# NetApp

# SnapMirror active sync

## Enterprise applications

NetApp
January 12, 2026

# Table of Contents

# SnapMirror active sync

## Overview

SnapMirror active sync allows you to build ultra high availability Oracle database environments where LUNs are available from two different storage clusters.

With SnapMirror active sync, there is no "primary" and "secondary" copy of the data. Each cluster can serve read IO from its local copy of the data, and each cluster will replicate a write to its partner. The result is symmetric IO behavior.

Among other options, this allows you to run Oracle RAC as an extended cluster with operational instances on both sites. Alternatively, you could build RPO=0 active-passive database clusters where single instance databases can be moved across sites during a site outage, and this process can be automated through products like Pacemaker or VMware HA. The foundation for all of these options is synchronous replication managed by SnapMirror active sync.

### Synchronous replication

In normal operation, SnapMirror active sync provides RPO=0 synchronous replica at all times, with one exception. If data cannot be replicated, ONTAP will release the requirement to replicate data and resume serving IO on one site while the LUNs on the other site are taken offline.

### Storage hardware

Unlike other storage disaster recovery solutions, SnapMirror active sync offers asymmetric platform flexibility. The hardware at each site does not need to be identical. This capability allows you to right-size the hardware used to support SnapMirror active sync. The remote storage system can be identical to the primary site if it needs to support a full production workload, but if a disaster results in reduced I/O, than a smaller system at the remote site might be more cost-effective.
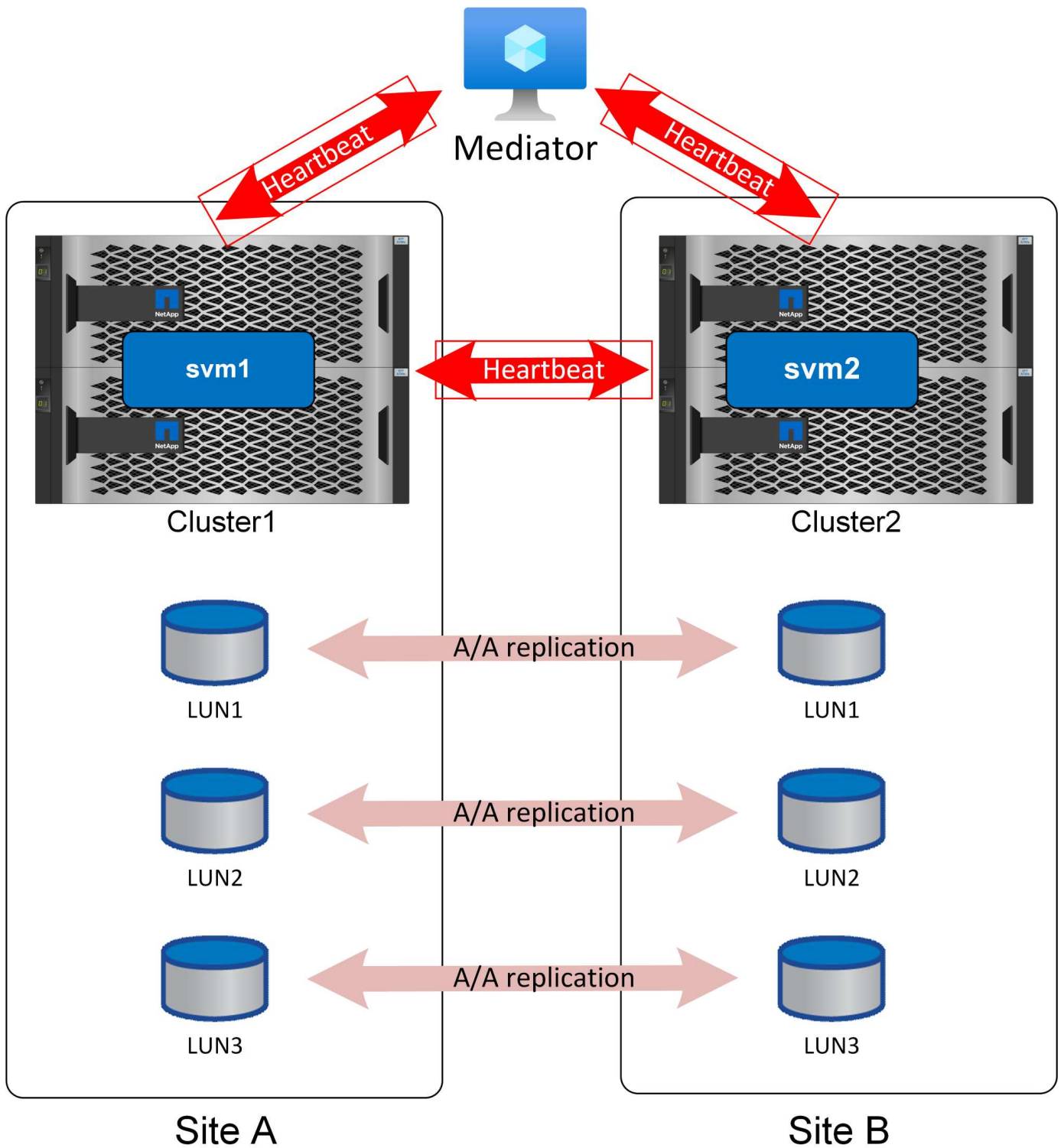
### ONTAP mediator

The ONTAP Mediator is a software application that is downloaded from NetApp support, and is typically deployed on a small virtual machine. The ONTAP Mediator is not a tiebreaker when used with SnapMirror active sync. It is an alternate communication channel for the two clusters that participate in SnapMirror active sync replication. Automated operations are driven by ONTAP based on the responses received from the partner via direct connections and via the mediator.

## ONTAP Mediator

The mediator is required for safely automating failover. Ideally, it would be placed on an independent 3rd site, but it can still function for most needs if colocated with one of the clusters participating in replication.

The mediator is not really a tiebreaker, although that is effectively the function it provides. The mediator helps in determining the state of the cluster nodes and assists in the automatic switchover process in the event of a site failure. Mediator does not transfer data under any circumstances.

The #1 challenge with automated failover is the split-brain problem, and that problem arises if your two sites lose connectivity with each other. What should happen? You do not want to have two different sites designate themselves as the surviving copies of the data, but how can a single site tell the difference between actual loss of the opposite site and an inability to communicate with the opposite site?

This is where the mediator enters the picture. If placed on a 3rd site, and each site has a separate network connection to that site, then you have an additional path for each site to validate the health of the other. Look at the picture above again and consider the following scenarios.

- What happens if the mediator fails or is unreachable from one or both sites?
  - The two clusters can still communicate with each other over the same link used for replication services.
  - Data is still served with RPO=0 protection
- What happens if Site A fails?
  - Site B will see both of the communication channels go down.
  - Site B will take over data services, but without RPO=0 mirroring
- What happens if Site B fails?
  - Site A will see both of the communication channels go down.
  - Site A will take over data services, but without RPO=0 mirroring

There is one other scenario to consider: Loss of the data replication link. If the replication link between sites is lost, RPO=0 mirroring will obviously be impossible. What should happen then?

