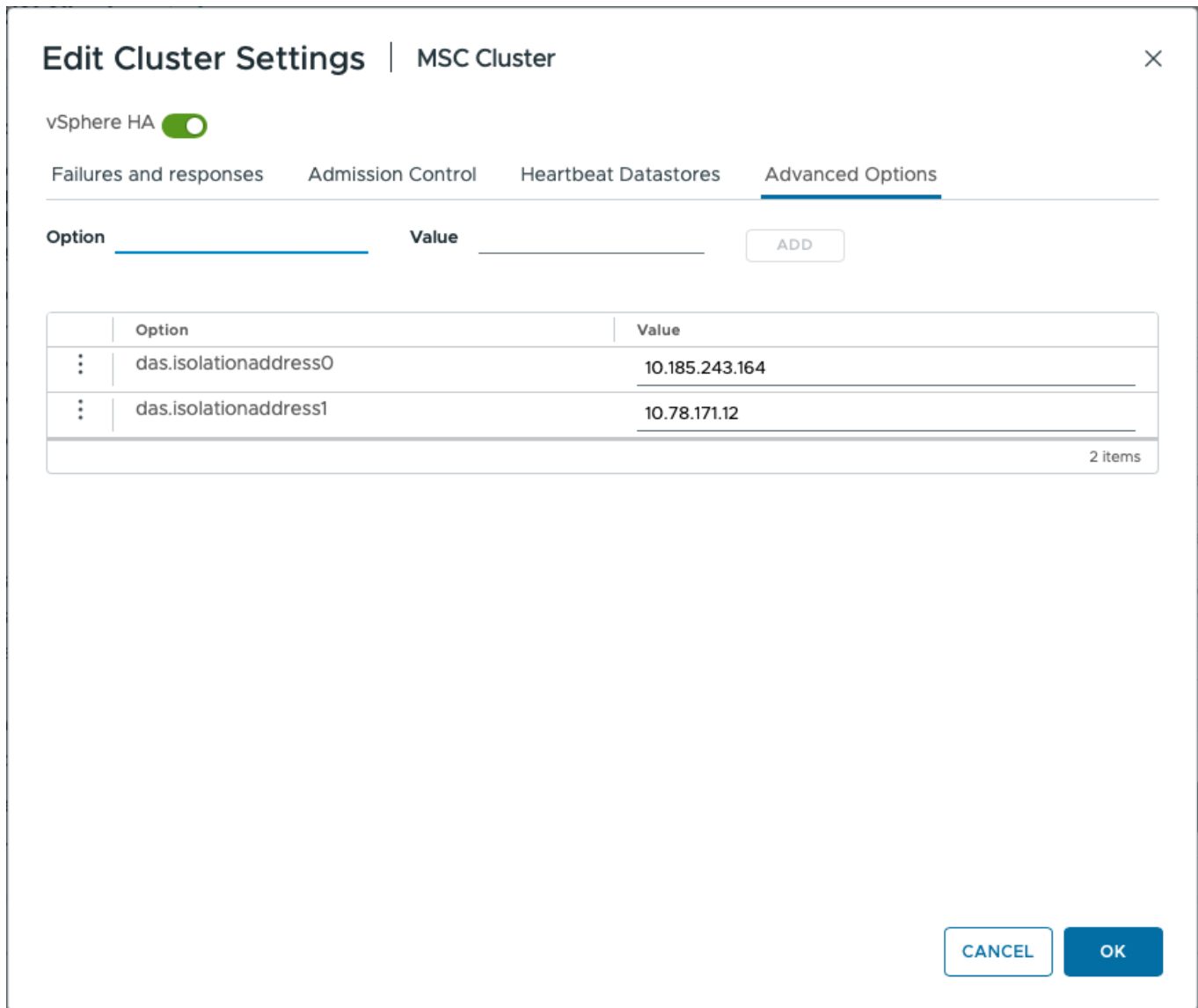


Figure 6 - vSphere HA Isolation Address Configuration

In a VMware vSphere stretched cluster environment, configure four heartbeat datastores instead of the minimum requirement of two, though the maximum remains five. Select two datastores from each site as preferred heartbeat datastores. This configuration ensures vSphere High Availability (HA) maintains datastore heartbeating capabilities during inter-site connection failures, allowing accurate host state determination in all scenarios.

To increase the number of heartbeat datastores beyond the default, add the `das.heartbeatDsPerHost` advanced setting, as illustrated in Figure 7. This setting enables full redundancy across both data center locations.

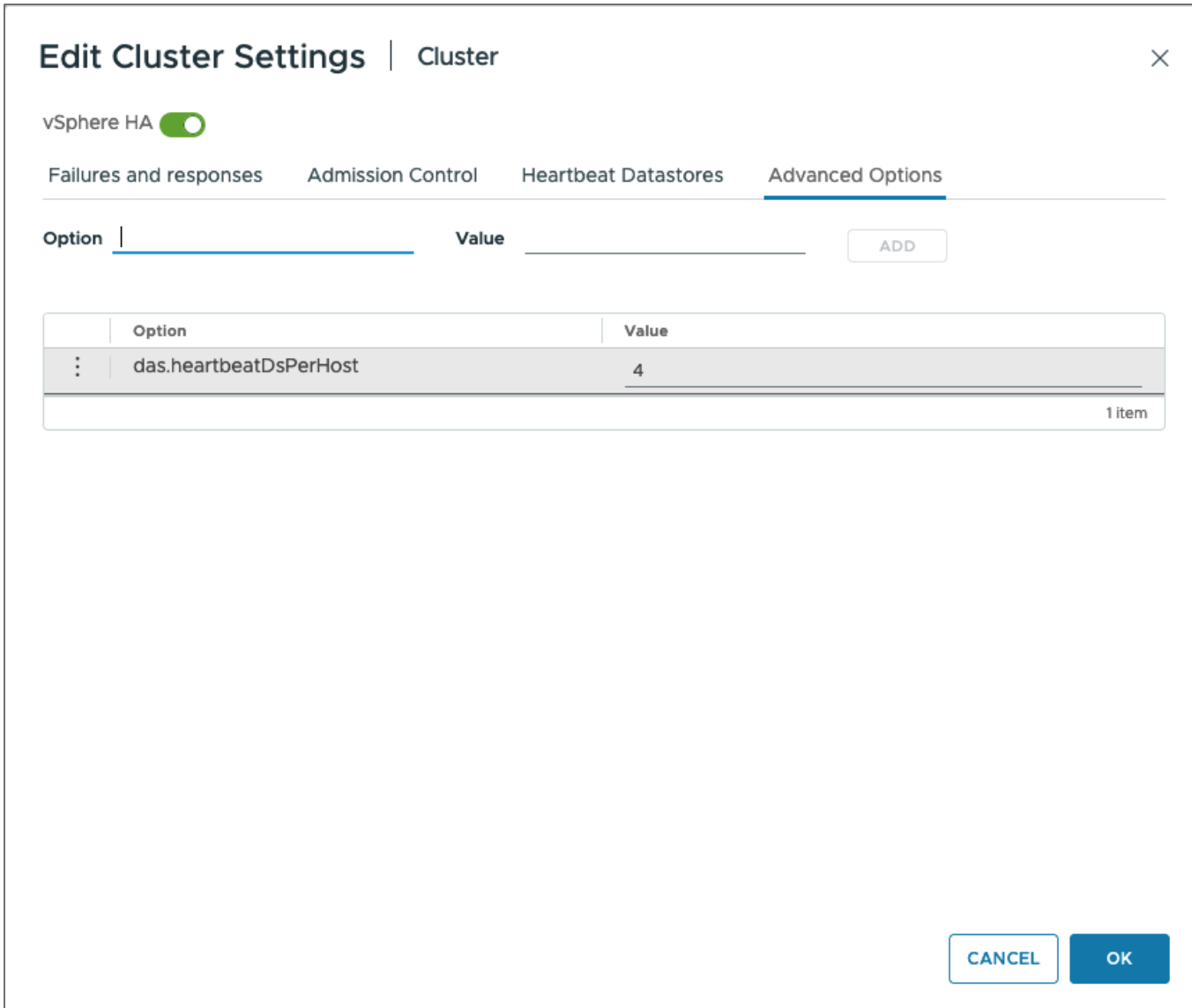


Figure 7 - vSphere HA Advanced Options

VMware recommends using the "Use datastores from the specified list and complement automatically if needed" option for heartbeat datastore selection in vSphere High Availability (HA). This setting allows vSphere HA to use alternative datastores if the four manually designated datastores become unavailable. For optimal redundancy, select two datastores in each physical location, ensuring datastore availability during site partitions.

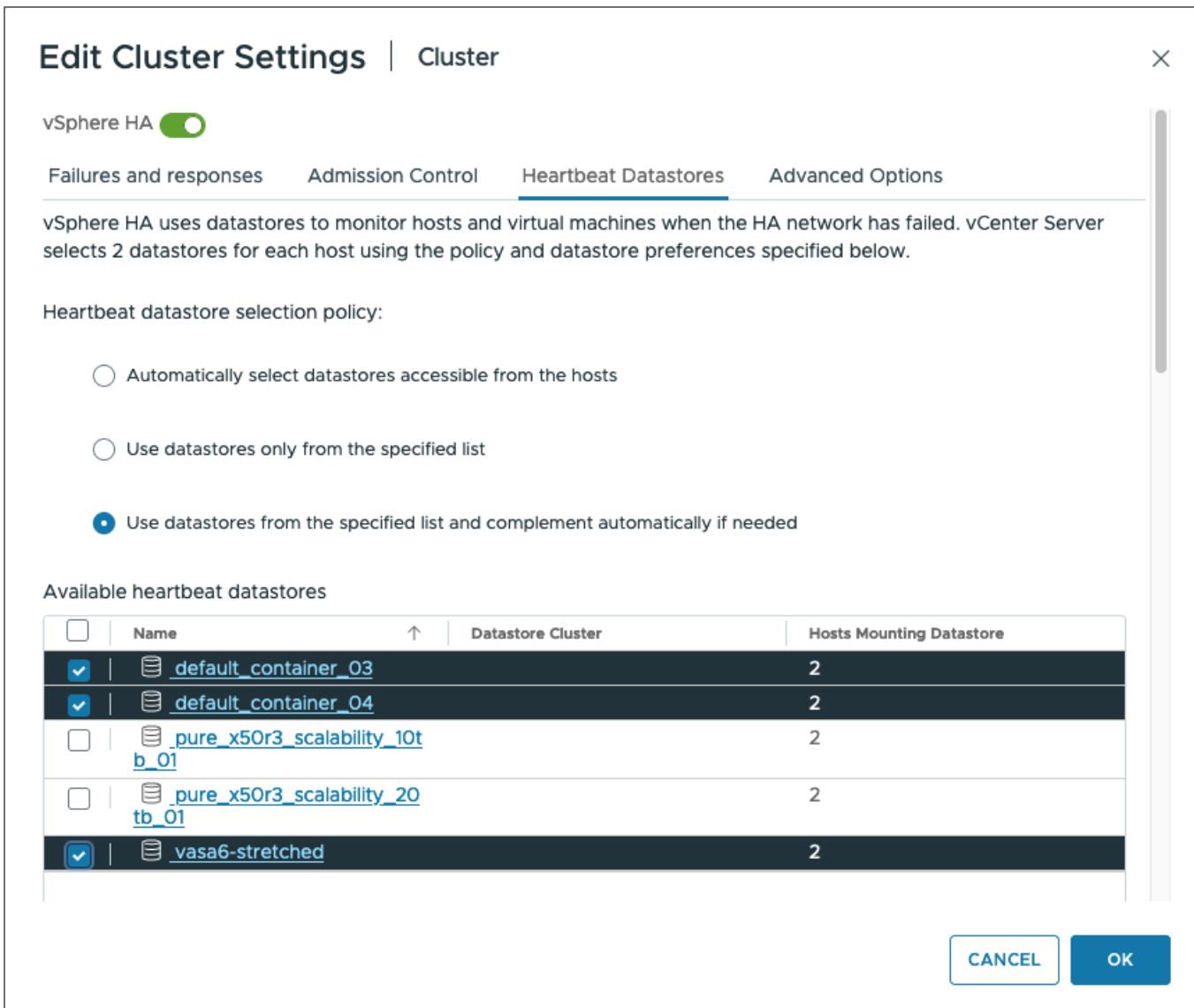


Figure 8 - Datastore Heartbeating

Permanent Device Loss and All Paths Down Scenarios

VMware Virtual Machine Component Protection (VMCP) enables automated VM failover during storage device failures through two scenarios: Permanent Device Loss (PDL) and All Paths Down (APD).

PDL occurs when a block storage array controller signals to the VMware vSphere host through a SCSI sense code that a device (LUN/PE) is permanently unavailable - for example, when a LUN/PE is set offline. In non-uniform models, this condition triggers appropriate host actions when device access is revoked.

During complete storage failures, PDL signaling becomes impossible due to lost array-host communication. The vSphere host identifies this as an APD condition. Storage network failures also result in APD conditions, as the host cannot determine the storage system's status.

To enable vSphere High Availability (HA) response to both conditions, configure the response settings for "Datastore with PDL" and "Datastore with APD" under Failures and Responses, as shown in Figure 9. Note that the current vSphere Client interface does not explicitly label these settings as VMCP.

