# VMware® High Availability (VMware HA): Deployment Best Practices

VMware® vSphere™ 4.1

**vm**ware®

This paper describes best practices and guidance for properly deploying VMware® High Availability (VMware HA) in VMware vSphere™ 4.1 ("vSphere"), including discussions on proper network and storage design, and recommendations on settings for host isolation response and admission control. This best-practice paper is not designed to replace existing product documentation. Instead it is designed to provide guidance based on the experience of VMware HA Product Management, Technical Marketing and Professional Services staff. Although there are many configuration options to choose from and architecture choices that can be made, the recommendations in this paper are based on the most reliable and resilient possible choices for the majority of cluster deployments. Customers with advanced needs, such as campus clustering and other nonstandard deployments, should review the *vSphere Availability Guide* for more information, or contact VMware Professional Services for assistance.

# Introduction

Downtime, whether planned or unplanned, brings with it considerable costs. Solutions to ensure higher levels of availability have traditionally been very costly, hard to implement and difficult to manage.

VMware vSphere makes it simpler and less expensive to provide higher levels of availability for important applications. With vSphere, organizations can easily increase the baseline level of availability provided for all applications, as well as provide higher levels of availability more easily and cost-effectively. vSphere makes it possible to reduce both planned and unplanned downtime. The revolutionary VMware vMotion™ (vMotion) capabilities in vSphere make it possible to perform planned maintenance with zero application downtime. VMware HA, a feature of vSphere, specifically reduces unplanned downtime by leveraging multiple VMware ESX® and VMware ESXi™ hosts configured as a cluster, to provide rapid recovery from outages as well as cost-effective high availability for applications running in virtual machines.

VMware HA protects application availability in the following ways:

• It protects against hardware failure by restarting the virtual machines on other hosts within the cluster.
• It protects against operating system failure by continuously monitoring a virtual machine and resetting it in the event that an operating system (OS) failure is detected.

Unlike other clustering solutions, VMware HA provides the infrastructure to protect all workloads within the cluster:

• You do not need to install additional software within the application or virtual machine. VMware HA protects all workloads. After VMware HA is configured, no actions are required to protect new virtual machines. They are automatically protected.
• You can combine VMware HA with VMware Distributed Resource Scheduler (VMware DRS) to protect against failures and to provide load balancing across the hosts within a cluster.

VMware HA has several advantages over traditional failover solutions, including:

• Minimal setup
• Reduced complexity (e.g., no need for quorum disks)
• Reduced hardware cost and setup
• Increased application availability without the expense of additional idle failover hosts or the complexity of maintaining identical hosts for failover pairs
• VMware DRS and vMotion integration

Refer to the *vSphere Availability Guide* for more information on VMware HA basics including how to create HA clusters, how VMware HA works, and explanations on configuration procedures.

## Design Principles for High Availability

The key to architecting a highly available computing environment is to eliminate single points of failure. With the potential of occurring anywhere in the environment, failures can affect both hardware and software. Building redundancy at vulnerable points will help mitigate downtime caused by failures, including failures at the following layers:

• Server components such as network adapters and host bus adapters (HBAs)

• Servers, including blades and blade chassis

• Networking components

• Storage arrays and storage networking

## Host Selection Considerations

Overall vSphere availability starts with proper host selection, including such items as redundant power supplies, error-correcting memory, remote monitoring and notification, and so on. Consideration also should be given to removing single points of failure in host locations, including distributing hosts across multiple racks or blade chassis to eliminate the ability for rack or chassis failure to impact an entire cluster.

When deploying a VMware HA cluster, it is a best practice to build the cluster out of identical server hardware. Using identical hardware provides a number of key advantages:

• It simplifies configuration and management of the servers using host profiles.