This is controlled by the preferred site status. In an SM-as relationship, one of the sites is secondary to the other. This has no effect on normal operations, and all data access is symmetric, but if replication is interrupted then the tie will have to be broken to resume operations. The result is the preferred site will continue operations without mirroring and the secondary site will halt IO processing until replication communication is restored.
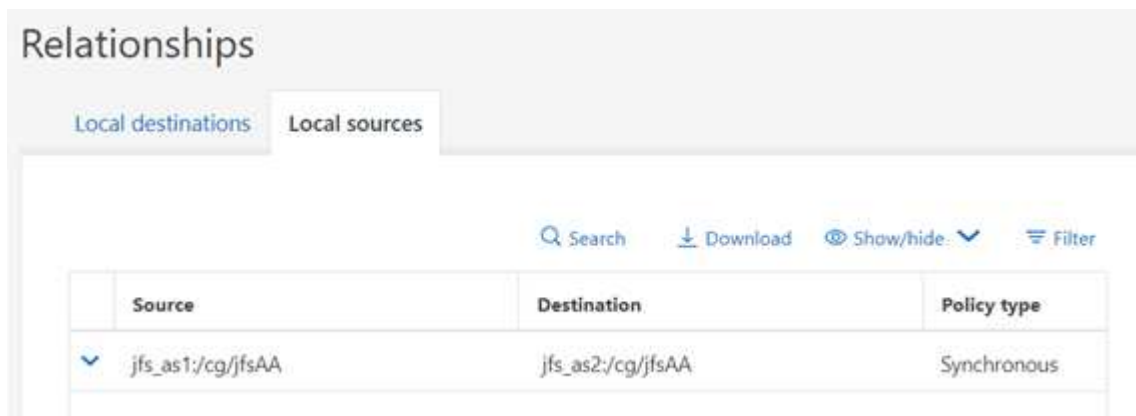
# SnapMirror active sync preferred site

SnapMirror active sync behavior is symmetric, with one important exception - preferred site configuration.

SnapMirror active sync will consider one site the "source" and the other the "destination". This implies a one-way replication relationship, but this does not apply to IO behavior. Replication is bidirectional and symmetric and IO response times are the same on either side of the mirror.

The `source` designation is controls the preferred site. If the replication link is lost, the LUN paths on the source copy will continue to serve data while the LUN paths on the destination copy will become unavailable until replication is reestablished and SnapMirror reenters a synchronous state. The paths will then resume serving data.

The sourced/destination configuration can be viewed via SystemManager:

## Relationships

| Local destinations | **Local sources** |

Q Search    ↓ Download    ⊘ Show/hide ✓    ≡ Filter

| Source | Destination | Policy type |
|---|---|---|
| ∨  jfs_as1:/cg/jfsAA | jfs_as2:/cg/jfsAA | Synchronous |

or at the CLI:

```
Cluster2::> snapmirror show -destination-path jfs_as2:/cg/jfsAA

                        Source Path: jfs_as1:/cg/jfsAA
                   Destination Path: jfs_as2:/cg/jfsAA
                  Relationship Type: XDP
            Relationship Group Type: consistencygroup
                 SnapMirror Schedule: -
              SnapMirror Policy Type: automated-failover-duplex
                   SnapMirror Policy: AutomatedFailOverDuplex
                         Tries Limit: -
                    Throttle (KB/sec): -
                        Mirror State: Snapmirrored
                Relationship Status: InSync
```

The key is that the source is the SVM on cluster1. As mentioned above, the terms "source" and "destination" don't describe the flow of replicated data. Both sites can process a write and replicate it to the opposite site. In effect, both clusters are sources and destinations. The effect of designating one cluster as a source simply controls which cluster survives as a read-write storage system if the replication link is lost.

# Network topology

## Uniform access

Uniform access networking means hosts are able to access paths on both sites (or failure domains within the same site).

An important feature of SM-as is the ability to configure the storage systems to know where the hosts are located. When you map the LUNs to a given host, you can indicate whether or not they are proximal to a given storage system.

### Proximity settings

Proximity refers to a per-cluster configuration that indicates a particular host WWN or iSCSI initiator ID belongs to a local host. It is a second, optional step for configuring LUN access.

The first step is the usual igroup configuration. Each LUN must be mapped to an igroup that contains the WWN/iSCSi IDs of the hosts that need access to that LUN. This controls which host has *access* to a LUN.
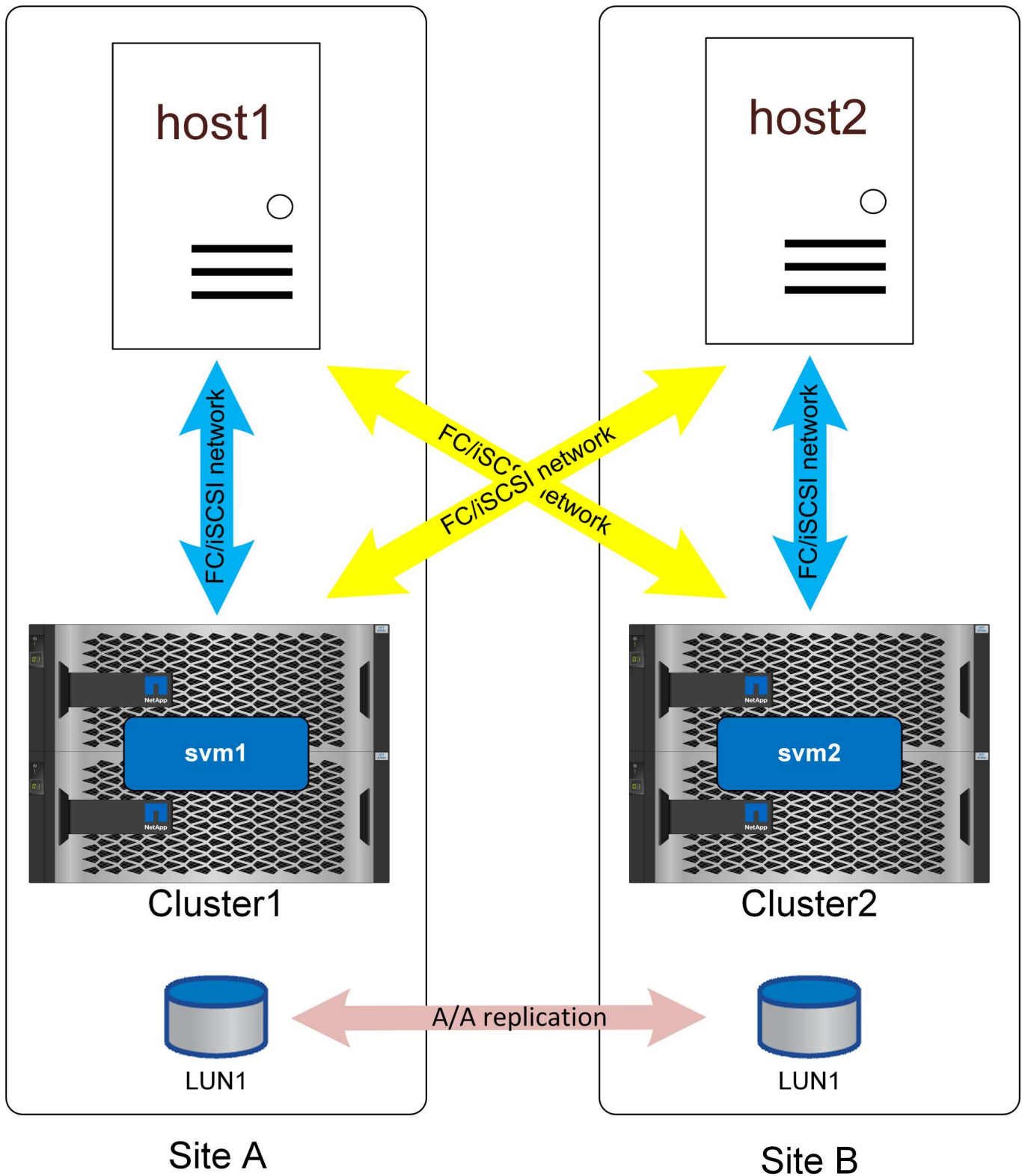
The second, optional step is to configure host proximity. This does not control access, it controls *priority*.

