



Deep dive

ONTAP Select

NetApp
December 11, 2020

This PDF was generated from https://docs.netapp.com/us-en/ontap-select/concept_stor_concepts_chars.html on December 11, 2020. Always check docs.netapp.com for the latest.



Table of Contents

- Deep dive 1
 - Storage 1
 - Networking 35
 - High availability architecture 60
 - Performance 69

Deep dive

Storage

Storage: General concepts and characteristics

Discover general storage concepts that apply to the ONTAP Select environment before exploring the specific storage components.

Phases of storage configuration

The major configuration phases of the ONTAP Select host storage include the following:

- Pre-deployment prerequisites
 - Make sure that each hypervisor host is configured and ready for an ONTAP Select deployment.
 - The configuration involves the physical drives, RAID controllers and groups, LUNs, as well as related network preparation.
 - This configuration is performed outside of ONTAP Select.
- Configuration using the hypervisor administrator utility
 - You can configure certain aspects of the storage using the hypervisor administration utility (for example, vSphere in a VMware environment).
 - This configuration is performed outside of ONTAP Select.
- Configuration using the ONTAP Select Deploy administration utility
 - You can use the Deploy administration utility to configure the core logical storage constructs.
 - This is performed either explicitly through CLI commands or automatically by the utility as part of a deployment.
- Post-deployment configuration
 - After an ONTAP Select deployment completes, you can configure the cluster using the ONTAP CLI or System Manager.
 - This configuration is performed outside of ONTAP Select Deploy.

Managed versus unmanaged storage

Storage that is accessed and directly controlled by ONTAP Select is managed storage. Any other storage on the same hypervisor host is unmanaged storage.

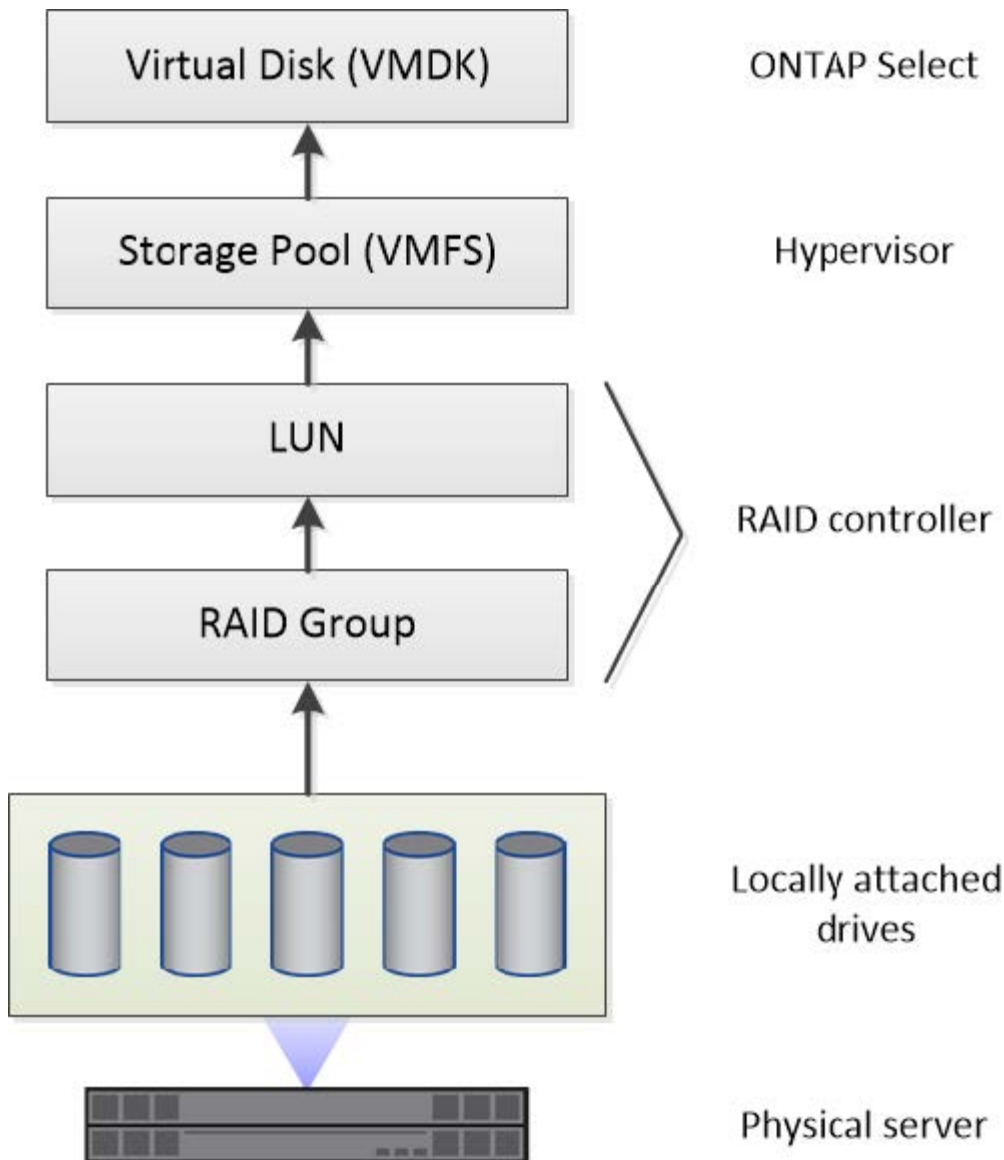
Homogeneous physical storage

All the physical drives comprising the ONTAP Select managed storage must be homogeneous. That is, all the hardware must be the same regarding the following characteristics:

- Type (SAS, NL-SAS, SATA, SSD)
- Speed (RPM)

Illustration of the local storage environment

Each hypervisor host contains local disks and other logical storage components that can be used by ONTAP Select. These storage components are arranged in a layered structure, from the physical disk.



Characteristics of the local storage components

There are several concepts that apply to the local storage components used in an ONTAP Select environment. You should be familiar with these concepts before preparing for an ONTAP Select deployment. These concepts are arranged according to category: RAID groups and LUNs, storage pools, and virtual disks.

Grouping physical drives into RAID groups and LUNs

One or more physical disks can be locally attached to the host server and available to ONTAP Select. The physical disks are assigned to RAID groups, which are then presented to the hypervisor host operating system as one or more LUNs. Each LUN is presented to the hypervisor host operating system as a physical hard drive.

When configuring an ONTAP Select host, you should be aware of the following:

- All managed storage must be accessible through a single RAID controller
- Depending on the vendor, each RAID controller supports a maximum number of drives per RAID group

One or more RAID groups

Each ONTAP Select host must have a single RAID controller. You should create a single RAID group for ONTAP Select. However, in certain situations you might consider creating more than one RAID group. Refer to [Best practices](#).

Storage pool considerations

There are several issues related to the storage pools that you should be aware of as part of preparing to deploy ONTAP Select.



In a VMware environment, a storage pool is synonymous with a VMware datastore.

Storage pools and LUNs

Each LUN is seen as a local disk on the hypervisor host and can be part of one storage pool. Each storage pool is formatted with a file system that the hypervisor host OS can use.

You must make sure that the storage pools are created properly as part of an ONTAP Select deployment. You can create a storage pool using the hypervisor administration tool. For example, with VMware you can use the vSphere client to create a storage pool. The storage pool is then passed in to the ONTAP Select Deploy administration utility.

Managing the virtual disks

There are several issues related to the virtual disks that you should be aware of as part of preparing to deploy ONTAP Select.

Virtual disks and file systems

The ONTAP Select virtual machine is allocated multiple virtual disk drives. Each virtual disk is actually a file contained in a storage pool and is maintained by the hypervisor. There are several types of disks used by ONTAP Select, primarily system disks and data disks.

You should also be aware of the following regarding virtual disks:

- The storage pool must be available before the virtual disks can be created.
- The virtual disks cannot be created before the virtual machine is created.
- You must rely on the ONTAP Select Deploy administration utility to create all virtual disks (that is, an administrator must never create a virtual disk outside of the Deploy utility).

Configuring the virtual disks

The virtual disks are managed by ONTAP Select. They are created automatically when you create a cluster using the Deploy administration utility.

Illustration of the external storage environment

The ONTAP Select vNAS solution enables ONTAP Select to use datastores residing on storage that is external to the hypervisor host. The datastores can be accessed through the network using VMware vSAN or directly at an external storage array.

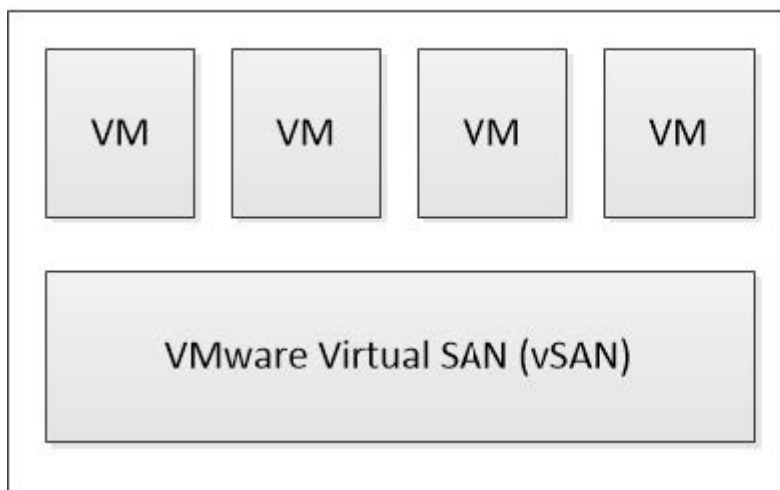
ONTAP Select can be configured to use the following types of VMware ESXi network datastores which are external to the hypervisor host:

- vSAN (Virtual SAN)
- VMFS
- NFS

vSAN datastores

Every ESXi host can have one or more local VMFS datastores. Normally these datastores are only accessible to the local host. However, VMware vSAN allows each of the hosts in an ESXi cluster to share all of the datastores in the cluster as if they were local. The following figure illustrates how vSAN creates a pool of datastores that are shared among the hosts in the ESXi cluster.

ESXi cluster



ONTAP Select
virtual machines

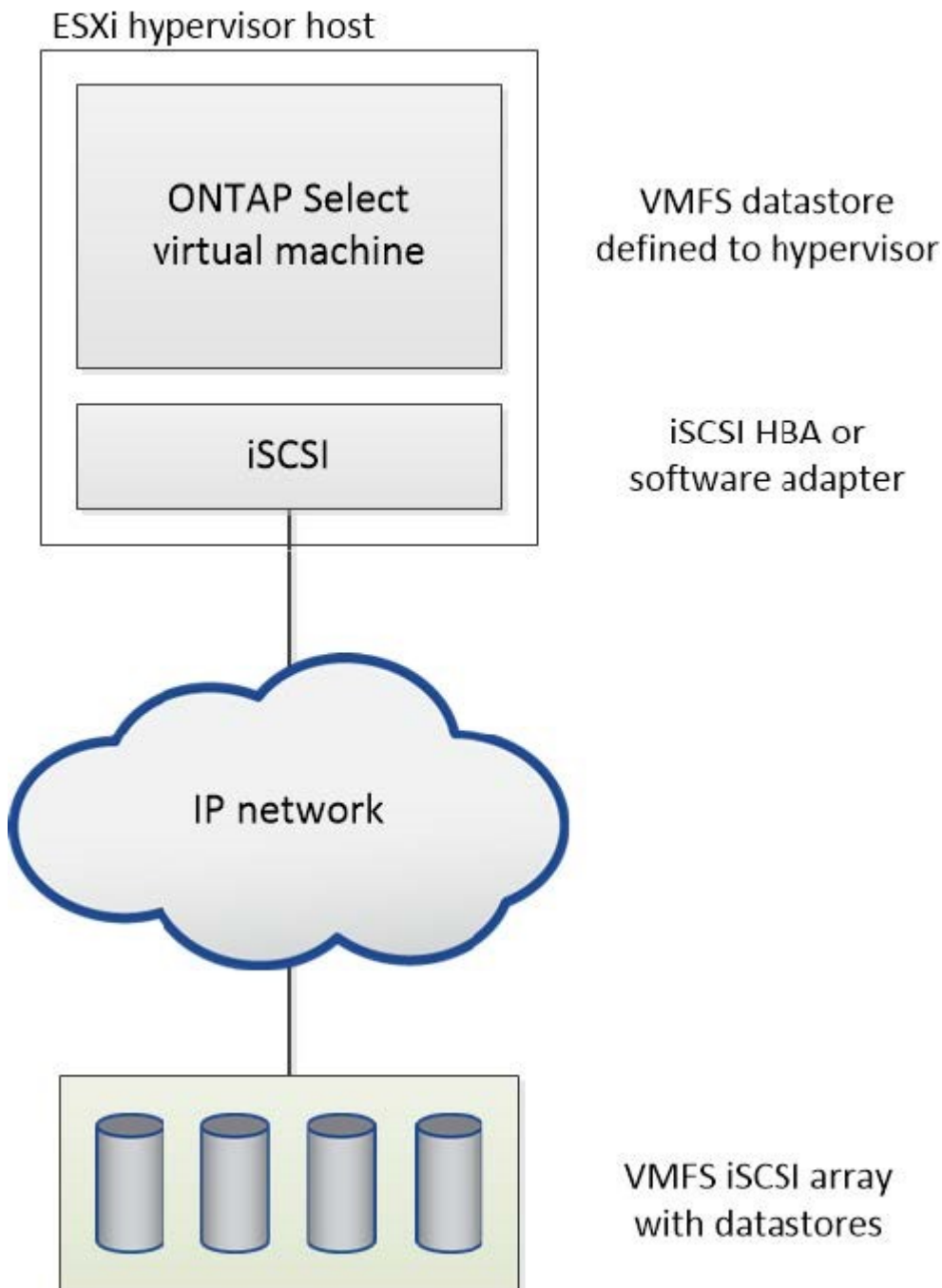
Shared datastores
accessed through vSAN

VMFS datastore on external storage array

You can create a VMFS datastore residing on an external storage array. The storage is accessed using one of several different network protocols. The following figure illustrates a VMFS datastore on an external storage array accessed using the iSCSI protocol.

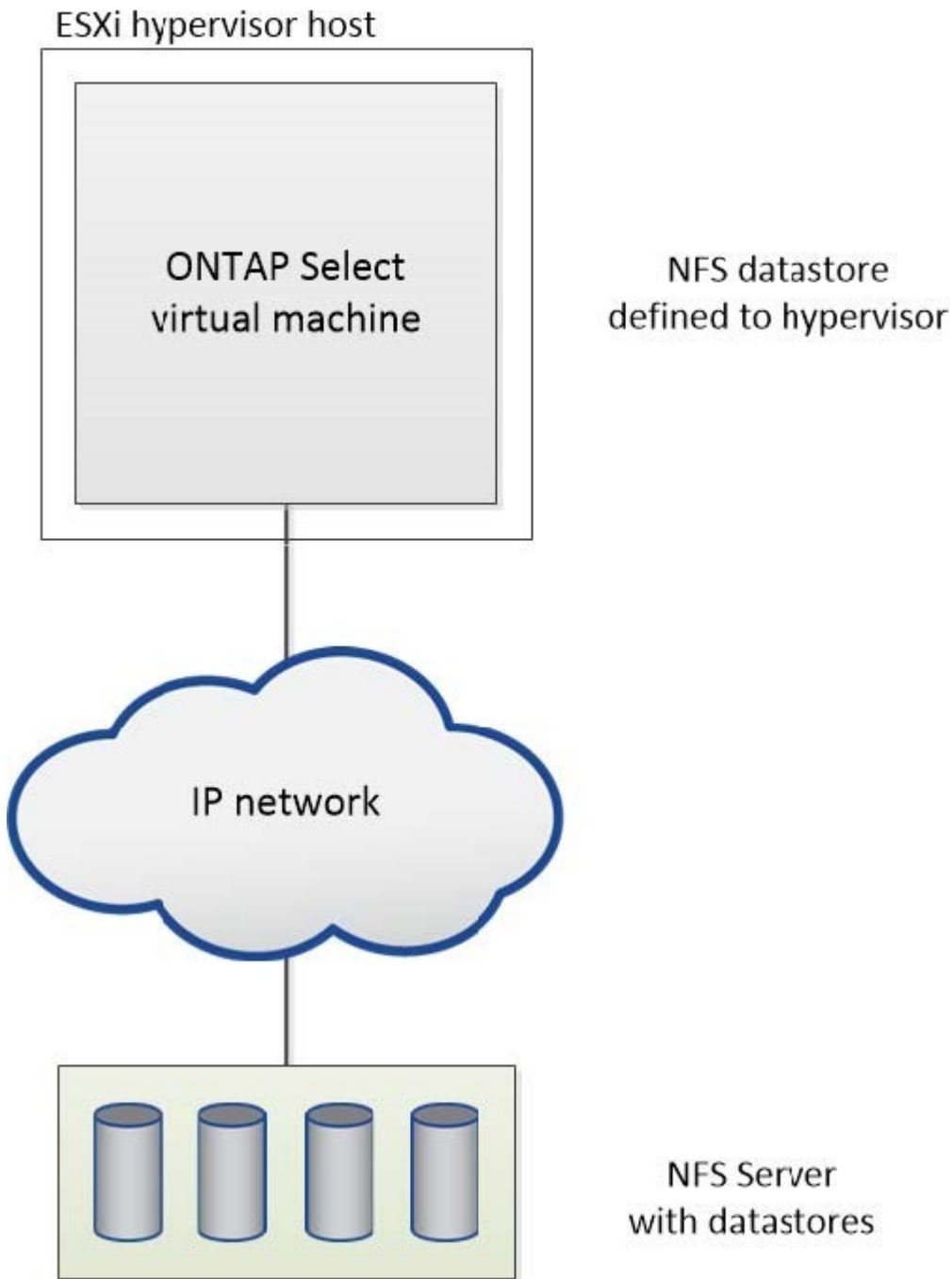


ONTAP Select supports all external storage arrays described in the VMware Storage/SAN Compatibility Guide, including iSCSI, Fiber Channel, and Fiber Channel over Ethernet.



NFS datastore on external storage array

You can create an NFS datastore residing on an external storage array. The storage is accessed using the NFS network protocol. The following figure illustrates an NFS datastore on external storage that is accessed through the NFS server appliance.



Hardware RAID services for local attached storage

When a hardware RAID controller is available, ONTAP Select can move RAID services to the hardware controller for both a write performance boost and protection against physical drive failures. As a result, RAID protection for all nodes

within the ONTAP Select cluster is provided by the locally attached RAID controller and not through ONTAP software RAID.



ONTAP Select data aggregates are configured to use RAID 0 because the physical RAID controller is providing RAID striping to the underlying drives. No other RAID levels are supported.

RAID controller configuration for local attached storage

All locally attached disks that provide ONTAP Select with backing storage must sit behind a RAID controller. Most commodity servers come with multiple RAID controller options across multiple price points, each with varying levels of functionality. The intent is to support as many of these options as possible, providing they meet certain minimum requirements placed on the controller.

The RAID controller that manages the ONTAP Select disks must meet the following requirements:

- The hardware RAID controller must have a battery backup unit (BBU) or flash-backed write cache (FBWC) and support 12Gbps of throughput.
- The RAID controller must support a mode that can withstand at least one or two disk failures (RAID 5 and RAID 6).
- The drive cache must be set to disabled.
- The write policy must be configured for writeback mode with a fallback to write through upon BBU or flash failure.
- The I/O policy for reads must be set to cached.

All locally attached disks that provide ONTAP Select with backing storage must be placed into RAID groups running RAID 5 or RAID 6. For SAS drives and SSDs, using RAID groups of up to 24 drives allows ONTAP to reap the benefits of spreading incoming read requests across a higher number of disks. Doing so provides a significant gain in performance. With SAS/SSD configurations, performance testing was performed against single-LUN versus multi-LUN configurations. No significant differences were found, so, for simplicity's sake, NetApp recommends creating the fewest number of LUNs necessary to support your configuration needs.

NL-SAS and SATA drives require a different set of best practices. For performance reasons, the minimum number of disks is still eight, but the RAID group size should not be larger than 12 drives. NetApp also recommends using one spare per RAID group; however, global spares for all RAID groups can be used. For example, you can use two spares for every three RAID groups, with each RAID group consisting of eight to 12 drives.



The maximum extent and datastore size for older ESX releases is 64TB, which can affect the number of LUNs necessary to support the total raw capacity provided by these large capacity drives.

RAID mode

Many RAID controllers support up to three modes of operation, each representing a significant difference in the data path taken by write requests. These three modes are as follows:

- **Writethrough.** All incoming I/O requests are written to the RAID controller cache and then immediately flushed to disk before acknowledging the request back to the host.
- **Writearound.** All incoming I/O requests are written directly to disk, circumventing the RAID controller cache.
- **Writeback.** All incoming I/O requests are written directly to the controller cache and immediately acknowledged back to the host. Data blocks are flushed to disk asynchronously using the controller.

Writeback mode offers the shortest data path, with I/O acknowledgment occurring immediately after the blocks enter cache. This mode provides the lowest latency and highest throughput for mixed read/write workloads. However, without the presence of a BBU or nonvolatile flash technology, users run the risk of losing data if the system incurs a power failure when operating in this mode.

ONTAP Select requires the presence of a battery backup or flash unit; therefore, we can be confident that cached blocks are flushed to disk in the event of this type of failure. For this reason, it is a requirement that the RAID controller be configured in writeback mode.

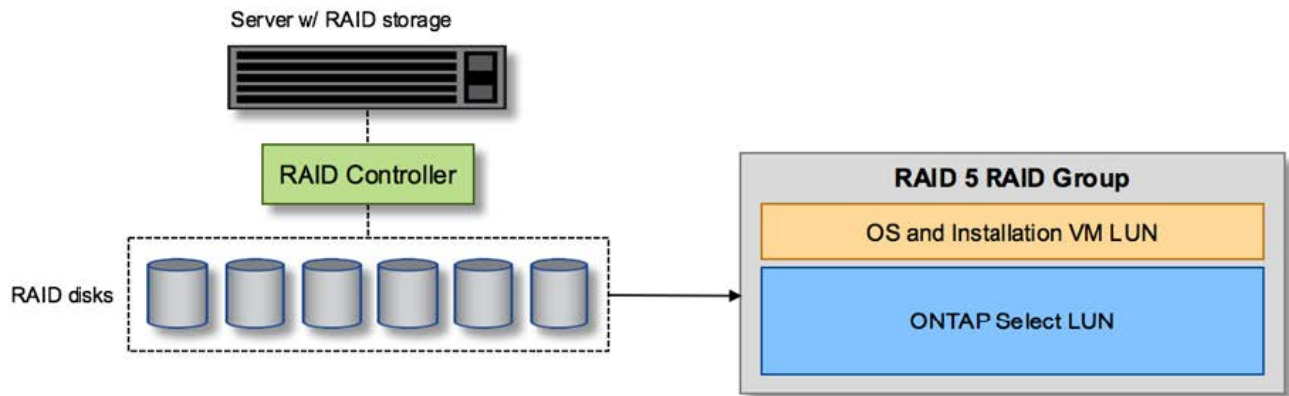
Local disks shared between ONTAP Select and OS

The most common server configuration is one in which all locally attached spindles sit behind a single RAID controller. You should provision a minimum of two LUNs: one for the hypervisor and one for the ONTAP Select VM.

For example, consider an HP DL380 g8 with six internal drives and a single Smart Array P420i RAID controller. All internal drives are managed by this RAID controller, and no other storage is present on the system.

The following figure shows this style of configuration. In this example, no other storage is present on the system; therefore, the hypervisor must share storage with the ONTAP Select node.

Server LUN configuration with only RAID-managed spindles



Provisioning the OS LUNs from the same RAID group as ONTAP Select allows the hypervisor OS (and any client VM that is also provisioned from that storage) to benefit from RAID protection. This configuration prevents a single-drive failure from bringing down the entire system.

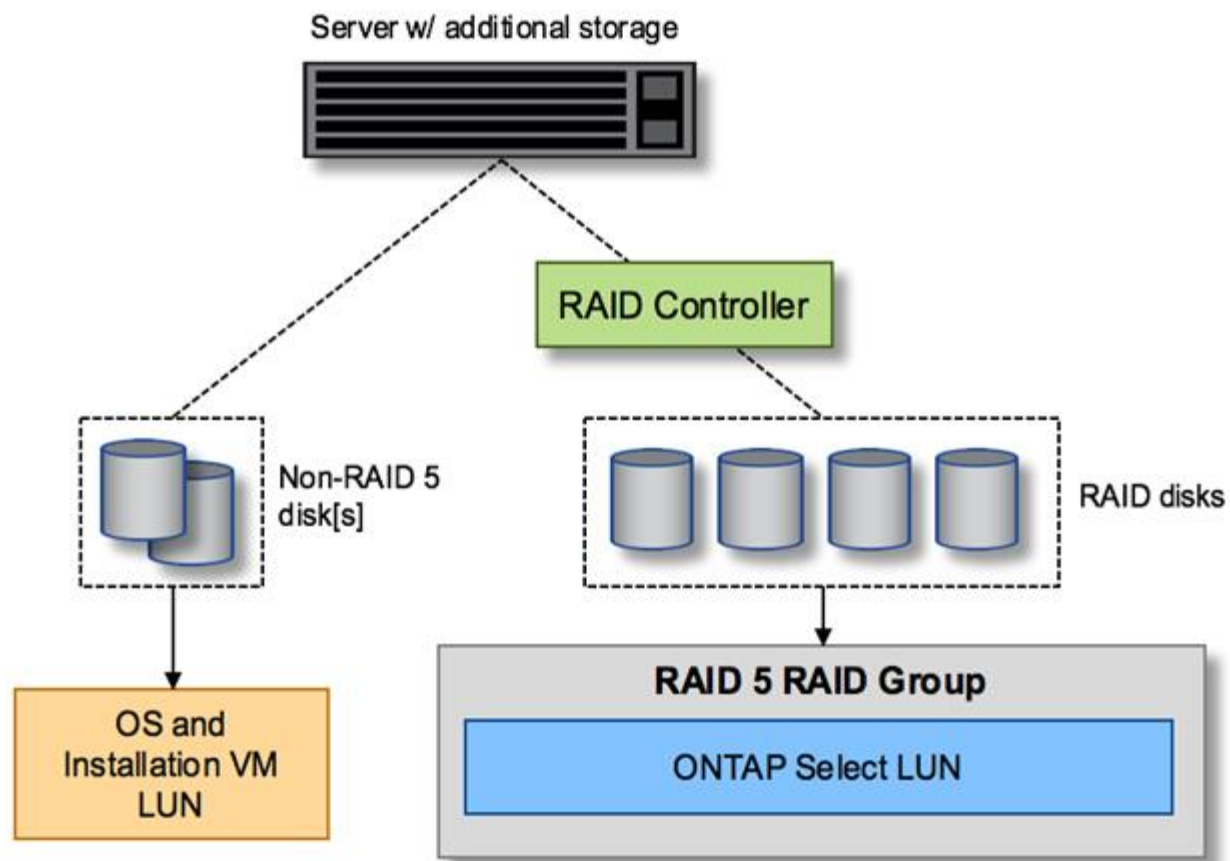
Local disks split between ONTAP Select and OS

The other possible configuration provided by server vendors involves configuring the system with multiple RAID or disk controllers. In this configuration, a set of disks is managed by one disk controller, which might or might not offer RAID services. A second set of disks is managed by a hardware RAID controller that is able to offer RAID 5/6 services.