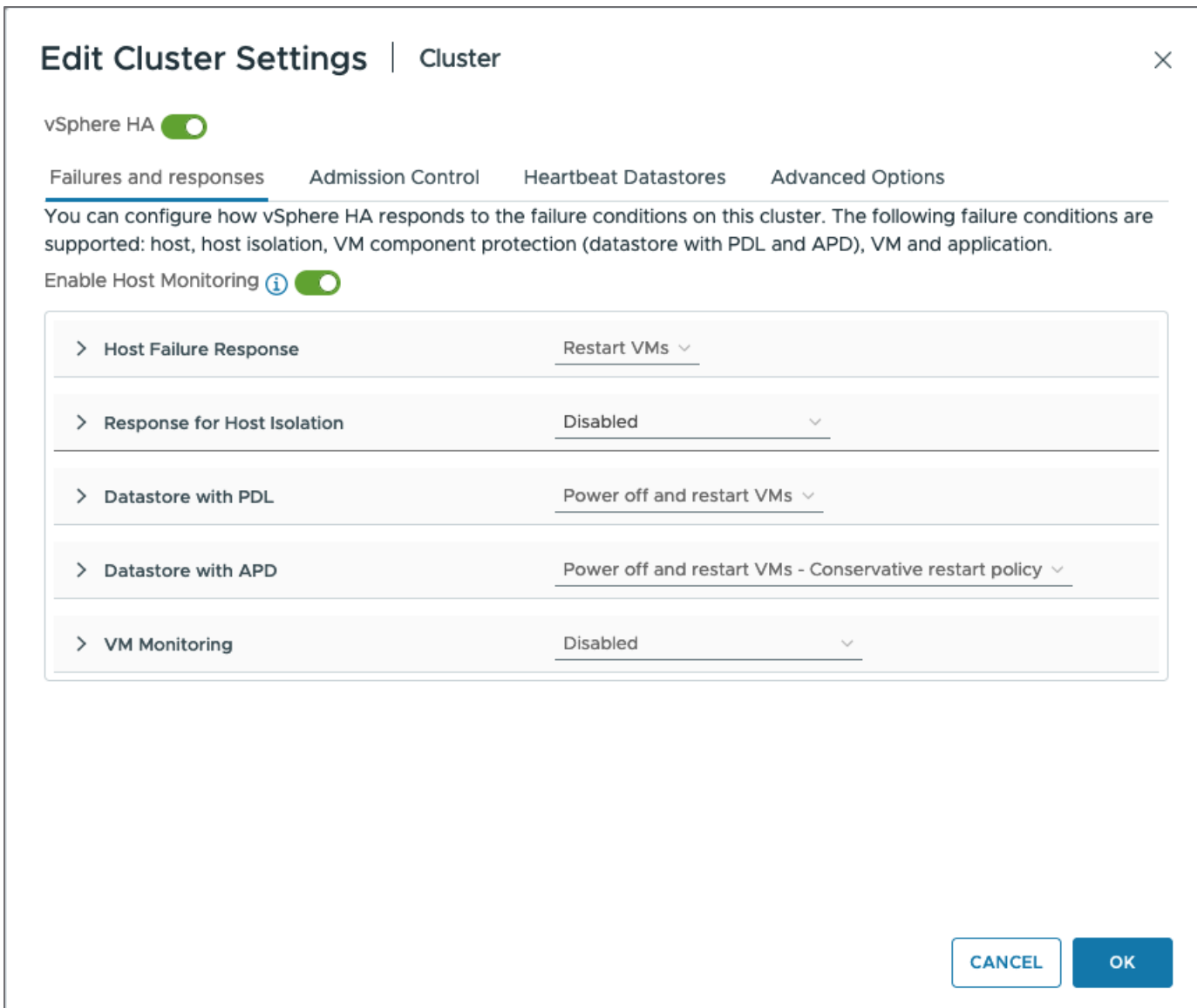


Figure 9 - VM Component Protection

The configuration screen can be found as follows:

- Log in to VMware vSphere Web Client.
- Click Hosts and Clusters.
- Click the cluster object.
- Click the Configure tab.
- Click vSphere Availability and then Edit.
- Select Failures and Responses.
- Select individual functionality, as displayed in figure 9.

For Permanent Device Loss (PDL) scenarios, configure the response in the Failures and Responses section. Select "Power off and restart VMs" as the recommended setting.

Note: Avoid using "Disabled" or "Issue Events" options, as they prevent VM termination during PDL conditions. This leaves the datastore in PDL state, requiring manual VM shutdown and blocking recovery for other VMs on the affected datastore.

For All Paths Down (APD) scenarios, configure settings in the same section shown in Figure 9. The interface enables response configuration for APD conditions, timing parameters, and system behaviour when failures resolve before the APD timeout, as shown in Figure 10.

Edit Cluster Settings | Cluster

All Paths Down (APD) Failure Response Allows you to configure the cluster to respond to APD Datastore failures

- Disabled
No action will be taken on the affected VMs.
- Issue events
No action will be taken on the affected VMs. Events will be generated.
- Power off and restart VMs - Conservative restart policy
A VM will be powered off, if HA determines the VM can be restarted on a different host.
- Power off and restart VMs - Aggressive restart policy
A VM will be powered off, if HA determines the VM can be restarted on a different host, or if HA cannot detect the resources on other hosts because of network connectivity loss (network partition).

Response for APD recovery Disabled

Delay for failure response 3 minutes

CANCEL OK

Figure 10 - VM Component Protection Detailed Configuration

When VMware vSphere detects an All Paths Down (APD) condition, it initiates a 140-second timer. After this period, the device enters APD timeout status, and vSphere High Availability (HA) begins its three-minute countdown. Once complete, vSphere HA typically restarts affected virtual machines (VMs), though VMware Virtual Machine Component Protection (VMCP) offers alternative responses. VMware recommends using "Power off and restart VMs – Conservative restart policy."

The Conservative and Aggressive policies differ in their VM handling. Conservative policy ensures a suitable host with necessary resources exists before powering off the affected VM. Aggressive policy powers off the VM first, then searches for a new host, risking VM downtime if no suitable host is found.

For vVols stretched storage, the Conservative policy provides particular value. vVols use both control (VASA) and data (I/O) paths. While losing either path triggers APD status, a VM with only control path loss may continue normal I/O operations. However, operations requiring control path access, such as snapshots or vMotion, become impossible. Therefore, restarting the VM on a functioning host represents best practice.

If storage access returns before the APD timeout, vSphere HA preserves VM state unless configured otherwise through "Response for APD recovery." VMware recommends leaving this setting inactive. When "Reset VMs" is selected, the system ignores restart priorities and dependencies, as these settings apply only to restarts, not resets.

Restart Ordering and Priority

In the VM Overrides section, VMware vSphere High Availability (HA) enables specific restart sequences and dependencies through VM Restart Priority settings. The priority levels range from Lowest to Highest, with Medium as the default setting.

For optimal failover operations, assign higher restart priorities to essential infrastructure services like DNS, Active Directory, and multi-tier applications. In multi-tier deployments, such as database, application, and web server configurations, the database tier typically requires the highest priority to ensure proper application startup sequence. Figure 11 demonstrates how to select multiple VMs for priority adjustment.

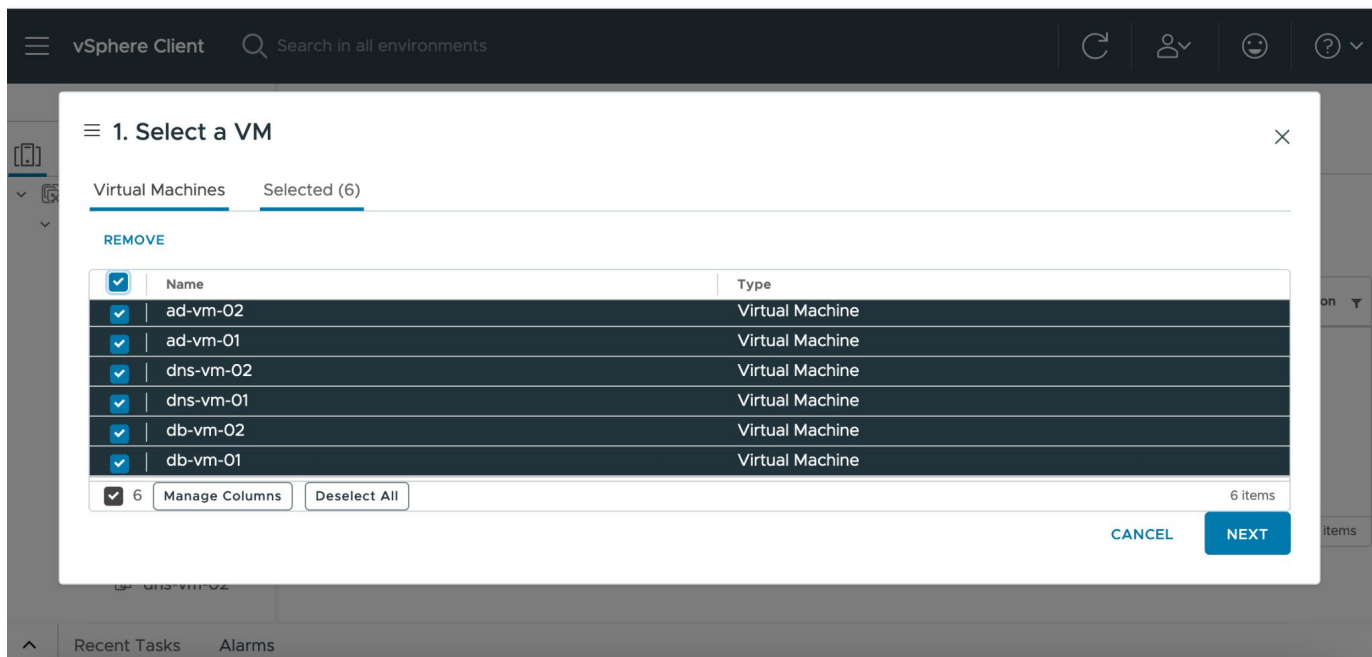


Figure 11 - VM Restart Priority

For these selected application tier VMs, set the restart priority to Highest, as shown in Figure 12. This ensures critical application components restart with top priority during failover events.

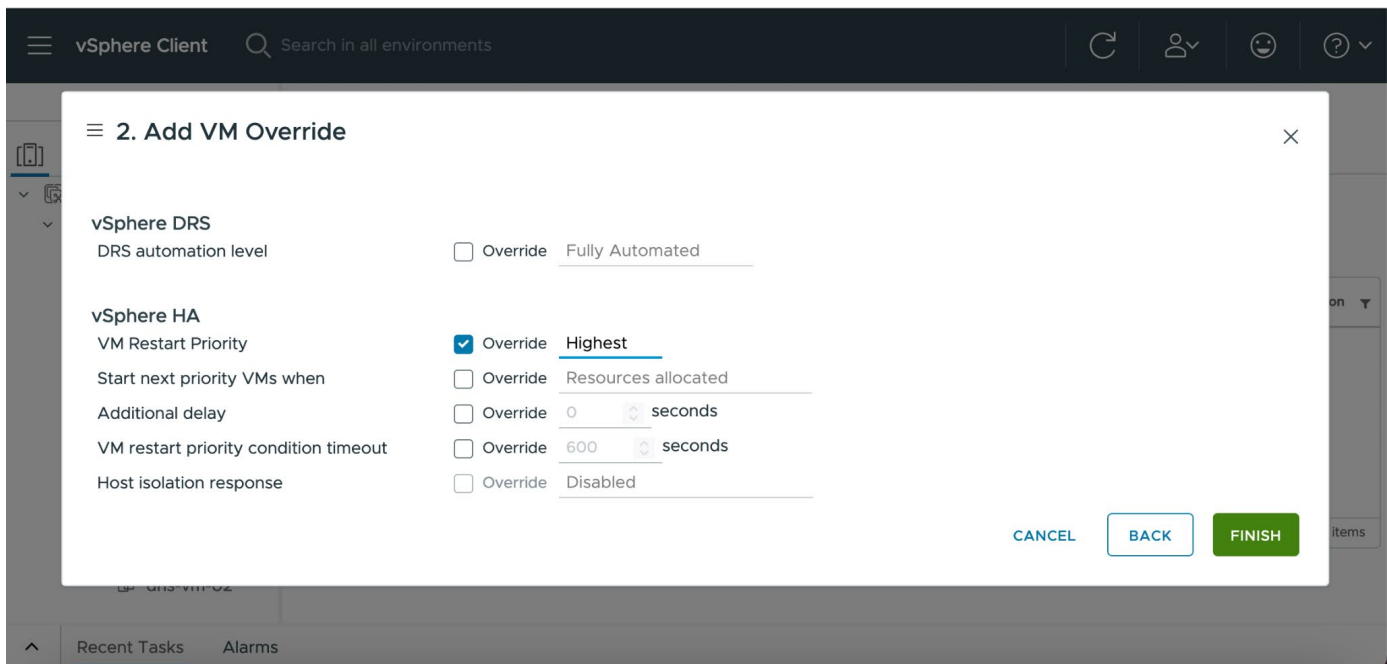


Figure 12 - Changing VM Restart Priority

vSphere HA allows configuration of VM startup sequence between priority groups through the "Start next priority VMs when" dropdown menu. Select "guest heartbeats detected" to wait for VMware Tools to signal VM readiness, as shown in Figure 13.

If no heartbeat is detected, the system waits 600 seconds before starting the next priority group. While this timeout value can be adjusted using the "or after timeout occurs at" setting, VMware recommends keeping the default 600-second interval.