For example, a host at site A might be configured to access a LUN that is protected by SnapMirror active sync, and since the SAN is extended across sites, paths are available to that LUN using storage on site A or storage on site B.

Without proximity settings, that host will use both storage systems equally because both storage systems will advertise active/optimized paths. If the SAN latency and/or bandwidth between sites is limited, this may not be desireable, and you may wish to ensure that during normal operation each host preferentially uses paths to the local storage system. This is configured by adding the host WWN/iSCSI ID to the local cluster as a proximal host. This can be done at the CLI or SystemManager.

**AFF**

With an AFF system, the paths would appear as shown below when host proximity has been configured.

host1

host2

FC/iSCSI network

FC/iSCSI network

FC/iSCSI network

FC/iSCSI network

svm1

svm2

Cluster1

Cluster2

LUN1

LUN1

A/A replication

Site A

Site B
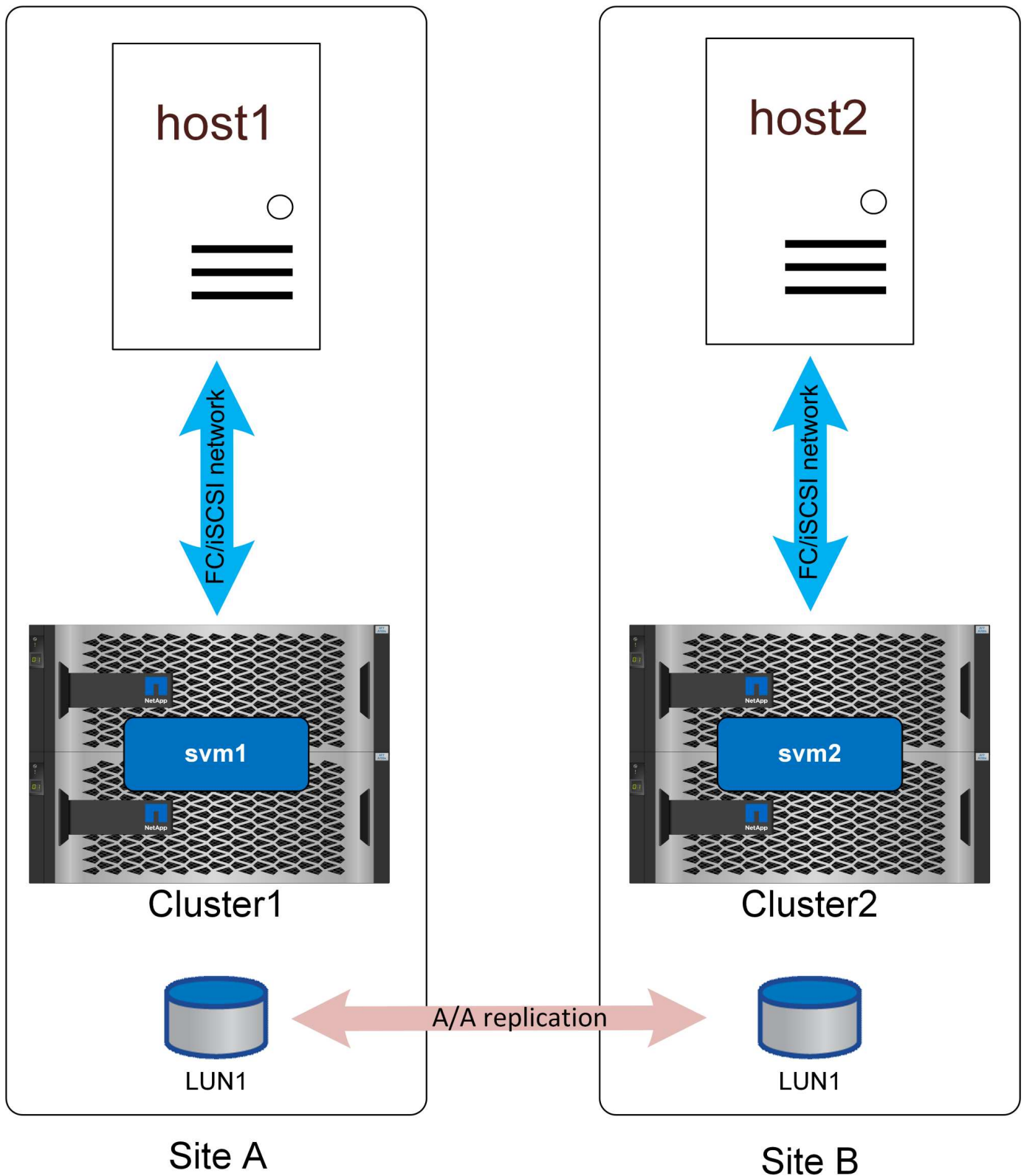
Active/Optimized Path

Active Path

In normal operation, all IO is local IO. Reads and writes are serviced from the local storage array. Write IO will, of course, also need to be replicated by the local controller to the remote system before being acknowledged, but all read IO will be serviced locally and will not incur extra latency by traversing the SAN link between sites.

The only time the nonoptimized paths will be used is when all active/optimized paths are lost. For example, if the entire array on site A lost power, the hosts at site A would still be able to access paths to the array on site B and therefore remain operational, although they would be experiencing higher latency.

There are redundant paths through the local cluster that are not shown on these diagrams for the sake of simplicity. ONTAP storage systems are HA themselves, so a controller failure should not result in site failure. It should merely result in a change in which local paths are used on the affected site.

## ASA

NetApp ASA systems offer active-active multipathing across all paths on a cluster. This also applies to SM-as configurations.