With this style of configuration, the set of spindles that sits behind the RAID controller that can provide RAID 5/6 services should be used exclusively by the ONTAP Select VM. Depending on the total storage capacity under management, you should configure the disk spindles into one or more RAID groups and one or more LUNs. These LUNs would then be used to create one or more datastores, with all datastores being protected by the RAID controller.

The first set of disks is reserved for the hypervisor OS and any client VM that is not using ONTAP storage, as shown in the following figure.

Server LUN configuration on mixed RAID/non-RAID system



Multiple LUNs

There are two cases for which single-RAID group/single-LUN configurations must change. When using NL-SAS or SATA drives, the RAID group size must not exceed 12 drives. In addition, a single LUN can become larger than the underlying hypervisor storage limits either individual file system extent maximum size or total storage pool maximum size. Then the underlying physical storage must be broken up into multiple LUNs to enable successful file system creation.

VMware vSphere virtual machine file system limits

The maximum size of a datastore on some versions of ESX is 64TB.

If a server has more than 64TB of storage attached, multiple LUNs might need to be provisioned, each smaller than 64TB. Creating multiple RAID groups to improve the RAID rebuild time for SATA/NL-SAS drives also results in multiple LUNs being provisioned.

When multiple LUNs are required, a major point of consideration is making sure that these LUNs have similar and consistent performance. This is especially important if all the LUNs are to be used in a single ONTAP aggregate. Alternatively, if a subset of one or more LUNs has a distinctly different performance profile, we strongly recommend isolating these LUNs in a separate ONTAP aggregate.

Multiple file system extents can be used to create a single datastore up to the maximum size of the

datastore. To restrict the amount of capacity that requires an ONTAP Select license, make sure to specify a capacity cap during the cluster installation. This functionality allows ONTAP Select to use (and therefore require a license for) only a subset of the space in a datastore.

Alternatively, one can start by creating a single datastore on a single LUN. When additional space requiring a larger ONTAP Select capacity license is needed, then that space can be added to the same datastore as an extent, up to the maximum size of the datastore. After the maximum size is reached, new datastores can be created and added to ONTAP Select. Both types of capacity extension operations are supported and can be achieved by using the ONTAP Deploy storage-add functionality. Each ONTAP Select node can be configured to support up to 400TB of storage. Provisioning capacity from multiple datastores requires a two-step process.

The initial cluster create can be used to create an ONTAP Select cluster consuming part or all of the space in the initial datastore. A second step is to perform one or more capacity addition operations using additional datastores until the desired total capacity is reached. This functionality is detailed in the section [Increasing storage capacity](#).



VMFS overhead is nonzero (see [VMware KB 1001618](#)), and attempting to use the entire space reported as free by a datastore has resulted in spurious errors during cluster create operations.

A 2% buffer is left unused in each datastore. This space does not require a capacity license because it is not used by ONTAP Select. ONTAP Deploy automatically calculates the exact number of gigabytes for the buffer, as long as a capacity cap is not specified. If a capacity cap is specified, that size is enforced first. If the capacity cap size falls within the buffer size, the cluster create fails with an error message specifying the correct maximum size parameter that can be used as a capacity cap:

```
"InvalidPoolCapacitySize: Invalid capacity specified for storage pool "ontap-select-storage-pool", Specified value: 34334204 GB. Available (after leaving 2% overhead space): 30948"
```

VMFS 6 is supported for both new installations and as the target of a Storage vMotion operation of an existing ONTAP Deploy or ONTAP Select VM.

VMware does not support in-place upgrades from VMFS 5 to VMFS 6. Therefore, Storage vMotion is the only mechanism that allows any VM to transition from a VMFS 5 datastore to a VMFS 6 datastore. However, support for Storage vMotion with ONTAP Select and ONTAP Deploy was expanded to cover other scenarios besides the specific purpose of transitioning from VMFS 5 to VMFS 6.

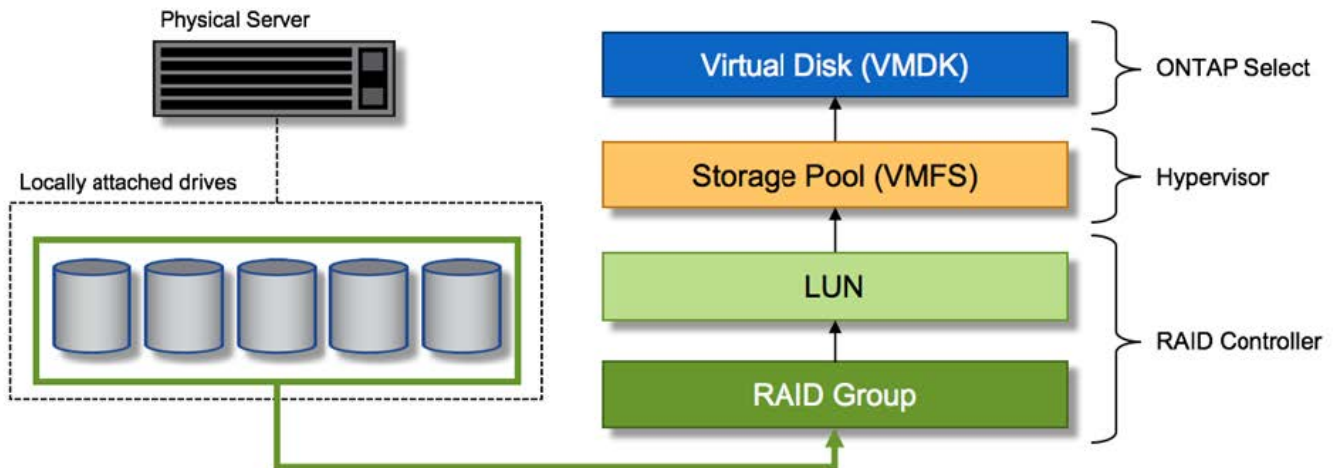
ONTAP Select virtual disks

At its core, ONTAP Select presents ONTAP with a set of virtual disks provisioned from one or more storage pools. ONTAP is presented with a set of virtual disks that it treats as physical, and the remaining portion of the storage stack is abstracted by the hypervisor. The following figure shows this relationship in more detail, highlighting the relationship between the physical RAID controller, the

hypervisor, and the ONTAP Select VM.

- RAID group and LUN configuration occur from within the server's RAID controller software. This configuration is not required when using VSAN or external arrays.
- Storage pool configuration occurs from within the hypervisor.
- Virtual disks are created and owned by individual VMs; in this example, by ONTAP Select.

Virtual disk to physical disk mapping



Virtual disk provisioning

To provide for a more streamlined user experience, the ONTAP Select management tool, ONTAP Deploy, automatically provisions virtual disks from the associated storage pool and attaches them to the ONTAP Select VM. This operation occurs automatically during both initial setup and during storage-add operations. If the ONTAP Select node is part of an HA pair, the virtual disks are automatically assigned to a local and mirror storage pool.

ONTAP Select breaks up the underlying attached storage into equal-sized virtual disks, each not exceeding 16TB. If the ONTAP Select node is part of an HA pair, a minimum of two virtual disks are created on each cluster node and assigned to the local and mirror plex to be used within a mirrored aggregate.

For example, an ONTAP Select can assigned a datastore or LUN that is 31TB (the space remaining after the VM is deployed and the system and root disks are provisioned). Then four ~7.75TB virtual disks are created and assigned to the appropriate ONTAP local and mirror plex.



Adding capacity to an ONTAP Select VM likely results in VMDKs of different sizes. For details, see the section [Increasing storage capacity](#). Unlike FAS systems, different sized VMDKs can exist in the same aggregate. ONTAP Select uses a RAID 0 stripe across these VMDKs, which results in the ability to fully use all the space in each VMDK regardless of its size.

Virtualized NVRAM

NetApp FAS systems are traditionally fitted with a physical NVRAM PCI card, a high-performing card containing nonvolatile flash memory. This card provides a significant boost in write performance by granting ONTAP with the ability to immediately acknowledge incoming writes back to the client. It can also schedule the movement of modified data blocks back to the slower storage media in a process known as destaging.

Commodity systems are not typically fitted with this type of equipment. Therefore, the functionality of this NVRAM card has been virtualized and placed into a partition on the ONTAP Select system boot disk. It is for this reason that placement of the system virtual disk of the instance is extremely important. This is also why the product requires the presence of a physical RAID controller with a resilient cache for local attached storage configurations.

NVRAM is placed on its own VMDK. Splitting the NVRAM in its own VMDK allows the ONTAP Select VM to use the vNVMe driver to communicate with its NVRAM VMDK. It also requires that the ONTAP Select VM uses hardware version 13, which is compatible with ESX 6.5 and newer.

Data path explained: NVRAM and RAID controller

The interaction between the virtualized NVRAM system partition and the RAID controller can be best highlighted by walking through the data path taken by a write request as it enters the system.

Incoming write requests to the ONTAP Select VM are targeted at the VM's NVRAM partition. At the virtualization layer, this partition exists within an ONTAP Select system disk, a VMDK attached to the ONTAP Select VM. At the physical layer, these requests are cached in the local RAID controller, like all block changes targeted at the underlying spindles. From here, the write is acknowledged back to the host.

At this point, physically, the block resides in the RAID controller cache, waiting to be flushed to disk. Logically, the block resides in NVRAM waiting for destaging to the appropriate user data disks.

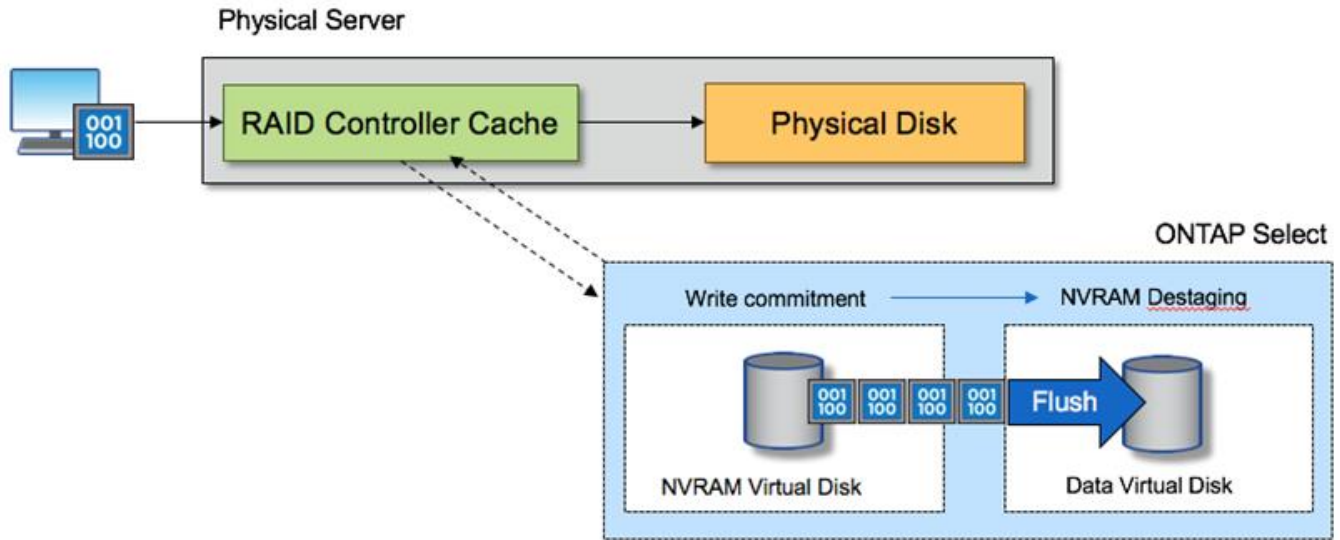
Because changed blocks are automatically stored within the RAID controller's local cache, incoming writes to the NVRAM partition are automatically cached and periodically flushed to physical storage media. This should not be confused with the periodic flushing of NVRAM contents back to ONTAP data disks. These two events are unrelated and occur at different times and frequencies.

The following figure shows the I/O path an incoming write takes. It highlights the difference between the physical layer (represented by the RAID controller cache and disks) and the virtual layer (represented by the VM's NVRAM and data virtual disks).



Although blocks changed on the NVRAM VMDK are cached in the local RAID controller cache, the cache is not aware of the VM construct or its virtual disks. It stores all changed blocks on the system, of which NVRAM is only a part. This includes write requests bound for the hypervisor, if it is provisioned from the same backing spindles.

Incoming writes to ONTAP Select VM



Note that the NVRAM partition is separated on its own VMDK. That VMDK is attached using the vNVME driver available in ESX versions of 6.5 or later. This change is most significant for ONTAP Select installations with software RAID, which do not benefit from the RAID controller cache.

Software RAID services for local attached storage

Software RAID is a RAID abstraction layer implemented within the ONTAP software stack. It provides the same functionality as the RAID layer within a traditional ONTAP platform such as FAS. The RAID layer performs drive parity calculations and provides protection against individual drive failures within an ONTAP Select node.

Independent of the hardware RAID configurations, ONTAP Select also provides a software RAID option. A hardware RAID controller might not be available or might be undesirable in certain environments, such as when ONTAP Select is deployed on a small form-factor commodity hardware. Software RAID expands the available deployment options to include such environments. To enable software RAID in your environment, here are some points to remember:

- It is available with a Premium or Premium XL license.
- It only supports SSD or NVMe (requires Premium XL license) drives for ONTAP root and data disks.
- It requires a separate system disk for the ONTAP Select VM boot partition.
 - Choose a separate disk, either an SSD or an NVMe drive, to create a datastore for the system disks (NVRAM, Boot/CF card, Coredump, and Mediator in a multi-node setup).

Notes

- The terms service disk and system disk are used interchangeably.
 - Service disks are the VMDKs that are used within the ONTAP Select VM to service various items such as clustering, booting, and so on.
 - Service disks are physically located on a single physical disk (collectively called the service/system physical disk) as seen from the host. That physical disk must contain a DAS datastore. ONTAP Deploy creates these service disks for the ONTAP Select VM during cluster deployment.
- It is not possible to further separate the ONTAP Select system disks across multiple datastores or across multiple physical drives.
- Hardware RAID is not deprecated.

Software RAID configuration for local attached storage

When using software RAID, the absence of a hardware RAID controller is ideal, but, if a system does have an existing RAID controller, it must adhere to the following requirements:

- The hardware RAID controller must be disabled such that disks can be presented directly to the system (a JBOD). This change can usually be made in the RAID controller BIOS
- Or the hardware RAID controller should be in the SAS HBA mode. For example, some BIOS configurations allow an “AHCI” mode in addition to RAID, which could be chosen to enable the JBOD mode. This enables a passthrough, so that the physical drives can be seen as is on the host.

Depending on maximum number of drives supported by the controller, an additional controller may be required. With the SAS HBA mode, ensure that the IO controller (SAS HBA) is supported with a minimum of 6Gb/s speed. However, NetApp recommends a 12Gbps speed.

No other hardware RAID controller modes or configurations is supported. For example, some controllers allow a RAID 0 support that can artificially enable disks to pass-through but the implications can be undesirable. The supported size of physical disks (SSD only) is between 200GB – 16TB.



Administrators need to keep track of which drives are in use by the ONTAP Select VM and prevent inadvertent use of those drives on the host.

ONTAP Select virtual and physical disks

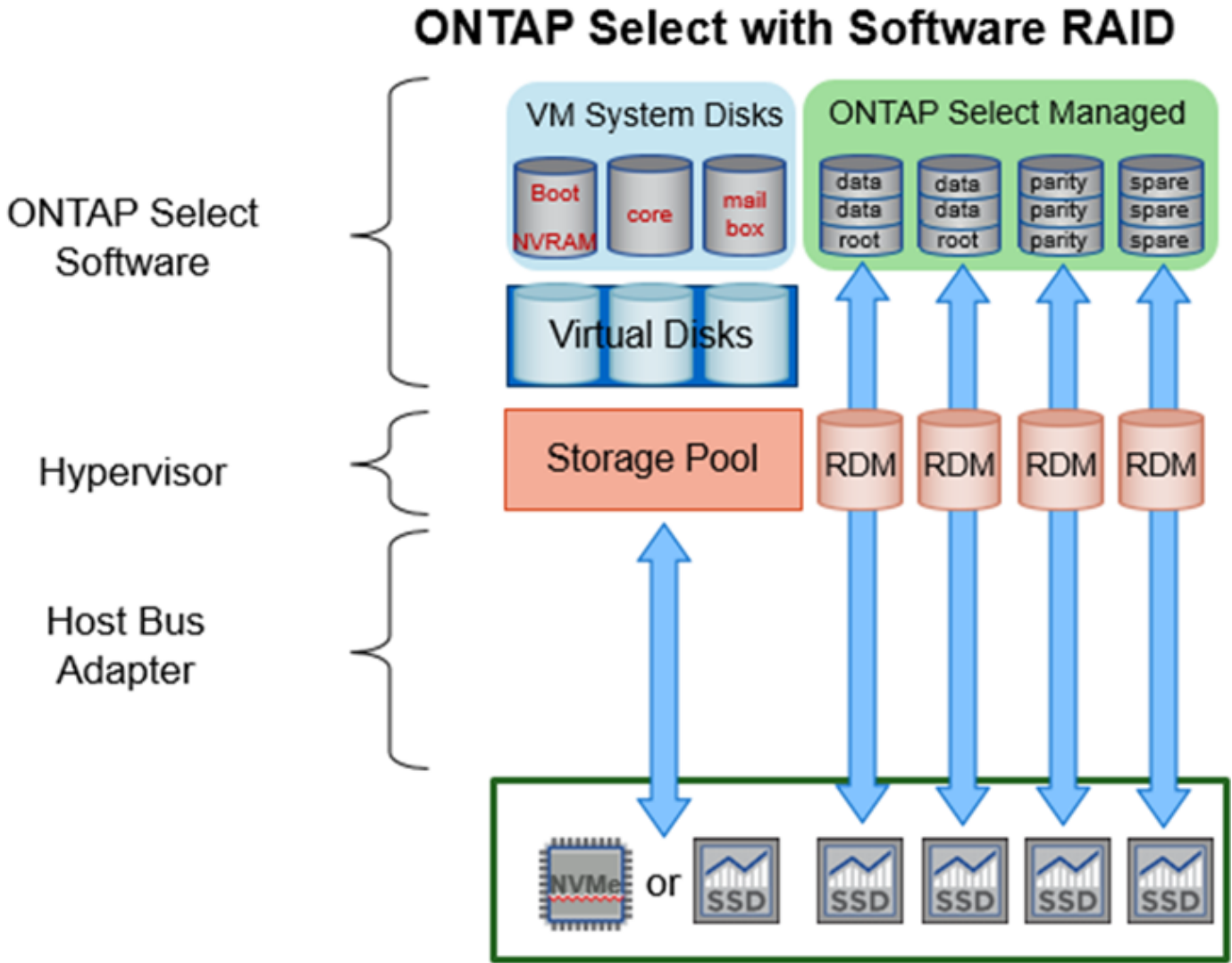
For configurations with hardware RAID controllers, physical disk redundancy is provided by the RAID controller. ONTAP Select is presented with one or more VMDKs from which the ONTAP admin can configure data aggregates. These VMDKs are striped in a RAID 0 format because using ONTAP software RAID is redundant, inefficient, and ineffective due to resiliency provided at the hardware level. Furthermore, the VMDKs used for system disks are in the same datastore as the VMDKs used to store user data.

When using software RAID, ONTAP Deploy presents ONTAP Select with a set of virtual disks (VMDKs)

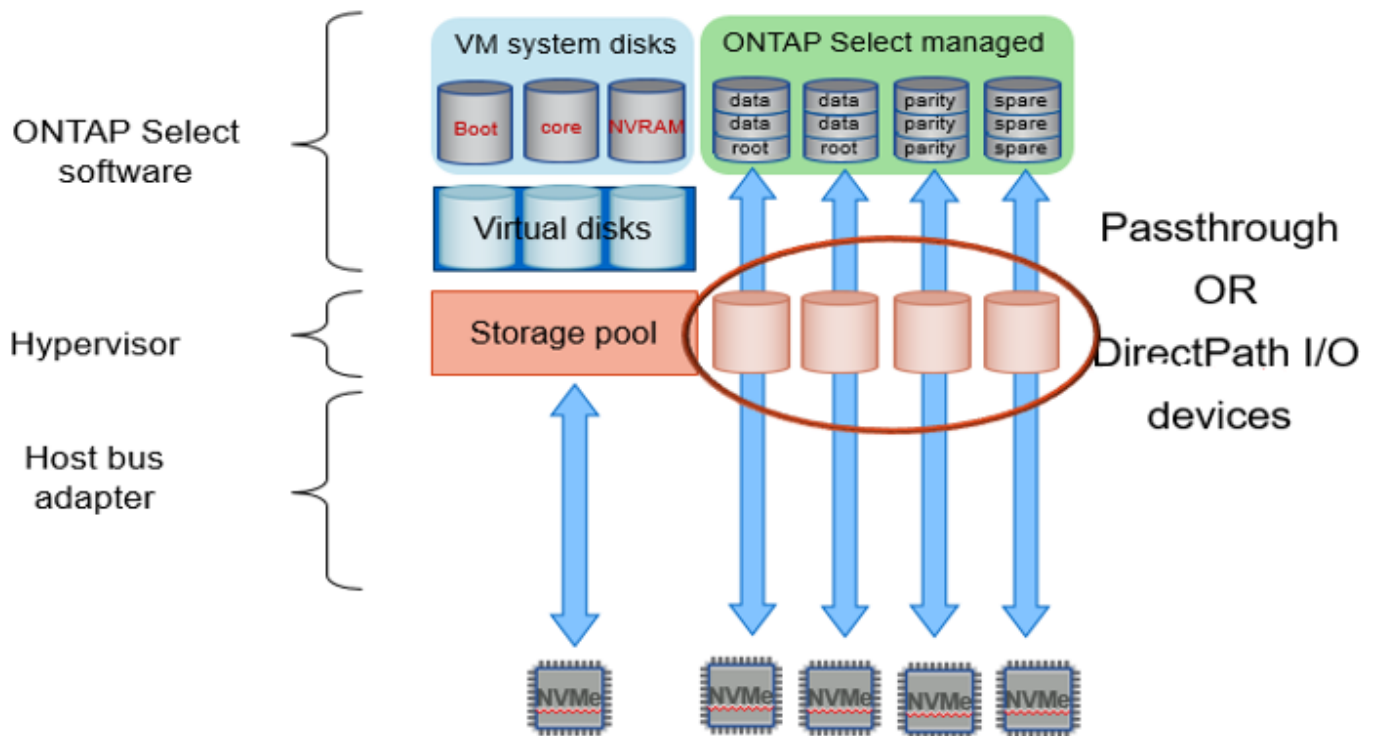
and physical disks Raw Device Mappings [RDMs] for SSDs and passthrough or DirectPath IO devices for NVMeS.

The following figures show this relationship in more detail, highlighting the difference between the virtualized disks used for the ONTAP Select VM internals and the physical disks used to store user data.

ONTAP Select software RAID: use of virtualized disks and RDMs



The system disks (VMDKs) reside in the same datastore and on the same physical disk. The virtual NVRAM disk requires a fast and durable media. Therefore, only NVMe and SSD-type datastores are supported.



The system disks (VMDKs) reside in the same datastore and on the same physical disk. The virtual NVRAM disk requires a fast and durable media. Therefore, only NVMe and SSD-type datastores are supported. When using NVMe drives for data, the system disk should also be an NVMe device for performance reasons. A good candidate for the system disk in an all NVMe configuration is an INTEL Optane card.

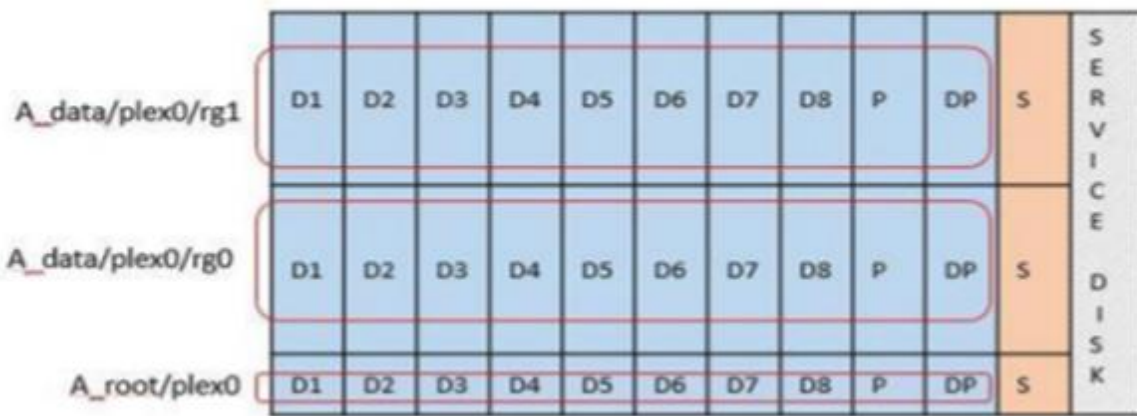


With the current release, it is not possible to further separate the ONTAP Select system disks across multiple datastores or multiple physical drives.

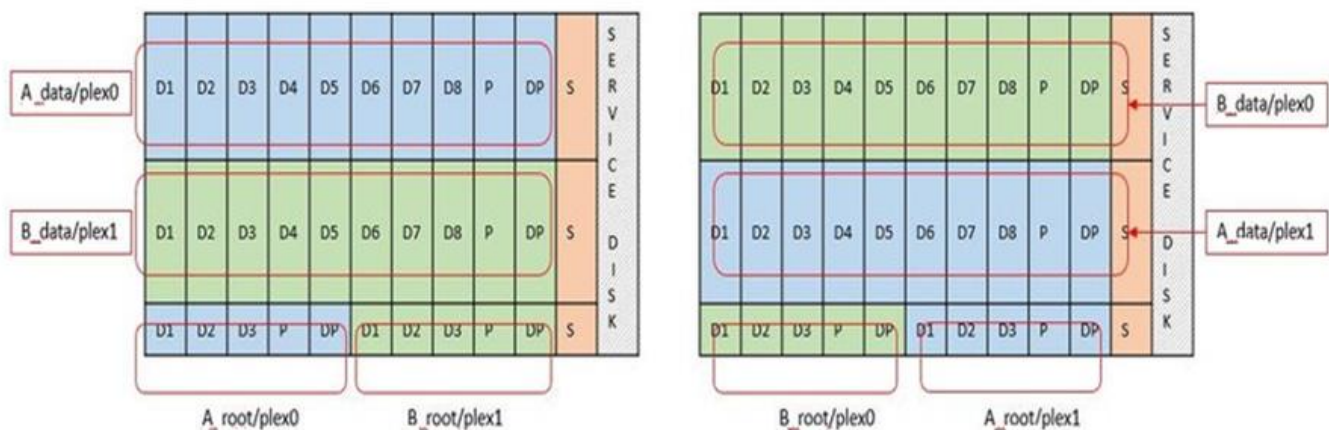
Each data disk is divided into three parts: a small root partition (stripe) and two equal-sized partitions to create two data disks seen within the ONTAP Select VM. Partitions use the Root Data Data (RD2) schema as shown in the following figures for a single node cluster and for a node in an HA pair.