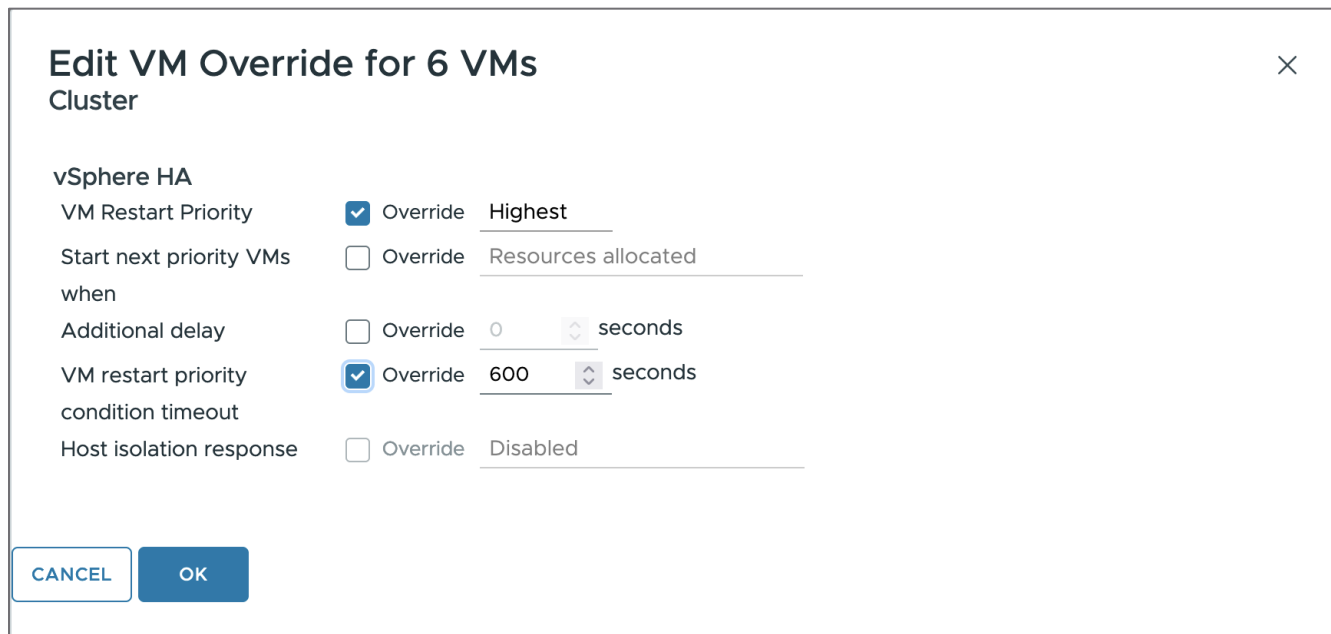


Figure 13 - Additional Delay

vSphere HA provides an additional method for controlling restart order through VM dependencies, which uses cluster rules. Create two VM groups (Figure 14), then establish a VM-to-VM rule specifying the dependency relationship (Figure 15).

The restart timing depends on the VM Dependency Restart Condition. With the default "Resources Allocated" cluster setting, the second group starts almost immediately after the first, functioning purely as a scheduling mechanism. For most scenarios, "Powered On" or "Guest Heartbeats Detected" conditions prove more practical.

These dependency rules are mandatory and cannot be violated. For example, if VMs in the first group fail to start and the condition is set to "Powered On," the second group remains powered off indefinitely.

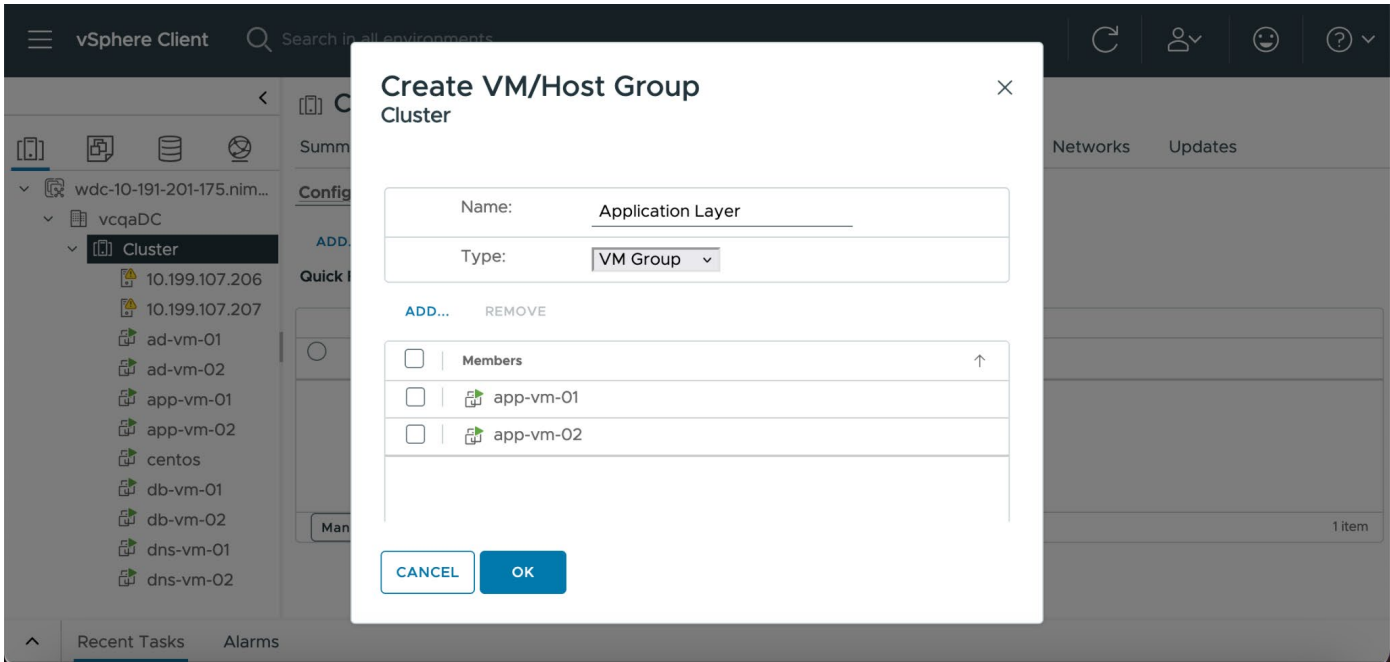


Figure 14 - Restart Dependency Group

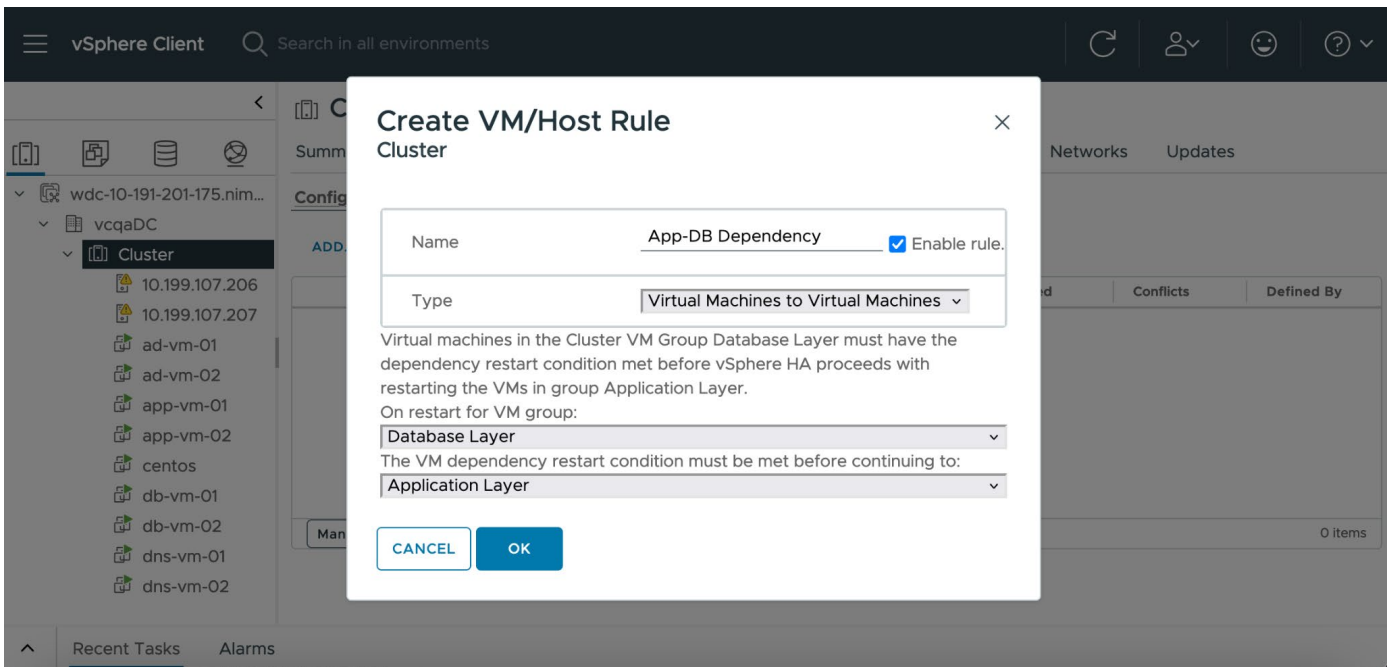


Figure 15 - Restart Dependency Rules

Proactive HA

vSphere DRS includes Proactive HA functionality, though it appears under the Availability settings in the interface. As downtime prevention remains crucial for stretched cluster configurations, VMware recommends enabling Proactive HA.

Proactive HA implementation requires installing a vendor-specific health provider plugin for the vSphere Client. Currently, health providers are available from HPE, Dell, and Cisco. After plugin installation, configure Proactive HA with the Automated automation level for immediate response to potential issues.

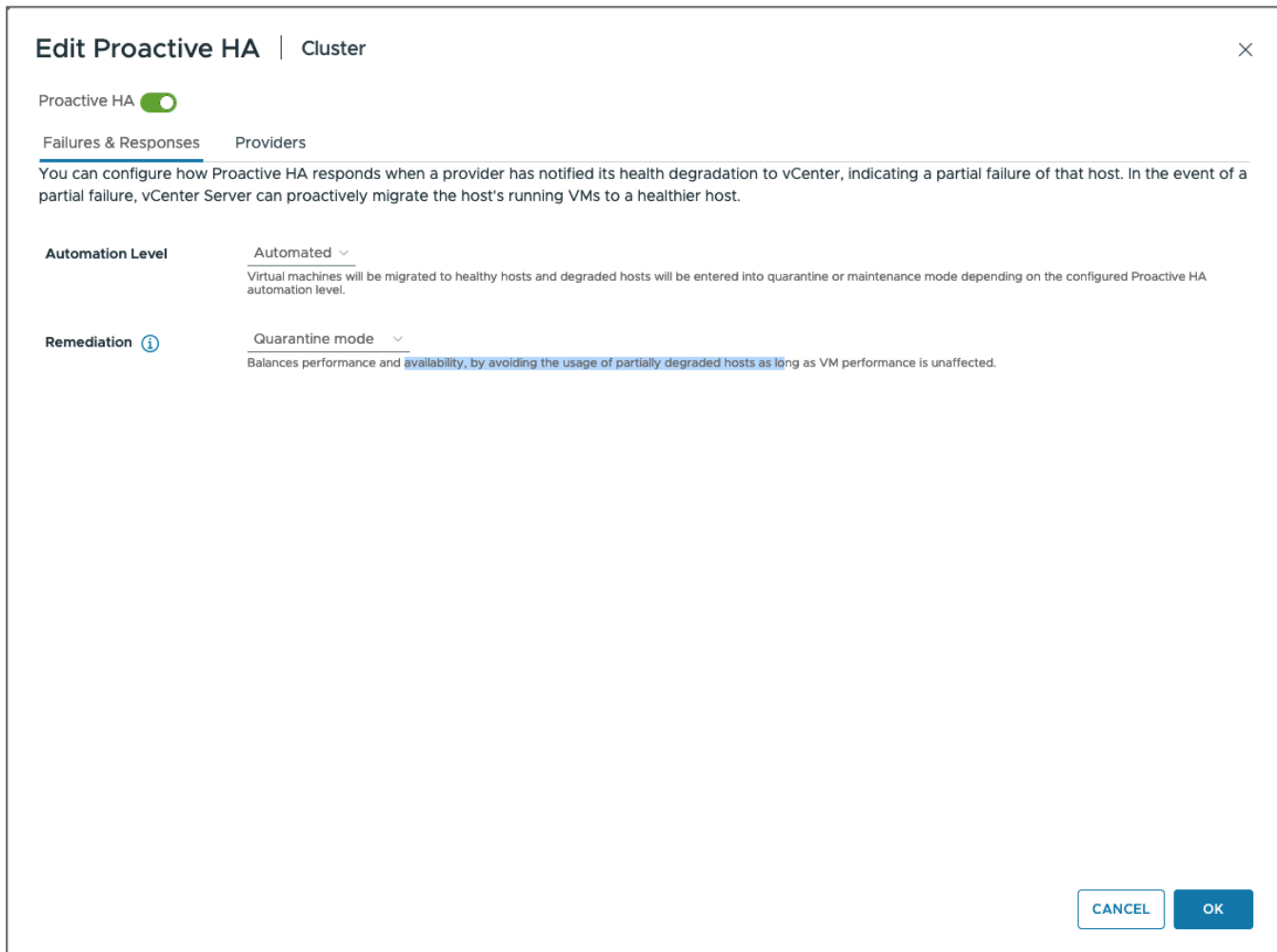


Figure 16 - Proactive HA Configuration

For Proactive HA remediation shown in Figure 16, VMware recommends using the mixed mode: Quarantine mode for moderate failures and Maintenance mode for severe failures. The available remediation options function as follows:

- “Quarantine Mode” prevents new VMs from being placed on the affected host while vSphere DRS attempts to migrate existing VMs. Migration success depends on resource availability and cluster rules.
- “Maintenance Mode” initiates complete host evacuation by migrating all VMs.
- “Quarantine for Moderate/Maintenance for Severe” applies different responses based on failure severity. This provides conservative evacuation for moderate issues and aggressive evacuation for severe ones.