• It increases the ability to handle server failures and reduces resource fragmentation. Using drastically different hardware will lead to an unbalanced cluster, as described in the "Admission Control" section. By default, VMware HA prepares for the worst-case scenario, in that the largest host in the cluster will fail. Therefore, in order to handle the worst case, more resources across all the hosts must be reserved, making them essentially unusable.

Cluster sizing also has a major impact on overall availability and level of consolidation possible. Smaller clusters require a larger percentage of resources available in order to facilitate their ability to handle failures (covered later in the "Admission Control" section). Larger clusters can get more complex to manage from a networking-and-storage perspective. Ideal cluster size for the majority of environments is between 6 and 10 nodes. This allows for ample spare capacity to handle failures, without becoming overly complex.

## Host Placement Considerations

VMware HA employs the concept of primary and secondary hosts. Primary hosts are responsible for making failover decisions for the cluster; secondary hosts only execute on these decisions. The first five ESX hosts that join a VMware HA cluster are designated as the primary VMware HA hosts. All subsequent hosts are designated as secondary hosts. Any host that joins the VMware HA cluster must communicate with an existing primary host to complete its configuration, except for when the first host is added to the cluster. At least one primary host must be functional for VMware HA to operate correctly. If all primary hosts are unavailable—that is, not responding—no hosts can successfully be configured for VMware HA, and no failovers can take place.

One of the primary hosts is also designated as the active primary host. Its responsibilities include:

• Deciding where to restart virtual machines

• Keeping track of failed restart attempts

• Determining when it is appropriate to keep trying to restart a virtual machine

See the *vSphere Availability Guide* for more details.

There are three scenarios that can cause a primary host in a VMware HA cluster to become unavailable: entering maintenance mode; powering off; and failure.

**Entering maintenance mode:** If a primary host enters maintenance mode, a secondary host, if available, will be chosen and promoted to become a new primary host.

**Failure and power-off:** If a primary host is powered off or fails, the total number of primary hosts is reduced by one. No secondary host is promoted to become a new primary host in this case. Therefore, VMware HA will be able to handle no more than four consecutive failures of primary hosts.

If all the primary hosts fail, the cluster loses VMware HA protection. In order to prevent loss of all primary hosts due to a single failure, it is highly recommended that the cluster be architected so no more than four hosts in a given cluster are placed in the same server rack or blade chassis. Although it would appear possible to add hosts to the cluster in a specific order so as to prevent more than four hosts from being designated as primary in a single rack or chassis, shifting a host to maintenance mode causes a new primary to be elected. For this reason it is critical to not allow more than four ESX hosts in the same cluster to exist in the same failure domain, as the primary role might shift around over time.

## Networking Design Considerations

Best practices network design falls into two specific areas: increasing resiliency of "client side" networking to ensure access from external systems to workloads running in vSphere; and increasing resiliency of communications used by VMware HA itself.

## General Networking Guidelines

The following suggestions are general best practices for configuring networking for improved availability.

• Configuring switches. If the physical network switches that connect your servers support the PortFast or an equivalent setting, enable it. This setting prevents a host from incorrectly determining that a network is isolated during the execution of lengthy Spanning Tree algorithms on boot. For more information on this option, refer to the documentation provided by your networking switch vendor.

• Disable host monitoring (using VMware vCenter™ Server, deselect the "Enable Host Monitoring" check box in the cluster's Settings dialog box, VMware HA->Host Monitoring Status) when performing any networking maintenance that might disable all heartbeat paths between hosts and cause isolation responses.

• Use DNS for name resolution rather than the error-prone method of manually editing the local `/etc/hosts` file on ESX hosts. If you do edit `/etc/hosts,` you must include both long and short names. VMware HA modifies the `/etc/hosts` file if the information in it is wrong or needs to be added.

• Use consistent port names on VLANs for virtual machine networks on all ESX hosts in the cluster. Port names are used to determine compatibility of the network by virtual machines. VMware HA will check whether a host is compatible before initiating a failover. If there are no hosts with matching port group names available, no failover is possible. Use of a documented naming scheme is highly recommended. Issues with port naming can be completely mitigated by using VMware distributed virtual switches.