P denotes a parity drive. **DP** denotes a dual parity drive and **S** denotes a spare drive.

RDD disk partitioning for single-node clusters



RDD disk partitioning for multinode clusters (HA pairs)



ONTAP software RAID supports the following RAID types: RAID 4, RAID-DP, and RAID-TEC. These are the same RAID constructs used by FAS and AFF platforms. For root provisioning ONTAP Select supports only RAID 4 and RAID-DP. When using RAID-TEC for the data aggregate, the overall protection is RAID-DP. ONTAP Select HA uses a shared-nothing architecture that replicates each node's configuration to the other node. That means each node must store its root partition and a copy of its peer's root partition. Since a data disk has a single root partition, that the minimum number of data disks will vary depending on whether the ONTAP Select node is part of an HA pair or not.

For single node clusters, all data partitions are used to store local (active) data. For nodes that are part of an HA pair, one data partition is used to store local (active) data for that node and the second data partition is used to mirror active data from the HA peer.

Passthrough (DirectPath IO) devices vs. Raw Device Maps (RDMs)

VMware ESX does not currently support NVMe disks as Raw Device Maps. For ONTAP Select to take direct control of NVMe disks, the NVMe drives must be configured in ESX as passthrough devices. Please note that configuring an NVMe device as a passthrough device requires support from the

server BIOS and it is a disruptive process, requiring an ESX host reboot. Furthermore, the maximum number of passthrough devices per ESX host is 16. However, ONTAP Deploy limits this to 14. This limit of 14 NVMe devices per ONTAP Select node means that an all NVMe configuration will provide a very high IOPs density (IOPs/TB) at the expense of total capacity. Alternatively, if a high performance configuration with larger storage capacity is desired, the recommended configuration is a large ONTAP Select VM size, an INTEL Optane card for the system disk, and a nominal number of SSD drives for data storage.



To take full advantage of NVMe performance, consider the large ONTAP Select VM size.

There is an additional difference between passthrough devices and RDMs. RDMs can be mapped to a running VM. Passthrough devices require a VM reboot. This means that any NVMe drive replacement or capacity expansion (drive addition) procedure will require an ONTAP Select VM reboot. The drive replacement and capacity expansion (drive addition) operation is driven by a workflow in ONTAP Deploy. ONTAP Deploy manages the ONTAP Select reboot for single node clusters and failover / failback for HA pairs. However it is important to note the difference between working with SSD data drives (no ONTAP Select reboot / failovers are required) and working with NVMe data drives (ONTAP Select reboot / failover is required).

Physical and virtual disk provisioning

To provide a more streamlined user experience, ONTAP Deploy automatically provisions the system (virtual) disks from the specified datastore (physical system disk) and attaches them to the ONTAP Select VM. This operation occurs automatically during the initial setup so that the ONTAP Select VM can boot. The RDMs are partitioned and the root aggregate is automatically built. If the ONTAP Select node is part of an HA pair, the data partitions are automatically assigned to a local storage pool and a mirror storage pool. This assignment occurs automatically during both cluster-creation operations and storage-add operations.

Because the data disks on the ONTAP Select VM are associated with the underlying physical disks, there are performance implications for creating configurations with a larger number of physical disks.



The root aggregate's RAID group type depends on the number of disks available. ONTAP Deploy picks the appropriate RAID group type. If it has sufficient disks allocated to the node, it uses RAID-DP, otherwise it creates a RAID-4 root aggregate.

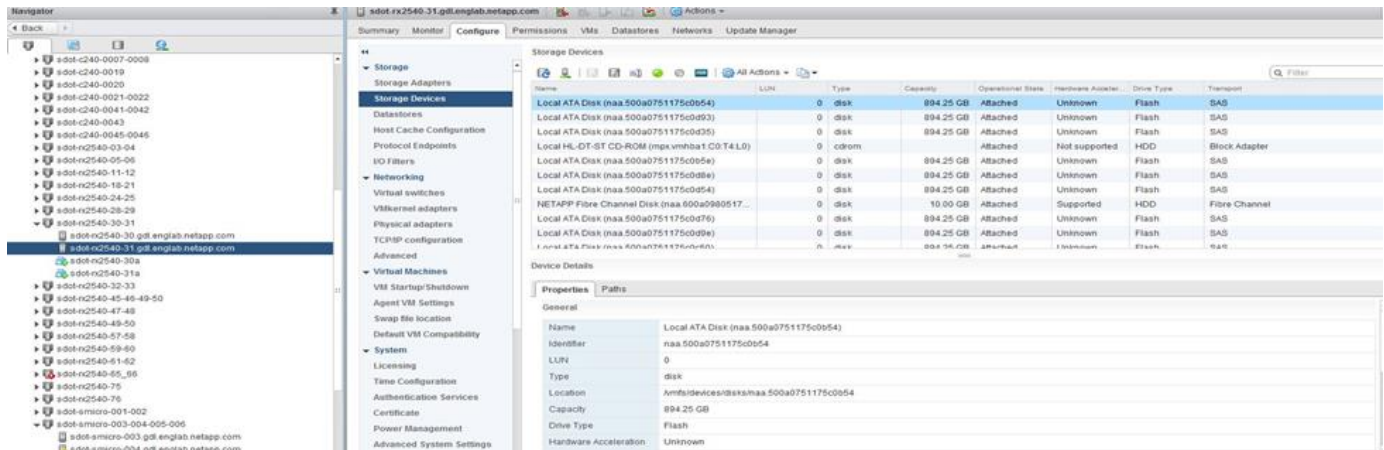
When adding capacity to an ONTAP Select VM using software RAID, the administrator must consider the physical drive size and the number of drives required. For details, see the section [Increasing storage capacity](#).

Similar to FAS and AFF systems, only drives with equal or larger capacities can be added to an existing RAID group. Larger capacity drives are right sized. If you are creating new RAID groups, the new RAID group size should match the existing RAID group size to make sure that the overall aggregate performance does not deteriorate.

Matching an ONTAP Select disk to the corresponding ESX disk

ONTAP Select disks are usually labeled NET x.y. You can use the following ONTAP command to obtain the disk UUID:

```
<system name>::> disk show NET-1.1
Disk: NET-1.1
Model: Micron_5100_MTFD
Serial Number: 1723175C0B5E
UID:
*500A0751:175C0B5E*:00000000:00000000:00000000:00000000:00000000:00000000:00000000:00000000
00
BPS: 512
Physical Size: 894.3GB
Position: shared
Checksum Compatibility: advanced_zoned
Aggregate: -
Plex: -This UID can be matched with the device UID displayed in the 'storage devices' tab
for the ESX host
```



In the ESXi shell, you can enter the following command to blink the LED for a given physical disk (identified by its naa.unique-id).

```
esxcli storage core device set -d <naa_id> -l=locator -L=<seconds>
```

Multiple drive failures when using software RAID

It is possible for a system to encounter a situation in which multiple drives are in a failed state at the same time. The behavior of the system depends on the aggregate RAID protection and the number of failed drives.

A RAID4 aggregate can survive one disk failure, a RAID-DP aggregate can survive two disk failures, and a RAID-TEC aggregate can survive three disks failures.

If the number of failed disks is less than the maximum number of failures that RAID type supports, and if a spare disk is available, the reconstruction process starts automatically. If spare disks are not available, the aggregate serves data in a degraded state until spare disks are added.

If the number of failed disks is more than the maximum number of failures that the RAID type supports, then the local plex is marked as failed, and the aggregate state is degraded. Data is served from the second plex residing on the HA partner. This means that any I/O requests for node 1 are sent through cluster interconnect port e0e (iSCSI) to the disks physically located on node 2. If the second plex also fails, then the aggregate is marked as failed and data is unavailable.

A failed plex must be deleted and recreated for the proper mirroring of data to resume. Note that a multi-disk failure resulting in a data aggregate being degraded also results in a root aggregate being degraded. ONTAP Select uses the root-data-data (RDD) partitioning schema to split each physical drive into a root partition and two data partitions. Therefore, losing one or more disks might impact multiple aggregates, including the local root or the copy of the remote root aggregate, as well as the local data aggregate and the copy of the remote data aggregate.

```
C3111E67::> storage aggregate plex delete -aggregate aggr1 -plex plex1
Warning: Deleting plex "plex1" of mirrored aggregate "aggr1" in a non-shared HA
configuration will disable its synchronous mirror protection and disable
        negotiated takeover of node "sti-rx2540-335a" when aggregate "aggr1" is online.
Do you want to continue? {y|n}: y
[Job 78] Job succeeded: DONE

C3111E67::> storage aggregate mirror -aggregate aggr1
Info: Disks would be added to aggregate "aggr1" on node "sti-rx2540-335a" in the
following manner:
    Second Plex
    RAID Group rg0, 5 disks (advanced_zoned checksum, raid_dp)

    Position  Disk                Type      Usable Physical
    -----  -
    shared    NET-3.2                SSD       -          -
    shared    NET-3.3                SSD       -          -
    shared    NET-3.4                SSD       208.4GB   208.4GB
    shared    NET-3.5                SSD       208.4GB   208.4GB
    shared    NET-3.12               SSD       208.4GB   208.4GB

    Aggregate capacity available for volume use would be 526.1GB.
    625.2GB would be used from capacity license.
Do you want to continue? {y|n}: y

C3111E67::> storage aggregate show-status -aggregate aggr1
Owner Node: sti-rx2540-335a
Aggregate: aggr1 (online, raid_dp, mirrored) (advanced_zoned checksums)
Plex: /aggr1/plex0 (online, normal, active, pool0)
```

RAID Group /aggr1/plex0/rg0 (normal, advanced_zoned checksums)

Position	Disk	Pool	Type	RPM	Usable Size	Physical Size	Status
shared	NET-1.1	0	SSD	-	205.1GB	447.1GB	(normal)
shared	NET-1.2	0	SSD	-	205.1GB	447.1GB	(normal)
shared	NET-1.3	0	SSD	-	205.1GB	447.1GB	(normal)
shared	NET-1.10	0	SSD	-	205.1GB	447.1GB	(normal)
shared	NET-1.11	0	SSD	-	205.1GB	447.1GB	(normal)

Plex: /aggr1/plex3 (online, normal, active, pool1)

RAID Group /aggr1/plex3/rg0 (normal, advanced_zoned checksums)

Position	Disk	Pool	Type	RPM	Usable Size	Physical Size	Status
shared	NET-3.2	1	SSD	-	205.1GB	447.1GB	(normal)
shared	NET-3.3	1	SSD	-	205.1GB	447.1GB	(normal)
shared	NET-3.4	1	SSD	-	205.1GB	447.1GB	(normal)
shared	NET-3.5	1	SSD	-	205.1GB	447.1GB	(normal)
shared	NET-3.12	1	SSD	-	205.1GB	447.1GB	(normal)

10 entries were displayed..



In order to test or simulate one or multiple drive failures, use the **storage disk fail -disk NET-x.y -immediate** command. If there is a spare in the system, the aggregate will begin to reconstruct. You can check the status of the reconstruction using the command **storage aggregate show**. You can remove the simulated failed drive using ONTAP Deploy. Note that ONTAP has marked the drive as **Broken**. The drive is not actually broken and can be added back using ONTAP Deploy. In order to erase the Broken label, enter the following commands in the ONTAP Select CLI:

```
set advanced
disk unfail -disk NET-x.y -spare true
disk show -broken
```

The output for the last command should be empty.

Virtualized NVRAM

NetApp FAS systems are traditionally fitted with a physical NVRAM PCI card. This card is a high-performing card containing nonvolatile flash memory that provides a significant boost in write performance. It does this by granting ONTAP the ability to immediately acknowledge incoming writes back to the client. It can also schedule the movement of modified data blocks back to slower storage media in a process known as destaging.

Commodity systems are not typically fitted with this type of equipment. Therefore, the functionality of the NVRAM card has been virtualized and placed into a partition on the ONTAP Select system boot disk.

It is for this reason that placement of the system virtual disk of the instance is extremely important.

VSAN and external array configurations

Virtual NAS (vNAS) deployments support ONTAP Select clusters on VSAN, some HCI products, NetApp HCI technology, and external array types of datastores. The underlying infrastructure of these configurations provide datastore resiliency.

The minimum requirement is that the underlying configuration is supported by VMware and should be listed on the respective VMware HCLs.

vNAS architecture

The vNAS nomenclature is used for all setups that do not use DAS. For multinode ONTAP Select clusters, this includes architectures for which the two ONTAP Select nodes in the same HA pair share a single datastore (including vSAN datastores). The nodes can also be installed on separate datastores from the same shared external array. This allows for array-side storage efficiencies to reduce the overall footprint of the entire ONTAP Select HA pair. The architecture of ONTAP Select vNAS solutions is very similar to that of ONTAP Select on DAS with a local RAID controller. That is to say that each ONTAP Select node continues to have a copy of its HA partner's data. ONTAP storage efficiency policies are node scoped. Therefore, array side storage efficiencies are preferable because they can potentially be applied across data sets from both ONTAP Select nodes.

It is also possible that each ONTAP Select node in an HA pair uses a separate external array. This is a common choice when using ONTAP Select Metrocluster SDS with external storage.

When using separate external arrays for each ONTAP Select node, it is very important that the two arrays provide similar performance characteristics to the ONTAP Select VM.

vNAS architectures versus local DAS with hardware RAID controllers

The vNAS architecture is logically most similar to the architecture of a server with DAS and a RAID controller. In both cases, ONTAP Select consumes datastore space. That datastore space is carved into VMDKs, and these VMDKs form the traditional ONTAP data aggregates. ONTAP Deploy makes sure that the VMDKs are properly sized and assigned to the correct plex (in the case of HA pairs) during cluster -create and storage-add operations.

There are two major differences between vNAS and DAS with a RAID controller. The most immediate difference is that vNAS does not require a RAID controller. vNAS assumes that the underlying external array provides the data persistence and resiliency that a DAS with a RAID controller setup would provide. The second and more subtle difference has to do with NVRAM performance.

vNAS NVRAM

The ONTAP Select NVRAM is a VMDK. In other words, ONTAP Select emulates a byte addressable space (traditional NVRAM) on top of a block addressable device (VMDK). However, the performance of the

NVRAM is absolutely critical to the overall performance of the ONTAP Select node.

For DAS setups with a hardware RAID controller, the hardware RAID controller cache acts as the de facto NVRAM cache, because all writes to the NVRAM VMDK are first hosted in the RAID controller cache.

For VNAS architectures, ONTAP Deploy automatically configures ONTAP Select nodes with a boot argument called Single Instance Data Logging (SIDL). When this boot argument is present, ONTAP Select bypasses the NVRAM and writes the data payload directly to the data aggregate. The NVRAM is only used to record the address of the blocks changed by the WRITE operation. The benefit of this feature is that it avoids a double write: one write to NVRAM and a second write when the NVRAM is destaged. This feature is only enabled for vNAS because local writes to the RAID controller cache have a negligible additional latency.

The SIDL feature is not compatible with all ONTAP Select storage efficiency features. The SIDL feature can be disabled at the aggregate level using the following command:

```
storage aggregate modify -aggregate aggr-name -single-instance-data-logging off
```

Note that write performance is affected if the SIDL feature is turned off. It is possible to re-enable the SIDL feature after all the storage efficiency policies on all the volumes in that aggregate are disabled:

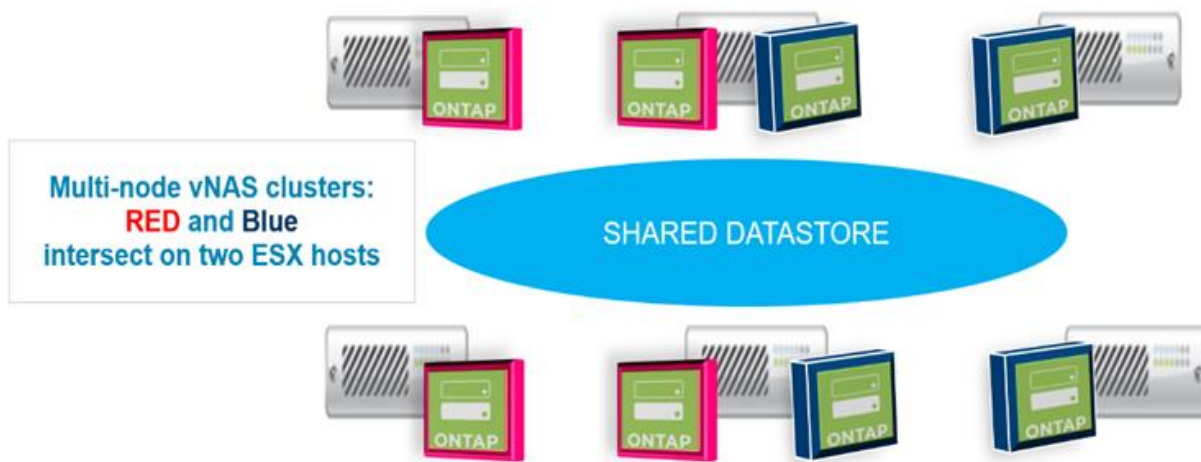
```
volume efficiency stop -all true -vserver * -volume * (all volumes in the affected aggregate)
```

Collocating ONTAP Select Nodes When Using vNAS

ONTAP Select includes support for multinode ONTAP Select clusters on shared storage. ONTAP Deploy enables the configuration of multiple ONTAP Select nodes on the same ESX host as long as these nodes are not part of the same cluster. Note that this configuration is only valid for VNAS environments (shared datastores). Multiple ONTAP Select instances per host are not supported when using DAS storage because these instances compete for the same hardware RAID controller.

ONTAP Deploy makes sure that the initial deployment of the multinode VNAS cluster does not place multiple ONTAP Select instances from the same cluster on the same host. The following figure shows for an example of a correct deployment of two four-node clusters that intersect on two hosts.

Initial deployment of multinode VNAS clusters



After deployment, the ONTAP Select nodes can be migrated between hosts. This could result in nonoptimal and unsupported configurations for which two or more ONTAP Select nodes from the same cluster share the same underlying host. NetApp recommends the manual creation of VM anti-affinity rules so that VMware automatically maintains physical separation between the nodes of the same cluster, not just the nodes from the same HA pair.



Anti-affinity rules require that DRS is enabled on the ESX cluster.

See the following example on how to create an anti-affinity rule for the ONTAP Select VMs. If the ONTAP Select cluster contains more than one HA pair, all nodes in the cluster must be included in this rule.



▼ Services

vSphere DRS
vSphere Availability

▼ vSAN

General
Disk Management
Fault Domains & Stretched Cluster
Health and Performance
iSCSI Targets
iSCSI Initiator Groups
Configuration Assist
Updates

▼ Configuration

General
Licensing
VMware EVC
VM/Host Groups

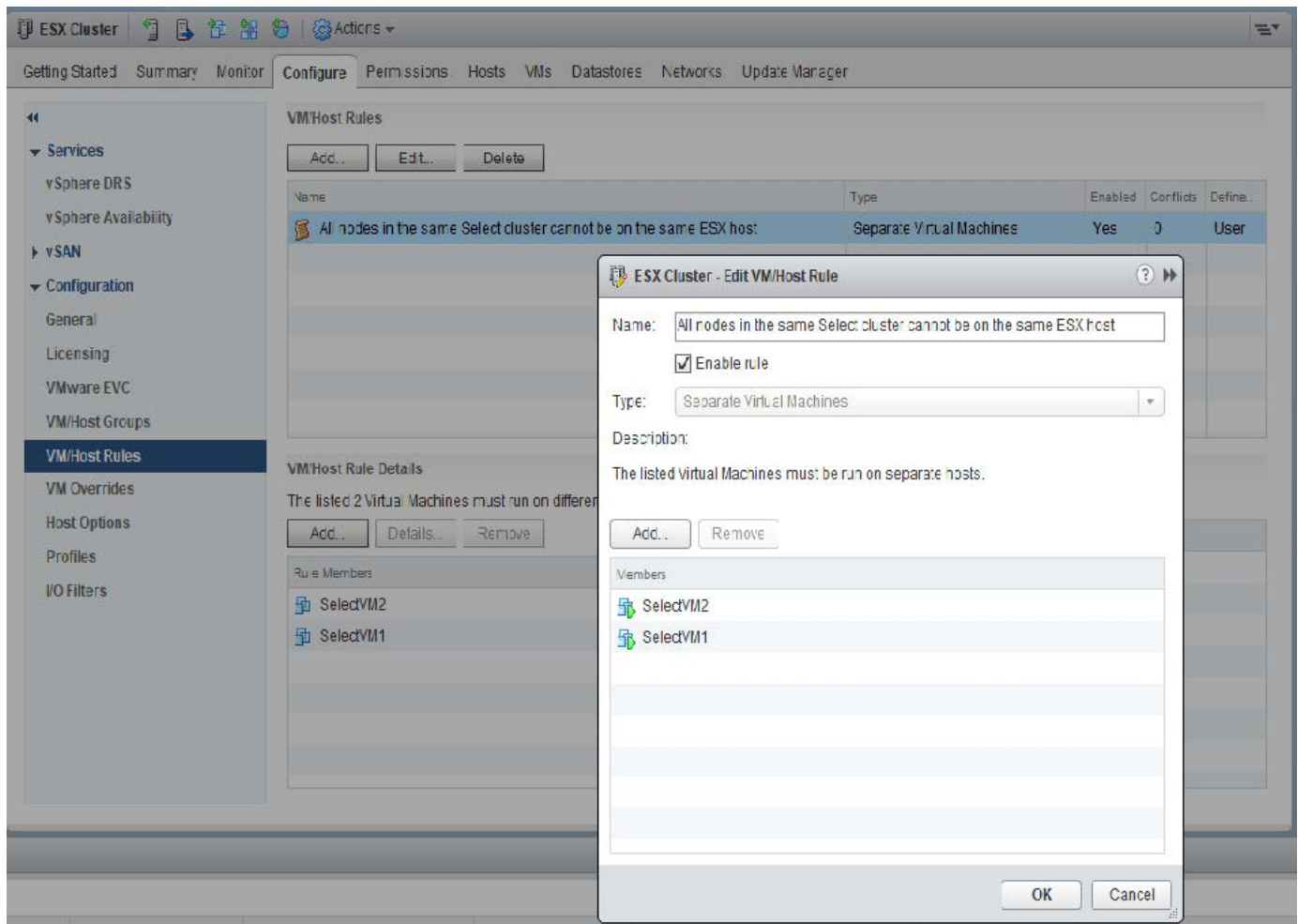
VM/Host Rules

VM Overrides
Host Options
Profiles
I/O Filters

VM/Host Rules

Name	Type	Enabled	Conflicts	Defined By
This list is empty.				

No VM/Host rule selected



Two or more ONTAP Select nodes from the same ONTAP Select cluster could potentially be found on the same ESX host for one of the following reasons:

- DRS is not present due to VMware vSphere license limitations or if DRS is not enabled.
- The DRS anti-affinity rule is bypassed because a VMware HA operation or administrator-initiated VM migration takes precedence.

Note that ONTAP Deploy does not proactively monitor the ONTAP Select VM locations. However, a cluster refresh operation reflects this unsupported configuration in the ONTAP Deploy logs:

 UnsupportedClusterConfiguration cluster 2018-05-16 11:41:19-04:00 ONTAP Select Deploy does not support multiple nodes within the same cluster sharing the same host:

Increasing storage capacity

ONTAP Deploy can be used to add and license additional storage for each node in an ONTAP Select cluster.

The storage-add functionality in ONTAP Deploy is the only way to increase the storage under management, and directly modifying the ONTAP Select VM is not supported. The following figure shows the “+” icon that initiates the storage-add wizard.

Cluster Details	
Name	onenode95IP15
ONTAP Image Version	9.5RC1
IPv4 Address	10.193.83.15
Netmask	255.255.255.128
Gateway	10.193.83.1
Last Refresh	-
Cluster Size	Single node cluster
Licensing	licensed
Domain Names	-
Server IP Addresses	-
NTP Server	216.239.35.0
Node Details	
Node	
Node	onenode95IP15-01 — 1.3 TB +
Host	10.193.39.54 — (Small (4 CPU, 16 GB Memory))

The following considerations are important for the success of the capacity-expansion operation. Adding capacity requires the existing license to cover the total amount of space (existing plus new). A storage-add operation that results in the node exceeding its licensed capacity fails. A new license with sufficient capacity should be installed first.

If the extra capacity is added to an existing ONTAP Select aggregate, then the new storage pool (datastore) should have a performance profile similar to that of the existing storage pool (datastore). Note that it is not possible to add non-SSD storage to an ONTAP Select node installed with an AFF-like personality (flash enabled). Mixing DAS and external storage is also not supported.

If locally attached storage is added to a system to provide for additional local (DAS) storage pools, you must build an additional RAID group and LUN (or LUNs). Just as with FAS systems, care should be taken to make sure that the new RAID group performance is similar to that of the original RAID group if you are adding new space to the same aggregate. If you are creating a new aggregate, the new RAID group layout could be different if the performance implications for the new aggregate are well understood.

The new space can be added to that same datastore as an extent if the total size of the datastore does not exceed the ESX-supported maximum datastore size. Adding a datastore extent to the datastore in which ONTAP Select is already installed can be done dynamically and does not affect the operations of the ONTAP Select node.

If the ONTAP Select node is part of an HA pair, some additional issues should be considered.

In an HA pair, each node contains a mirror copy of the data from its partner. Adding space to node 1 requires that an identical amount of space is added to its partner, node 2, so that all the data from node 1 is replicated to node 2. In other words, the space added to node 2 as part of the capacity-add operation for node 1 is not visible or accessible on node 2. The space is added to node 2 so that node 1 data is fully protected during an HA event.

There is an additional consideration with regard to performance. The data on node 1 is synchronously replicated to node 2. Therefore, the performance of the new space (datastore) on node 1 must match the performance of the new space (datastore) on node 2. In other words, adding space on both nodes, but using different drive technologies or different RAID group sizes, can lead to performance issues.