Proactive HA responses vary by health provider plugin version and monitored components. Basic plugins monitor power supplies, fans, and storage (SD cards), while advanced versions may include memory and network monitoring.

Response severity follows a color-coded system (Table 2), with Yellow and Red states triggering Proactive HA actions. Though some plugins allow manual severity adjustment, VMware recommends keeping vendor-defined severity settings.

State	Color
Unknown	Gray
OK	Green
Moderately Degraded	Yellow
Severely Degraded	Red

Table 2 - Proactive HA Status

For detailed configuration instructions for health providers and their specific settings, consult your server vendor's documentation.

Distributed Resource Scheduler (DRS)

vSphere DRS provides load distribution and additional features valuable for stretched cluster environments. VMware recommends enabling DRS to balance CPU and memory workloads across cluster hosts. However, storage and network resource management requires careful consideration.

In stretched cluster environments, implement vSphere DRS affinity rules to create logical VM separation, reducing storage and network traffic overhead while improving availability. This separation proves particularly valuable for infrastructure services like Microsoft Active Directory and DNS, ensuring these critical services remain distributed across sites.

Site Affinity

vSphere DRS affinity rules help prevent unnecessary downtime and reduce storage and network traffic by enforcing site preferences. VMware recommends aligning VM-to-host affinity rules with storage configurations, ensuring VMs run on hosts at the same site as their primary read/write storage array.

For example, VMs on the Frimley01 datastore should prefer hosts in the Frimley data center. This maintains VM connection to primary storage during inter-site network failures and keeps read I/O operations local. Note that storage vendors use different terms for LUN/PE array relationships. This document uses "storage site affinity" to indicate preferred LUN/PE access location.

VMware recommends implementing "should rules" rather than "must rules." During site failures, vSphere HA can violate "should rules" to maintain service availability. "Must rules" prevent vSphere HA from restarting VMs on alternate data center hosts during failures. vSphere DRS shares these rules with vSphere HA through a compatibility list for startup governance.

While rare, vSphere DRS may violate "should rules" during extreme host saturation with aggressive settings. Monitor for rule violations as they can affect availability and performance.

Two key advanced settings control this behavior:

- `das.respectVmVmAntiAffinityRules` (defaults to true, set to false if you want HA to not respect VM-VM affinity rules)
- `das.respectVmHostSoftAffinityRules` (defaults to true, set to false if want HA to not respect VM-Host affinity rules)

Keep VM-VM anti-affinity rules enabled for guest clustering scenarios. Note that vSphere 8.0 Update 3 does not yet support guest clustering on stretched containers. Maintain default VM-Host affinity settings for proper site affinity.

For site definition, create host groups per site and assign VMs based on their datastore affinity. Automate this process using VMware Aria Orchestrator or VMware PowerCLI. If automation isn't possible, use consistent naming conventions and regularly validate group assignments.

Figures 17 through 20 demonstrate this configuration, showing VMs assigned to the Bluefin data center group.

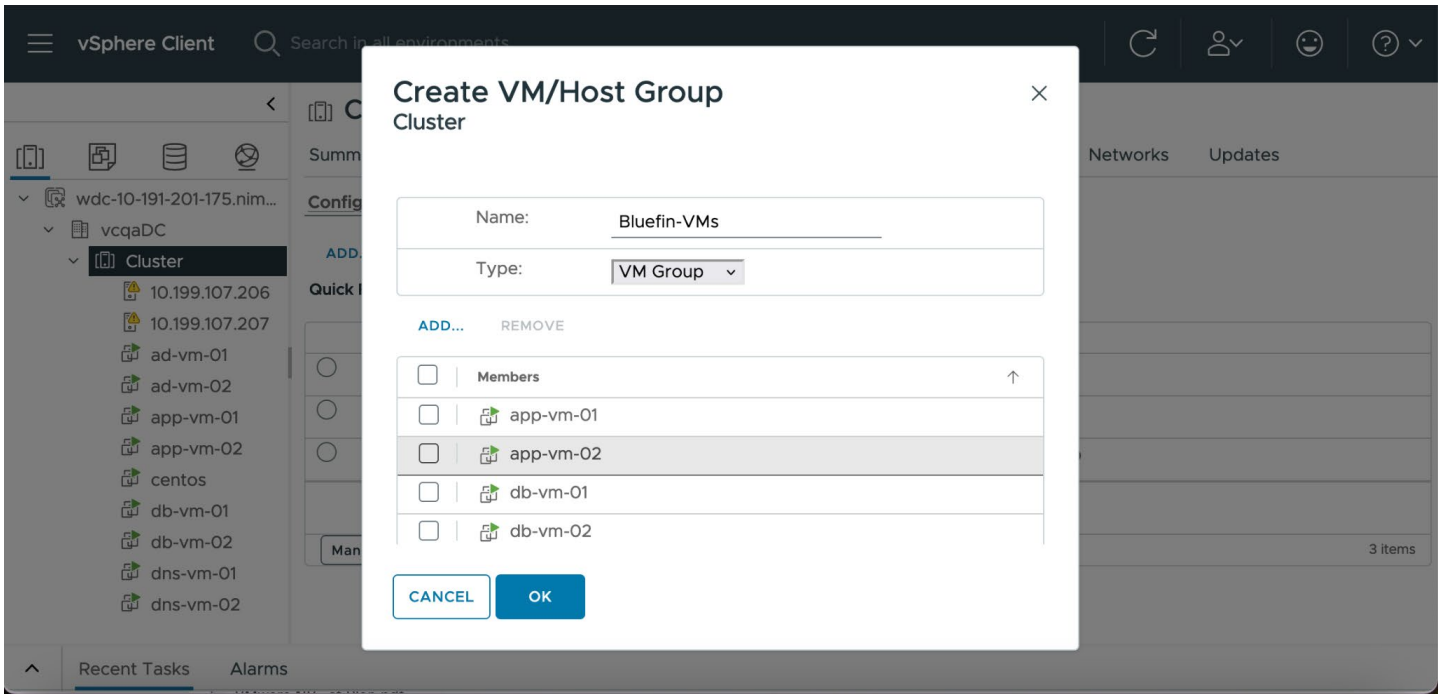


Figure 17 - VM Group in DRS

Create a Bluefin host group containing all hosts located at the Bluefin site.

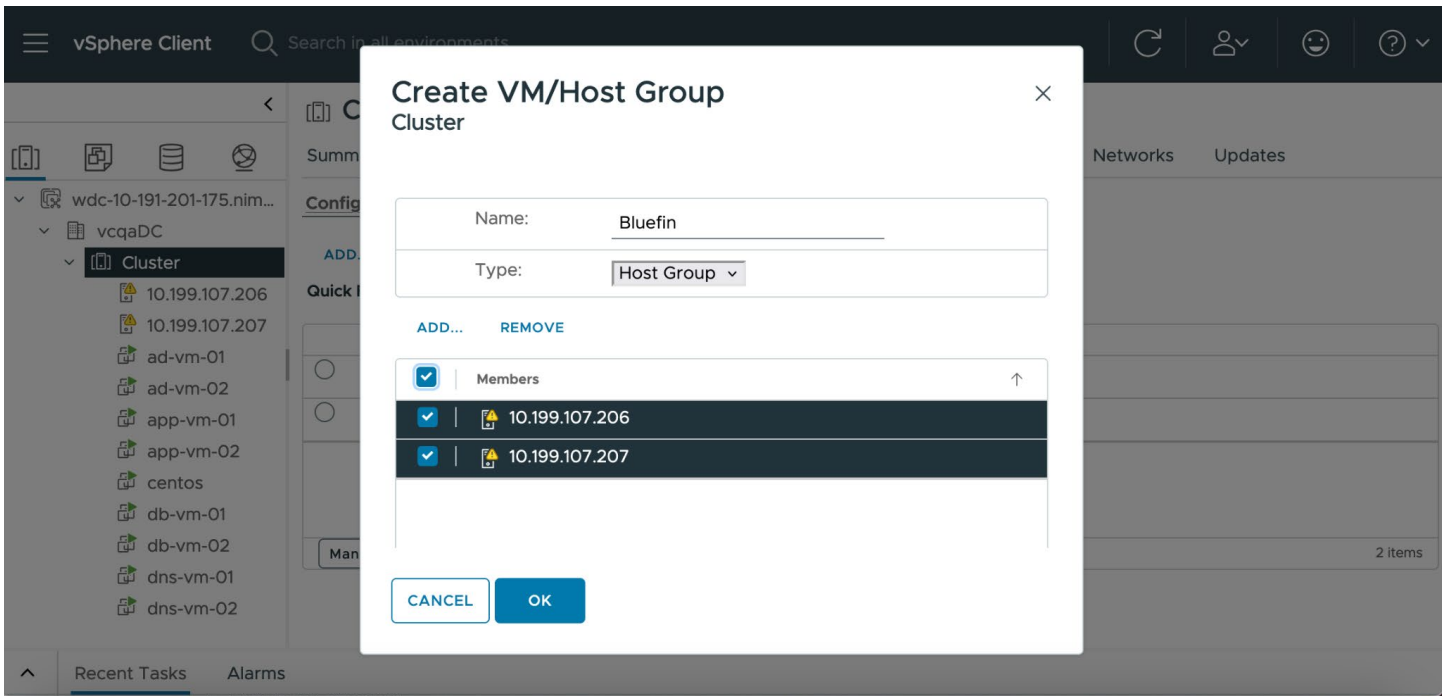


Figure 18 - Host Group in DRS

Create a "should run on" rule linking the Bluefin host and VM groups to establish site affinity.

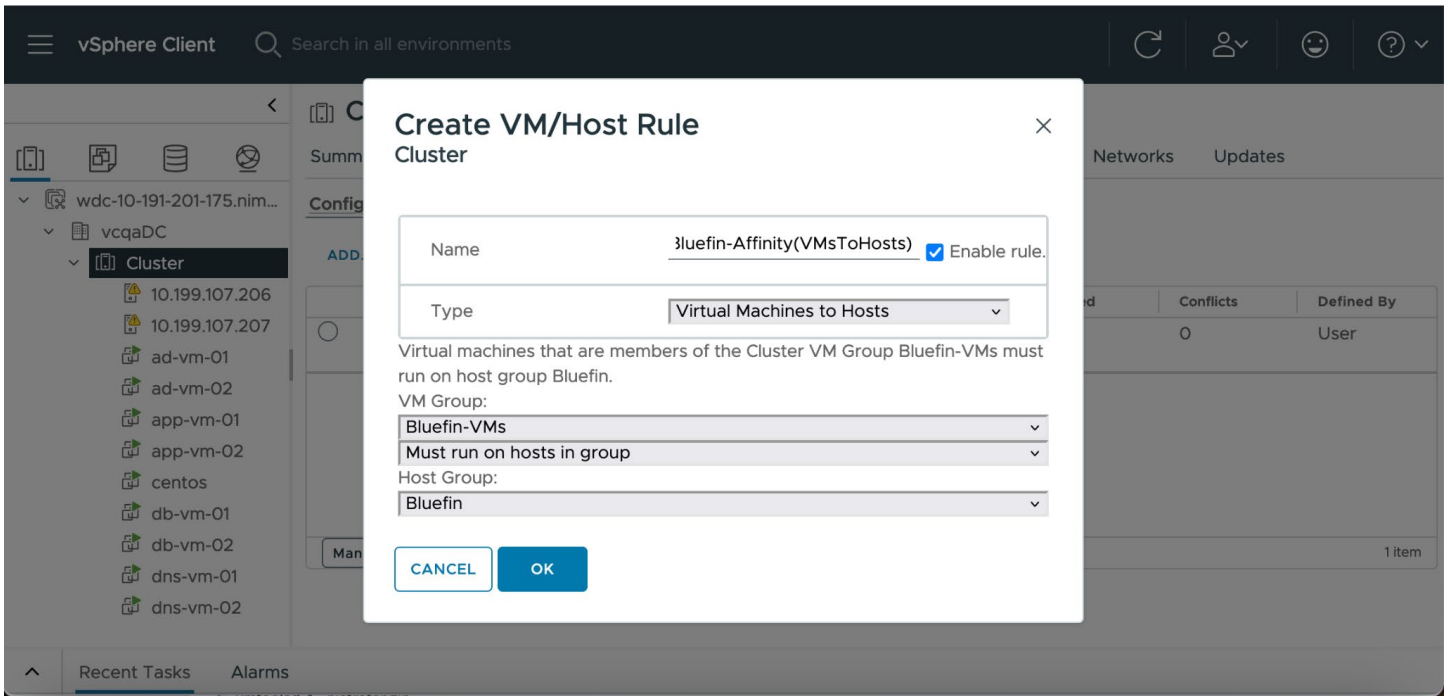


Figure 19 - Rule Definition in DRS

Configure separate "should run on" rules for each site to establish proper affinity relationships across both locations.