host1

host2

FC/iSCSI network

FC/iSCSI network

svm1

svm2

Cluster1

Cluster2

LUN1 ← A/A replication → LUN1

Site A
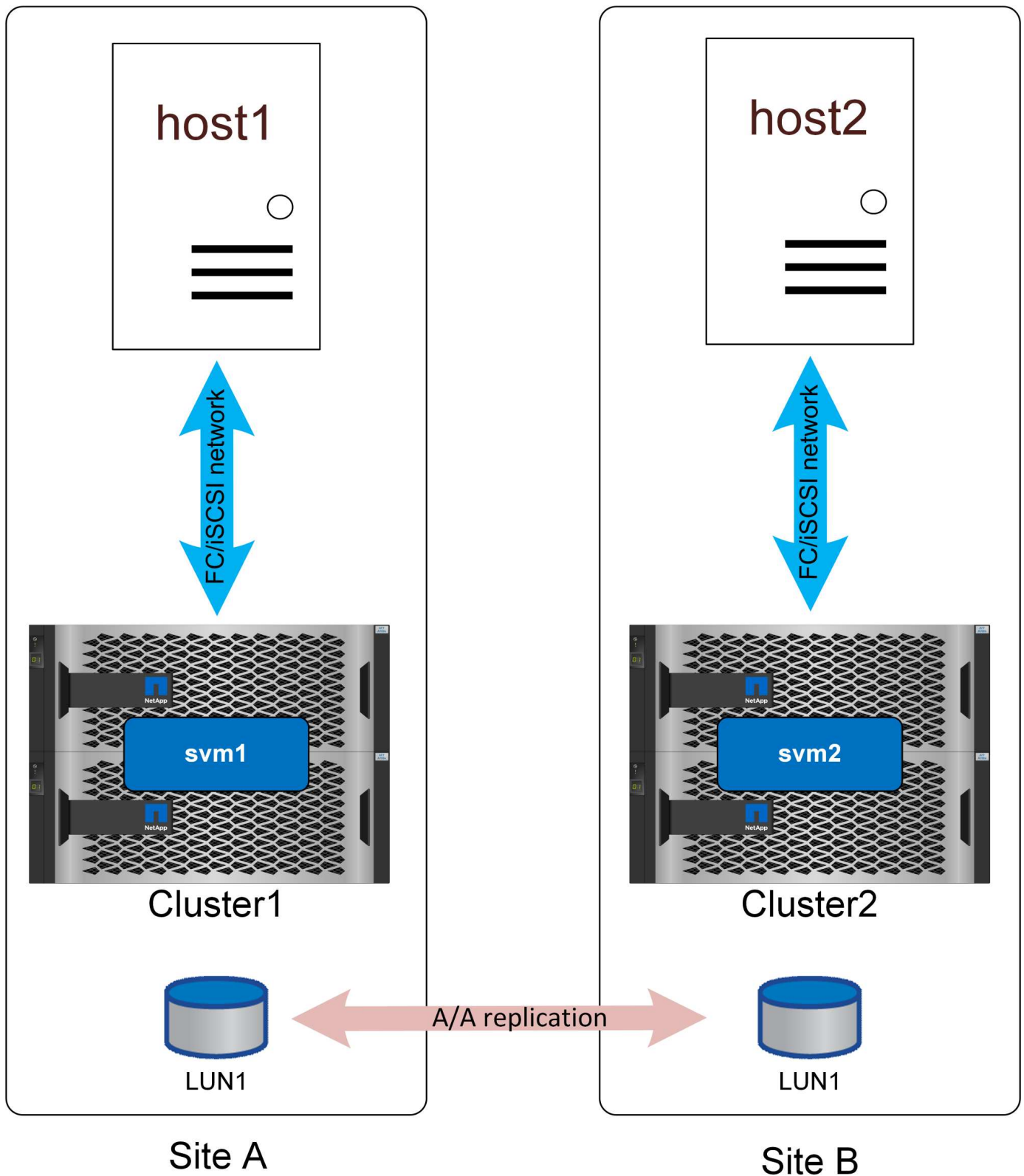
Site B

## Active/Optimized Path

An ASA configuration with non-uniform access would work largely the same as it would with AFF. With uniform access, IO would be crossing the WAN. This may or may not be desirable.

If the two sites were 100 meters apart with fiber connectivity there should be no detectable additional latency crossing the WAN, but if the sites were a long distance apart then read performance would suffer on both sites. In contrast, with AFF those WAN-crossing paths would only be used if there were no local paths available and day-to-day performance would be better because all IO would be local IO. ASA with nonuniform access network would be an option to gain the cost and feature benefits of ASA without incurring a cross-site latency access penalty.

ASA with SM-as in a low-latency configuration offers two interesting benefits. First, it essentially **doubles** the performance for any single host because IO can be serviced by twice as many controllers using twice as many paths. Second, in a single-site environment it offers extreme availability because an entire storage system could be lost without interrupting host access.

## Nonuniform access

Nonuniform access networking means each host only has access to ports on the local storage system. The SAN is not extended across sites (or failure domains within the same site).

The primary benefit to this approach is SAN simplicity - you remove the need to stretch a SAN over the network. Some customers don't have sufficiently low-latency connectivity between sites or lack the

infrastructure to tunnel FC SAN traffic over an intersite network.

The disadvantage to nonuniform access is that certain failure scenarios, including loss of the replication link, will result some hosts losing access to storage. Applications that run as single instances, such as a non-clustered database that is inherently only running on a single host at any given mount would fail if local storage connectivity was lost. The data would still be protected, but the database server would no longer have access. It would need to be restarted on a remote site, preferably through an automated process. For example, VMware HA can detect an all-paths-down situation on one server and restart a VM on another server where paths are available.

In contrast, a clustered application such as Oracle RAC can deliver a service that is simultaneously available at two different sites. Losing a site doesn't mean loss of the application service as a whole. Instances are still available and running at the surviving site.

In many cases, the additional latency overhead of an application accessing storage across a site-to-site link would be unacceptable. This means that the improved availability of uniform networking is minimal, since loss of storage on a site would lead to the need to shut down services on that failed site anyway.

> (i) There are redundant paths through the local cluster that are not shown on these diagrams for the sake of simplicity. ONTAP storage systems are HA themselves, so a controller failure should not result in site failure. It should merely result in a change in which local paths are used on the affected site.

# Oracle Configurations

## Overview

The use of SnapMirror active sync does not necessarily add to or change any best practices for operating a database.
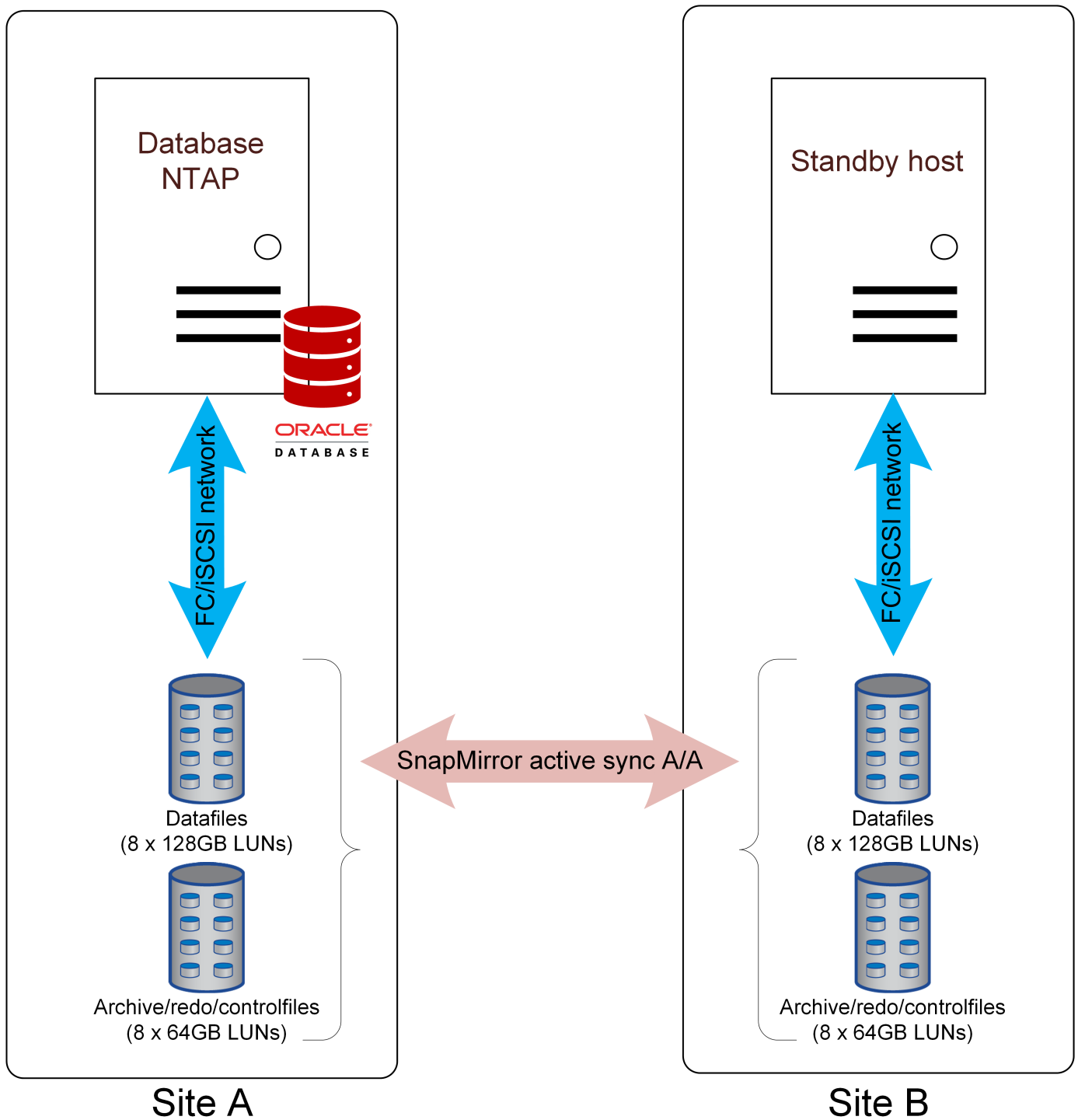
The best architecture depends on the business requirements. For example, if the goal is to have RPO=0 protection against data loss, but the RTO is relaxed, then using Oracle Single Instance databases and replicating the LUNs with SM-as might be sufficient as well as less expensive from an Oracle licensing standpoing. Failure of the remote site would not interrupt operations, and loss of the primary site would result in LUNs at the surviving site that are online and ready to be used.