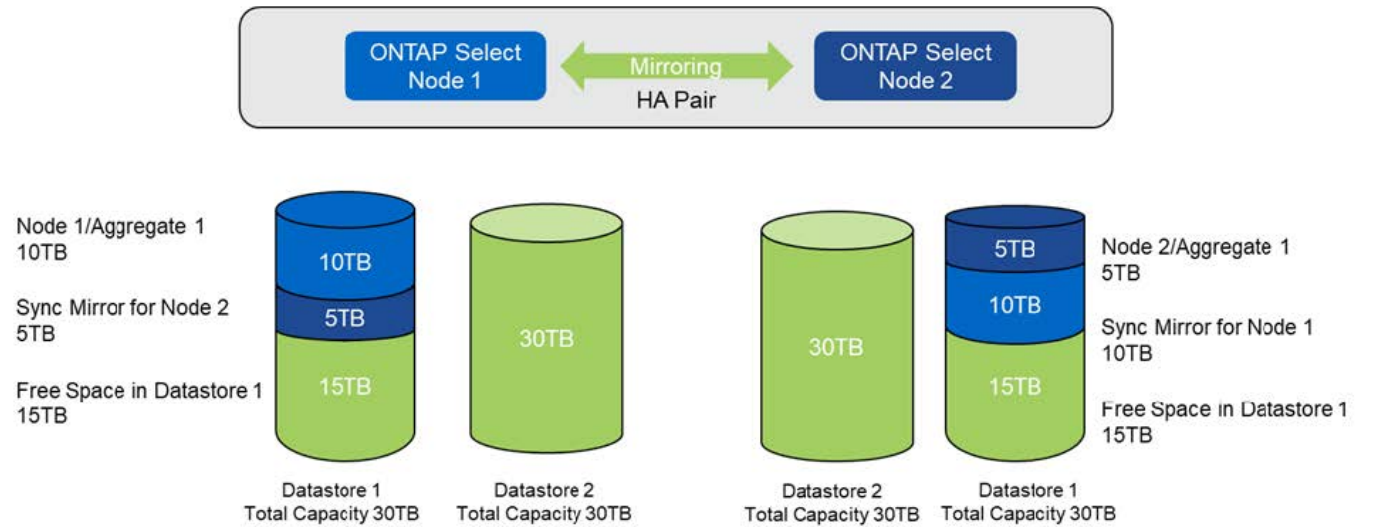
This is due to the RAID SyncMirror operation used to maintain a copy of the data on the partner node.

To increase user-accessible capacity on both nodes in an HA pair, two storage-add operations must be performed, one for each node. Each storage-add operation requires additional space on both nodes. The total space required on each node is equal to the space required on node 1 plus the space required on node 2.

Initial setup is with two nodes, each node having two datastores with 30TB of space in each datastore. ONTAP Deploy creates a two-node cluster, with each node consuming 10TB of space from datastore 1. ONTAP Deploy configures each node with 5TB of active space per node.

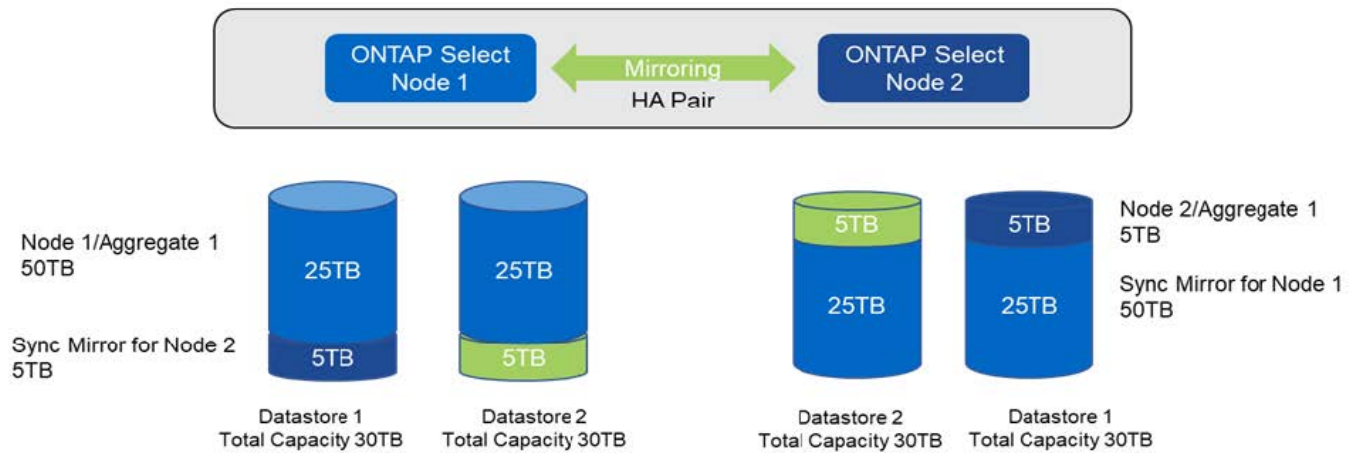
The following figure shows the results of a single storage-add operation for node 1. ONTAP Select still uses an equal amount of storage (15TB) on each node. However, node 1 has more active storage (10TB) than node 2 (5TB). Both nodes are fully protected as each node hosts a copy of the other node’s data. There is additional free space left in datastore 1, and datastore 2 is still completely free.

Capacity distribution: allocation and free space after a single storage-add operation



Two additional storage-add operations on node 1 consume the rest of datastore 1 and a part of datastore 2 (using the capacity cap). The first storage-add operation consumes the 15TB of free space left in datastore 1. The following figure shows the result of the second storage-add operation. At this point, node 1 has 50TB of active data under management, while node 2 has the original 5TB.

Capacity distribution: allocation and free space after two additional storage-add operation for node 1



The maximum VMDK size used during capacity add operations is 16TB. The maximum VMDK size used during cluster create operations continues to be 8TB. ONTAP Deploy creates correctly sized VMDKs depending on your configuration (a single-node or multinode cluster) and the amount of capacity being added. However, the maximum size of each VMDK should not exceed 8TB during the cluster create operations and 16TB during the storage-add operations.

Increasing capacity for ONTAP Select with Software RAID

The storage-add wizard can similarly be used to increase capacity under management for ONTAP Select nodes using software RAID. The wizard only presents those DAS SDD drives that are available and can be mapped as RDMs to the ONTAP Select VM.

Though it is possible to increase the capacity license by a single TB, when working with software RAID, it is not possible to physically increase the capacity by a single TB. Similar to adding disks to a FAS or AFF array, certain factors dictate the minimum amount of storage that can be added in a single operation.

Note that in an HA pair, adding storage to node 1 requires that an identical number of drives is also available on the node's HA pair (node 2). Both the local drives and the remote disks are used by one storage-add operation on node 1. That is to say, the remote drives are used to make sure that the new storage on node 1 is replicated and protected on node 2. In order to add locally usable storage on node 2, a separate storage-add operation and a separate and equal number of drives must be available on both nodes.

ONTAP Select partitions any new drives into the same root, data, and data partitions as the existing drives. The partitioning operation takes place during the creation of a new aggregate or during the expansion on an existing aggregate. The size of the root partition stripe on each disk is set to match the existing root partition size on the existing disks. Therefore, each one of the two equal data partition sizes can be calculated as the disk total capacity minus the root partition size divided by two. The root partition stripe size is variable, and it is computed during the initial cluster setup as follows. Total root space required (68GB for a single-node cluster and 136GB for HA pairs) is divided across the initial number of disks minus any spare and parity drives. The root partition stripe size is maintained to be constant on all the drives being added to the system.

If you are creating a new aggregate, the minimum number of drives required varies depending on the RAID type and whether the ONTAP Select node is part of an HA pair.

If adding storage to an existing aggregate, some additional considerations are necessary. It is possible to add drives to an existing RAID group, assuming that the RAID group is not at the maximum limit already. Traditional FAS and AFF best practices for adding spindles to existing RAID groups also apply here, and creating a hot spot on the new spindle is a potential concern. In addition, only drives of equal or larger data partition sizes can be added to an existing RAID group. As explained above, the data partition size is not the same as drive raw size. If the data partitions being added are larger than the existing partitions, the new drives is right-sized. In other words, a portion of capacity of each new drive remains unutilized.

It is also possible to use the new drives to create a new RAID group as part of an existing aggregate. In this case, the RAID group size should match the existing RAID group size.

Storage efficiency support

ONTAP Select provides storage efficiency options that are similar to the storage efficiency options present on FAS and AFF arrays.

ONTAP Select virtual NAS (vNAS) deployments using all-flash VSAN or generic flash arrays should follow the best practices for ONTAP Select with non-SSD DAS storage.

An AFF-like personality is automatically enabled on new installations as long as you have DAS storage with SSD drives and a Premium license.

With an AFF-like personality, the following inline SE features are automatically enabled during installation:

- Inline zero pattern detection
- Volume inline deduplication
- Volume background deduplication
- Adaptive inline compression
- Inline data compaction
- Aggregate inline deduplication
- Aggregate background deduplication

To verify that ONTAP Select has enabled all the default storage efficiency policies, run the following command on a newly created volume:

```

<system name>::> set diag
Warning: These diagnostic commands are for use by NetApp personnel only.
Do you want to continue? {y|n}: y
twonode95IP15::~*> sis config
Vserver:                                SVM1
Volume:                                _export1_NFS_volume
Schedule:                               -
Policy:                                auto
Compression:                            true
Inline Compression:                      true
Compression Type:                        adaptive
Application IO Si                        8K
Compression Algorithm:                   lzopro
Inline Dedupe:                           true
Data Compaction:                         true
Cross Volume Inline Deduplication:       true
Cross Volume Background Deduplication:   true

```



For ONTAP Select upgrades from 9.4, you must install ONTAP Select 9.4 on DAS SSD storage with a Premium license. In addition, the **Enable Storage Efficiencies** check box must be checked during initial cluster installation with ONTAP Deploy. Enabling an AFF-like personality post-ONTAP upgrade when prior conditions have not been met requires the manual creation of a boot argument and a node reboot. Contact technical support for further details.

The following table summarizes the various storage efficiency options available, enabled by default, or not enabled by default but recommended, depending on the ONTAP Select version and media type.

ONTAP Select storage efficiency configurations

ONTAP Select Features	9.6 / 9.5 Premium or Premium XL ⁴ (DAS SSD)	9.4 ¹ / 9.3 ² Premium (DAS SSD)	9.6 / 9.5 / 9.4 ¹ / 9.3 ² All Licenses (DAS HDD)	9.6 All Licenses (vNAS)	9.5 / 9.4 ¹ / 9.3 ² Premium or Standard (vNAS) ³
Inline zero detection	Yes (default)	Yes Enabled by user on a per-volume basis	Yes Enabled by user on a per-volume basis	Yes Enabled by user on a per-volume basis	Not supported
Volume inline deduplication	Yes (default)	Yes (recommended) Enabled by user on a per-volume basis	Not available	Not supported	Not supported

ONTAP Select Features	9.6 / 9.5 Premium or Premium XL⁴ (DAS SSD)	9.4¹ / 9.3² Premium (DAS SSD)	9.6 / 9.5 / 9.4¹ / 9.3² All Licenses (DAS HDD)	9.6 All Licenses (vNAS)	9.5 / 9.4¹ / 9.3² Premium or Standard (vNAS)³
32K inline compression (secondary compression)	Yes Enabled by user on a per volume basis.	Yes Enabled by user on a per-volume basis	Yes Enabled by user on a per-volume basis	Not supported	Not supported
8K inline compression (adaptive compression)	Yes (default)	Yes (recommended) Enabled by user on a per-volume basis	Yes Enabled by user on a per volume basis	Not supported	Not supported
Background compression	Not supported	Not supported	Yes Enabled by user on a per volume basis	Yes Enabled by user on a per-volume basis	Not supported
Compression scanner	Yes	Yes Enabled by user on a per-volume basis	Yes	Yes Enabled by user on a per-volume basis	Not supported
Inline data compaction	Yes (default)	Yes (recommended) Enabled by user on a per-volume basis	Yes Enabled by user on a per volume basis	Not supported	Not supported
Compaction scanner	Yes	Yes Enabled by user on a per-volume basis	Yes	Not supported	Not supported
Aggregate inline deduplication	Yes (default)	Yes (recommended) Enabled by user on a per volume basis with space guarantee = none)	N/A	Not supported	Not supported
Volume background deduplication	Yes (default)	Yes (recommended)	Yes Enabled by user on a per volume basis	Yes Enabled by user on a per-volume basis	Not supported

ONTAP Select Features	9.6 / 9.5 Premium or Premium XL ⁴ (DAS SSD)	9.4 ¹ / 9.3 ² Premium (DAS SSD)	9.6 / 9.5 / 9.4 ¹ / 9.3 ² All Licenses (DAS HDD)	9.6 All Licenses (vNAS)	9.5 / 9.4 ¹ / 9.3 ² Premium or Standard (vNAS) ³
Aggregate background deduplication	Yes (default)	Yes (recommended) Enabled by user on a per volume basis with space guarantee = none)	N/A	Not supported	Not supported

¹ONTAP Select 9.4 on DAS SSDs (requires Premium license) allows existing data in an aggregate to be deduped using aggregate-level background cross volume scanners. This one-time operation is performed manually for volumes created before 9.4.

²ONTAP Select 9.3 on DAS SSDs (requires Premium license) supports aggregate-level background deduplication; however, this feature must be enabled after creating the aggregate.

³ONTAP Select 9.5 vNAS by default does not support any storage efficiency policies. Review the vNAS section for details on Single Instance Data Logging (SIDL).

⁴ONTAP Select 9.6 supports a new license (Premium XL) and a new VM size (large). However, the large VM is only supported for DAS configurations using software RAID. Hardware RAID and vNAS configurations are not supported with the large ONTAP Select VM in the current release.

Notes on upgrade behavior for DAS SSD configurations

After upgrading to ONTAP Select 9.5 or later, wait for the `system node upgrade-revert show` command to indicate that the upgrade has completed before verifying the storage efficiency values for existing volumes.

On a system upgraded to ONTAP Select 9.5 or later, a new volume created on an existing aggregate or a newly created aggregate has the same behavior as a volume created on a fresh deployment. Existing volumes that undergo the ONTAP Select code upgrade have most of the same storage efficiency policies as a newly created volume with some variations:

Scenario 1 If no storage efficiency policies were enabled on a volume prior to the upgrade, then:

- Volumes with `space guarantee = volume` do not have inline data-compaction, aggregate inline deduplication, and aggregate background deduplication enabled. These options can be enabled post-upgrade.
- Volumes with `space guarantee = none` do not have background compression enabled. This option can be enabled post upgrade.
- Storage efficiency policy on the existing volumes is set to auto after upgrade.

Scenario 2 If some storage efficiencies are already enabled on a volume prior to the upgrade, then:

- Volumes with `space guarantee = volume` do not see any difference after upgrade.
- Volumes with `space guarantee = none` have aggregate background deduplication turned on.
- Volumes with `storage policy inline-only` have their policy set to auto.
- Volumes with user defined storage efficiency policies have no change in policy, with the exception of volumes with `space guarantee = none`. These volumes have aggregate background deduplication enabled.

Notes on Upgrade Behavior for DAS HDD Configuration

Storage efficiency features enabled prior to the upgrade are retained after the upgrade to ONTAP Select 9.5 or later. If no storage efficiencies were enabled prior to the upgrade, no storage efficiencies are enabled post-upgrade.

Networking

Networking: General concepts and characteristics

First become familiar with general networking concepts that apply to the ONTAP Select environment. Then explore the specific characteristics and options available with the single-node and multi-node clusters.

Physical networking

The physical network supports an ONTAP Select cluster deployment primarily by providing the underlying layer two switching infrastructure. The configuration related to the physical network includes both the hypervisor host and the broader switched network environment.

Host NIC options

Each ONTAP Select hypervisor host must be configured with either two or four physical ports. The exact configuration you choose depends on several factors, including:

- Whether the cluster contains one or multiple ONTAP Select hosts
- What hypervisor operating system is used
- How the virtual switch is configured
- Whether LACP is used with the links or not

Physical switch configuration

You must make sure that the configuration of the physical switches supports the ONTAP Select deployment. The physical switches are integrated with the hypervisor-based virtual switches. The exact configuration you choose depends on several factors. The primary considerations include the following:

- How will you maintain separation between the internal and external networks?
- Will you maintain a separation between the data and management networks?
- How will the layer two VLANs be configured?

Logical networking

ONTAP Select uses two different logical networks, separating the traffic according to type. Specifically, traffic can flow among the hosts within the cluster as well as to the storage clients and other machines outside of the cluster. The virtual switches managed by the hypervisors help support the logical network.

Internal network

With a multi-node cluster deployment, the individual ONTAP Select nodes communicate using an isolated “internal” network. This network is not exposed or available outside of the nodes in the ONTAP Select cluster.



The internal network is only present with a multi-node cluster.

The internal network has the following characteristics:

- Used to process ONTAP intra-cluster traffic including:
 - Cluster
 - High Availability Interconnect (HA-IC)
 - RAID Synch Mirror (RSM)
- Single layer-two network based on a VLAN
- Static IP addresses are assigned by ONTAP Select:
 - IPv4 only
 - DHCP not used
 - Link-local address
- The MTU size is 9000 bytes by default and can be adjusted within 7500-9000 range (inclusive)

External network

The external network processes traffic between the nodes of an ONTAP Select cluster and the external storage clients as well as the other machines. The external network is a part of every cluster deployment and has the following characteristics:

- Used to process ONTAP traffic including:
 - Data (NFS, CIFS, iSCSI)
 - Management (cluster and node; optionally SVM)

- Intercluster (optional)
- Optionally supports VLANs:
 - Data port group
 - Management port group
- IP addresses that are assigned based on the configuration choices of the administrator:
 - IPv4 or IPv6
- MTU size is 1500 bytes by default (can be adjusted)

The external network is present with clusters of all sizes.

Virtual machine networking environment

The hypervisor host provides several networking features.

ONTAP Select relies on the following capabilities exposed through the virtual machine:

Virtual machine ports

There are several ports available for use by ONTAP Select. They are assigned and used based on several factors, including the size of the cluster.

Virtual switch

The virtual switch software within the hypervisor environment, whether vSwitch (VMware) or Open vSwitch (KVM), joins the ports exposed by the virtual machine with the physical Ethernet NIC ports. You must configure a vSwitch for every ONTAP Select host, as appropriate for your environment.

Single and multiple node network configurations

ONTAP Select supports both single node and multinode network configurations.

Single node network configuration

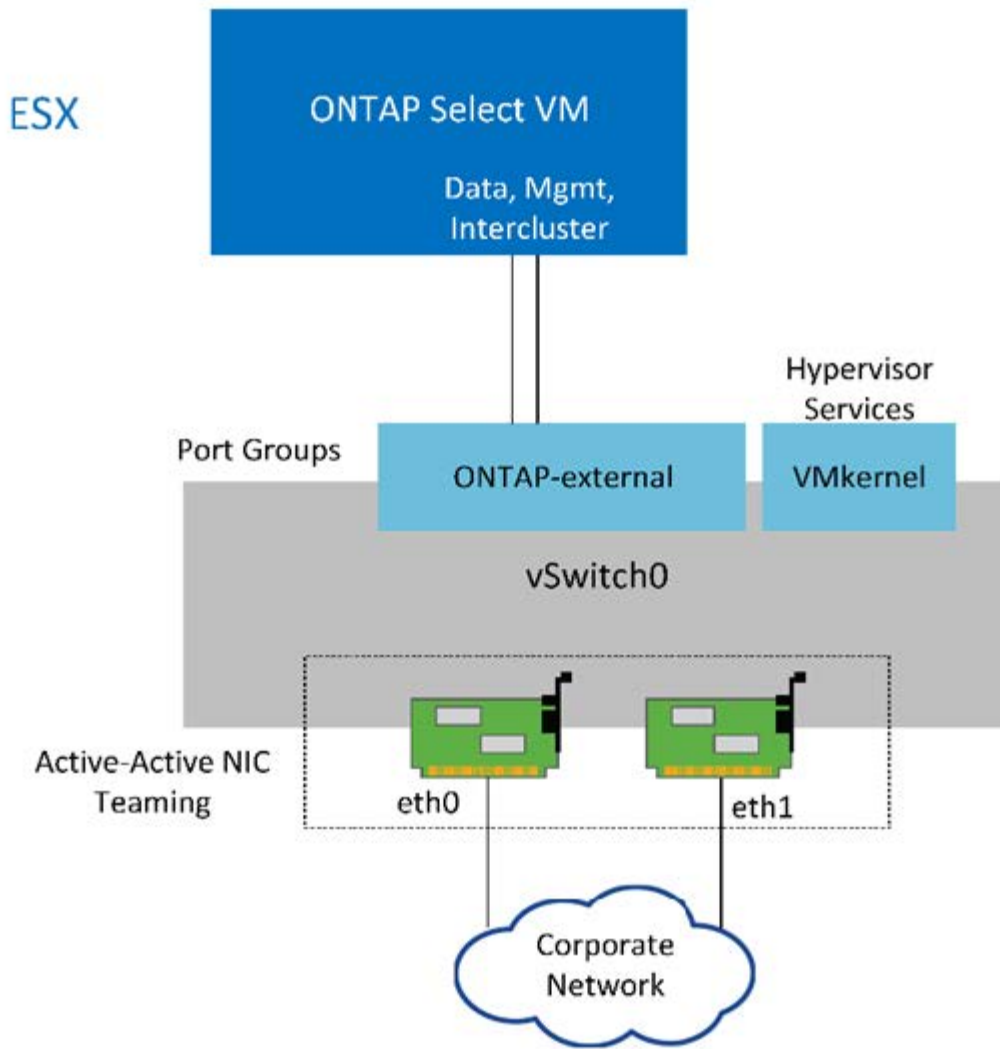
Single-node ONTAP Select configurations do not require the ONTAP internal network, because there is no cluster, HA, or mirror traffic.

Unlike the multinode version of the ONTAP Select product, each ONTAP Select VM contains three virtual network adapters, presented to ONTAP network ports e0a, e0b, and e0c.

These ports are used to provide the following services: management, data, and intercluster LIFs.

The relationship between these ports and the underlying physical adapters can be seen in the following figure, which depicts one ONTAP Select cluster node on the ESX hypervisor.

Network configuration of single-node ONTAP Select cluster



Even though two adapters are sufficient for a single-node cluster, NIC teaming is still required.

LIF assignment

As explained in the multinode LIF assignment section of this document, IPspaces are used by ONTAP Select to keep cluster network traffic separate from data and management traffic. The single-node variant of this platform does not contain a cluster network. Therefore, no ports are present in the cluster IPspace.



Cluster and node management LIFs are automatically created during ONTAP Select cluster setup. The remaining LIFs can be created post deployment.

Management and data LIFs (e0a, e0b, and e0c)

ONTAP ports e0a, e0b, and e0g are delegated as candidate ports for LIFs that carry the following types of traffic:

- SAN/NAS protocol traffic (CIFS, NFS, and iSCSI)
- Cluster, node, and SVM management traffic
- Intercluster traffic (SnapMirror and SnapVault)

Multinode network configuration

The multinode ONTAP Select network configuration consists of two networks.

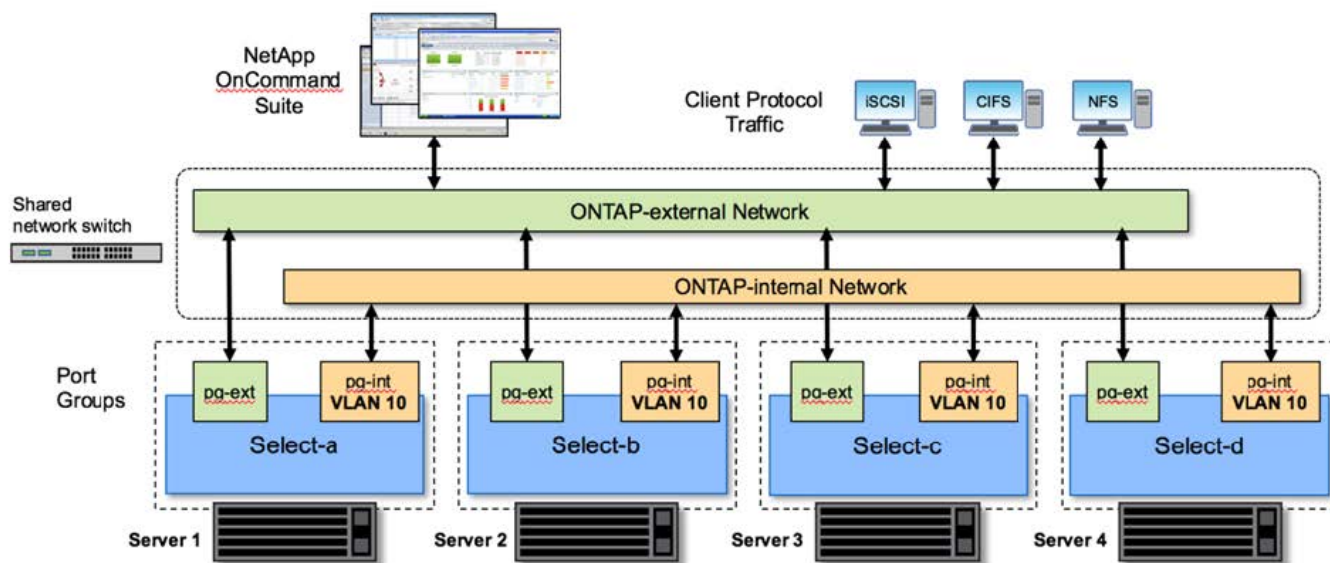
These are an internal network, responsible for providing cluster and internal replication services, and an external network, responsible for providing data access and management services. End-to-end isolation of traffic that flows within these two networks is extremely important in allowing you to build an environment that is suitable for cluster resiliency.

These networks are represented in the following figure, which shows a four-node ONTAP Select cluster running on a VMware vSphere platform. Six- and eight-node clusters have a similar network layout.



Each ONTAP Select instance resides on a separate physical server. Internal and external traffic is isolated using separate network port groups, which are assigned to each virtual network interface and allow the cluster nodes to share the same physical switch infrastructure.

Overview of an ONTAP Select multinode cluster network configuration



Each ONTAP Select VM contains seven virtual network adapters presented to ONTAP as a set of seven network ports, e0a through e0g. Although ONTAP treats these adapters as physical NICs, they are in fact virtual and map to a set of physical interfaces through a virtualized network layer. As a result, each hosting server does not require six physical network ports.



Adding virtual network adapters to the ONTAP Select VM is not supported.

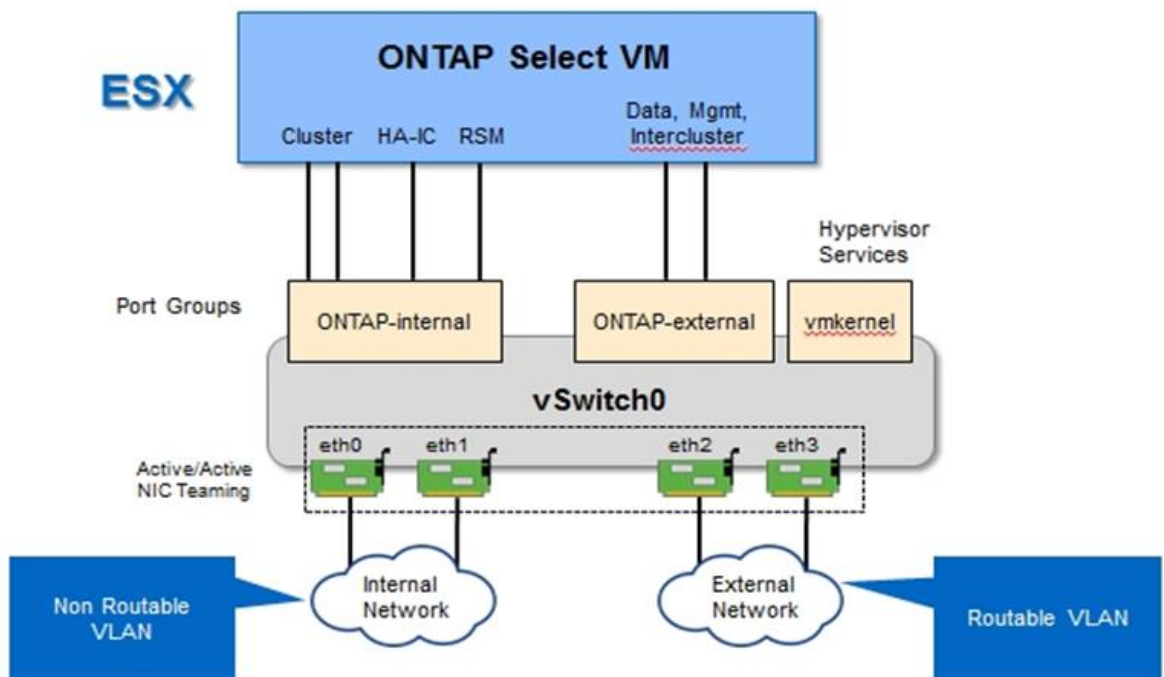
These ports are preconfigured to provide the following services:

- e0a, e0b, and e0g. Management and data LIFs
- e0c, e0d. Cluster network LIFs
- e0e. RSM
- e0f. HA interconnect

Ports e0a, e0b, and e0g reside on the external network. Although ports e0c through e0f perform several different functions, collectively they compose the internal Select network. When making network design decisions, these ports should be placed on a single layer-2 network. There is no need to separate these virtual adapters across different networks.