• Host firewalls. On ESX/ESXi hosts, VMware HA needs and automatically opens the following firewall ports:

 – Incoming port: TCP/UDP 8042-8045

 – Outgoing port: TCP/UDP 2050-2250

Configure redundant networking from ESX hosts to network switching hardware if possible. Using network adapter teaming can enhance overall network availability as well as increase overall throughput.

## Setting Up Redundancy for VMware HA Networking

Networking redundancy between cluster hosts is absolutely critical for VMware HA reliability. Redundant management networking enables the reliable detection of failures and prevents isolation conditions from occurring. For this document, "management network" refers to the service console network on ESX 4.1 or the VMkernel network selected for use as a management network on ESXi 4.1.

You can implement network redundancy at the network adapter level or at the management network level. In most implementations, network adapter teaming provides sufficient redundancy, and it is explained here. If you wish to add additional redundancy, see the *ESX/ESXi Configuration Guide* for more information on setting up redundant management networks.

### Network Adapter Teaming and Management Networks

Using a team of two network adapters connected to separate physical switches can improve the reliability of the management network. Because servers connected to each other through two network adapters—and through separate switches—have two independent paths for sending and receiving heartbeats, the cluster is more resilient. It should be noted, however, that in this setup, only one network adapter is active at one time. As such, a heartbeat can only be sent or received through a single path.

To configure a network adapter team for the management network, configure the vNICs in the vSwitch configuration for the ESX/ESXi host, for active/standby configuration.

Requirements:

• Two physical network adapters
• VLAN trunking

Recommended:

• Two physical switches

The vSwitch should be configured as follows:

• Load balancing = route based on the originating virtual port ID (default)
• Failback = "No"
• vSwitch0: two physical network adapters (for example: vmnic0 & vmnic2)
• Two port groups (for example, vMotion and management)

In this example, the management network runs on vSwitch0 as active on vmnic0 and standby on vmnic2. The vMotion network runs on vSwitch0 as active on vmnic2 and standby on vmnic0.

Each port group has a VLAN ID assigned and it runs dedicated on its own physical network adapter; only in the case of a failure is it switched over to the standby network adapter. Failback is set to "No" because in the case of physical switch failure and restart, ESX might incorrectly recognize that the switch is back online when its ports first come online. In reality, the switch might not be forwarding on any packets until it is fully online. However, when Failback is set to "No" and an issue arises, both your management network and vMotion network will be running on the same network adapter and will keep on running until you manually intervene.

The following diagram depicts this scenario:

## Storage Design Considerations

To maintain a constant connection between an ESX/ESXi host and its storage, ESX/ESXi supports multipathing. Multipathing is a technique that enables you to use more than one physical path that transfers data between the host and an external storage device.

In case of a failure of any element in the SAN network, such as an adapter, switch or cable, ESX/ESXi can switch to another physical path that does not use the failed component. This process of path switching to avoid failed components is known as *path failover*.

In addition to path failover, multipathing provides load balancing. Load balancing is the process of distributing I/O loads across multiple physical paths. Load balancing reduces or removes potential bottlenecks.

For Fibre Channel SAN configurations, multipathing setup is very specific to the HBA, switch and array components chosen. Please see the *Fibre Channel Configuration Guide* for more information.

For iSCSI configurations, both ESX and ESXi support creating a second iSCSI initiator to enable multipathing configurations. See the *iSCSI SAN Configuration Guide* for more details on setting up multiple iSCSI initiators.

VMware strongly recommends multiple paths to storage for maximum resiliency for all block storage configurations.