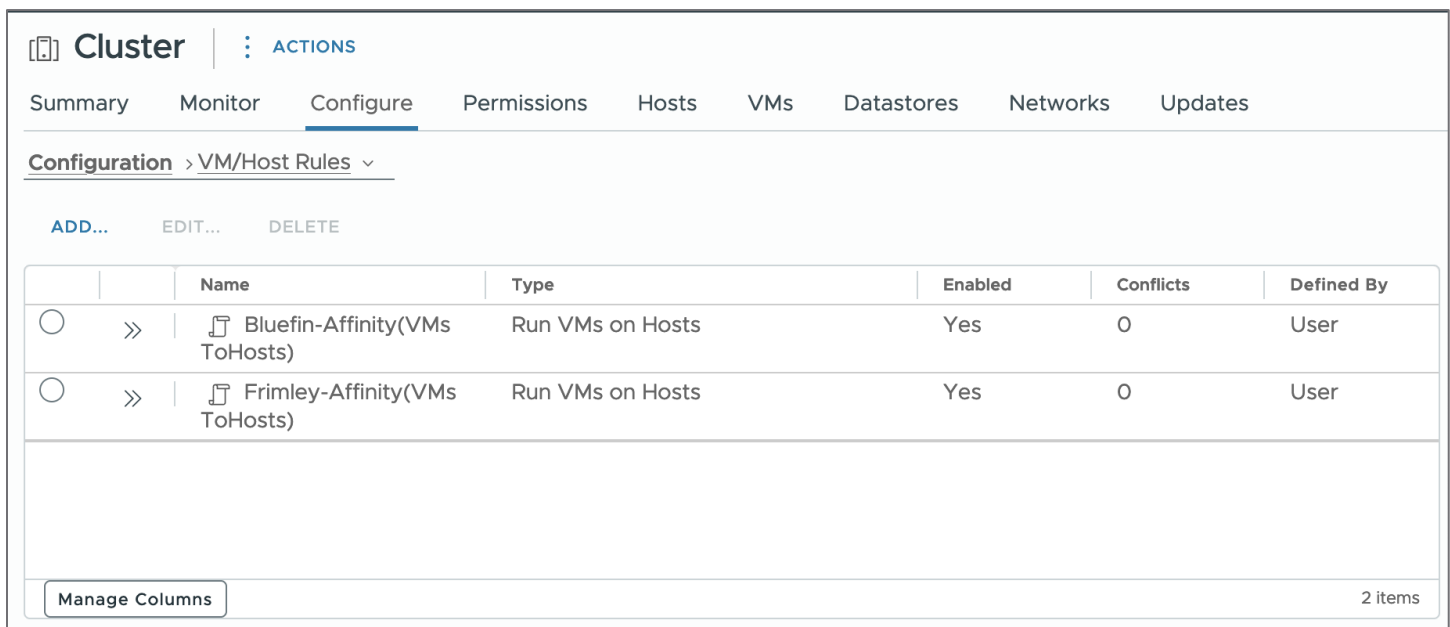


Figure 20 - VM to Host Rules in DRS

Advanced Settings

When configuring vSphere DRS advanced settings in stretched cluster environments, maintain default values for options not specifically discussed in VMware documentation. The VM Distribution setting, which evenly distributes VMs across hosts based on quantity rather than resources, may override non-mandatory VM-to-Host rules. This could result in unexpected VM placement.

Before implementing any advanced options, test all failure scenarios to validate outcomes match expectations. Document and verify results for each configuration change.

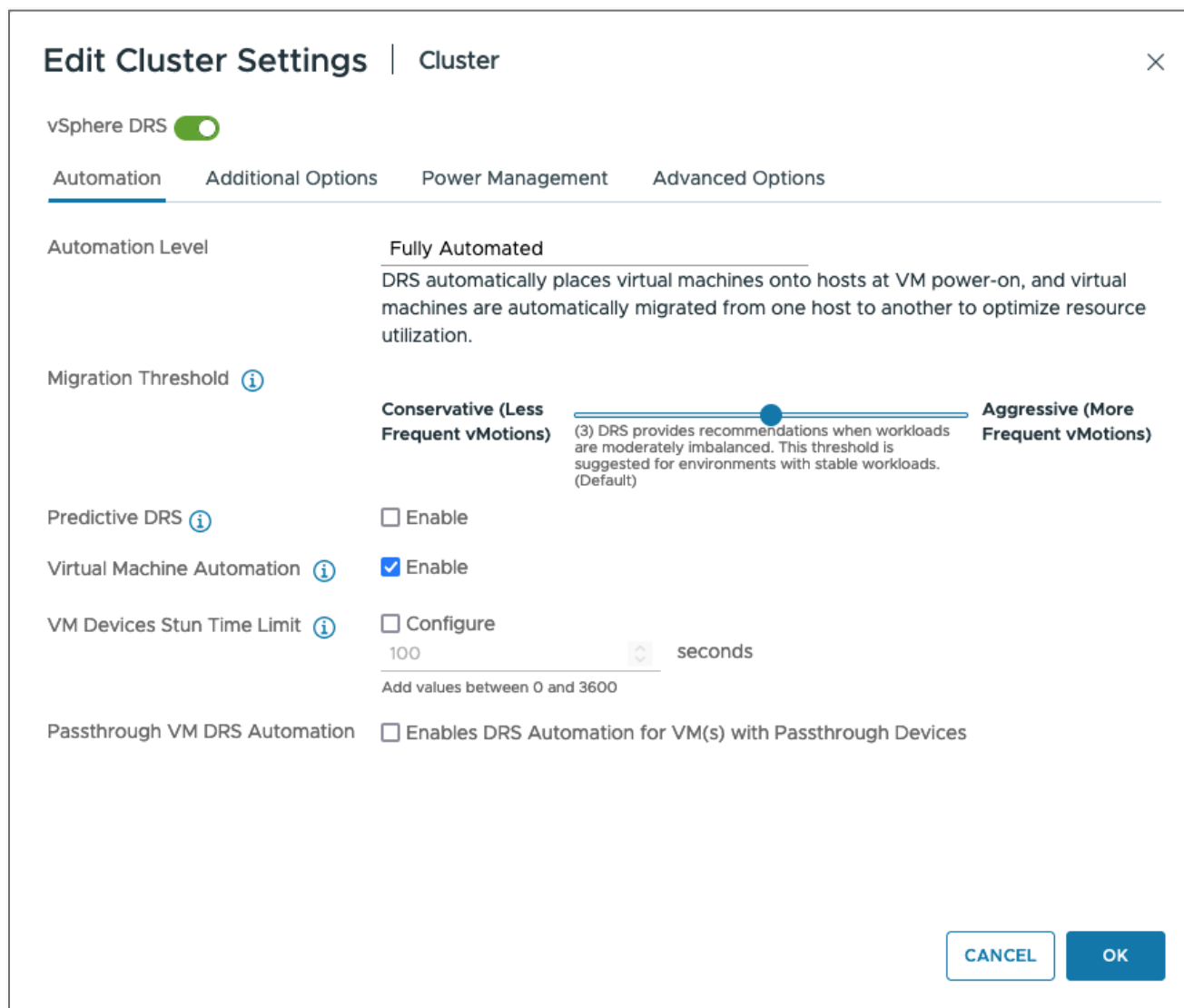


Figure 21 - Advanced DRS Options

Correcting Affinity Rule Violations

After a rule violation, vSphere DRS prioritises correcting affinity rules over load balancing. The system generates VM migration recommendations to align with host group assignments.

vSphere DRS activates every 5 minutes by default as well as when cluster conditions change, such as host reconnection. Testing confirms DRS generates corrective recommendations within 30 seconds of host reconnection. The vSphere vMotion network capacity may require multiple DRS cycles to resolve all violations.

Storage DRS

Storage DRS does not support vVols due to fundamental architectural differences. In VMFS environments, a datastore typically connects to a single LUN, providing consistent performance across all VMDKs, modified only by VM storage shares.

With vVols, each virtual volume links to a separate LUN (in SCSI-based implementations) within the same datastore. These LUNs can have different performance characteristics based on RAID levels and storage policies. Storage DRS assumptions

about datastore-level performance metrics and space constraints do not apply to vVol architecture, making its recommendations invalid for vVol environments.

Known Issues and Considerations

Known issues and considerations for VMware vVols in vSphere 8.0 Update 3 stretched storage clusters:

- VASA Provider downtime causes esxcli vVol commands to delay up to one minute due to network interaction attempts. This will be resolved in VCF 9.0. Consider unregistering the VASA Provider during extended outages.
- Do not unregister the only VASA Provider supporting mounted datastores, as this causes datastore inaccessibility. Stop workloads first if unregistering is necessary. Use the non-disruptive upgrade workflow for VASA Provider updates.
- VASA version upgrades require manual administrator initiation, even when vSphere detects new versions. Schedule upgrades during low-activity periods.
- ESXi hosts cannot use mixed storage protocols (FC, iSCSI) to access the same storage. Configure array access using a single protocol.
- Linked Clones must remain within the same vVol datastore. This requirement helps prevent issues when datastores fail over to different sites during mirror link failures.
- Keep all VM components within one datastore, except during Storage vMotion. Multiple datastores risk VM unavailability if sites lose connectivity.
- The vvolTerminateVMOnPdl setting applies only to non-stretched containers.
- Leave the TerminateVMOnPDL advanced option disabled for vVol stretched storage. Use HA/VMCP settings to manage PDL responses.
- Unmount vVol datastores from all hosts before removing stretch configuration. This requirement will be addressed in a future release.

Failure Scenarios

A properly architected environment prevents many potential cluster failures through VMware vSphere HA, vSphere DRS, and storage redundancy. While zero-impact failures like single network cable issues are documented in storage vendor materials, this section examines these failure scenarios:

Common failure scenarios:

- Single-host failure in Frimley data center
- Single-host isolation in Frimley data center
- Storage partition
- Data center partition
- Disk shelf failure in Frimley data center
- Full storage failure in Frimley data center
- Full compute failure in Frimley data center
- Full compute failure in Frimley data center with full storage failure in Bluefin data center
- Complete Frimley data center loss
- Loss of VASA Providers in Frimley data center
- Loss of VASA Providers at both sites

Single Host Failure in Frimley Data Center

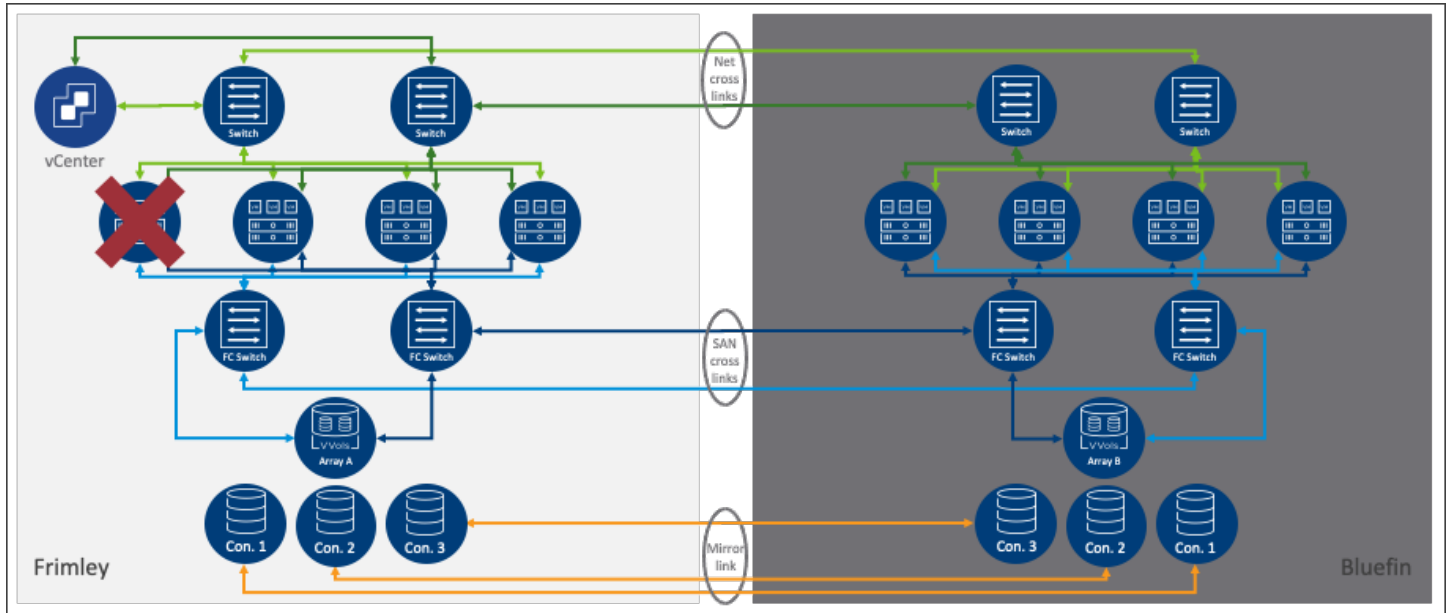


Figure 22 - Single Host Failure Scenario

Figure 22 illustrates a complete host failure scenario in the Frimley data center.

Result

vSphere HA successfully restarted all VMs while maintaining VM-to-host affinity rules.

Explanation

When a host fails, the vSphere HA primary node detects the failure through the absence of network heartbeats from the affected host. The primary node then begins monitoring for datastore heartbeats. Since the host has experienced complete failure, it cannot generate datastore heartbeats, which the vSphere HA primary node also detects as missing. During this validation period, the primary node conducts an additional availability check by pinging the management addresses of the failed host. When all these verification methods prove unsuccessful, the primary node declares the host dead and initiates restart procedures for all protected VMs that were running on the host before contact was lost.

The VM-to-host affinity rules configured at the cluster level are designated as "should rules." vSphere HA respects these VM-to-host affinity rules during the restart process, ensuring VMs restart within their correct site when resources are available.