The relationship between these ports and the underlying physical adapters is illustrated in the following figure, which depicts one ONTAP Select cluster node on the ESX hypervisor.

Network configuration of a single node that is part of a multinode ONTAP Select cluster



Segregating internal and external traffic across different physical NICs prevents latencies from being introduced into the system due to insufficient access to network resources. Additionally, aggregation through NIC teaming makes sure that failure of a single network adapter does not prevent the ONTAP Select cluster node from accessing the respective network.

Note that both the external network and internal network port groups contain all four NIC adapters in a symmetrical manner. The active ports in the external network port group are the standby ports in the internal network. Conversely, the active ports in the internal network port group are the standby ports in the external network port group.

LIF assignment

With the introduction of IPspaces, ONTAP port roles have been deprecated. Like FAS arrays, ONTAP Select clusters contain both a default IPspace and a cluster IPspace. By placing network ports e0a, e0b, and e0g into the default IPspace and ports e0c and e0d into the cluster IPspace, those ports have essentially been walled off from hosting LIFs that do not belong. The remaining ports within the ONTAP Select cluster are consumed through the automatic assignment of interfaces providing internal services. They are not exposed through the ONTAP shell, as is the case with the RSM and HA interconnect interfaces.



Not all LIFs are visible through the ONTAP command shell. The HA interconnect and RSM interfaces are hidden from ONTAP and are used internally to provide their respective services.

The network ports and LIFs are explained in detail in the following sections.

Management and data LIFs (e0a, e0b, and e0g)

ONTAP ports e0a, e0b, and e0g are delegated as candidate ports for LIFs that carry the following types of traffic:

- SAN/NAS protocol traffic (CIFS, NFS, and iSCSI)
- Cluster, node, and SVM management traffic
- Intercluster traffic (SnapMirror and SnapVault)



Cluster and node management LIFs are automatically created during ONTAP Select cluster setup. The remaining LIFs can be created post deployment.

Cluster network LIFs (e0c, e0d)

ONTAP ports e0c and e0d are delegated as home ports for cluster interfaces. Within each ONTAP Select cluster node, two cluster interfaces are automatically generated during ONTAP setup using link local IP addresses (169.254.x.x).



These interfaces cannot be assigned static IP addresses, and additional cluster interfaces should not be created.

Cluster network traffic must flow through a low-latency, nonrouted layer-2 network. Due to cluster throughput and latency requirements, the ONTAP Select cluster is expected to be physically located within proximity (for example, multipack, single data center). Building four-node, six-node, or eight-node stretch cluster configurations by separating HA nodes across a WAN or across significant geographical distances is not supported. A stretched two-node configuration with a mediator is supported.

For details, see the section [Two-node stretched HA \(MetroCluster SDS\) best practices](#).



To make sure of maximum throughput for cluster network traffic, this network port is configured to use jumbo frames (7500 to 9000 MTU). For proper cluster operation, verify that jumbo frames are enabled on all upstream virtual and physical switches providing internal network services to ONTAP Select cluster nodes.

RAID SyncMirror traffic (e0e)

Synchronous replication of blocks across HA partner nodes occurs using an internal network interface residing on network port e0e. This functionality occurs automatically, using network interfaces configured by ONTAP during cluster setup, and requires no configuration by the administrator.



Port e0e is reserved by ONTAP for internal replication traffic. Therefore, neither the port nor the hosted LIF is visible in the ONTAP CLI or in System Manager. This interface is configured to use an automatically generated link local IP address, and the reassignment of an alternate IP address is not supported. This network port requires the use of jumbo frames (7500 to 9000 MTU).

HA interconnect (e0f)

NetApp FAS arrays use specialized hardware to pass information between HA pairs in an ONTAP cluster. Software-defined environments, however, do not tend to have this type of equipment available (such as InfiniBand or iWARP devices), so an alternate solution is needed. Although several possibilities were considered, ONTAP requirements placed on the interconnect transport required that this functionality be emulated in software. As a result, within an ONTAP Select cluster, the functionality of the HA interconnect (traditionally provided by hardware) has been designed into the OS, using Ethernet as a transport mechanism.

Each ONTAP Select node is configured with an HA interconnect port, e0f. This port hosts the HA interconnect network interface, which is responsible for two primary functions:

- Mirroring the contents of NVRAM between HA pairs
- Sending/receiving HA status information and network heartbeat messages between HA pairs

HA interconnect traffic flows through this network port using a single network interface by layering remote direct memory access (RDMA) frames within Ethernet packets.



In a manner similar to the RSM port (e0e), neither the physical port nor the hosted network interface is visible to users from either the ONTAP CLI or from System Manager. As a result, the IP address of this interface cannot be modified, and the state of the port cannot be changed. This network port requires the use of jumbo frames (7500 to 9000 MTU).

ONTAP Select internal and external network

Characteristics of ONTAP Select internal and external networks.

ONTAP Select internal network

The internal ONTAP Select network, which is only present in the multinode variant of the product, is responsible for providing the ONTAP Select cluster with cluster communication, HA interconnect, and synchronous replication services. This network includes the following ports and interfaces:

- **e0c, e0d.** Hosting cluster network LIFs
- **e0e.** Hosting the RSM LIF
- **e0f.** Hosting the HA interconnect LIF

The throughput and latency of this network are critical in determining the performance and resiliency of the ONTAP Select cluster. Network isolation is required for cluster security and to make sure that system interfaces are kept separate from other network traffic. Therefore, this network must be used exclusively by the ONTAP Select cluster.



Using the Select internal network for traffic other than Select cluster traffic, such as application or management traffic, is not supported. There can be no other VMs or hosts on the ONTAP internal VLAN.

Network packets traversing the internal network must be on a dedicated VLAN-tagged layer-2 network. This can be accomplished by completing one of the following tasks:

- Assigning a VLAN-tagged port group to the internal virtual NICs (e0c through e0f) (VST mode)
- Using the native VLAN provided by the upstream switch where the native VLAN is not used for any other traffic (assign a port group with no VLAN ID, that is, EST mode)

In all cases, VLAN tagging for internal network traffic is done outside of the ONTAP Select VM.



Only ESX standard and distributed vSwitches are supported. Other virtual switches or direct connectivity between ESX hosts are not supported. The internal network must be fully opened; NAT or firewalls are not supported.

Within an ONTAP Select cluster, internal traffic and external traffic are separated using virtual layer-2 network objects known as port groups. Proper vSwitch assignment of these port groups is extremely important, especially for the internal network, which is responsible for providing cluster, HA interconnect, and mirror replication services. Insufficient network bandwidth to these network ports can cause performance degradation and even affect the stability of the cluster node. Therefore, four-node, six-node, and eight-node clusters require that the internal ONTAP Select network use 10Gb connectivity; 1Gb NICs are not supported. Tradeoffs can be made to the external network, however, because limiting the flow of incoming data to an ONTAP Select cluster does not affect its ability to operate reliably.

A two-node cluster can use either four 1Gb ports for internal traffic or a single 10Gb port instead of the two 10Gb ports required by the four-node cluster. In an environment in which conditions prevent the server from being fit with four 10Gb NIC cards, two 10Gb NIC cards can be used for the internal

network and two 1Gb NICs can be used for the external ONTAP network.

Internal network validation and troubleshooting

The internal network in a multinode cluster can be validated by using the network connectivity checker functionality. This function can be invoked from the Deploy CLI running the `network connectivity-check start` command.

Run the following command to view the output of the test:

```
network connectivity-check show --run-id X (X is a number)
```

This tool is only useful for troubleshooting the internal network in a multinode Select cluster. The tool should not be used to troubleshoot single-node clusters (including vNAS configurations), ONTAP Deploy to ONTAP Select connectivity, or client-side connectivity issues.

The cluster create wizard (part of the ONTAP Deploy GUI) includes the internal network checker as an optional step available during the creation of multinode clusters. Given the important role that the internal network plays in multinode clusters, making this step part of the cluster create workflow improves the success rate of cluster create operations.

Starting with ONTAP Deploy 2.10, the MTU size used by the internal network can be set between 7,500 and 9,000. The network connectivity checker can also be used to test MTU size between 7,500 and 9,000. The default MTU value is set to the value of the virtual network switch. That default would have to be replaced with a smaller value if a network overlay like VXLAN is present in the environment.

ONTAP Select external network

The ONTAP Select external network is responsible for all outbound communications by the cluster and, therefore, is present on both the single-node and multinode configurations. Although this network does not have the tightly defined throughput requirements of the internal network, the administrator should be careful not to create network bottlenecks between the client and ONTAP VM, because performance issues could be mischaracterized as ONTAP Select problems.



In a manner similar to internal traffic, external traffic can be tagged at the vSwitch layer (VST) and at the external switch layer (EST). In addition, the external traffic can be tagged by the ONTAP Select VM itself in a process known as VGT. See the section [Data and management traffic separation](#) for further details.

The following table highlights the major differences between the ONTAP Select internal and external networks.

Internal versus external network quick reference

Description	Internal Network	External Network
Network services	Cluster HA/IC RAID SyncMirror (RSM)	Data management Intercluster (SnapMirror and SnapVault)
Network isolation	Required	Optional
Frame size (MTU)	7,500 to 9,000	1,500 (default) 9,000 (supported)
IP address assignment	Autogenerated	User-defined
DHCP support	No	No

NIC teaming

To make sure that the internal and external networks have both the necessary bandwidth and resiliency characteristics required to provide high performance and fault tolerance, physical network adapter teaming is recommended. Two-node cluster configurations with a single 10Gb link are supported. However, the NetApp recommended best practice is to make use of NIC teaming on both the internal and the external networks of the ONTAP Select cluster.

MAC address generation

The MAC addresses assigned to all ONTAP Select network ports are generated automatically by the included deployment utility. The utility uses a platform-specific, organizationally unique identifier (OUI) specific to NetApp to make sure there is no conflict with FAS systems. A copy of this address is then stored in an internal database within the ONTAP Select installation VM (ONTAP Deploy), to prevent accidental reassignment during future node deployments. At no point should the administrator modify the assigned MAC address of a network port.

Supported network configurations

Select the best hardware and configure your network to optimize performance and resiliency.

Server vendors understand that customers have different needs and choice is critical. As a result, when purchasing a physical server, there are numerous options available when making network connectivity decisions. Most commodity systems ship with various NIC choices that provide single-port and multiport options with varying permutations of speed and throughput. Starting with ONTAP Select 9.8, 25Gb/s and 40Gb/s NIC adapters are supported with VMWare ESX.

Because the performance of the ONTAP Select VM is tied directly to the characteristics of the underlying hardware, increasing the throughput to the VM by selecting higher-speed NICs results in a higher-performing cluster and a better overall user experience. Four 10Gb NICs or two higher-speed NICs (25/40 Gb/s) can be used to achieve a high performance network layout. There are a number of other configurations that are also supported. For two-node clusters, 4 x 1Gb ports or 1 x 10Gb ports are

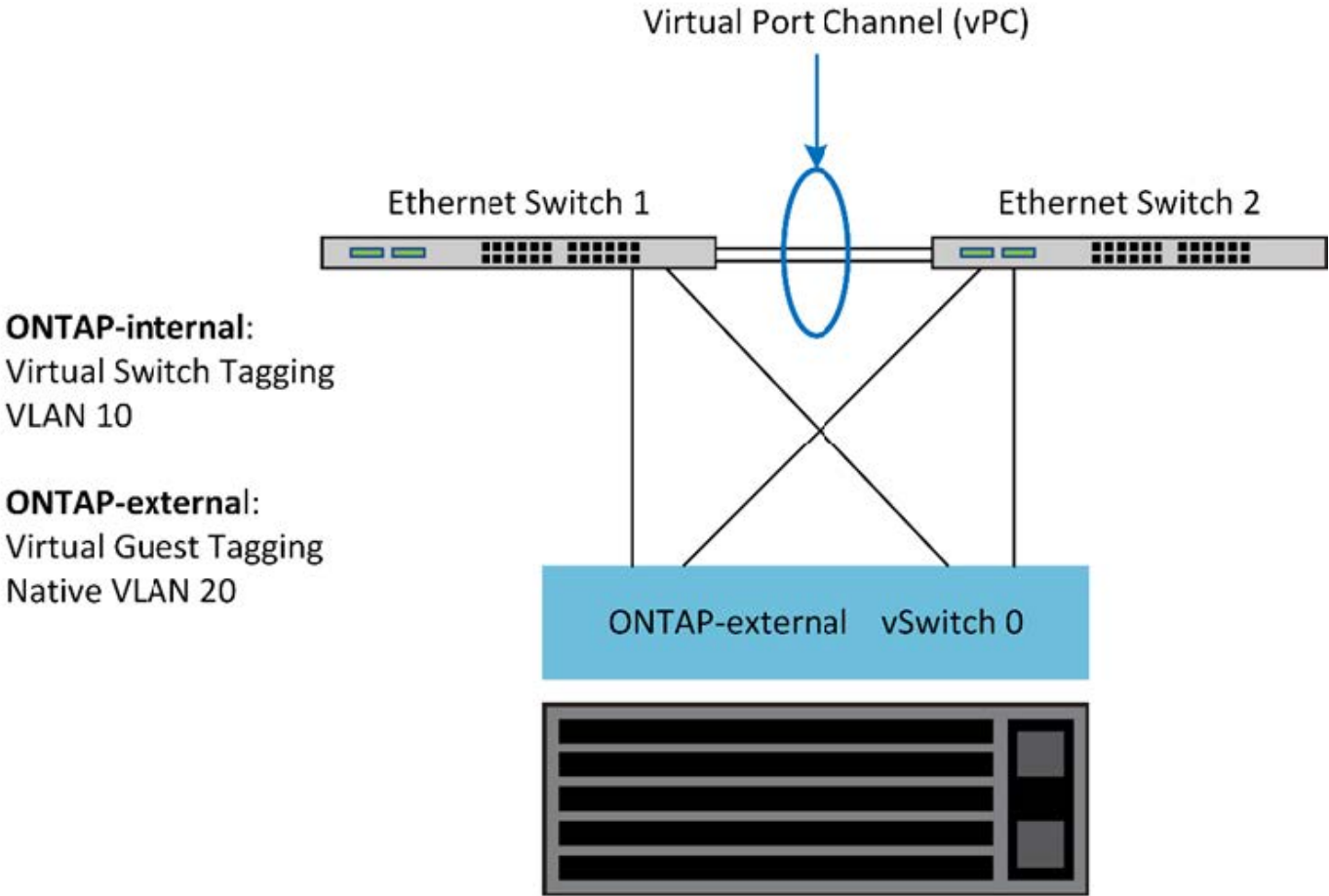
supported. For single node clusters, 2 x 1Gb ports are supported.

Network minimum and recommended configurations

	Minimum Requirements	Recommendations
Single node clusters	2 x 1Gb	2 x 10Gb
Two node clusters/MetroCluster SDS	4 x 1Gb or 1 x 10Gb	2 x 10Gb
4/6/8 node clusters	2 x 10Gb	4 x 10Gb or 2 x 25/40Gb

Network configuration using multiple physical switches

When sufficient hardware is available, NetApp recommends using the multiswitch configuration shown in the following figure, due to the added protection against physical switch failures.



VMWare vSphere vSwitch Configuration

ONTAP Select vSwitch configuration and load-balancing policies for two-NIC and four-NIC configurations.

ONTAP Select supports the use of both standard and distributed vSwitch configurations. Distributed vSwitches support link aggregation constructs (LACP). Link aggregation is a common network

construct used to aggregate bandwidth across multiple physical adapters. LACP is a vendor-neutral standard that provides an open protocol for network endpoints that bundle groups of physical network ports into a single logical channel. ONTAP Select can work with port groups that are configured as a Link Aggregation Group (LAG). However, NetApp recommends using the individual physical ports as simple uplink (trunk) ports to avoid the LAG configuration. In these cases, the best practices for standard and distributed vSwitches are identical.

This section describes the vSwitch configuration and load-balancing policies that should be used in both two-NIC and four-NIC configurations.

When configuring the port groups to be used by ONTAP Select, the following best practices should be followed; the load-balancing policy at the port-group level is Route Based on Originating Virtual Port ID. VMware recommends that STP be set to Portfast on the switch ports connected to the ESXi hosts.

All vSwitch configurations require a minimum of two physical network adapters bundled into a single NIC team. ONTAP Select supports a single 10Gb link for two-node clusters. However, it is a NetApp best practice to make sure of hardware redundancy through NIC aggregation.

On a vSphere server, NIC teams are the aggregation construct used to bundle multiple physical network adapters into a single logical channel, allowing the network load to be shared across all member ports. It's important to remember that NIC teams can be created without support from the physical switch. Load-balancing and failover policies can be applied directly to a NIC team, which is unaware of the upstream switch configuration. In this case, policies are only applied to outbound traffic.



Static port channels are not supported with ONTAP Select. LACP-enabled channels are supported with distributed vSwitches but using LACP LAGs may result in un-even load distribution across the LAG members.

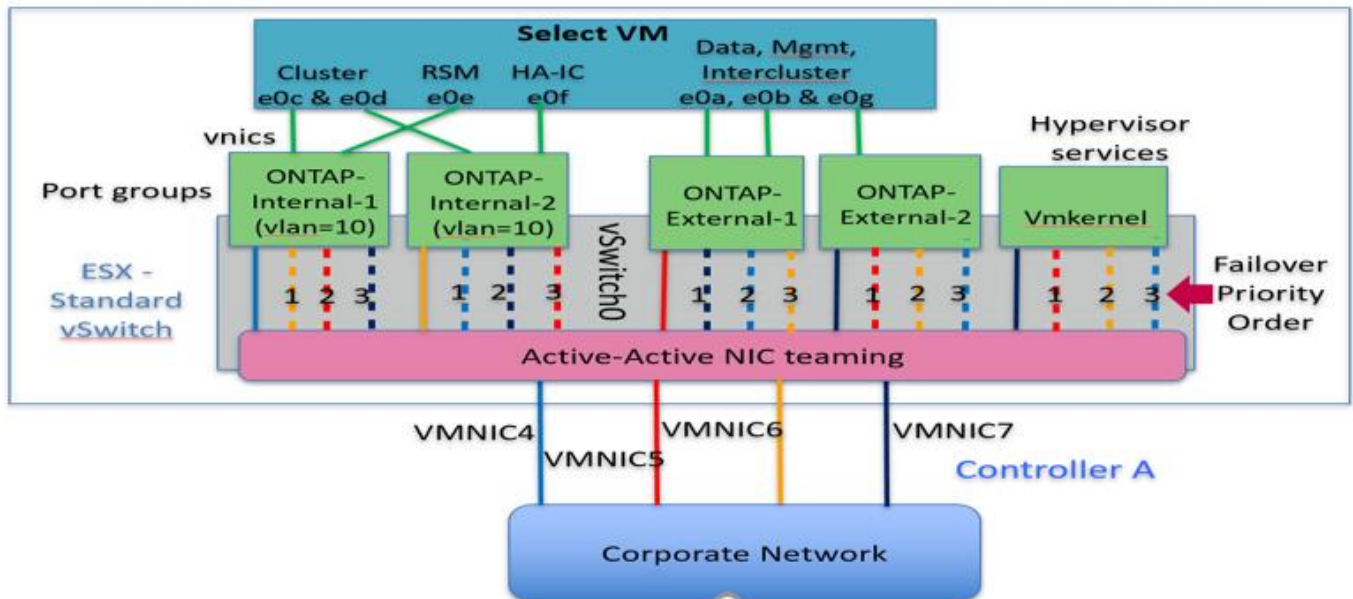
For single node clusters, ONTAP Deploy configures the ONTAP Select VM to use a port group for the external network and either the same port group or, optionally, a different port group for the cluster and node management traffic. For single node clusters, the desired number of physical ports can be added to the external port group as active adapters.

For multinode clusters, ONTAP Deploy configures each ONTAP Select VM to use one or two port groups for the internal network and separately, one or two port groups for the external network. The cluster and node management traffic can either use the same port group as the external traffic, or optionally a separate port group. The cluster and node management traffic cannot share the same port group with internal traffic.

Standard or distributed vSwitch and four physical ports per Node

Four port groups can be assigned to each node in a multinode cluster. Each port group has a single active physical port and three standby physical ports as in the following figure.

vSwitch with four physical ports per node



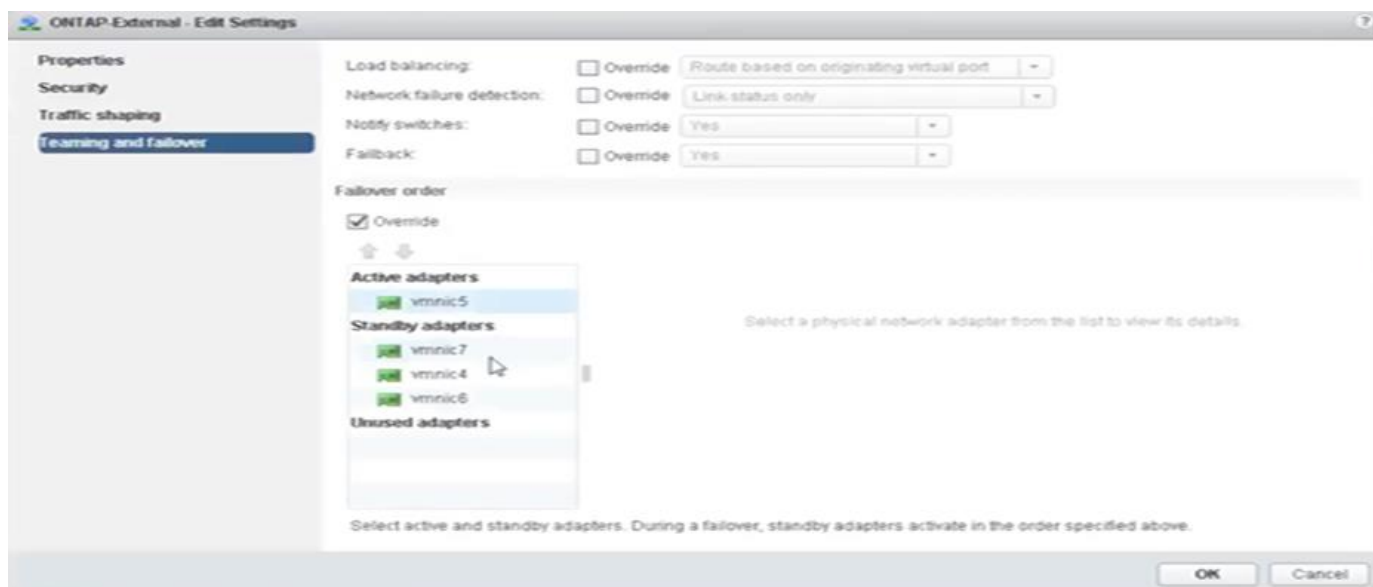
The order of the ports in the standby list is important. The following table provides an example of the physical port distribution across the four port groups.

Network minimum and recommended configurations

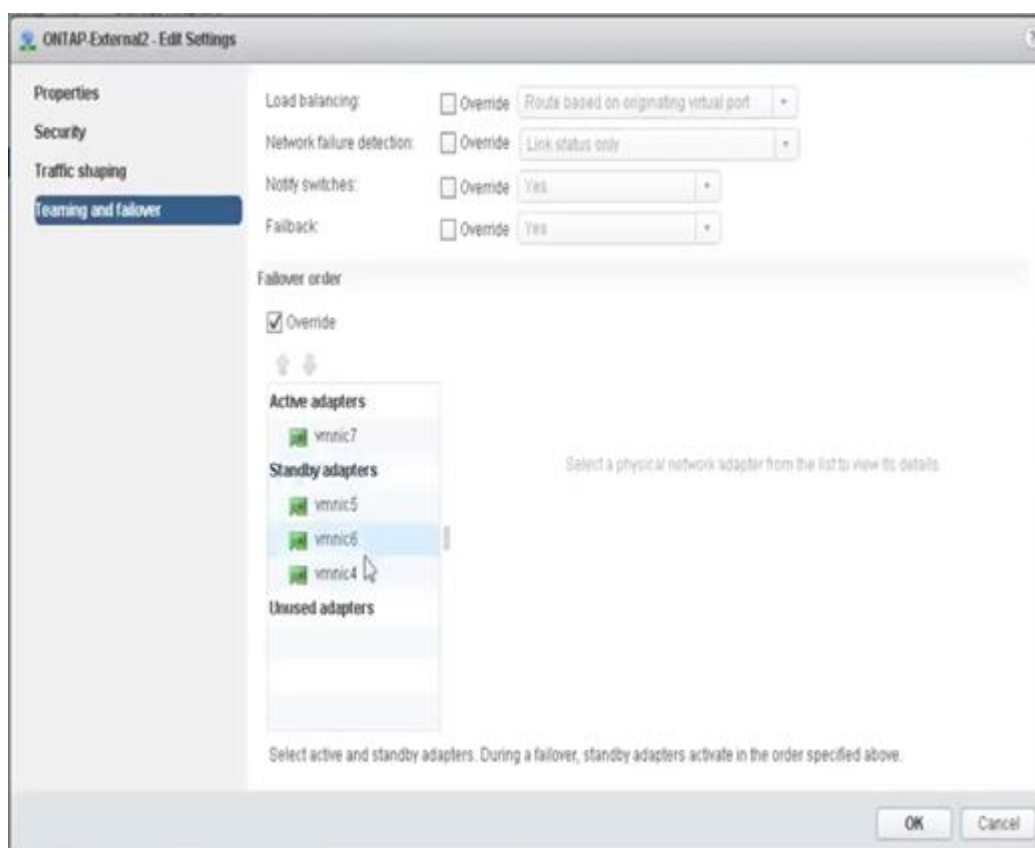
Port Group	External 1	External 2	Internal 1	Internal 2
Active	vmnic0	vmnic1	vmnic2	vmnic3
Standby 1	vmnic1	vmnic0	vmnic3	vmnic2
Standby 2	vmnic2	vmnic3	vmnic0	vmnic1
Standby 3	vmnic3	vmnic2	vmnic1	vmnic0

The following figures show the configurations of the external network port groups from the vCenter GUI (ONTAP-External and ONTAP-External2). Note that the active adapters are from different network cards. In this setup, vmnic 4 and vmnic 5 are dual ports on the same physical NIC, while vmnic 6 and vmnic 7 are similarly dual ports on a separate NIC (vnmics 0 through 3 are not used in this example). The order of the standby adapters provides a hierarchical fail over with the ports from the internal network being last. The order of internal ports in the standby list is similarly swapped between the two external port groups.