If the RTO was more strict, basic active-passive automation through scripts or clusterware such as Pacemaker or Ansible would improve failover time. For example, VMware HA could be configured to detect VM failure on the primary site and active the VM on the remote site.

Finally, for extremely rapid failover, Oracle RAC could be deployed across sites. The RTO would essentially be zero because the database would be online and available on both sites at all times.

## Oracle Single-Instance

The examples explained below show some of the many options for to deploying Oracle Single Instance databases with SnapMirror active sync replication.

Site A
Site B

**Failover with a preconfigured OS**

SnapMirror active sync delivers a synchronous copy of the data at the disaster recovery site, but making that data available requires an operating system and the associated applications. Basic automation can dramatically improve the failover time of the overall environment. Clusterware products such as Pacemaker are often used to create a cluster across the sites, and in many cases the failover process can be driven with simple scripts.

If the primary nodes are lost, the clusterware (or scripts) will bring the databases online at the alternate site. One option is to create standby servers that are preconfigured for the SAN resources that make up the database. If the primary site fails, the clusterware or scripted alternative performs a sequence of actions similar

to the following:

1. Detect failure of primary site
2. Perform discovery of FC or iSCSI LUNs
3. Mounting file systems and/or mounting ASM disk groups
4. Starting the database

The primary requirement of this approach is a running OS in place on the remote site. It must be preconfigured with Oracle binaries, which also means that tasks such as Oracle patching must be performed on the primary and standby site. Alternatively, the Oracle binaries can be mirrored to the remote site and mounted if a disaster is declared.

The actual activation procedure is simple. Commands such as LUN discovery require just a few commands per FC port. File system mounting is nothing more than a `mount` command, and both databases and ASM can be started and stopped at the CLI with a single command.

**Failover with a virtualized OS**

Failover of database environments can be extended to include the operating system itself. In theory, this failover can be done with boot LUNs, but most often it is done with a virtualized OS. The procedure is similar to the following steps:

1. Detect failure of primary site
2. Mounting the datastores hosting the database server virtual machines
3. Starting the virtual machines
4. Starting databases manually or configuring the virtual machines to automatically start the databases.

For example, an ESX cluster could span sites. In the event of disaster, the virtual machines can be brought online at the disaster recovery site after the switchover.

**Storage failure protection**

The diagram above shows the use of nonuniform access, where the SAN is not stretched across sites. This may be simpler to configure, and in some cases may be the only option given the current SAN capabilities, but it also means that failure of the primary storage system would cause a database outage until the application was failed over.

For additional resilience, the solution could be deployed with uniform access. This would allow the applications to continue operating using the paths advertised from the opposite site.

## Oracle Extended RAC

Many customers optimize their RTO by stretching an Oracle RAC cluster across sites, yielding a fully active-active configuration. The overall design becomes more complicated because it must include quorum management of Oracle RAC.
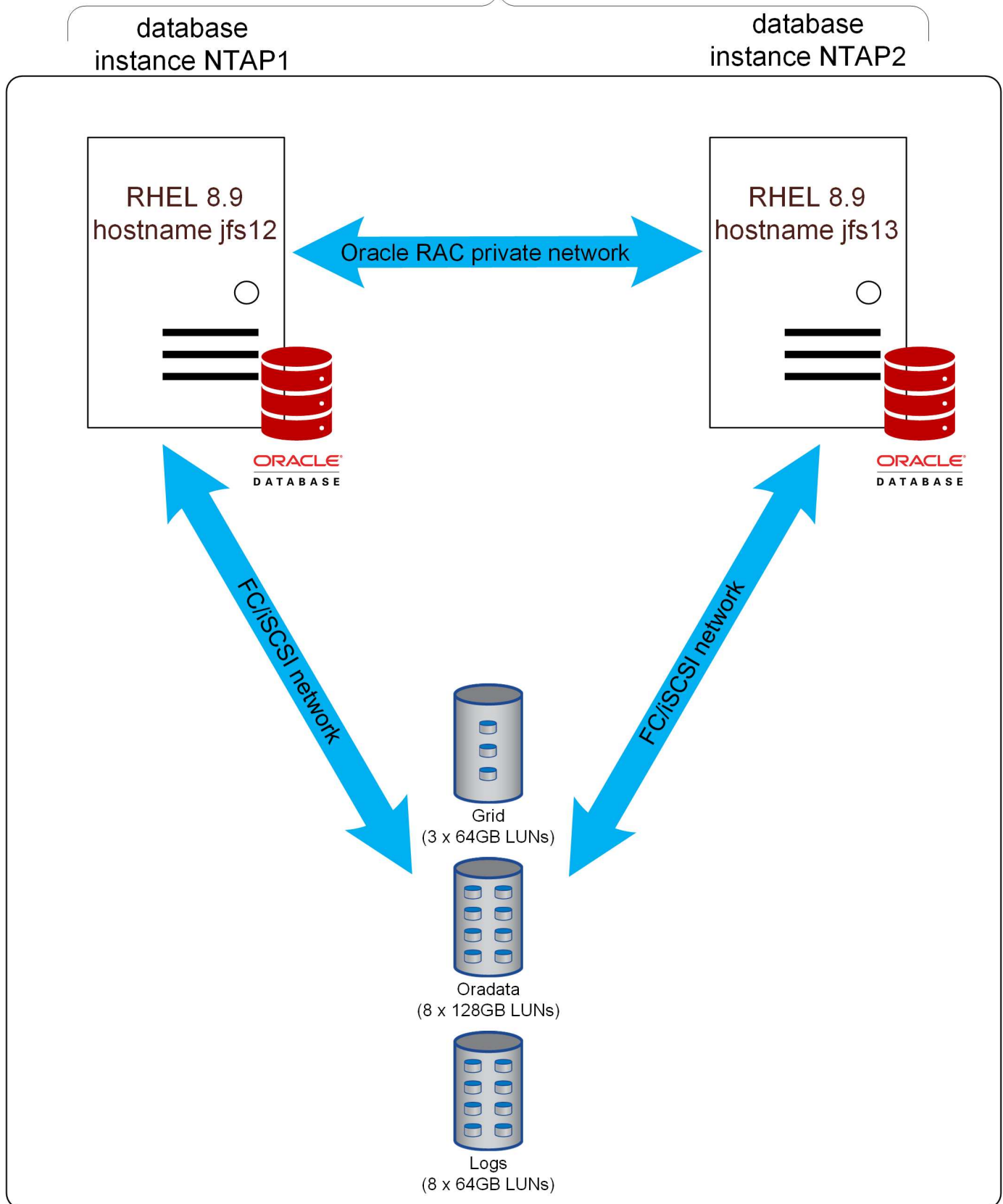
Traditional extended RAC clustered relied on ASM mirroring to provide data protection. This approach works, but it also requires a lot of manual configuration steps and imposes overhead on the network infrastructure. In contrast, allowing SnapMirror active sync to take responsibility for data replication dramatically simplifies the solution. Operations such as synchronization, resynchronization after disruptions, failovers, and quorum management are easier, plus the SAN does not need to be distributed across sites which simplifies SAN

design and management.