## Understanding Host Isolation

A key mechanism within VMware HA is the ability for a host to detect when it has been isolated from the rest of the cluster. An ESX host is considered to be isolated when it loses the ability to exchange heartbeat or status information with other hosts in the VMware HA cluster via the service console network—or management network on ESXi 4.1—and it cannot communicate with its configured "isolation addresses." The key concept to understand about host isolation is that the host in question has detected that it is likely dealing with a communication issue, and that other healthy hosts might very well exist. However, other hosts in the cluster cannot detect whether a host has been isolated from the network or has been powered off, failed or crashed. They only detect that the host is unreachable, and therefore must execute the standard response of failing over virtual machines to other healthy hosts. By default, a host detecting that it is isolated will shut down running virtual machines on the premise that other healthy hosts remain in the cluster, and that they can better provide service for the virtual machines. This response can be tuned in specific cases, as explained later.

Host isolation happens only in cases where management network connectivity between hosts in the cluster has failed but a host remains running. Although host isolation response enables some level of tuning to allow for proper failover of virtual machines to other, more viable hosts, an outage will still occur during failover. It is strongly recommended that reliable, redundant networking be configured between hosts for VMware HA use, to prevent the possibility of this occurrence.

### Host Isolation Detection

The mechanism that detects host isolation is based on the presence of heartbeats. When a host receives no heartbeats from any of the other nodes for 13 seconds (default setting), it will ping the configured isolation address(es).

The default isolation address is the gateway specified for the management network, but there is a possibility of specifying one or more additional isolation addresses with the advanced settings: das.isolationaddress[X] (where X is 1–10) and das.usedefaultisolationaddress (set to "False" if default gateway address is nonpingable). We recommend setting at least one additional isolation address, such as a router close to the hosts.

When isolation has been confirmed, meaning no heartbeats have been received and VMware HA was unable to ping any of the isolation addresses, VMware HA will execute the host isolation response. The host isolation response options are: 1) shut down, 2) power off, and  3) remain powered on. These are described in the following section.

If only one heartbeat is received or only a single isolation address can be pinged, the isolation response will not be triggered. Keep in mind that primary VMware HA hosts send heartbeats to primary and secondary hosts, whereas secondary hosts send heartbeats only to primary hosts. This means that in rare cases where all secondary nodes are on one side of a split of the management network, the secondary nodes could all detect themselves to be isolated. This is another good example of the need for redundant management networking.

### Host Isolation Response

The action taken by VMware HA for virtual machines running on a host when the host has lost its ability to communicate with other hosts over the management network and cannot ping the isolation addresses is called *host isolation response*. Host isolation does not necessarily mean that the virtual machine network is down—only that the management network, and possibly others, are down.

Host isolation response configures how an ESX/ESXi host will respond when it detects that it is isolated. The restarting of virtual machines by VMware HA on other hosts in the cluster in the event of a host isolation or host fault is dependent on the "host monitoring" setting for that particular host. If host monitoring is disabled, restart of virtual machines is also disabled on other hosts following a host failure or isolation. Essentially, a host will always perform the programmed host isolation response when it detects that it is isolated. The host monitoring setting determines whether virtual machines will be restarted elsewhere following this event.

There are three settings for a VMware HA cluster's host isolation response: leave powered on; power off; and shut down. For vSphere 4.0 and 4.1, the default host isolation response setting for the VMware HA cluster is "shut down." This setting can also be configured independently for each virtual machine in the cluster.

### Shut down

We recommend that you use the default "shut down" setting because it gracefully shuts down virtual machines in the event of host isolation. VMware HA uses VMware Tools to initiate this shutdown; as such, VMware Tools should be installed within the guest operating system if the "shut down" isolation response is selected. When a guest shutdown is initiated, VMware HA will wait up to 5 minutes—by default, configurable, using advanced options—for the guest OS to shut down before it powers off the virtual machine. If the guest shuts down cleanly, no power-off will occur.

VMware strongly recommends the "shut down" setting for all environments. Assuming a reliable, redundant service console network, cases of host isolation should be nearly nonexistent. In the unlikely event of an isolation occurrence, "shut down" makes sure that the virtual machine is properly shut down, with minimal time added to the failover.