Part 1: ONTAP Select external port group configurations



Part 2: ONTAP Select external port group configurations



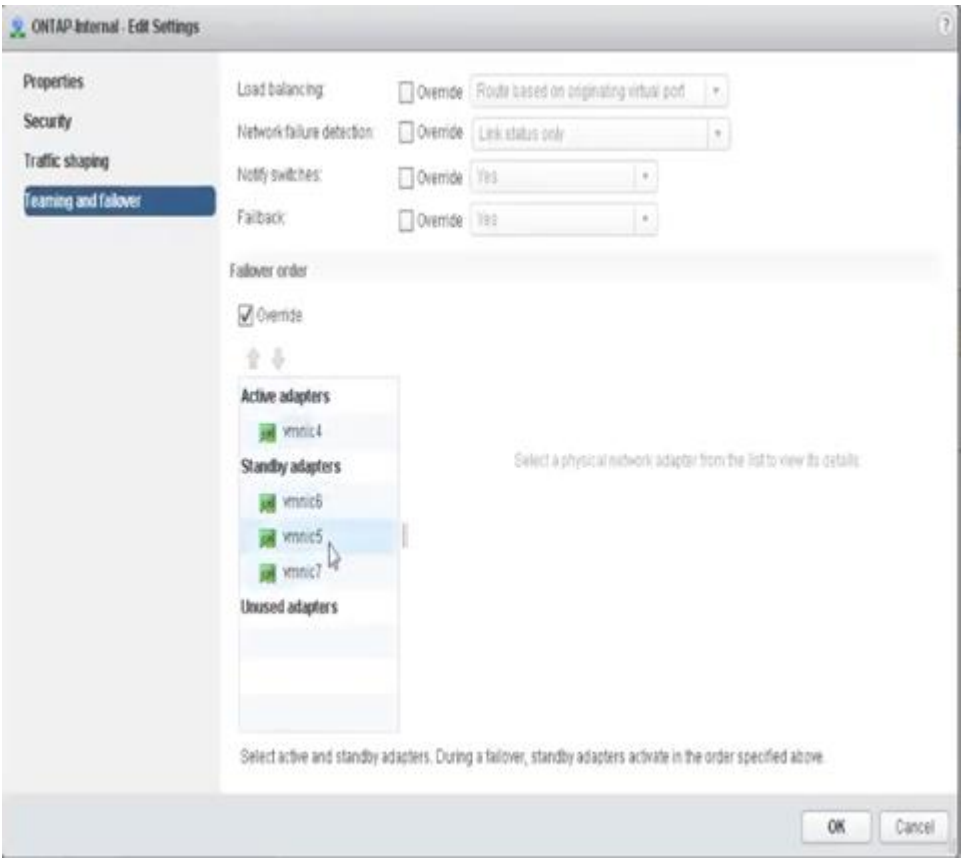
For readability, the assignments are as follows:

ONTAP-External	ONTAP-External2
Active adapters: vmnic5	Active adapters: vmnic7
Standby adapters: vmnic7, vmnic4, vmnic6	Standby adapters: vmnic5, vmnic6, vmnic4

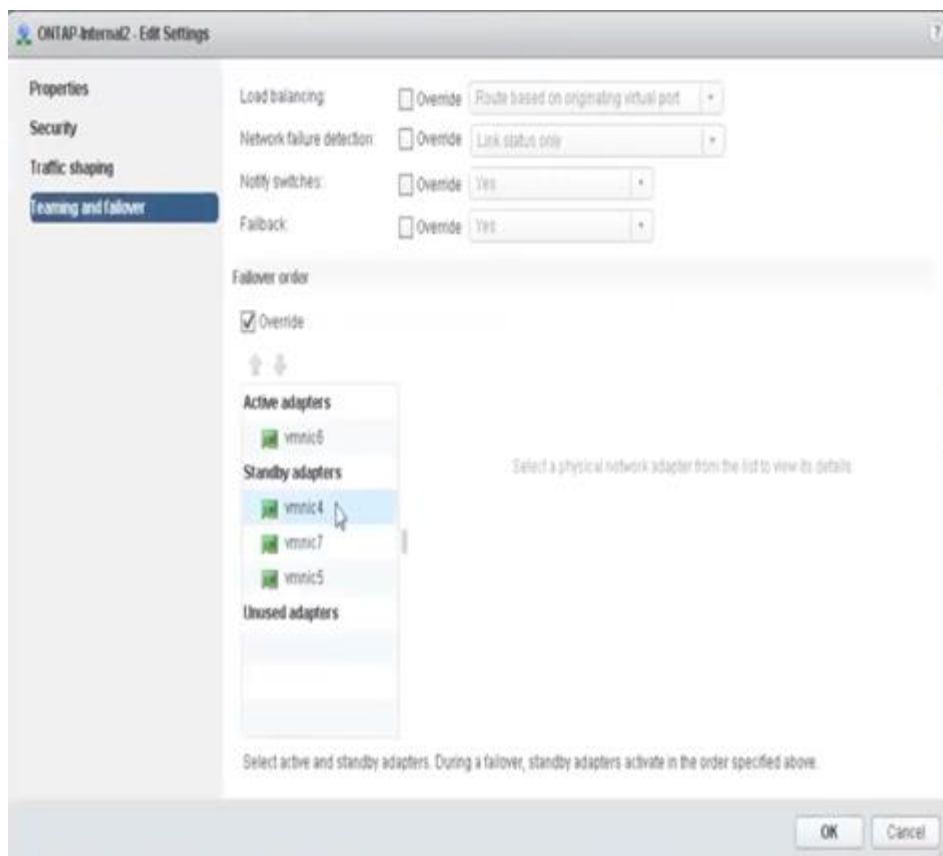
The following figures show the configurations of the internal network port groups (ONTAP-Internal

and ONTAP-Internal2). Note that the active adapters are from different network cards. In this setup, vmnic 4 and vmnic 5 are dual ports on the same physical ASIC, whereas vmnic 6 and vmnic 7 are similarly dual ports on a separate ASIC. The order of the standby adapters provides a hierarchical fail over with the ports from the external network being last. The order of external ports in the standby list is similarly swapped between the two internal port groups.

Part 1: ONTAP Select internal port group configurations



Part 2: ONTAP Select internal port groups



For readability, the assignments are as follows:

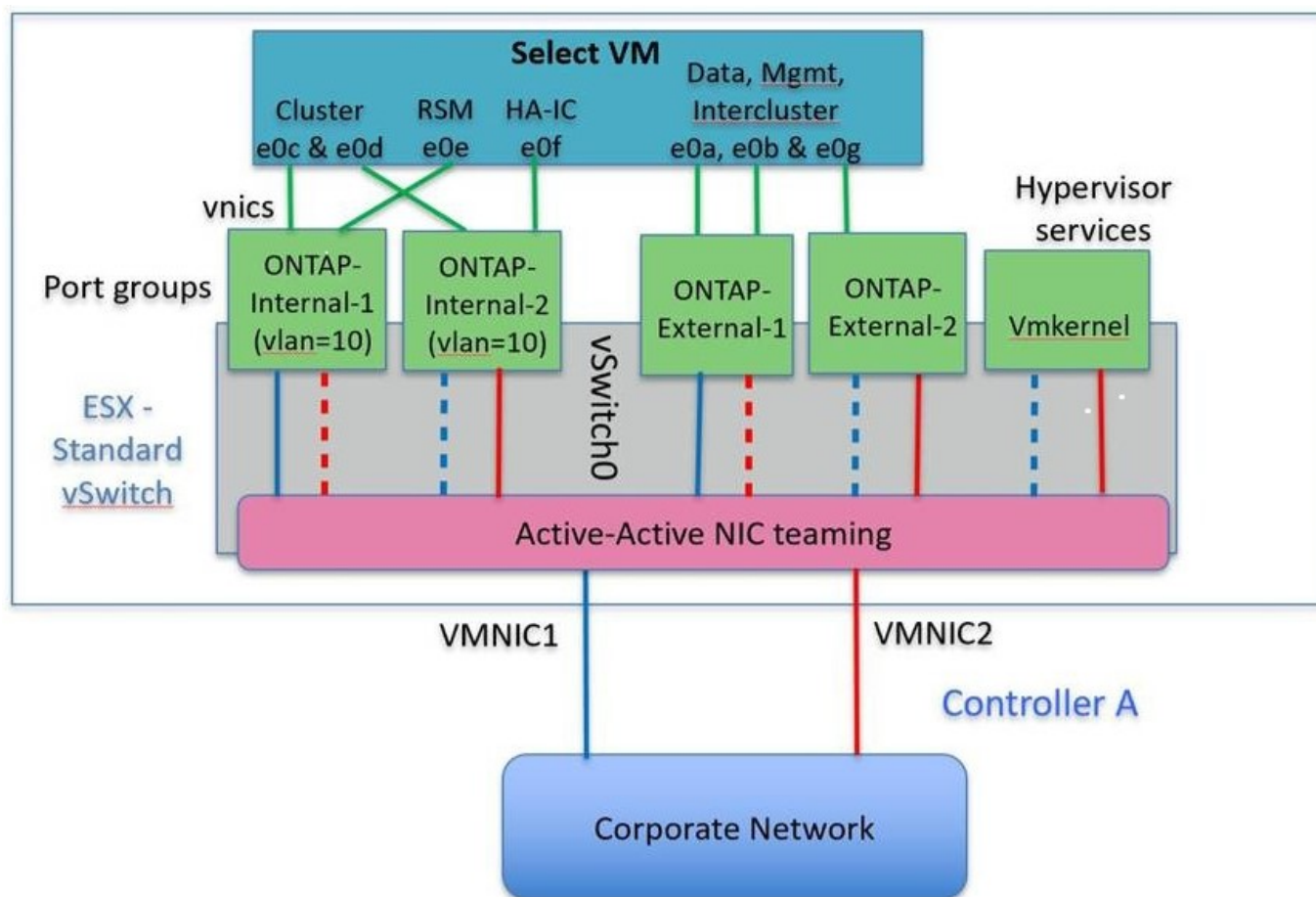
ONTAP-Internal	ONTAP-Internal2
Active adapters: vmnic4	Active adapters: vmnic6
Standby adapters: vmnic6, vmnic5, vmnic7	Standby adapters: vmnic4, vmnic7, vmnic5

Standard or distributed vSwitch and two physical ports per node

When using two high speed (25/40Gb) NICs, the recommended port group configuration is conceptually very similar to the configuration with four 10Gb adapters. Four port groups should be used even when using only two physical adapters. The port group assignments are as follows:

Port Group	External 1 (e0a,e0b)	Internal 1 (e0c,e0e)	Internal 2 (e0d,e0f)	External 2 (e0g)
Active	vmnic0	vmnic0	vmnic1	vmnic1
Standby	vmnic1	vmnic1	vmnic0	vmnic0

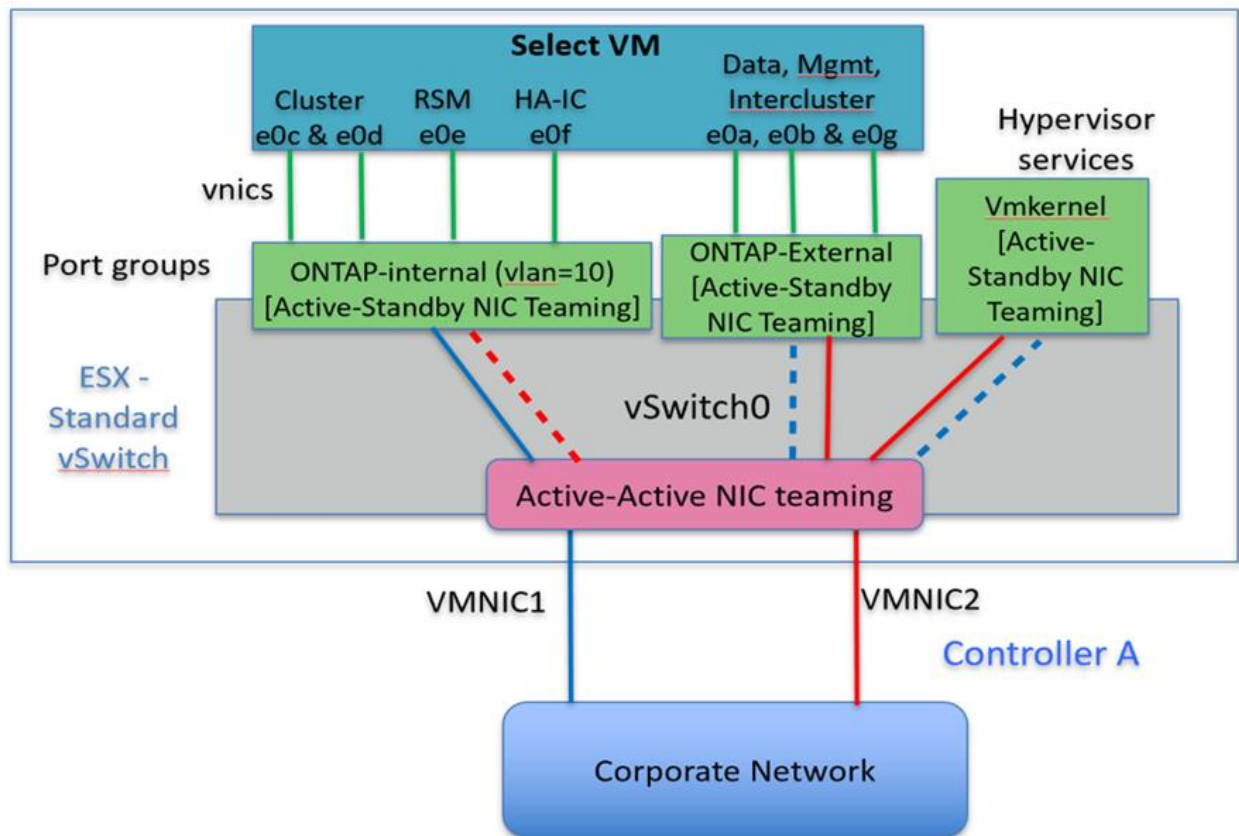
vSwitch with two high speed (25/40Gb) physical ports per node



When using two physical ports (10Gb or less), each port group should have an active adapter and a standby adapter configured opposite to each other. The internal network is only present for multinode ONTAP Select clusters. For single-node clusters, both adapters can be configured as active in the external port group.

The following example shows the configuration of a vSwitch and the two port groups responsible for handling internal and external communication services for a multinode ONTAP Select cluster. The external network can use the internal network VMNIC in the event of a network outage because the internal network VMNICs are part of this port group and configured in standby mode. The opposite is the case for the external network. Alternating the active and standby VMNICs between the two port groups is critical for the proper failover of the ONTAP Select VMs during network outages.

vSwitch with two physical ports (10Gb or less) per node

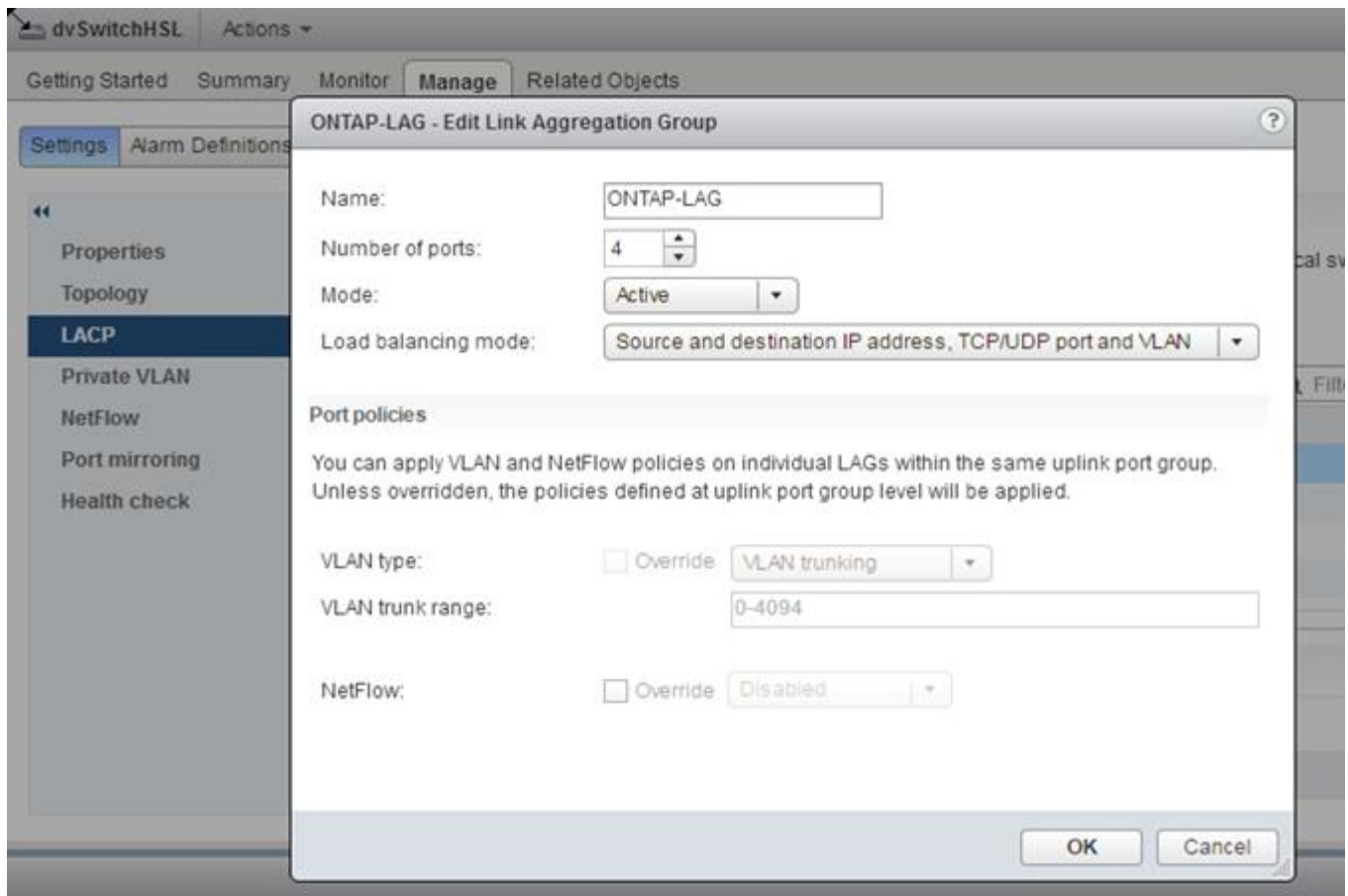


Distributed vSwitch with LACP

When using distributed vSwitches in your configuration, LACP can be used (though it is not a best practice) in order to simplify the network configuration. The only supported LACP configuration requires that all the VMNICs are in a single LAG. The uplink physical switch must support an MTU size between 7,500 to 9,000 on all the ports in the channel. The internal and external ONTAP Select networks should be isolated at the port group level. The internal network should use a nonroutable (isolated) VLAN. The external network can use either VST, EST, or VGT.

The following examples show the distributed vSwitch configuration using LACP.

LAG properties when using LACP




External port group configurations using a distributed vSwitch with LACP enabled

ONTAP-External Settings

General
Policies
Security
Traffic Shaping
VLAN
Teaming and Failover
Resource Allocation
Monitoring
Miscellaneous
Advanced

Policies

Teaming and Failover

Load Balancing:  Route based on IP hash



Network Failover Detection: Link status only

Notify Switches: Yes

Failback: Yes

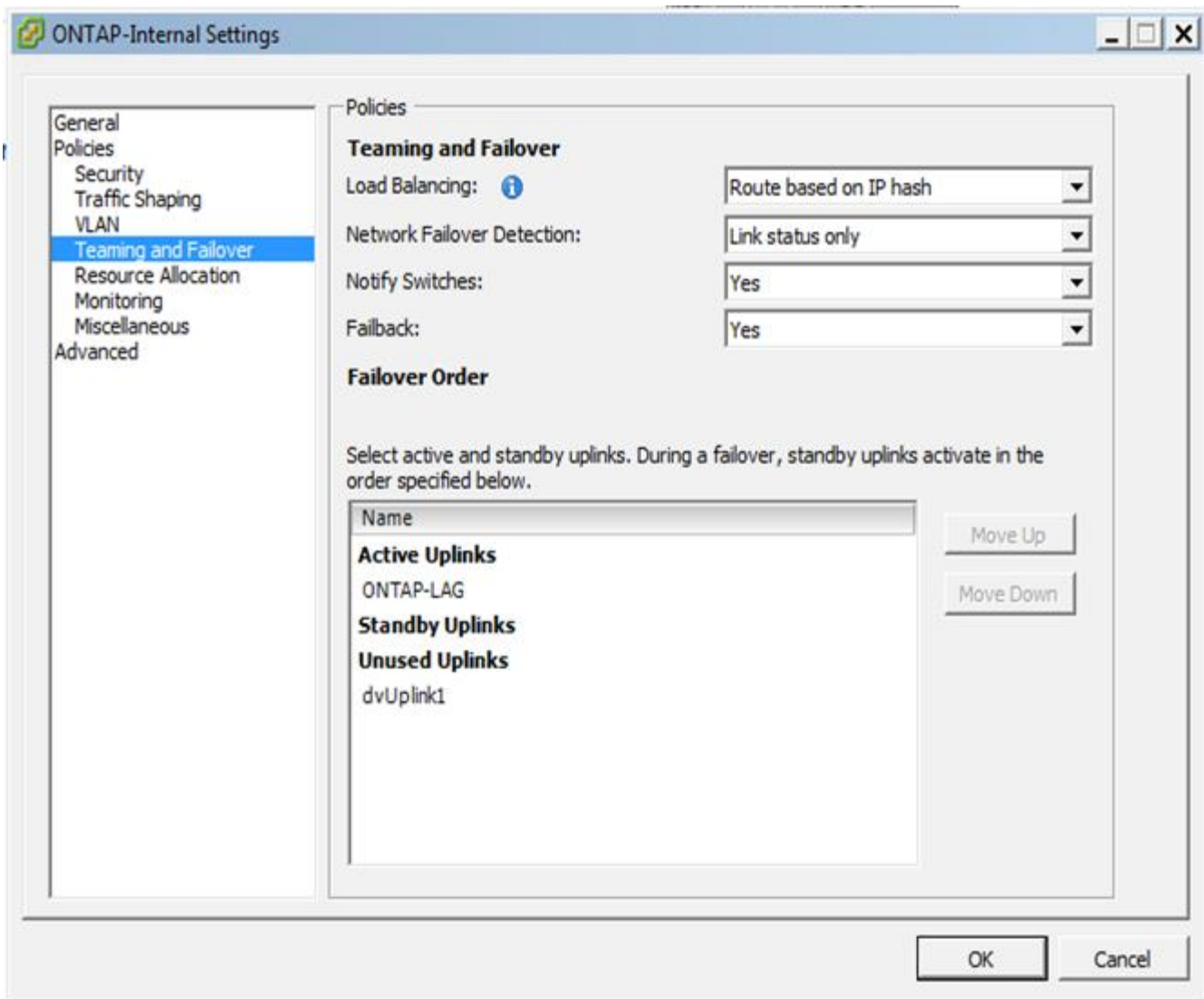
Failover Order

Select active and standby uplinks. During a failover, standby uplinks activate in the order specified below.

Name	
Active Uplinks	
ONTAP-LAG	
Standby Uplinks	
Unused Uplinks	
dvUplink1	

OK Cancel

Internal port group configurations using a distributed vSwitch with LACP enabled



LACP requires that you configure the upstream switch ports as a port channel. Prior to enabling this on the distributed vSwitch, make sure that an LACP-enabled port channel is properly configured.

Physical switch configuration

Upstream physical switch configuration details based on single-switch and multi-switch environments.

Careful consideration should be taken when making connectivity decisions from the virtual switch layer to physical switches. Separation of internal cluster traffic from external data services should extend to the upstream physical networking layer through isolation provided by layer-2 VLANs.

Physical switch ports should be configured as trunkports. ONTAP Select external traffic can be separated across multiple layer-2 networks in one of two ways. One method is by using ONTAP VLAN-tagged virtual ports with a single port group. The other method is by assigning separate port groups in VST mode to management port e0a. You must also assign data ports to e0b and e0c/e0g depending on

the ONTAP Select release and the single-node or multinode configuration. If the external traffic is separated across multiple layer-2 networks, the uplink physical switch ports should have those VLANs in its allowed VLAN list.

ONTAP Select internal network traffic occurs using virtual interfaces defined with link local IP addresses. Because these IP addresses are nonroutable, internal traffic between cluster nodes must flow across a single layer-2 network. Route hops between ONTAP Select cluster nodes are unsupported.

Best Practice

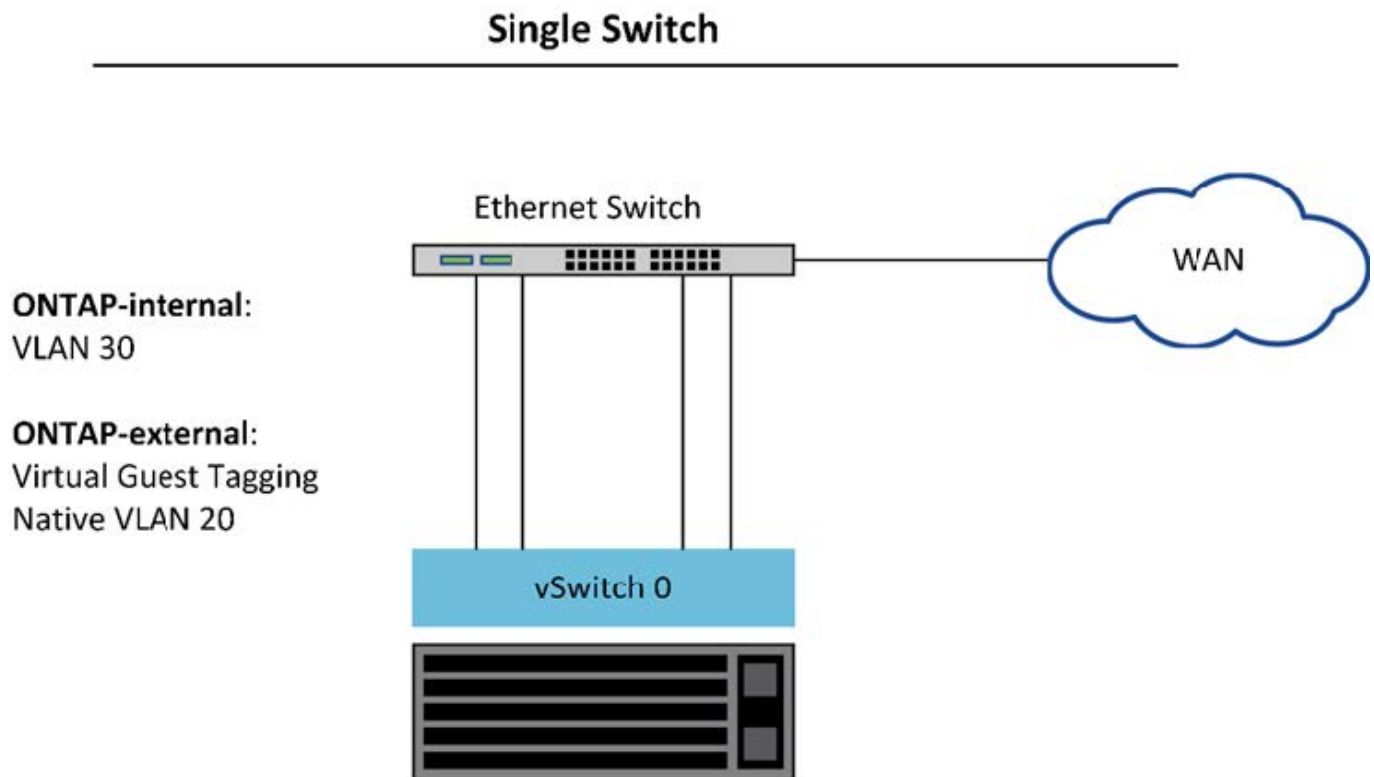
Shared physical switch

The following figure depicts a possible switch configuration used by one node in a multinode ONTAP Select cluster. In this example, the physical NICs used by the vSwitches hosting both the internal and external network port groups are cabled to the same upstream switch. Switch traffic is kept isolated using broadcast domains contained within separate VLANs.