**Replication**

They key to understanding RAC functionality on SnapMirror active sync is to view storage as a single set of LUNs which are hosted on mirrored storage. For example:

# Database NTAP

## database instance NTAP1

## database instance NTAP2



RHEL 8.9
hostname jfs12

Oracle RAC private network

RHEL 8.9
hostname jfs13

ORACLE®
DATABASE

ORACLE®
DATABASE

FC/iSCSI network

FC/iSCSI network

Grid
(3 x 64GB LUNs)

Oradata
(8 x 128GB LUNs)

Logs
(8 x 64GB LUNs)

There is no primary copy or mirror copy. Logically, there is only a single copy of each LUN, and that LUN is available on SAN paths that are located on two different storage systems. From a host point of view, there are no storage failovers; instead there are path changes. Various failure events might lead to loss of certain paths

to the LUN while other paths remain online. SnapMirror active sync ensures the same data is available across all operational paths.

**Storage configuration**

In this example configuration, the ASM disks are configured the same as they would be in any single-site RAC configuration on enterprise storage. Since the storage system provides data protection, ASM external redundancy would be used.

**Uniform vs nonuninform access**

The most important consideration with Oracle RAC on SnapMirror active sync is whether to use uniform or nonuniform access.
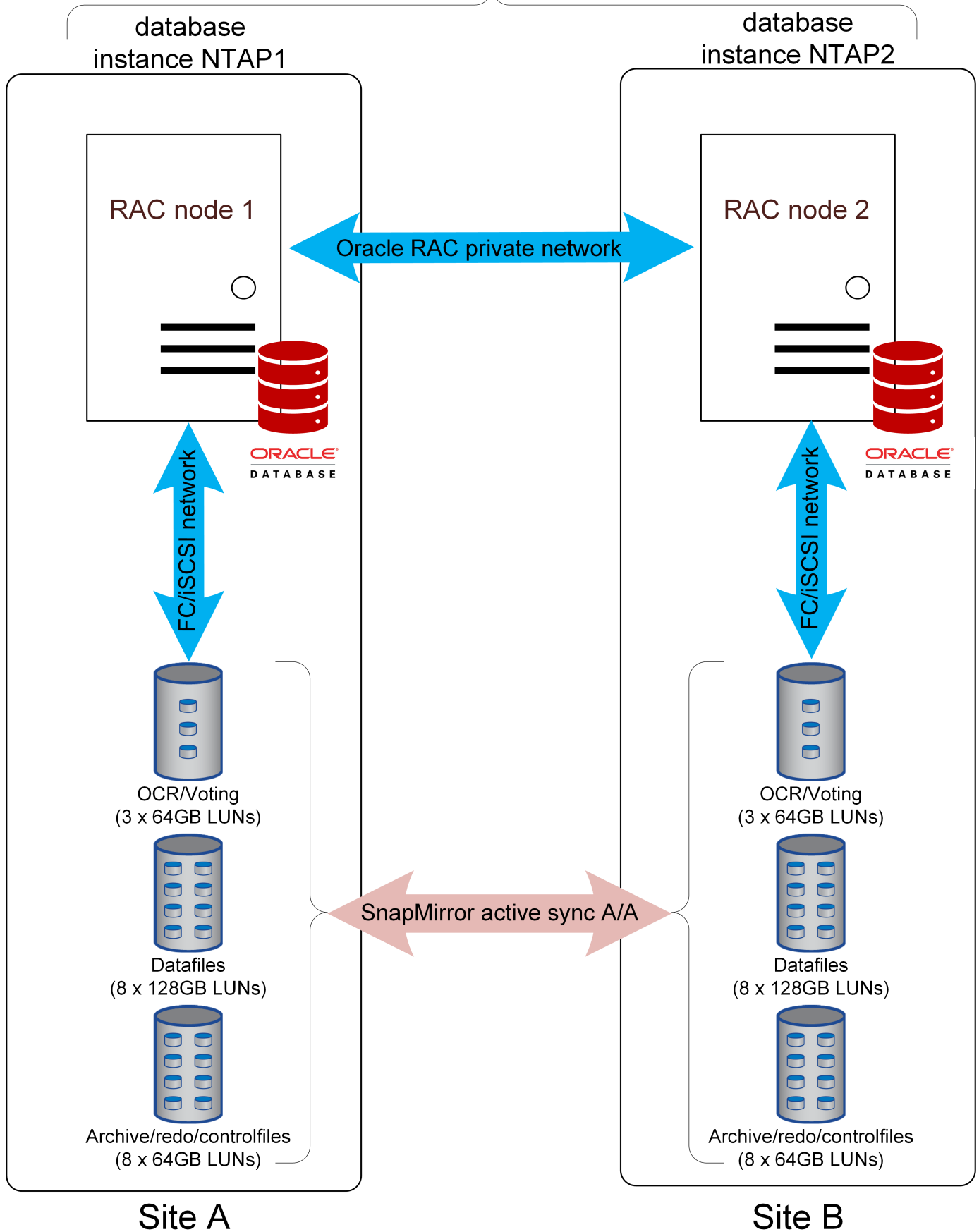
Uniform access means each host can see paths on both clusters. Nonuniform access means hosts can only see paths to the local cluster.

Neither option is specifically recommended or discouraged. Some customers have dark fibre readily available to connect sites, others either do not have such connectivity or their SAN infrastructure doesn't support a long-distance ISL.

**Nonuniform access**

Nonuniform access is simpler to configure from a SAN perspective.

Database NTAP

database instance NTAP1

database instance NTAP2

RAC node 1

RAC node 2

Oracle RAC private network

ORACLE® DATABASE

FC/iSCSI network

OCR/Voting
(3 x 64GB LUNs)

Datafiles
(8 x 128GB LUNs)

Archive/redo/controlfiles
(8 x 64GB LUNs)

SnapMirror active sync A/A

Site A

Site B

The primary downside of the nonuniform access approach is that loss of site-to-site ONTAP connectivity or loss of a storage system will result in loss of the database instances at one site. This obviously is not desirable, but it may be an acceptable risk in exchange for a simpler SAN configuration.

**Uniform access**

Uniform access requires extending the SAN across sites. The primary benefit is that loss of a storage system will not result in loss of a database instance. Instead, it would result in a multipathing change in which paths are currently in use.