### Power off

Setting the host isolation response to "power off" enables the user to trade a hard power-off of the virtual machine for reduced failover time. Since "power off" is similar to pulling the plug, there is also a chance—higher for Windows-based guest operating systems—that the guest operating systems might be inconsistent and that crash recovery within the guests might be invoked on restart.

Using the "power off" option essentially trades restart time for possible problems with integrity in the virtual machine. With a properly designed service console network, actual host isolation, while other machines remain viable in the cluster, is very small. It is strongly recommended that this option not be used unless the administrator completely understands the trade-off. In other words, if you have enough isolation problems that you feel the need to tune for faster restart, you should consider enhancing your service console networking rather than performing isolation tuning

### Leave powered on

The "leave powered on" setting does exactly as the name implies. In a host isolation situation, the virtual machine is left powered on. This option is typically used only on rare configurations where the network configuration might lead to a loss of the management network while retaining access to storage, and the user desires the system to continue processing compute operations, such as a financial calculation environment where uptime of the virtual machine is more important than client/customer access. Configuring "leave powered on" should be considered only after carefully examining the likelihood of failure of a properly designed management network, and the environment-specific needs of the virtual machines in question.

### How VMware HA Avoids Split-Brain Scenarios

In any host isolation response setting, the remaining hosts in the VMware HA cluster will always try to restart the isolated host's virtual machines. In situations where virtual machines still hold their VMFS file locks, such as a break of the service console network leaving hosts running and heartbeating each other on both sides. VMFS locking prevents any true split-brain scenario. However, in the case where "leave powered on" is selected as the isolation response and both the service console network and the storage network are isolated, a split-brain scenario could occur. In these situations, ESX will detect that a virtual machine is active but the VMFS file lock is owned by another host, and it will act upon this. ESX will automatically power-off the virtual machine that has lost access to its VMFS files.

### Host Isolation and VMware Fault Tolerance

Host isolation responses are not performed on virtual machines enabled with VMware Fault Tolerance (VMware FT). The rationale is that the primary and secondary fault-tolerant virtual machine pairs are already communicating via the fault tolerance logging network. They either will continue to have network connectivity and continue to function, or they have lost network and are not heartbeating over the fault tolerance logging network; one of them will then take over as the primary fault-tolerant virtual machine. Because VMware HA does not offer better protection than that, it skips fault-tolerant virtual machines when initiating host isolation response.

## Admission Control

VMware vCenter Server uses VMware HA admission control to ensure that sufficient resources in the cluster are reserved for virtual machine recovery in the event of host failure. Admission control will prevent the following if they encroach into the resources reserved for virtual machines restarted due to failure:

• The power-on of new virtual machines

• Changes of virtual machine memory or CPU reservations

• A vMotion of a virtual machine into the cluster from another cluster

This mechanism is highly recommended in order to guarantee the availability of virtual machines. Beginning with vSphere 4.0, VMware HA offers three configuration options for choosing your admission control strategy:

• Host Failures Cluster Tolerates policy (default): VMware HA ensures that a specified number of hosts can fail and sufficient resources remain in the cluster to fail over all the virtual machines from those hosts. VMware HA uses a concept called slots to calculate both available and needed resources for a failing over virtual machines from a failed host.

• Percentage of cluster resources reserved as failover spare capacity: VMware HA ensures that a specified percentage of aggregate cluster resources are reserved for failover. This policy is recommended for situations where virtual machines must be hosted with significantly different CPU and memory reservations in the same cluster, or where they have differently sized hosts in terms of CPU and memory capacity.

• Specify a failover host: VMware HA designates a specific host as the failover host. When a host fails, VMware HA attempts to restart its virtual machines on the specified failover host.

The best-practices recommendation from VMware staff for admission control is as follows.