For the ONTAP Select internal network, tagging is done at the port group level. While the following example uses VGT for the external network, both VGT and VST are supported on that port group.

Network configuration using shared physical switch

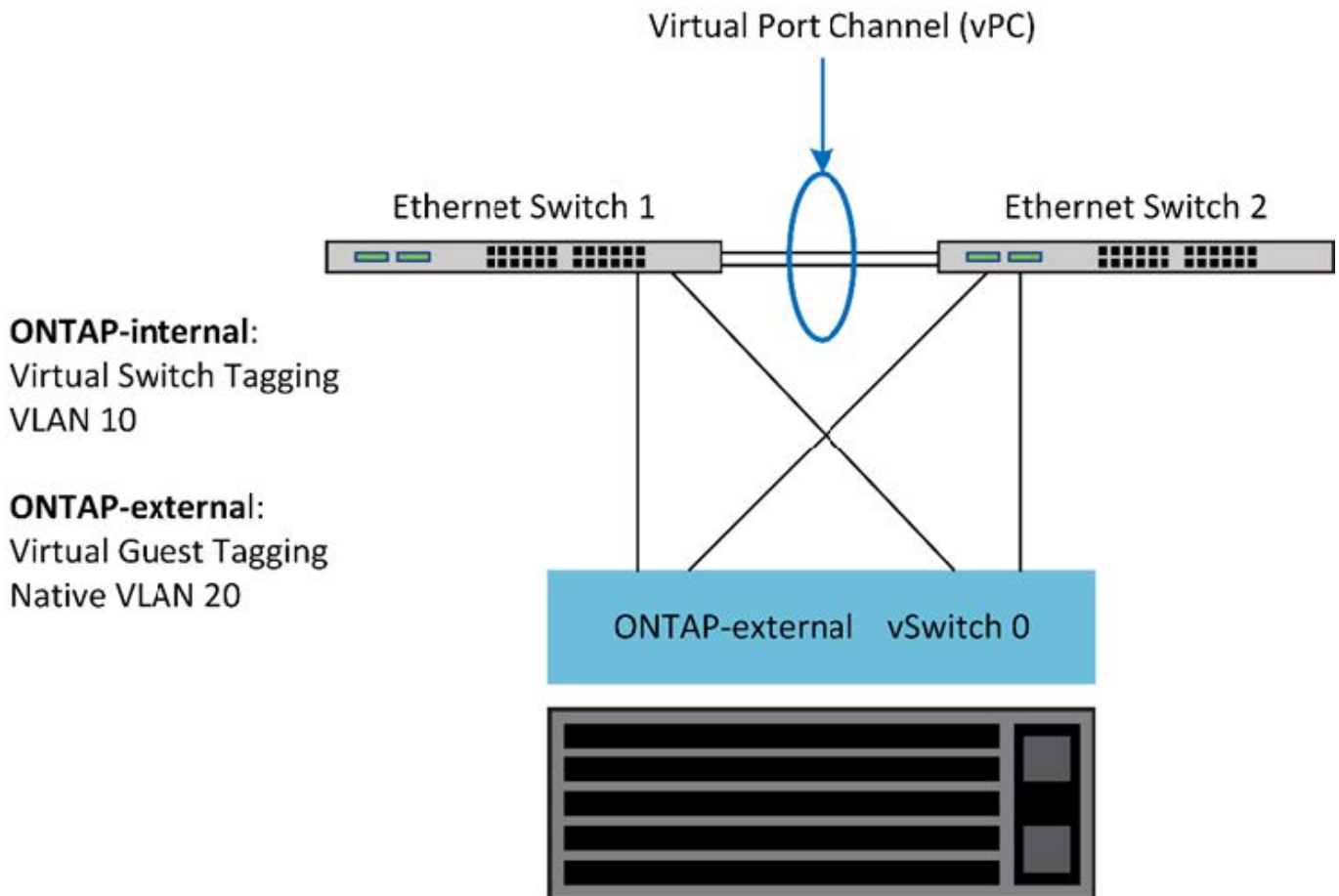


In this configuration, the shared switch becomes a single point of failure. If possible, multiple switches should be used to prevent a physical hardware failure from causing a cluster network outage.

Multiple physical switches

When redundancy is needed, multiple physical network switches should be used. The following figure shows a recommended configuration used by one node in a multinode ONTAP Select cluster. NICs from both the internal and external port groups are cabled into different physical switches, protecting the user from a single hardware-switch failure. A virtual port channel is configured between switches to prevent spanning tree issues.

Network configuration using multiple physical switches



Data and management traffic separation

Isolate data traffic and management traffic into separate layer-2 networks.

ONTAP Select external network traffic is defined as data (CIFS, NFS, and iSCSI), management, and replication (SnapMirror) traffic. Within an ONTAP cluster, each style of traffic uses a separate logical interface that must be hosted on a virtual network port. On the multinode configuration of ONTAP Select, these are designated as ports e0a and e0b/e0g. On the single node configuration, these are designated as e0a and e0b/e0c, while the remaining ports are reserved for internal cluster services.

NetApp recommends isolating data traffic and management traffic into separate layer-2 networks. In the ONTAP Select environment, this is done using VLAN tags. This can be achieved by assigning a

VLAN-tagged port group to network adapter 1 (port e0a) for management traffic. Then you can assign a separate port group(s) to ports e0b and e0c (single-node clusters) and e0b and e0g (multinode clusters) for data traffic.

If the VST solution described earlier in this document is not sufficient, collocating both data and management LIFs on the same virtual port might be required. To do so, use a process known as VGT, in which VLAN tagging is performed by the VM.

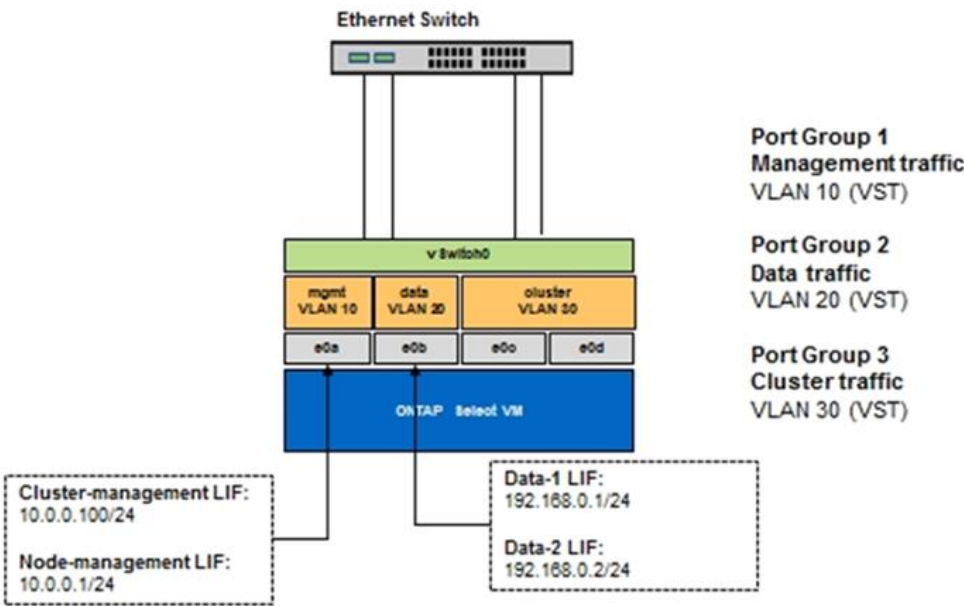


Data and management network separation through VGT is not available when using the ONTAP Deploy utility. This process must be performed after cluster setup is complete.

There is an additional caveat when using VGT and two-node clusters. In two-node cluster configurations, the node management IP address is used to establish connectivity to the mediator before ONTAP is fully available. Therefore, only EST and VST tagging is supported on the port group mapped to the node management LIF (port e0a). Furthermore, if both the management and the data traffic are using the same port group, only EST/VST are supported for the entire two-node cluster.

Both configuration options, VST and VGT, are supported. The following figure shows the first scenario, VST, in which traffic is tagged at the vSwitch layer through the assigned port group. In this configuration, cluster and node management LIFs are assigned to ONTAP port e0a and tagged with VLAN ID 10 through the assigned port group. Data LIFs are assigned to port e0b and either e0c or e0g and given VLAN ID 20 using a second port group. The cluster ports use a third port group and are on VLAN ID 30.

Data and management separation using VST



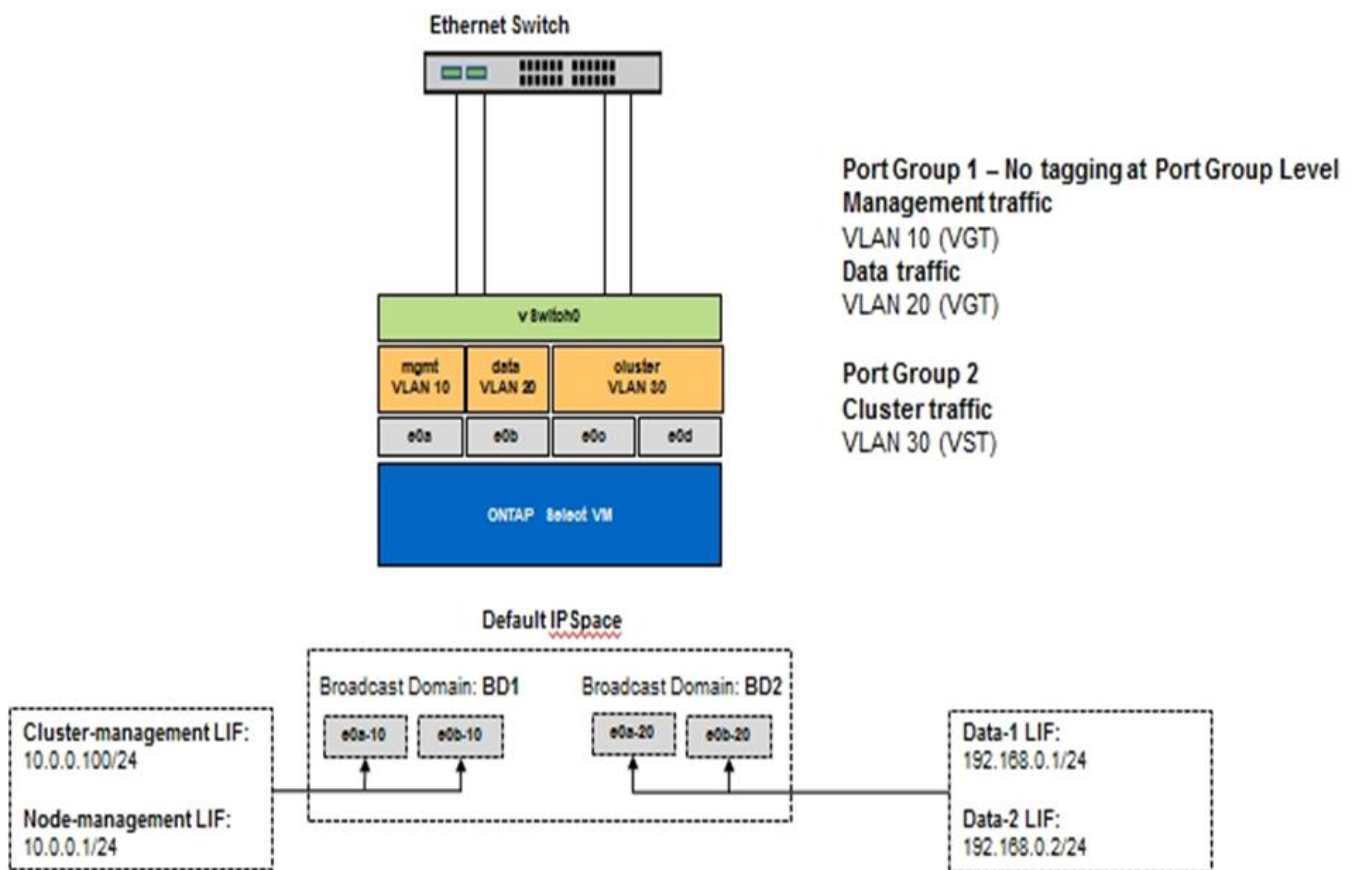
The following figure shows the second scenario, VGT, in which traffic is tagged by the ONTAP VM using VLAN ports that are placed into separate broadcast domains. In this example, virtual ports e0a-10/e0b-

10/(e0c or e0g)-10 and e0a-20/e0b-20 are placed on top of VM ports e0a and e0b. This configuration allows network tagging to be performed directly within ONTAP, rather than at the vSwitch layer. Management and data LIFs are placed on these virtual ports, allowing further layer-2 subdivision within a single VM port. The cluster VLAN (VLAN ID 30) is still tagged at the port group.

Notes:

- This style of configuration is especially desirable when using multiple IPspaces. Group VLAN ports into separate custom IPspaces if further logical isolation and multitenancy are desired.
- To support VGT, the ESXi/ESX host network adapters must be connected to trunk ports on the physical switch. The port groups connected to the virtual switch must have their VLAN ID set to 4095 to enable trunking on the port group.

Data and management separation using VGT



High availability architecture

High availability configurations

Discover high availability options to select the best HA configuration for your environment.

Although customers are starting to move application workloads from enterprise-class storage appliances to software-based solutions running on commodity hardware, the expectations and needs around resiliency and fault tolerance have not changed. An HA solution providing a zero recovery point objective (RPO) protects the customer from data loss due to a failure from any component in the infrastructure stack.

A large portion of the SDS market is built on the notion of shared-nothing storage, with software replication providing data resiliency by storing multiple copies of user data across different storage silos. ONTAP Select builds on this premise by using the synchronous replication features (RAID SyncMirror) provided by ONTAP to store an extra copy of user data within the cluster. This occurs within the context of an HA pair. Every HA pair stores two copies of user data: one on storage provided by the local node, and one on storage provided by the HA partner. Within an ONTAP Select cluster, HA and synchronous replication are tied together, and the functionality of the two cannot be decoupled or used independently. As a result, the synchronous replication functionality is only available in the multinode offering.

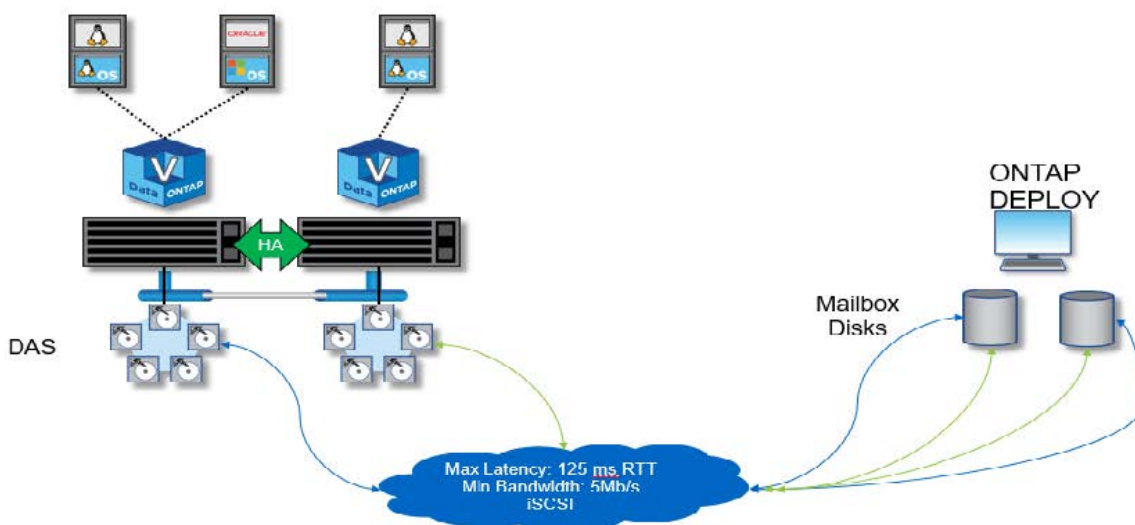


In an ONTAP Select cluster, synchronous replication functionality is a function of HA implementation, not a replacement for the asynchronous SnapMirror or SnapVault replication engines. Synchronous replication cannot be used independently from HA.

There are two ONTAP Select HA deployment models: the multinode clusters (four, six, or eight nodes) and the two-node clusters. The salient feature of a two-node ONTAP Select cluster is the use of an external mediator service to resolve split-brain scenarios. The ONTAP Deploy VM serves as the default mediator for all the two-node HA pairs that it configures.

The two architectures are represented in the following figures.

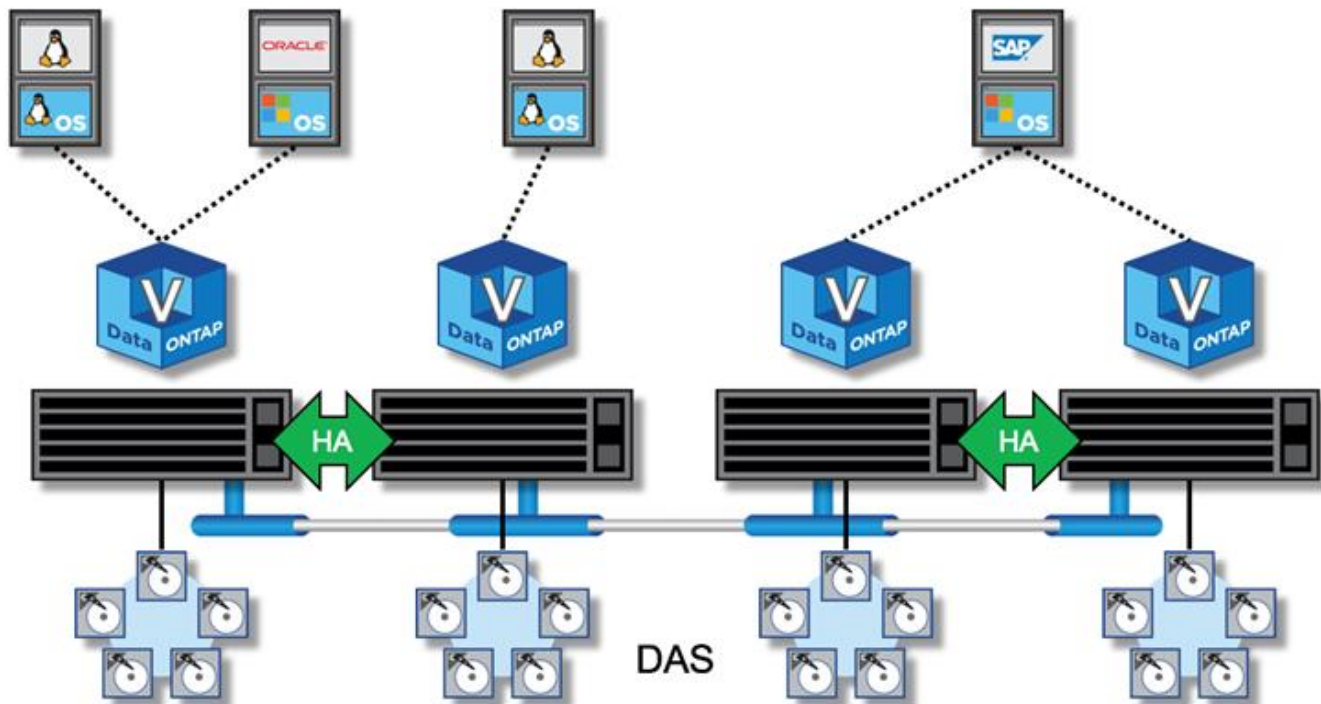
Two-node ONTAP Select cluster with remote mediator and using local-attached storage





The two-node ONTAP Select cluster is composed of one HA pair and a mediator. Within the HA pair, data aggregates on each cluster node are synchronously mirrored, and, in the event of a failover, there is no loss of data.

Four-node ONTAP Select cluster using local-attached storage



- The four-node ONTAP Select cluster is composed of two HA pairs. Six-node and eight-node clusters are composed of three and four HA pairs, respectively. Within each HA pair, data aggregates on each cluster node are synchronously mirrored, and, in the event of a failover, there is no loss of data.
- Only one ONTAP Select instance can be present on a physical server when using DAS storage. ONTAP Select requires unshared access to the local RAID controller of the system and is designed to manage the locally attached disks, which would be impossible without physical connectivity to the storage.

Two-node HA versus multi-node HA

Unlike FAS arrays, ONTAP Select nodes in an HA pair communicate exclusively over the IP network. That means that the IP network is a single point of failure (SPOF), and protecting against network partitions and split-brain scenarios becomes an important aspect of the design. The multi-node cluster can sustain single-node failures because the cluster quorum can be established by the three or more surviving nodes. The two-node cluster relies on the mediator service hosted by the ONTAP Deploy VM to achieve the same result.

The heartbeat network traffic between the ONTAP Select nodes and the ONTAP Deploy mediator service is minimal and resilient so that the ONTAP Deploy VM can be hosted in a different data center

than the ONTAP Select two-node cluster.



The ONTAP Deploy VM becomes an integral part of a two-node cluster when serving as the mediator for that cluster. If the mediator service is not available, the two-node cluster continues serving data, but the storage failover capabilities of the ONTAP Select cluster are disabled. Therefore, the ONTAP Deploy mediator service must maintain constant communication with each ONTAP Select node in the HA pair. A minimum bandwidth of 5Mbps and a maximum round-trip time (RTT) latency of 125ms are required to allow proper functioning of the cluster quorum.

If the ONTAP Deploy VM acting as a mediator is temporarily or potentially permanently unavailable, a secondary ONTAP Deploy VM can be used to restore the two-node cluster quorum. This results in a configuration in which the new ONTAP Deploy VM is unable to manage the ONTAP Select nodes, but it successfully participates in the cluster quorum algorithm. The communication between the ONTAP Select nodes and the ONTAP Deploy VM is done by using the iSCSI protocol over IPv4. The ONTAP Select node management IP address is the initiator, and the ONTAP Deploy VM IP address is the target. Therefore, it is not possible to support IPv6 addresses for the node management IP addresses when creating a two-node cluster. The ONTAP Deploy hosted mailbox disks are automatically created and masked to the proper ONTAP Select node management IP addresses at the time of two-node cluster creation. The entire configuration is automatically performed during setup, and no further administrative action is required. The ONTAP Deploy instance creating the cluster is the default mediator for that cluster.

An administrative action is required if the original mediator location must be changed. It is possible to recover a cluster quorum even if the original ONTAP Deploy VM is lost. However, NetApp recommends that you back up the ONTAP Deploy database after every two-node cluster is instantiated.

Two-node HA versus two-node stretched HA (MetroCluster SDS)

It is possible to stretch a two-node, active/active HA cluster across larger distances and potentially place each node in a different data center. The only distinction between a two-node cluster and a two-node stretched cluster (also referred to as MetroCluster SDS) is the network connectivity distance between nodes.

The two-node cluster is defined as a cluster for which both nodes are located in the same data center within a distance of 300m. In general, both nodes have uplinks to the same network switch or set of interswitch link (ISL) network switches.

Two-node MetroCluster SDS is defined as a cluster for which nodes are physically separated (different rooms, different buildings, and different data centers) by more than 300m. In addition, each node's uplink connections are connected to separate network switches. The MetroCluster SDS does not require dedicated hardware. However, the environment should adhere to requirements for latency (a maximum of 5ms for RTT and 5ms for jitter, for a total of 10ms) and physical distance (a maximum of 10km).

MetroCluster SDS is a premium feature and requires a Premium license. The Premium license supports

the creation of both small and medium VMs, as well as HDD and SSD media.



MetroCluster SDS is supported with both local attached storage (DAS) and shared storage (vNAS). Note that vNAS configurations usually have a higher innate latency because of the network between the ONTAP Select VM and shared storage. MetroCluster SDS configurations must provide a maximum of 10ms of latency between the nodes, including the shared storage latency. In other words, only measuring the latency between the Select VMs is not adequate because shared storage latency is not negligible for these configurations.

HA RSM and mirrored aggregates

Prevent data loss using RAID SyncMirror (RSM), mirrored aggregates, and the write path.

Synchronous replication

The ONTAP HA model is built on the concept of HA partners. ONTAP Select extends this architecture into the nonshared commodity server world by using the RAID SyncMirror (RSM) functionality that is present in ONTAP to replicate data blocks between cluster nodes, providing two copies of user data spread across an HA pair.

A two-node cluster with a mediator can span two data centers. For more information, see the section [Two-node stretched HA \(MetroCluster SDS\) best practices](#).

Mirrored aggregates

An ONTAP Select cluster is composed of two to eight nodes. Each HA pair contains two copies of user data, synchronously mirrored across nodes over an IP network. This mirroring is transparent to the user, and it is a property of the data aggregate, automatically configured during the data aggregate creation process.

All aggregates in an ONTAP Select cluster must be mirrored for data availability in the event of a node failover and to avoid an SPOF in case of hardware failure. Aggregates in an ONTAP Select cluster are built from virtual disks provided from each node in the HA pair and use the following disks:

- A local set of disks (contributed by the current ONTAP Select node)
- A mirrored set of disks (contributed by the HA partner of the current node)



The local and mirror disks used to build a mirrored aggregate must be the same size. These aggregates are referred to as plex 0 and plex 1 (to indicate the local and remote mirror pairs, respectively). The actual plex numbers can be different in your installation.