If hosts within the VM-to-host group lack resources or become unavailable for restarts, vSphere HA may override these "should rules" and restart VMs on any available cluster host, regardless of location or rule constraints. When this occurs, vSphere DRS attempts to correct rule violations during its next activation by migrating VMs to align with their affinity rules. VMware recommends manually triggering vSphere DRS after resolving the failure cause to ensure proper VM placement and prevent performance issues, particularly in active/passive configurations where both reads and writes route through the active storage node.

Single Host Isolation in Frimley Data Center

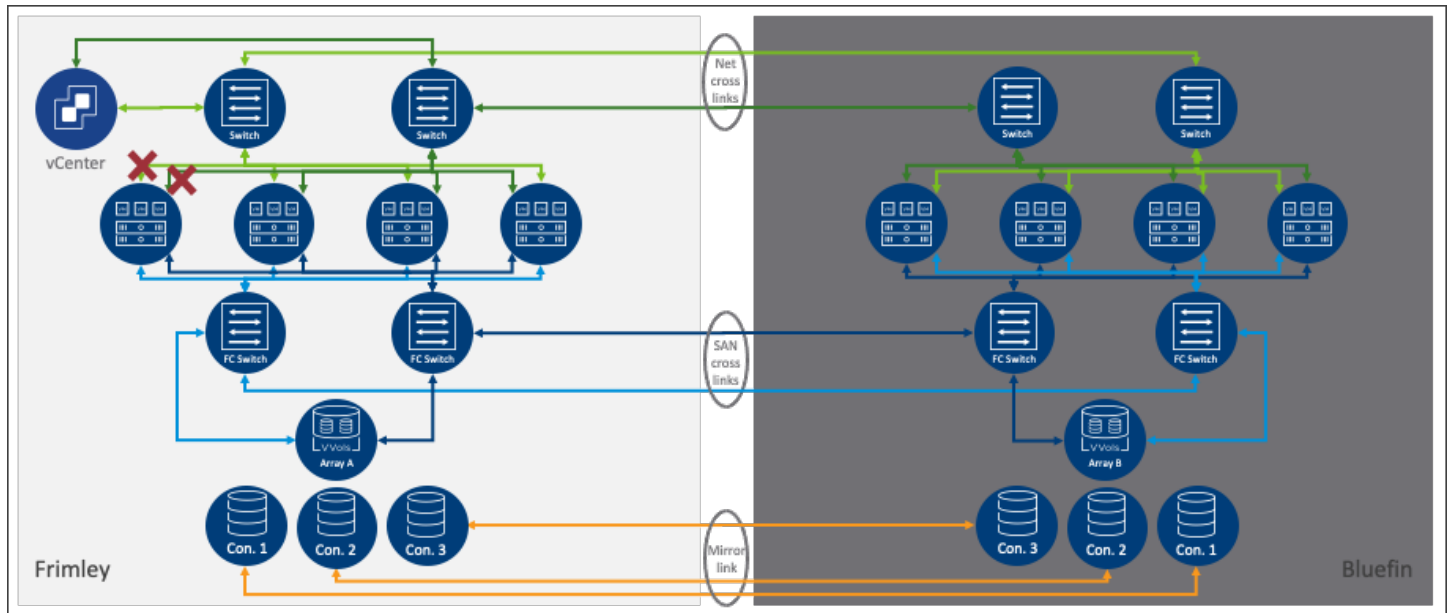


Figure 23 - Single Host Isolation Scenario

Figure 23 illustrates a single host isolation scenario in the Frimley data center.

Result

VMs remain running as isolation response is configured to maintain powered-on state.

Explanation

During host isolation, the vSphere HA primary node first detects the loss of network heartbeats from the affected host. The primary node then monitors datastore heartbeats. Though network-isolated, the host continues generating datastore heartbeats, allowing the primary node to determine the host remains operational but lacks network connectivity. Based on isolation response settings, the affected host can power off VMs, shut them down, or maintain their powered-on state. This isolation response triggers 30 seconds after isolation detection.

VMware recommends configuring isolation response based on business needs and infrastructure design. For most environments, maintaining powered-on VMs represents best practice. Host isolation rarely occurs in properly designed environments due to modern infrastructure redundancy. However, in environments using network storage protocols like iSCSI and NFS with converged networks, VMware recommends a power-off response since network outages may affect both host isolation and datastore communication.

If configured for "power off" or "shut down" instead of the recommended "leave powered on" setting, the vSphere HA primary node restarts VMs on available cluster nodes. While VM-to-host affinity rules use "should rules" configuration, vSphere HA typically respects these preferences, restarting VMs within their designated sites under normal conditions.

Storage Partition

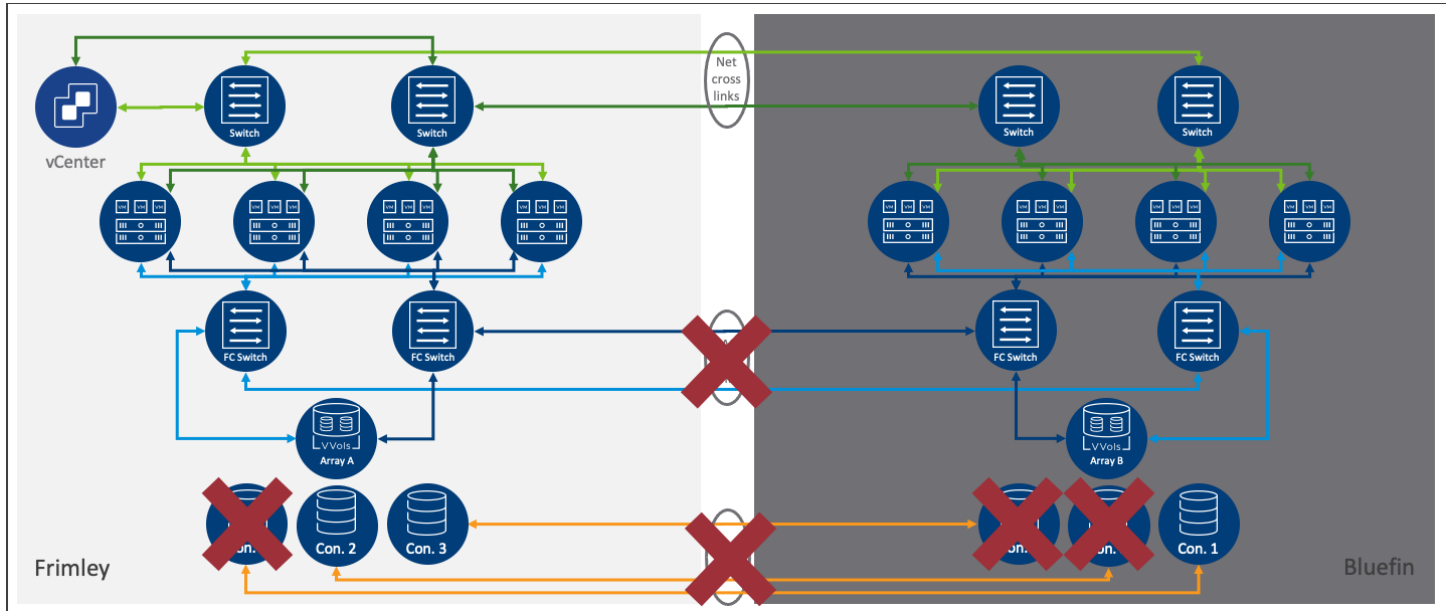


Figure 24 - Storage Partition

Figure 24 illustrates a failure on the storage network between data centers.

Result

VMs continue running without interruption.

Explanation

Each container/datastore maintains defined storage site affinity, aligned with vSphere DRS rules. When storage remains available within a site, VMs experience no impact. For example, container "Con. 1" in Figure 24 has Bluefin affinity and fails over to the Bluefin site. Since VMs with Bluefin affinity should already run on Bluefin-affiliated containers, they continue operating normally.

If a VM runs contrary to its affinity rules (for example, operating from Frimley while its disk resides on a Bluefin-affiliated datastore), it loses I/O capability during inter-site storage partition. This triggers a Permanent Device Loss (PDL) condition, prompting vSphere HA to restart the VM according to configured PDL response settings.

Data Center Partition

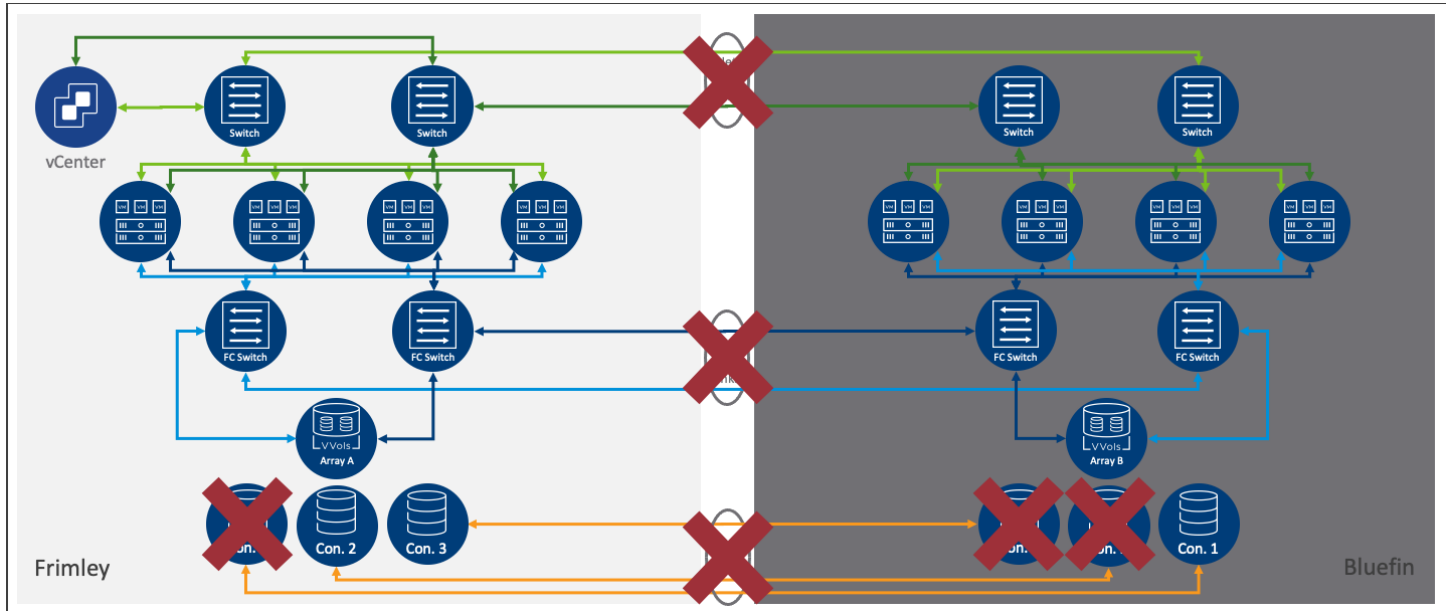


Figure 25 – Data Center Partition

Figure 25 illustrates a scenario where the Frimley data center has become isolated from the Bluefin data center.

Result

VMs continue running without interruption.

Explanation

During complete data center isolation, the scenario combines storage partition and host isolation behaviors. Properly configured vSphere DRS rules prevent VM impact.

vSphere HA follows this process during cluster partition:

The Frimley vSphere HA primary node detects Bluefin host unreachability through missing network heartbeats and storage heartbeats, as inter-site storage connections have failed. VMs with Frimley affinity continue running normally. The system cannot restart other VMs since Bluefin-affiliated datastores remain inaccessible to Frimley hosts.

Meanwhile, Bluefin hosts initiate primary election after losing connection to the Frimley primary. The new Bluefin primary evaluates VM status and restart needs. Since Bluefin-affiliated VMs continue running and Frimley datastores remain inaccessible, no restarts occur.

If VM-to-host affinity rules were violated, this sequence occurs:

1. A Frimley-affiliated VM running in Bluefin loses datastore access, triggering PDL and lock file updates cease, stopping all vVol datastore I/O for the VM. This results in the VM being unable to write to or read from disk as the datastore enters PDL state.
2. Frimley vSphere HA restarts the VM since hosts in Frimley cannot detect its Bluefin instance due to inability to heartbeat to shared storage.
3. A Frimley host acquires the VMDK lock and powers on the VM since the datastore is only available to Frimley.
4. With proper PDL response configuration (Power off and restart VMs), the VM powers off in Bluefin after PDL detection.

Incorrect PDL configuration may cause duplicate running VMs for these reasons:

- Network heartbeat from the VM's host is missing due to site connection loss
- Datastore heartbeat is absent due to lost shared storage access between sites
- Management address ping fails from lost site connectivity
- The Frimley primary attempts restart after losing Bluefin communication, unable to detect actual VM state

Upon site reconnection, this creates a VM split-brain scenario with duplicate MAC addresses. However, only the VM accessing its files continues running after vSphere HA detection.

VMware recommends monitoring vSphere HA and DRS cluster rules alignment with datastore site affinity to prevent unnecessary downtime during these scenarios.