• Select the "percentage of cluster resources reserved" for admission control. This policy offers the most flexibility in terms of host and virtual machine sizing. In most cases, a simple calculation of 1/N, where N = total nodes in the cluster will yield adequate sparing.

• Ensure that all cluster hosts are sized equally. An "unbalanced" cluster results in excess capacity being reserved to handle failure of the largest possible node.

- Attempt to keep virtual machine sizing requirements similar across all configured virtual machines. The Host Failures Cluster Tolerates policy uses a notion of *slot sizes* to calculate the amount of capacity needed to be reserved for each virtual machine. The slot size is based on the largest reserved memory and CPU needed for any virtual machine. Mixing virtual machines of largely different CPU and memory requirements will cause the slot size calculation to default to the largest possible for all virtual machines, limiting consolidation. See the *vSphere Availability Guide* for more information on slot size calculation and overriding slot size calculation in cases where you must configure different size virtual machines in the same cluster.

With the release of vSphere 4.1, VMware HA now has the added capability to balance virtual machine loading on failover, reducing "resource fragmentation" in a cluster when virtual machines of different CPU and memory requirements are configured. VMware HA will now migrate virtual machines by invoking VMware DRS to create more contiguous slots on a host to increase the chance for larger virtual machines to be restarted. This does not guarantee enough contiguous resources to restart all the failed virtual machines. It simply means that vSphere will make the best effort to restart all virtual machines with the host resources remaining after a failure.

Based on available resources, admission control calculates the capacity required for a failover. In other words, if a host is placed into maintenance mode or it is disconnected, it is taken out of the equation. This might cause issues with powering up virtual machines when systems are down or in maintenance mode. In reality, this is not a problem. Rather, it is vSphere warning you that you do not have adequate capacity to handle an additional host failure in the current cluster configuration.

There are certain implications to using strict admission control to guarantee availability of failover capacity when using VMware Distributed Power Management (VMware DPM). See VMware Knowledge Base article http://kb.vmware.com/kb/1007006 for more information on admission control and VMware DPM.

## Affinity Rules

A virtual machine–host affinity rule specifies whether or not the members of a selected virtual machine VMware DRS group can run on the members of a specific host VMware DRS group. Unlike a virtual machine–virtual machine affinity rule, which specifies affinity, or anti-affinity, between individual virtual machines, a virtual machine–host affinity rule specifies an affinity relationship between a group of virtual machines and a group of hosts. There are "required" rules, designated by "must," and "preferential" rules, designated by "should." See the *vSphere Resource Management Guide* for more details on setting up virtual machine–host affinity rules.

When preferential rules are used to specify affinity of virtual machines to hosts, VMware HA will restart virtual machines from a failed ESX host on any available ESX host in a VMware HA cluster. Then VMware DRS will migrate the virtual machines back to their host groups to preserve the preferential rule if possible.

However, when required rules are used, VMware HA will restart virtual machines only on an ESX host in the same host VMware DRS group. If no available hosts are in the host VMware DRS group, or if the hosts are resource constrained, the restart will fail.

**Recommendation:** Avoid widespread use of virtual machine–host affinity rules. These rules, especially required rules, restrict placement of virtual machines, so take extra care in deploying virtual machine–host affinity rules.

# Summary

VMware HA—along with new VMware vSphere 4.1 capabilities of clusters and resource pools and tighter integration with other VMware tools such as VMware vCenter Server, vMotion, and VMware DRS—greatly simplifies virtual machine provisioning, resource allocation, load balancing and migration, while also providing an easy-to-use, cost-effective, high-availability and failover solution for applications running in virtual machines. Using VMware vSphere 4.1 and VMware HA helps to eliminate single points of failure in the deployment of business-critical applications in virtual machines. At the same time, it enables you to maintain other inherent virtualization benefits such as higher system utilization, closer alignment of IT resources with business goals and priorities, and more streamlined, simplified and automated administration of larger infrastructure installations and systems.