This approach is fundamentally different from the way standard ONTAP clusters work. This applies to

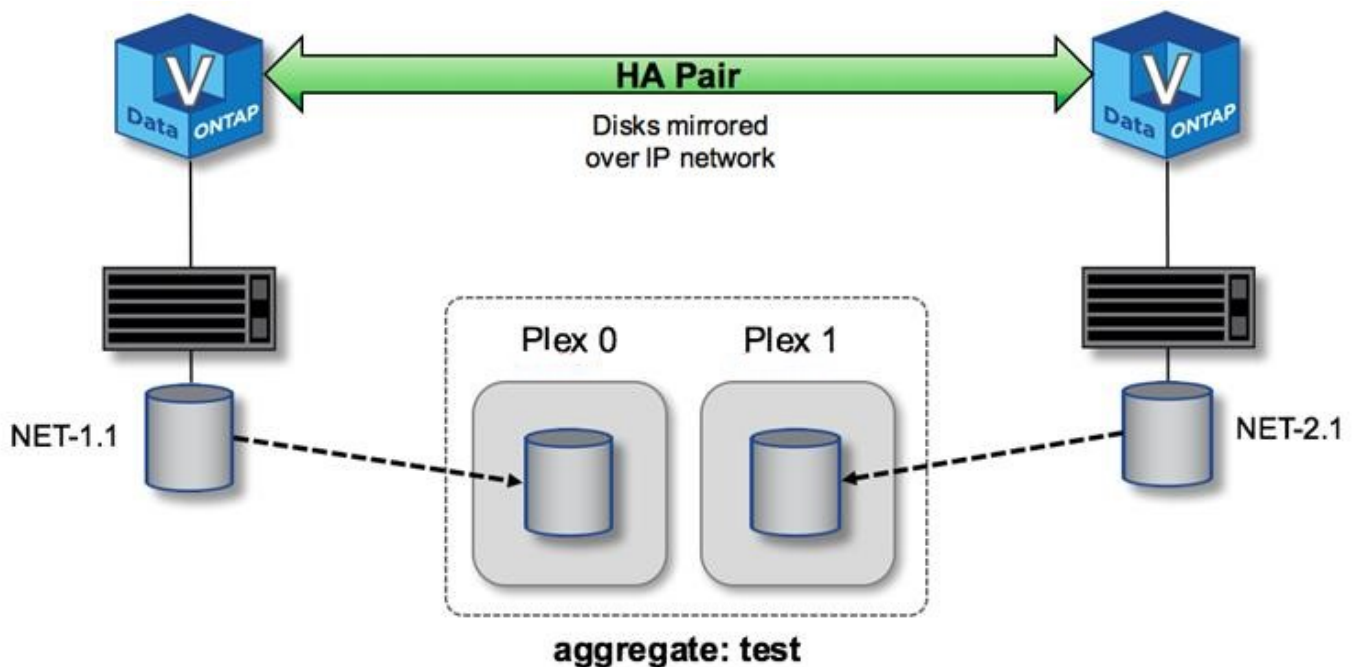
all root and data disks within the ONTAP Select cluster. The aggregate contains both local and mirror copies of data. Therefore, an aggregate that contains N virtual disks offers N/2 disks' worth of unique storage, because the second copy of data resides on its own unique disks.

The following figure shows an HA pair within a four-node ONTAP Select cluster. Within this cluster is a single aggregate (test) that uses storage from both HA partners. This data aggregate is composed of two sets of virtual disks: a local set, contributed by the ONTAP Select owning cluster node (Plex 0), and a remote set, contributed by the failover partner (Plex 1).

Plex 0 is the bucket that holds all local disks. Plex 1 is the bucket that holds mirror disks, or disks responsible for storing a second replicated copy of user data. The node that owns the aggregate contributes disks to Plex 0, and the HA partner of that node contributes disks to Plex 1.

In the following figure, there is a mirrored aggregate with two disks. The contents of this aggregate are mirrored across our two cluster nodes, with local disk NET-1.1 placed into the Plex 0 bucket and remote disk NET-2.1 placed into the Plex 1 bucket. In this example, aggregate test is owned by the cluster node to the left and uses local disk NET-1.1 and HA partner mirror disk NET-2.1.

ONTAP Select mirrored aggregate



When an ONTAP Select cluster is deployed, all virtual disks present on the system are automatically assigned to the correct plex, requiring no additional step from the user regarding disk assignment. This prevents the accidental assignment of disks to an incorrect plex and provides optimal mirror disk configuration.

Write Path

Synchronous mirroring of data blocks between cluster nodes and the requirement for no data loss with a system failure have a significant impact on the path an incoming write takes as it propagates

through an ONTAP Select cluster. This process consists of two stages:

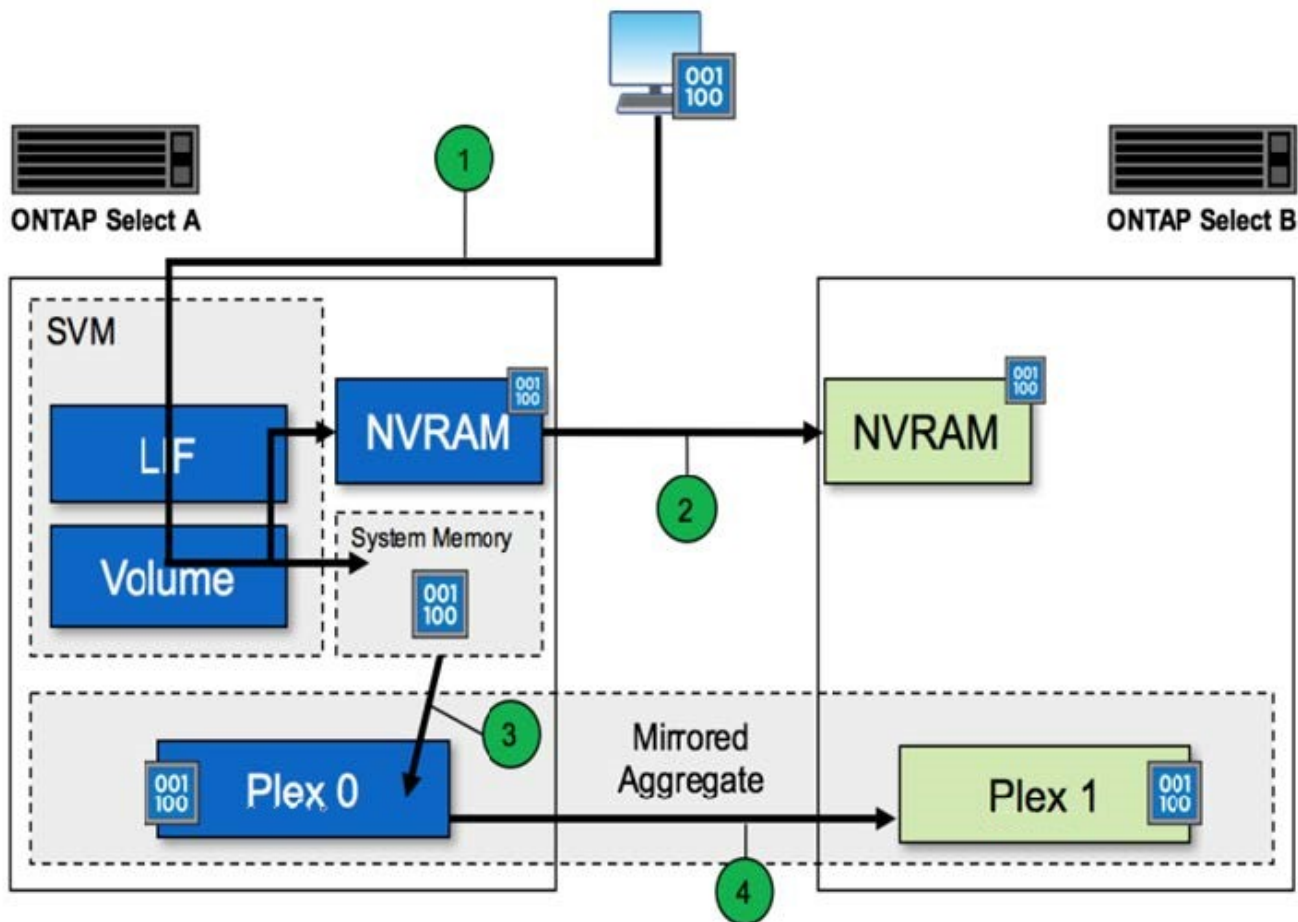
- Acknowledgment
- Destaging

Writes to a target volume occur over a data LIF and are committed to the virtualized NVRAM partition, present on a system disk of the ONTAP Select node, before being acknowledged back to the client. On an HA configuration, an additional step occurs, because these NVRAM writes are immediately mirrored to the HA partner of the target volume's owner before being acknowledged. This process makes sure of the file system consistency on the HA partner node, if there is a hardware failure on the original node.

After the write has been committed to NVRAM, ONTAP periodically moves the contents of this partition to the appropriate virtual disk, a process known as destaging. This process only happens once, on the cluster node owning the target volume, and does not happen on the HA partner.

The following figure shows the write path of an incoming write request to an ONTAP Select node.

ONTAP Select write path workflow



Incoming write acknowledgment includes the following steps:

- Writes enter the system through a logical interface owned by ONTAP Select node A.

- Writes are committed to the NVRAM of node A and mirrored to the HA partner, node B.
- After the I/O request is present on both HA nodes, the request is then acknowledged back to the client.

ONTAP Select destaging from NVRAM to the data aggregate (ONTAP CP) includes the following steps:

- Writes are destaged from virtual NVRAM to virtual data aggregate.
- Mirror engine synchronously replicates blocks to both plexes.

HA additional details

HA disk heartbeating, HA mailbox, HA heartbeating, HA Failover, and Giveback work to enhance data protection.

Disk heartbeating

Although the ONTAP Select HA architecture leverages many of the code paths used by the traditional FAS arrays, some exceptions exist. One of these exceptions is in the implementation of disk-based heartbeating, a nonnetwork-based method of communication used by cluster nodes to prevent network isolation from causing split-brain behavior. A split-brain scenario is the result of cluster partitioning, typically caused by network failures, whereby each side believes the other is down and attempts to take over cluster resources.

Enterprise-class HA implementations must gracefully handle this type of scenario. ONTAP does this through a customized, disk-based method of heartbeating. This is the job of the HA mailbox, a location on physical storage that is used by cluster nodes to pass heartbeat messages. This helps the cluster determine connectivity and therefore define quorum in the event of a failover.

On FAS arrays, which use a shared storage HA architecture, ONTAP resolves split-brain issues in the following ways:

- SCSI persistent reservations
- Persistent HA metadata
- HA state sent over HA interconnect

However, within the shared-nothing architecture of an ONTAP Select cluster, a node is only able to see its own local storage and not that of the HA partner. Therefore, when network partitioning isolates each side of an HA pair, the preceding methods of determining cluster quorum and failover behavior are unavailable.

Although the existing method of split-brain detection and avoidance cannot be used, a method of mediation is still required, one that fits within the constraints of a shared-nothing environment. ONTAP Select extends the existing mailbox infrastructure further, allowing it to act as a method of mediation in the event of network partitioning. Because shared storage is unavailable, mediation is accomplished through access to the mailbox disks over NAS. These disks are spread throughout the

cluster, including the mediator in a two-node cluster, using the iSCSI protocol. Therefore, intelligent failover decisions can be made by a cluster node based on access to these disks. If a node can access the mailbox disks of other nodes outside of its HA partner, it is likely up and healthy.



The mailbox architecture and disk-based heartbeating method of resolving cluster quorum and split-brain issues are the reasons the multinode variant of ONTAP Select requires either four separate nodes or a mediator for a two-node cluster.

HA mailbox posting

The HA mailbox architecture uses a message post model. At repeated intervals, cluster nodes post messages to all other mailbox disks across the cluster, including the mediator, stating that the node is up and running. Within a healthy cluster at any point in time, a single mailbox disk on a cluster node has messages posted from all other cluster nodes.

Attached to each Select cluster node is a virtual disk that is used specifically for shared mailbox access. This disk is referred to as the mediator mailbox disk, because its main function is to act as a method of cluster mediation in the event of node failures or network partitioning. This mailbox disk contains partitions for each cluster node and is mounted over an iSCSI network by other Select cluster nodes. Periodically, these nodes post health statuses to the appropriate partition of the mailbox disk. Using network-accessible mailbox disks spread throughout the cluster allows you to infer node health through a reachability matrix. For example, cluster nodes A and B can post to the mailbox of cluster node D, but not to the mailbox of node C. In addition, cluster node D cannot post to the mailbox of node C, so it is likely that node C is either down or network isolated and should be taken over.

HA heartbeating

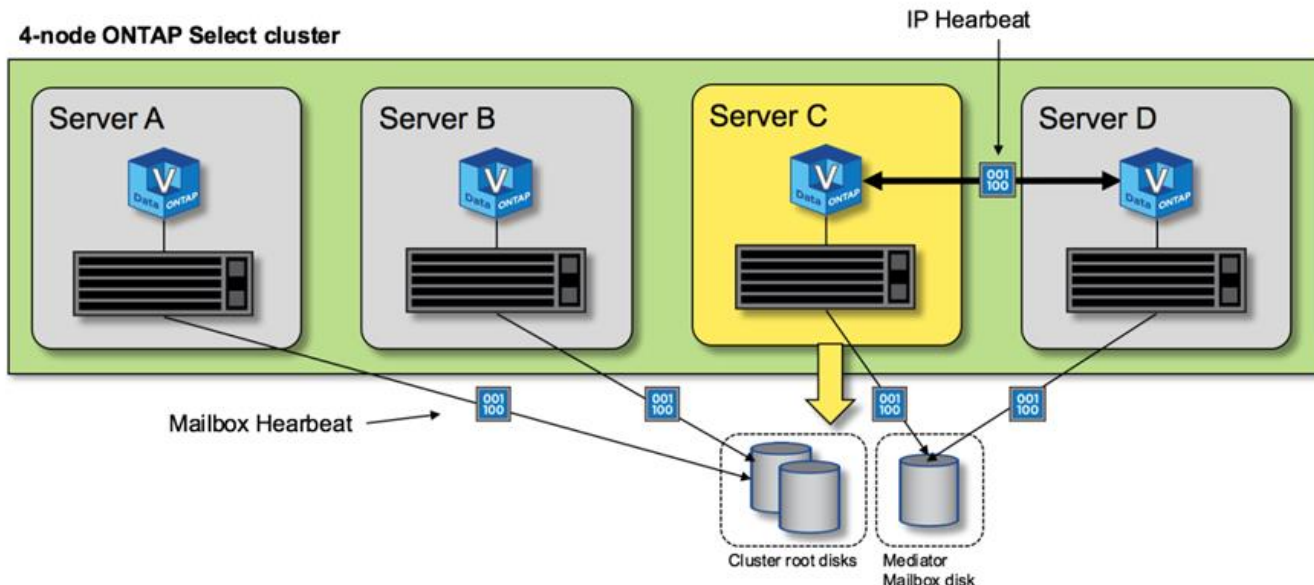
Like with NetApp FAS platforms, ONTAP Select periodically sends HA heartbeat messages over the HA interconnect. Within the ONTAP Select cluster, this is performed over a TCP/IP network connection that exists between HA partners. Additionally, disk-based heartbeat messages are passed to all HA mailbox disks, including mediator mailbox disks. These messages are passed every few seconds and read back periodically. The frequency with which these are sent and received allows the ONTAP Select cluster to detect HA failure events within approximately 15 seconds, the same window available on FAS platforms. When heartbeat messages are no longer being read, a failover event is triggered.

The following figure shows the process of sending and receiving heartbeat messages over the HA interconnect and mediator disks from the perspective of a single ONTAP Select cluster node, node C.



Network heartbeats are sent over the HA interconnect to the HA partner, node D, while disk heartbeats use mailbox disks across all cluster nodes, A, B, C, and D.

HA heartbeating in a four-node cluster: steady state



HA failover and giveback

During a failover operation, the surviving node assumes the data serving responsibilities for its peer node using the local copy of its HA partner's data. Client I/O can continue uninterrupted, but changes to this data must be replicated back before giveback can occur. Note that ONTAP Select does not support a forced giveback because this causes changes stored on the surviving node to be lost.

The sync back operation is automatically triggered when the rebooted node rejoins the cluster. The time required for the sync back depends on several factors. These factors include the number of changes that must be replicated, the network latency between the nodes, and the speed of the disk subsystems on each node. It is possible that the time required for sync back will exceed the auto giveback window of 10 minutes. In this case, a manual giveback after the sync back is required. The progress of the sync back can be monitored using the following command:

```
storage aggregate status -r -aggregate <aggregate name>
```

Performance

Performance general considerations

Performance varies based on hardware configuration.

The performance numbers described in this section are intended as a rough estimate of the performance of an ONTAP Select cluster and are not a performance guarantee.

The performance of an ONTAP Select cluster can vary considerably due to the characteristics of the underlying hardware and configuration. As a matter of fact, the specific hardware configuration is the biggest factor in the performance of a particular ONTAP Select instance. Here are some of the factors

that affect the performance of a specific ONTAP Select instance:

- **Core frequency.** In general, a higher frequency is preferable.
- **Single socket versus multsocket.** ONTAP Select does not use multsocket features, but the hypervisor overhead for supporting multsocket configurations accounts for some amount of deviation in total performance.
- **RAID card configuration and associated hypervisor driver.** The default driver provided by the hypervisor might need to be replaced by the hardware vendor driver.
- **Drive type and number of drives in the RAID group(s).**
- **Hypervisor version and patch level.**

This section includes performance comparisons only when the testing was performed on the exact same test bed to highlight the impact of a specific feature. In general, we document the hardware environment and run the highest performance configuration possible on that platform.

ONTAP Select 9.4 performance: Premium HA direct-attached SSD storage

ONTAP Select 9.4 performance with premium HA direct-attached SSD storage.

Reference platform

ONTAP Select 9.4 (Premium) hardware (per node):

- Cisco UCS C240 M4S2:
 - Intel Xeon CPU E5-2697 at 2.60GHz
 - 2 x sockets; 14 x CPUs per socket
 - 56 x logical CPUs (HT enabled)
 - 256GB RAM
 - VMware ESXi 6.5
 - Drives per host: 24 X371A NetApp 960GB SSD

Client hardware:

- 4 x NFSv3 IBM 3550m4 clients

Configuration information:

- 1,500 MTU for data path between clients and Select cluster
- No storage efficiency features in use (compression, deduplication, Snapshot copies, SnapMirror, and so on)

The following table lists the throughput measured against read/write workloads on an HA pair of

ONTAP Select Premium nodes. Performance measurements were taken using the SIO load-generating tool.

Performance results for a single node (part of a four-node medium instance) ONTAP Select 9.4 cluster on DAS (SSD)

Description	Sequential Read 64KiB	Sequential Write 64KiB	Random Read 8KiB	Random Write 8KiB	Random WR/ RD (50/50) 8KiB
ONTAP 9.4 Select Medium instance with DAS (SSD)	1045MBps 16,712 IOPS	251MBps 4016 IOPS	492MBps 62,912 IOPS	141MBps 18,048 IOPS	218MBps 27,840 IOPS

64K sequential read

Details:

- SIO direct I/O enabled
- 2 x data NIC
- 1 x data aggregate (2TB)
- 64 volumes; 64 SIO procs/threads
- 32 volumes per node (64 total)
- 1 x SIO procs per volume; 1 x SIO thread per file
- 1 x files per volume; files are 12000MB each

64K sequential write

Details:

- SIO direct I/O enabled
- 2 x data NIC
- 1 x data aggregate (2TB):
- 64 volumes; 128 SIO procs/threads
- 32 volumes per node (64 total)
- 2 x SIO procs per volume; 1 x SIO thread per file
- 2 x files per volume; files are 30720MB each

8K random read

Details:

- SIO direct I/O enabled

- 2 x data NIC
- 1 x data aggregate (2TB):
- 64 volumes; 64 SIO procs/threads
- 32 volumes per node (64 total)
- 1 x SIO procs per volume; 8 x SIO thread per file
- 1 x files per volume; files are 12228MB each

8K random write

Details:

- SIO direct I/O enabled
- 2 x data NIC
- 1 x data aggregate (2TB)
- 64 volumes; 64 SIO procs/threads
- 32 volumes per node (64 total)
- 1 x SIO procs per volume; 8 x SIO thread per file
- 1 x files per volume; files are 8192MB each

8K random 50% write 50% read

Details:

- SIO direct I/O enabled
- 2 x data NIC
- 1 x data aggregate (2TB)
- 64 volumes; 64 SIO procs/threads
- 32 volumes per node (64 total)
- 1 x SIO procs per volume; 20 x SIO thread per file
- 1 x files per volume; files are 12228MB each

ONTAP Select 9.5 performance: Premium HA direct-attached SSD storage

ONTAP Select 9.5 performance with premium HA direct-attached SSD storage.

Reference platform

ONTAP Select 9.5 (Premium) hardware (per node):

- Cisco UCS C240 M4SX:

- Intel Xeon CPU E5-2620 at 2.1GHz
- 2 x sockets; 16 x CPUs per socket
- 128GB RAM
- VMware ESXi 6.5
- Drives per host: 24 900GB SSD

Client hardware:

- 5 x NFSv3 IBM 3550m4 clients

Configuration information:

- 1,500 MTU for data path between clients and Select cluster
- No storage efficiency features in use (compression, deduplication, Snapshot copies, SnapMirror, and so on)

The following table lists the throughput measured against read/write workloads on an HA pair of ONTAP Select Premium nodes using both software RAID and hardware RAID. Performance measurements were taken using the SIO load-generating tool.

Performance results for a single node (part of a four-node medium instance) ONTAP Select 9.5 cluster on DAS (SSD) with software RAID and hardware RAID

Description	Sequential Read 64KiB	Sequential Write 64KiB	Random Read 8KiB	Random Write 8KiB	Random WR/ RD (50/50) 8KiB
ONTAP 9.5 Select Medium instance with DAS (SSD) hardware RAID	1,714MiBps	412MiBps	391MiBps	1251MiBps	309MiBps
ONTAP 9.5 Select Medium instance with DAS (SSD) software RAID	1,674MiBps	360MiBps	451MiBps	223MiBps	293MiBps

64K sequential read

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node

- 1 x data aggregate per node (2TB hardware RAID), (8TB software RAID)
- 64 SIO procs, 1 thread per proc
- 32 volumes per node
- 1 x files per proc; files are 12000MB each

64K sequential write

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node
- 1 x data aggregate per node (2TB hardware RAID), (4TB software RAID)
- 128 SIO procs, 1 thread per proc
- volumes per node 32 (hardware RAID), 16 (software RAID)
- 1 x files per proc; files are 30720MB each

8K random read

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node
- 1 x data aggregate per node (2TB hardware RAID), (4TB software RAID)
- 64 SIO procs, 8 threads per proc
- volumes per node 32
- 1 file per proc; files are 12228MB each

8K random write

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node
- 1 x data aggregate per node (2TB hardware RAID), (4TB software RAID)
- 64 SIO procs, 8 threads per proc
- volumes per node 32

- 1 file per proc; files are 8192MB each

8K random 50% write 50% read

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node
- 1 x data aggregate per node (2TB hardware RAID), (4TB software RAID)
- 64 SIO procs, 20 threads per proc
- volumes per node 32
- 1 file per proc; files are 12228MB each

ONTAP Select 9.6 performance: Premium HA direct-attached SSD storage

Performance information for the reference platform.

Reference platform

ONTAP Select 9.6 (Premium XL) Hardware (per Node)

- FUJITSU PRIMERGY RX2540 M4:
 - Intel® Xeon® Gold 6142b CPU at 2.6 GHz
 - 32 physical cores (16 x 2 sockets), 64 logical
 - 256 GB RAM
 - Drives per host: 24 960GB SSD
 - ESX 6.5U1

Client hardware

- 5 x NFSv3 IBM 3550m4 clients

Configuration information

- SW RAID 1 x 9 + 2 RAID-DP (11 drives)
- 22+1 RAID-5 (RAID-0 in ONTAP) / RAID cache NVRAM
- No storage efficiency features in use (compression, deduplication, Snapshot copies, SnapMirror, and so on)

The following table lists the throughput measured against read/write workloads on an HA pair of ONTAP Select Premium nodes using both software RAID and hardware RAID. Performance

measurements were taken using the SIO load-generating tool.

Performance results for a single node (part of a four-node medium instance) ONTAP Select 9.5 cluster on DAS (SSD) with software RAID and hardware RAID

Description	Sequential Read 64KiB	Sequential Write 64KiB	Random Read 8KiB	Random Write 8KiB	Random WR/ RD (50/50) 8KiB
ONTAP 9.6 Select Large instance with DAS (SSD) software RAID	2171 MiBps	559 MiBps	954 MiBps	394 MiBps	564 MiBps
ONTAP 9.6 Select Medium instance with DAS (SSD) software RAID	2090 MiBps	592 MiBps	677 MiBps	335 MiBps	441 3MiBps
ONTAP 9.6 Select Medium instance with DAS (SSD) hardware RAID	2038 MiBps	520 MiBps	578 MiBps	325 MiBps	399 MiBps

64K sequential read

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node
- 1 x data aggregate per node (2TB hardware RAID), (8TB software RAID)
- 64 SIO procs, 1 thread per proc
- 32 volumes per node
- 1 x files per proc; files are 12000MB each

64K sequential write

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node

- 1 x data aggregate per node (2TB hardware RAID), (4TB software RAID)
- 128 SIO procs, 1 thread per proc
- volumes per node 32 (hardware RAID), 16 (software RAID)
- 1 x files per proc; files are 30720MB each

8K random read

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node
- 1 x data aggregate per node (2TB hardware RAID), (4TB software RAID)
- 64 SIO procs, 8 threads per proc
- volumes per node 32
- 1 x files per proc; files are 12228MB each

8K random write

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node
- 1 x data aggregate per node (2TB hardware RAID), (4TB software RAID)
- 64 SIO procs, 8 threads per proc
- volumes per node 32
- 1 x files per proc; files are 8192MB each

8K random 50% write 50% read

Details:

- SIO direct I/O enabled
- 2 nodes
- 2 x data NIC per node
- 1 x data aggregate per node (2TB hardware RAID), (4TB software RAID)
- 64 SIO proc208 threads per proc
- volumes per node 32

- 1 x files per proc; files are 12228MB each

Copyright Information

Copyright © 2020 NetApp, Inc. All rights reserved. Printed in the U.S. No part of this document covered by copyright may be reproduced in any form or by any means-graphic, electronic, or mechanical, including photocopying, recording, taping, or storage in an electronic retrieval system-without prior written permission of the copyright owner.

Software derived from copyrighted NetApp material is subject to the following license and disclaimer:

THIS SOFTWARE IS PROVIDED BY NETAPP “AS IS” AND WITHOUT ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE, WHICH ARE HEREBY DISCLAIMED. IN NO EVENT SHALL NETAPP BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

NetApp reserves the right to change any products described herein at any time, and without notice. NetApp assumes no responsibility or liability arising from the use of products described herein, except as expressly agreed to in writing by NetApp. The use or purchase of this product does not convey a license under any patent rights, trademark rights, or any other intellectual property rights of NetApp.

The product described in this manual may be protected by one or more U.S. patents, foreign patents, or pending applications.

RESTRICTED RIGHTS LEGEND: Use, duplication, or disclosure by the government is subject to restrictions as set forth in subparagraph (c)(1)(ii) of the Rights in Technical Data and Computer Software clause at DFARS 252.277-7103 (October 1988) and FAR 52-227-19 (June 1987).

Trademark Information

NETAPP, the NETAPP logo, and the marks listed at <http://www.netapp.com/TM> are trademarks of NetApp, Inc. Other company and product names may be trademarks of their respective owners.