Disk Shelf Failure in Frimley Data Center

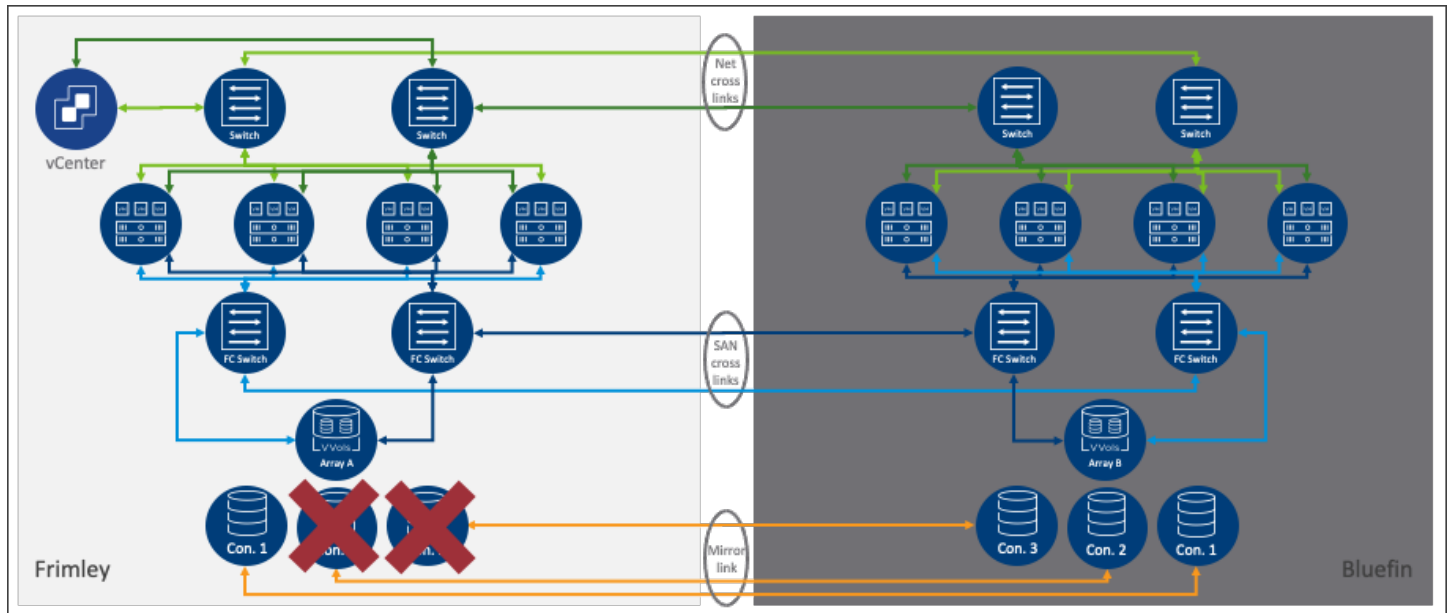


Figure 26 - Disk Shelf Failure

Figure 26 illustrates the loss of a disk shelf in the Frimley data center where “Con 2” and “Con 3” are impacted.

Result

VMs continue running without interruption.

Explanation

When a disk shelf fails in the Frimley data center, the storage processor detects the failure, marks the container unavailable, and flags array paths from Frimley as Permanent Device Loss (PDL). The hosts seamlessly transition from these paths to operational paths connected to the Bluefin array, accessing the mirror copy in the Bluefin data center. VMs experience only a brief I/O response time increase during this transition.

Full Storage Failure in Frimley Data Center

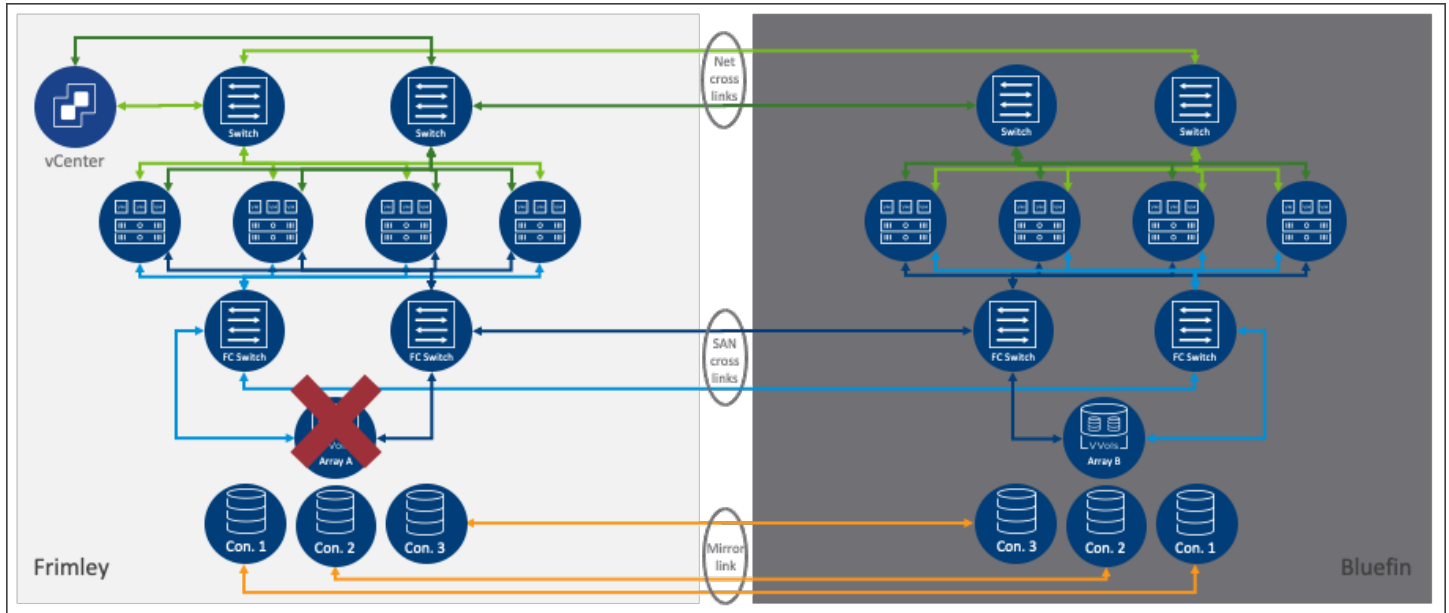


Figure 27 - Full Storage Failure

Figure 27 illustrates a complete storage failure scenario in the Frimley data center.

Result

VMs continue running without interruption.

Explanation

During complete storage failure in the Frimley data center, automated failover depends on array configuration. Without an external witness, manual 'take over' commands may be required to make containers/datastores available in the Bluefin site, as the Bluefin array cannot differentiate between mirror link loss and array failure. The array must assume containers with Frimley affinity continue operating there unless an external witness confirms losing connection to the Frimley array. This scenario assumes witness usage, enabling automatic failover within 15 seconds.

The failover appears seamless to VMs, with only a brief I/O and control path pause up to 15 seconds. Operations resume immediately when the Bluefin array brings containers online, allowing both I/O and control path operations to continue with the operational Bluefin array.

vSphere HA remains inactive during this failure type. Though datastore heartbeats briefly pause, the vSphere HA primary agent only checks these heartbeats after three seconds of missing network heartbeats. Since network heartbeats continue throughout the storage failure, vSphere HA requires no restart intervention.

Permanent Device Loss

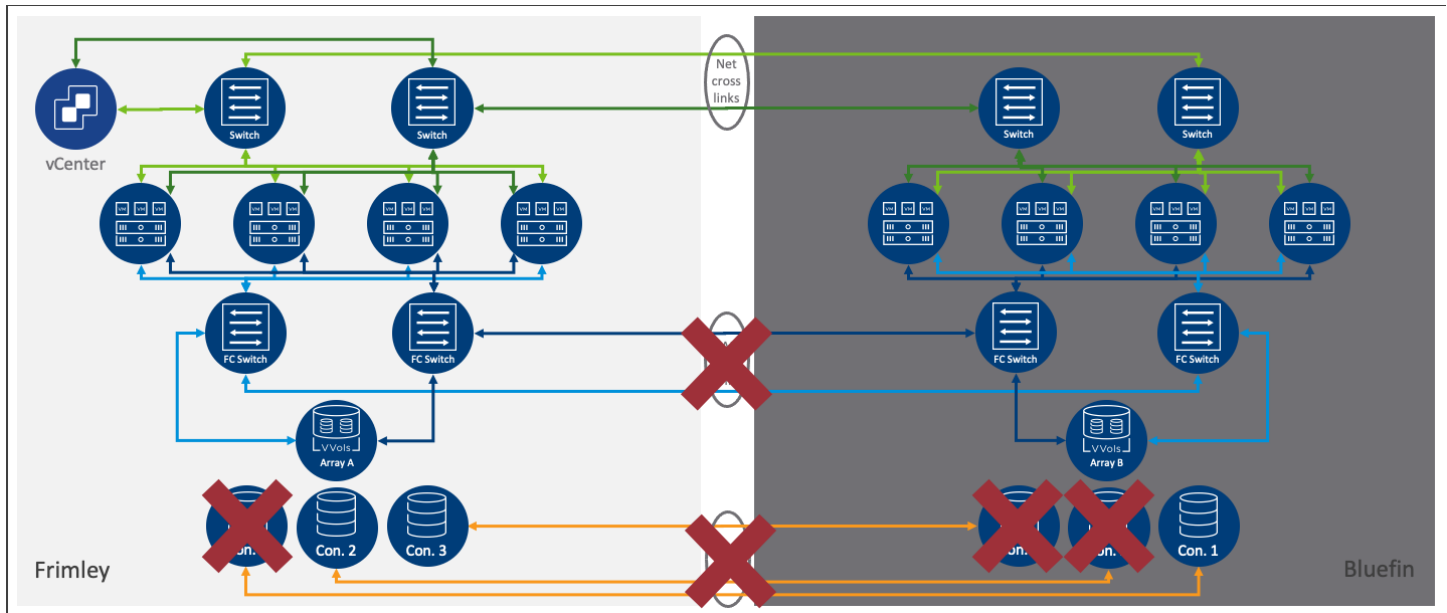


Figure 28 - Permanent Device Loss

Figure 28 illustrates a scenario in the Frimley data center where one or more datastores enter PDL conditions due to loss of SAN and mirroring links. Hosts in Frimley see "Con 1" in a PDL state, and hosts in Bluefin see "Con 2" and "Con 3" in PDL.

Result

VMs either continue running or restart through VM Component Protection (VMCP) as needed.

Explanation

In stretched storage configurations, Permanent Device Loss (PDL) typically occurs from misconfigurations or when VMs run outside their optimal site during errors. When PDL occurs, VMCP halts VMs running on Frimley hosts using datastore "Con 1," then vSphere HA restarts them on Bluefin hosts where the datastore remains available. Similarly, VMs using "Con 2" or "Con 3" on Bluefin hosts undergo power off and restart through vSphere HA.

As shown in Figure 28, vVol datastores should fail over to their affinity-aligned arrays. This may require no restarts if VMs run on their expected sites. However, during combined array and SAN cross-link failures, VMCP must restart PDL-affected VMs.

VMware recommends configuring "Response for Datastore with Permanent Device Loss (PDL)" to "Power off and restart VMs," as noted in the vSphere HA section. This setting ensures appropriate PDL response and enables datastores to exit PDL state when conditions resolve. Figure 29 demonstrates this configuration.

Edit Cluster Settings | Cluster ✕

vSphere HA

Failures and responses | Admission Control | Heartbeat Datastores | Advanced Options

You can configure how vSphere HA responds to the failure conditions on this cluster. The following failure conditions are supported: host, host isolation, VM component protection (datastore with PDL and APD), VM and application.

Enable Host Monitoring i

> Host Failure Response	Restart VMs v
> Response for Host Isolation	Disabled v
> Datastore with PDL	Power off and restart VMs v
> Datastore with APD	Power off and restart VMs - Conservative restart policy v
> VM Monitoring	Disabled v

CANCEL
OK

Figure 29 - VMCP PDL Configuration

Full Compute Failure in Frimley Data Center

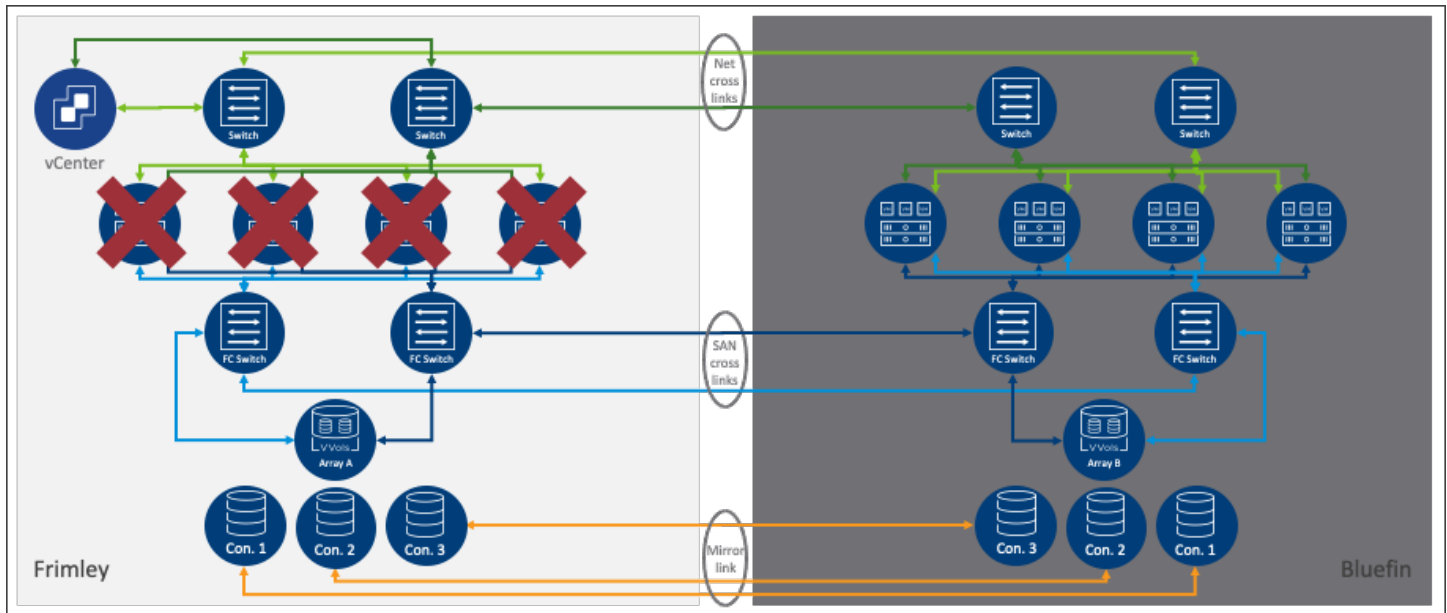


Figure 30 - Full Compute Failure

Figure 30 illustrates a complete compute failure scenario in the Frimley data center.