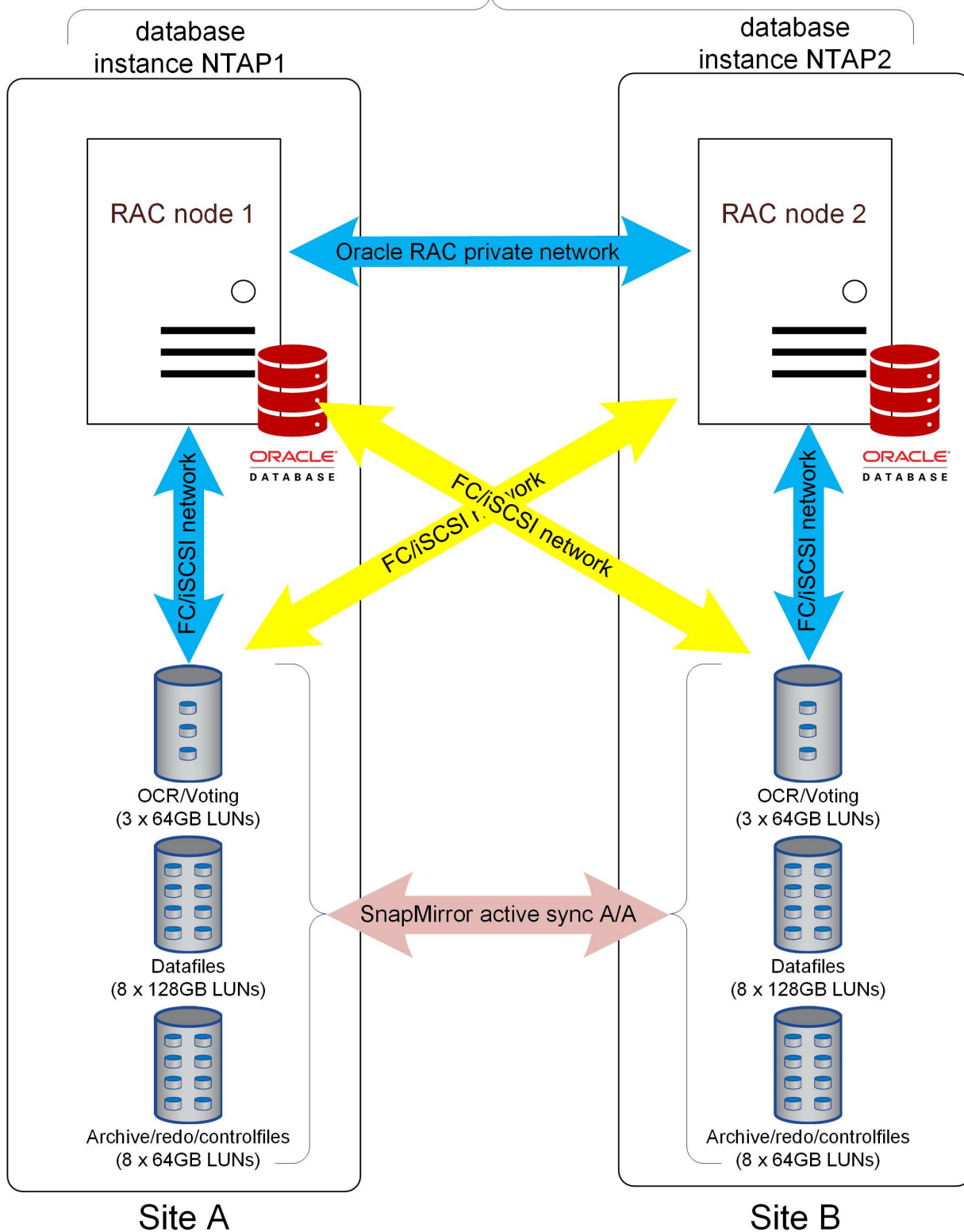There are several ways to configure nonuniform access.

> ℹ️ In the diagrams below, there are also active but nonoptimized paths present that would be used during simple controller failures, but those paths are not shown in the interest of simplifying the diagrams.

**AFF with proximity settings**

If there is significant latency between sites, then AFF systems can be configured with host proximity settings. This allows each storage system to be aware of which hosts are local and which are remote and assign path priorities appropriately.

Database NTAP

database instance NTAP1

database instance NTAP2

RAC node 1

RAC node 2

Oracle RAC private network

ORACLE
DATABASE

ORACLE
DATABASE

FC/iSCSI network

FC/iSCSI network
FC/iSCSI network

FC/iSCSI network

OCR/Voting
(3 x 64GB LUNs)

OCR/Voting
(3 x 64GB LUNs)

SnapMirror active sync A/A

Datafiles
(8 x 128GB LUNs)

Datafiles
(8 x 128GB LUNs)

Archive/redo/controlfiles
(8 x 64GB LUNs)

Archive/redo/controlfiles
(8 x 64GB LUNs)
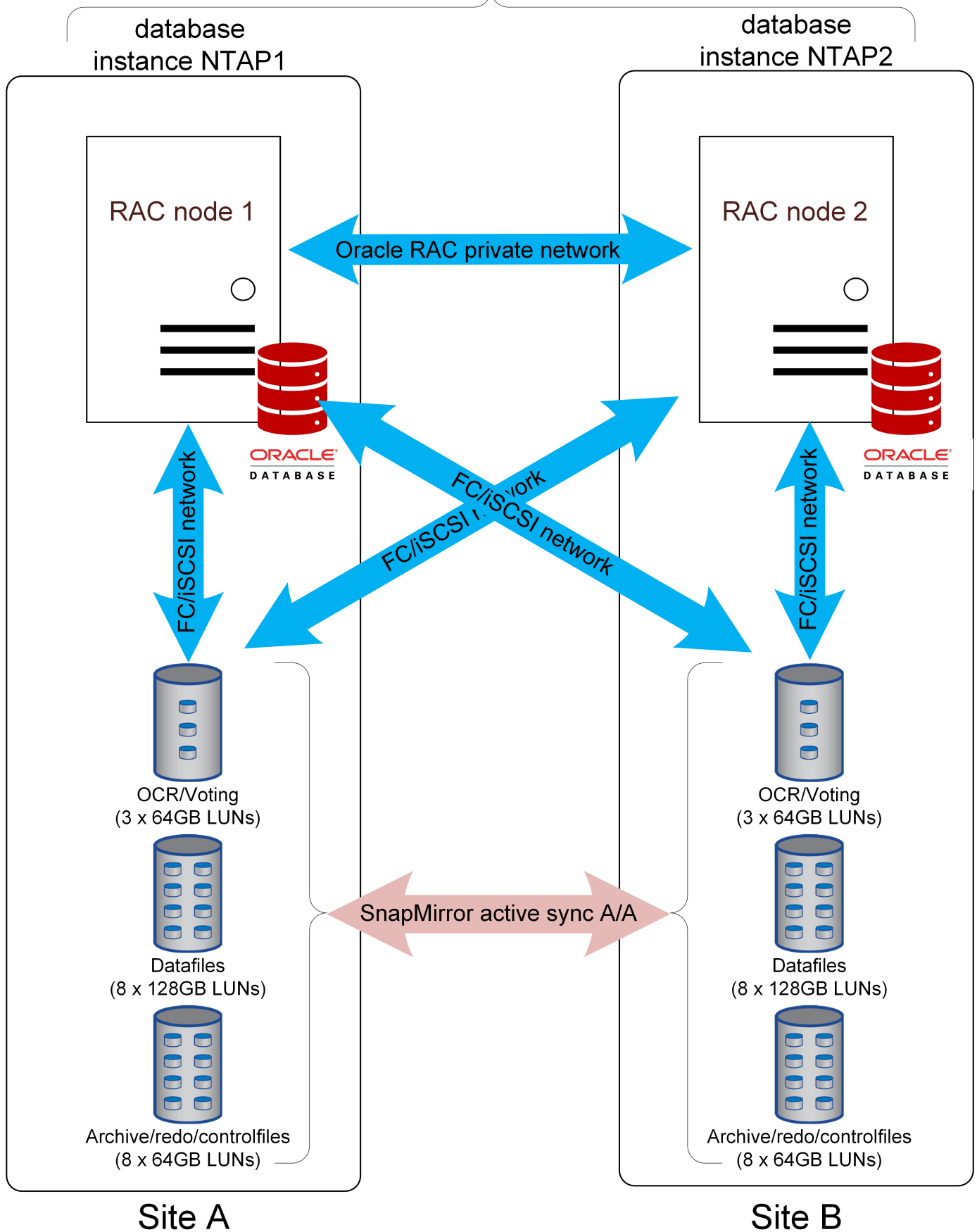
Site A

Site B

Active/Optimized Path

Active Path

In normal operation, each Oracle instance would preferentially use the local active/optimized paths. The result is that all reads would be serviced by the local copy of the blocks. This yields the lowest possible latency. Write IO is similarly sent down paths to the local controller. The IO must still be replicated before being acknowledged and therefor would still incur the additional latency of crossing the site-to-site network, but this cannot be avoided in a synchronous replication solution.