Result

All VMs successfully restart in the Bluefin data center.

Explanation

When complete compute failure occurs in Frimley, the vSphere HA primary node location affects recovery. After Bluefin hosts detect missing network heartbeats from Frimley, they initiate primary election, establishing a new primary within approximately 20 seconds. This new primary evaluates host failures and impacted VMs, then initiates VM restarts.

The restart schedule succeeds due to vSphere HA admission control reserving 50% of CPU and memory capacity, as recommended. This configuration ensures sufficient resources for restarting VMs from half the hosts.

vSphere HA can process 32 concurrent restarts per host, minimizing restart latency. VM Overrides enables restart sequencing through five priority levels: lowest, low, medium, high, and highest. VMs restart in descending priority order, with highest-priority VMs starting first.

When Frimley hosts return online, vSphere DRS detects their availability. The initial DRS run corrects only affinity rule violations, while resource imbalance resolves during the next full invocation. Though DRS runs automatically every 5 minutes or during power state changes through vSphere Client, VMware recommends manual DRS activation after problem resolution to verify proper affinity rule compliance.

Loss of Frimley Data Center

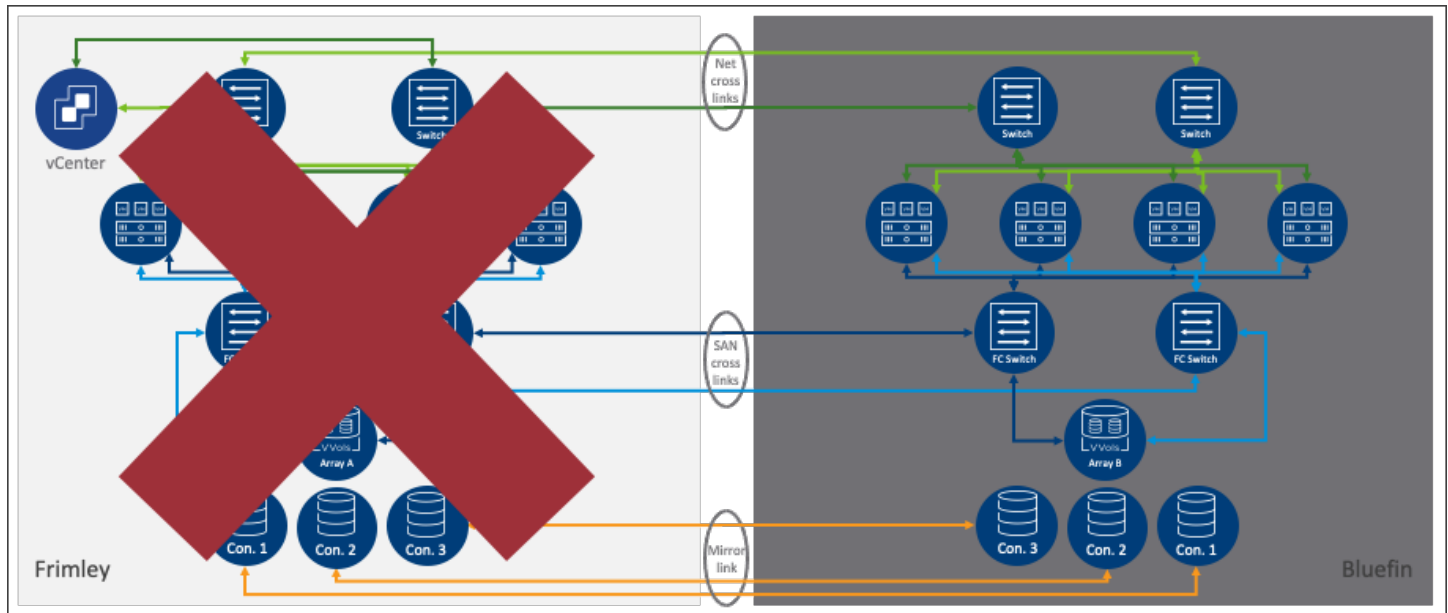


Figure 31 - Full Data Center Failure

Figure 31 illustrates a scenario where the Frimley data center is offline.

Result

All VMs successfully restart in the Bluefin data center.

Explanation

After losing contact with the Frimley vSphere HA primary, Bluefin hosts elect a new primary node. Since the Frimley storage system has failed, recovery requires a 'take over' command at the surviving site, though this may occur automatically with an external witness.

The new vSphere HA primary accesses per-datastore files to identify protected VMs, then attempts to restart VMs not running on surviving Bluefin hosts.

Note: vSphere HA discontinues VM restart attempts after 30 minutes by default. If storage administrators do not execute the 'take over' command within this window, vSphere administrators must manually restart VMs once storage becomes available.

Loss of VASA Providers in Frimley Data Center

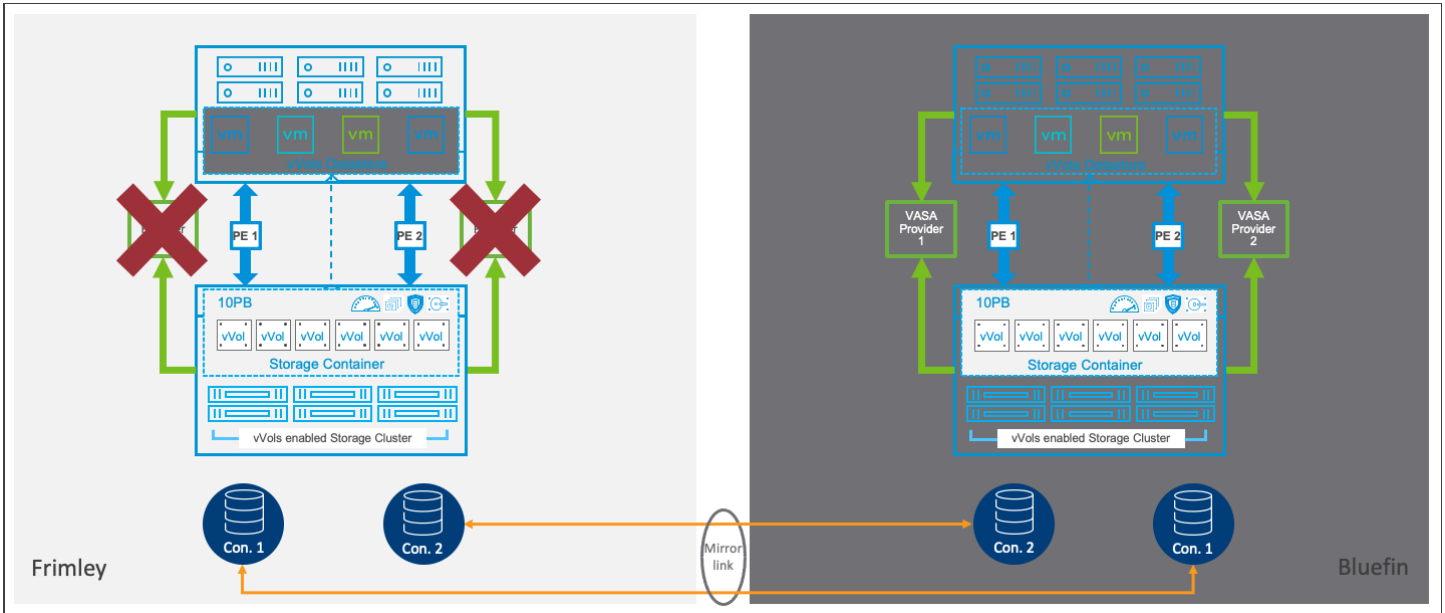


Figure 32 - Loss of VASA Providers in One Site

Figure 32 illustrates a scenario where the Frimley data center VASA providers are not available.

Result

VMs continue running without interruption.

Explanation

The Frimley data center experiences loss of both VASA Providers for its array. Stretched containers/datstores maintain functionality since hosts across both sites retain access to the Bluefin VASA Providers, which handle VASA operations equally well for stretched containers.

VMware recommends dual VASA Provider support for all containers, stretched or not. This requires either two embedded VASA Providers per array (for legacy and unstretched containers) or two external (Common) VASA Providers for arrays. Stretched containers specifically require two to four VASA Providers per container.

Figure 32 illustrates two VASA Providers per array, though SAN cross-links remain hidden. Network cross-links, while not shown, operate to enable ESXi host access to all available VASA Providers. The arrays present two containers, each with a Protocol Endpoint (PE) maintaining I/O path accessibility. These PEs stretch across sites for vVol Stretched Storage Containers.

Loss of VASA Providers in Both Data Centers

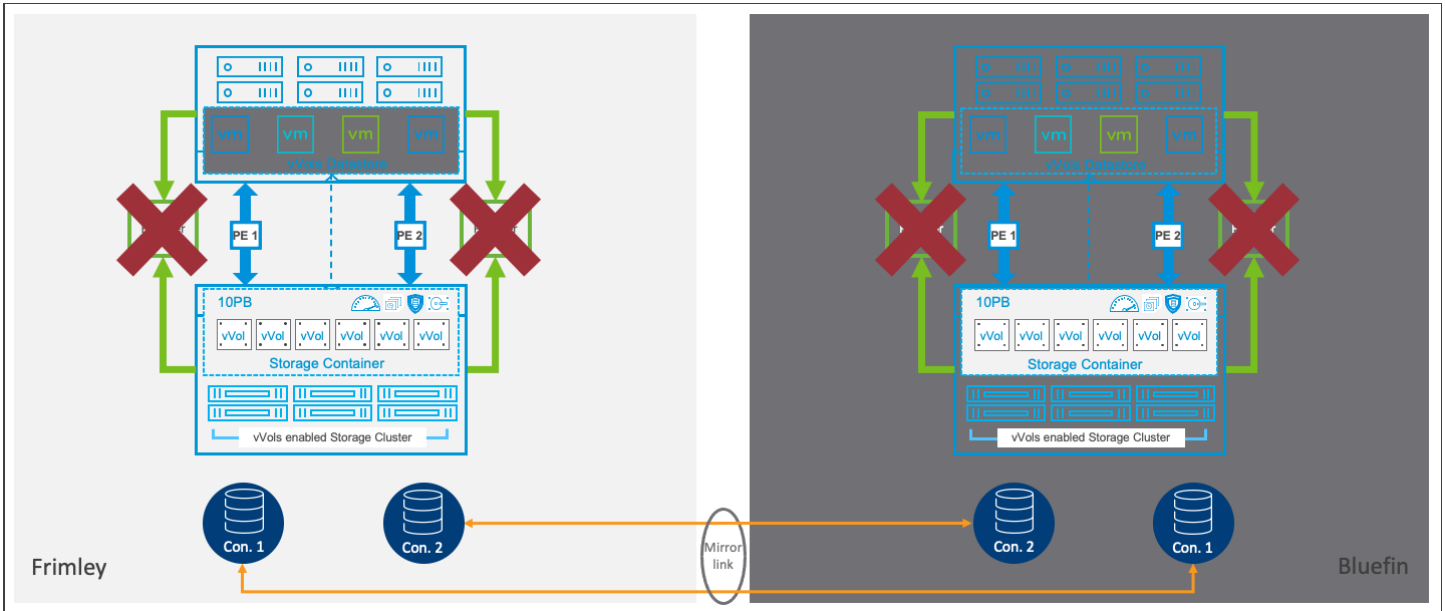


Figure 33 - Loss of VASA Providers in Both Sites

Figure 33 illustrates a scenario where all VASA providers are unavailable.

Result

VMs continue running with potential operational impact.

Explanation

When both data centers lose all VASA Providers, stretched containers/datastores maintain I/O paths but lose VASA operations capability, affecting workflows like snapshots, VM power operations, and vMotion.

VMs engaged in VASA-dependent operations (like partial snapshot creation) may become non-operational, while other VMs continue running. Though I/O paths remain functional, all hosts mark datastores as All Paths Down (APD) to prevent vSphere HA from starting new VMs.

VMCP detects VMs running on APD storage and attempts relocation. However, since no hosts can access the datastores, VMCP's conservative APD restart policy keeps VMs running in their current location. VMs only power off once VMCP identifies a host with datastore access and sufficient CPU and memory resources.

This scenario, with all stretched container VASA Providers unavailable but functioning I/O paths, represents a critical but not complete failure state.

Summary & Conclusion

Properly configured stretched clusters enhance resiliency and enable cross-site workload movement. This paper clarifies how VMware vSphere HA and vSphere DRS handle various failure scenarios in stretched cluster environments, providing configuration guidance for optimal performance.

The discussion emphasizes site affinity importance, vSphere HA and DRS cluster rules and groups functionality, and their interactions. Maintaining these configurations over time ensures cluster reliability and predictable behavior during failure events.