**ASA / AFF without proximity settings**

If there is no significant latency between sites, then AFF systems can be configured without host proximity settings, or ASA can be used.

Database NTAP

database instance NTAP1

database instance NTAP2

RAC node 1

RAC node 2

Oracle RAC private network

ORACLE DATABASE

ORACLE DATABASE

FC/iSCSI network

FC/iSCSI network

FC/iSCSI network

FC/iSCSI network

OCR/Voting
(3 x 64GB LUNs)

OCR/Voting
(3 x 64GB LUNs)

SnapMirror active sync A/A

Datafiles
(8 x 128GB LUNs)

Datafiles
(8 x 128GB LUNs)

Archive/redo/controlfiles
(8 x 64GB LUNs)

Archive/redo/controlfiles
(8 x 64GB LUNs)

Site A

Site B

Each host will be able to use all operational paths on both storage systems. This potentially improves performance significantly by allowing each host to draw upon the performance potential of two clusters, not just one.

With ASA, not only would all paths to both clusters be considered active and optimized, but the paths on partner controllers would also be active. The result would be all-active SAN paths on the entire cluster, all the time.

> ⓘ ASA systems may also be used in a nonuniform access configuration. Since no cross-site paths exist, there would be no impact on performance resulting from IO crossing the ISL.

## RAC tiebreaker

While extended RAC using SnapMirror active sync is a symmetric architecture with respect to IO, there is one exception that is connected to split-brain management.

What happens if the replication link is lost and neither site has quorum? What should happen? This question applies to both the Oracle RAC and the ONTAP behavior. If changes cannot be replicated across sites, and you want to resume operations, one of the sites will have to survive and the other site will have to become unavailable.

The ONTAP Mediator addresses this requirement at the ONTAP layer. There are multiple options for RAC tiebreaking.

### Oracle tiebreakers

The best method to manage split-brain Oracle RAC risks is to use an odd number of RAC nodes, preferably by use of a 3rd site tiebreaker. If a 3rd site is unavailable, the tiebreaker instance could be placed on one site of the two sites, effectively designating it a preferred survivor site.

### Oracle and css_critical

With an even number of nodes, the default Oracle RAC behavior is that one of the nodes in the cluster will be deemed more important than the other nodes. The site with that higher priority node will survive site isolation while the nodes on the other site will evict. The prioritization is based on multiple factors, but you can also control this behavior using the css_critical setting.

In the example architecture, the hostnames for the RAC nodes are jfs12 and jfs13. The current settings for css_critical are as follows:

```
[root@jfs12 ~]# /grid/bin/crsctl get server css_critical
CRS-5092: Current value of the server attribute CSS_CRITICAL is no.

[root@jfs13 trace]# /grid/bin/crsctl get server css_critical
CRS-5092: Current value of the server attribute CSS_CRITICAL is no.
```

If you want the site with jfs12 to be the preferred site, change this value to yes on a site A node and restart services.

```
[root@jfs12 ~]# /grid/bin/crsctl set server css_critical yes
CRS-4416: Server attribute 'CSS_CRITICAL' successfully changed. Restart
Oracle High Availability Services for new value to take effect.

[root@jfs12 ~]# /grid/bin/crsctl stop crs
CRS-2791: Starting shutdown of Oracle High Availability Services-managed
resources on 'jfs12'
CRS-2673: Attempting to stop 'ora.crsd' on 'jfs12'
CRS-2790: Starting shutdown of Cluster Ready Services-managed resources on
server 'jfs12'
CRS-2673: Attempting to stop 'ora.ntap.ntappdb1.pdb' on 'jfs12'
…
CRS-2673: Attempting to stop 'ora.gipcd' on 'jfs12'
CRS-2677: Stop of 'ora.gipcd' on 'jfs12' succeeded
CRS-2793: Shutdown of Oracle High Availability Services-managed resources
on 'jfs12' has completed
CRS-4133: Oracle High Availability Services has been stopped.

[root@jfs12 ~]# /grid/bin/crsctl start crs
CRS-4123: Oracle High Availability Services has been started.
```

# Failure scenarios

## Overview

Planning a complete SnapMirror active sync application architecture requires understanding how SM-as will respond in various planned and unplanned failover scenarios.

For the following examples, assume that site A is configured as the preferred site.

### Loss of replication connectivity

If SM-as replication is interrupted, write IO cannot be completed because it would be impossible for a cluster to replicate changes to the opposite site.

### Site A (Preferred site)

The result of replication link failure on the preferred site will be an approximate 15 second pause in write IO processing as ONTAP retries replicated write operations before it determines that the replication link is genuinely unreachable. After the 15 seconds elapses, the site A system resumes read and write IO processing. The SAN paths will not change, and the LUNs will remain online.

### Site B

Since site B is not the SnapMirror active sync preferred site, its LUN paths will become unavailable after about 15 seconds.

**Storage system failure**

The result of a storage system failure is nearly identical to the result of losing the replication link. The surviving site should experience a roughly 15 second IO pause. Once that 15 second period elapses, IO will resume on that site as usual.

**Loss of the mediator**

The mediator service does not directly control storage operations. It functions as an alternate control path between clusters. It exists primarily to automate failover without the risk of a split-brain scenario. In normal operation, each cluster is replicating changes to its partner, and each cluster therefore can verify that the partner cluster is online and serving data. If the replication link failed, replication would cease.

The reason a mediator is required for safe automated failover is because it would otherwise be impossible for a storage cluster to be able to determine whether loss of bidirectional communication was the result of a network outage or actual storage failure.

The mediator provides an alternate path for each cluster to verify the health of its partner. The scenarios are as follows:

- If a cluster can contact its partner directly, replication services are operational. No action required.
- If a preferred site cannot contact its partner directly or via the mediator, it will assume the partner is either actually unavailable or was isolated and has taken its LUN paths offline. The preferred site will then proceed to release the RPO=0 state and continue processing both read and write IO.
- If a non-preferred site cannot contact its partner directly, but can contact it via the mediator, it will take its paths offline and await the return of the replication connection.
- If a non-preferred site cannot contact its partner directly or via an operational mediator, it will assume the partner is either actually unavailable or was isolated and has taken its LUN paths offline. The non-preferred site will then proceed to release the RPO=0 state and continue processing both read and write IO. It will assume the role of the replication source and will become the new preferred site.

If the mediator is wholly unavailable:

- Failure of replication services for any reason, including failure of the nonpreferred site or storage system, will result in the preferred site releasing the RPO=0 state and resuming read and write IO processing. The non-preferred site will take its paths offline.
- Failure of the preferred site will result in an outage because the non-preferred site will be unable to verify that the opposite site is truly offline and therefore it would not be safe for the nonpreferred site to resume services.

**Restoring services**

After a failure is resolved, such as restoring site-to-site connectivity or powering on a failed system, the SnapMirror active sync endpoints will automatically detect the presence of a faulty replication relationship and bring it back to an RPO=0 state. Once synchronous replication is reestablished, the failed paths will come online again.

In many cases, clustered applications will automatically detect the return of failed paths, and those applications will also come back online. In other cases, a host-level SAN scan may be required, or applications may need to be brought back online manually. It depends on the application and how it is configured, and in general such tasks can be easily automated. ONTAP itself is self-healing and should not require any user intervention to resume RPO=0 storage operations.

**Manual failover**

Changing the preferred site requires a simple operation. IO will pause for a second or two as authority over replication behavior switches between clusters, but IO is otherwise unaffected.
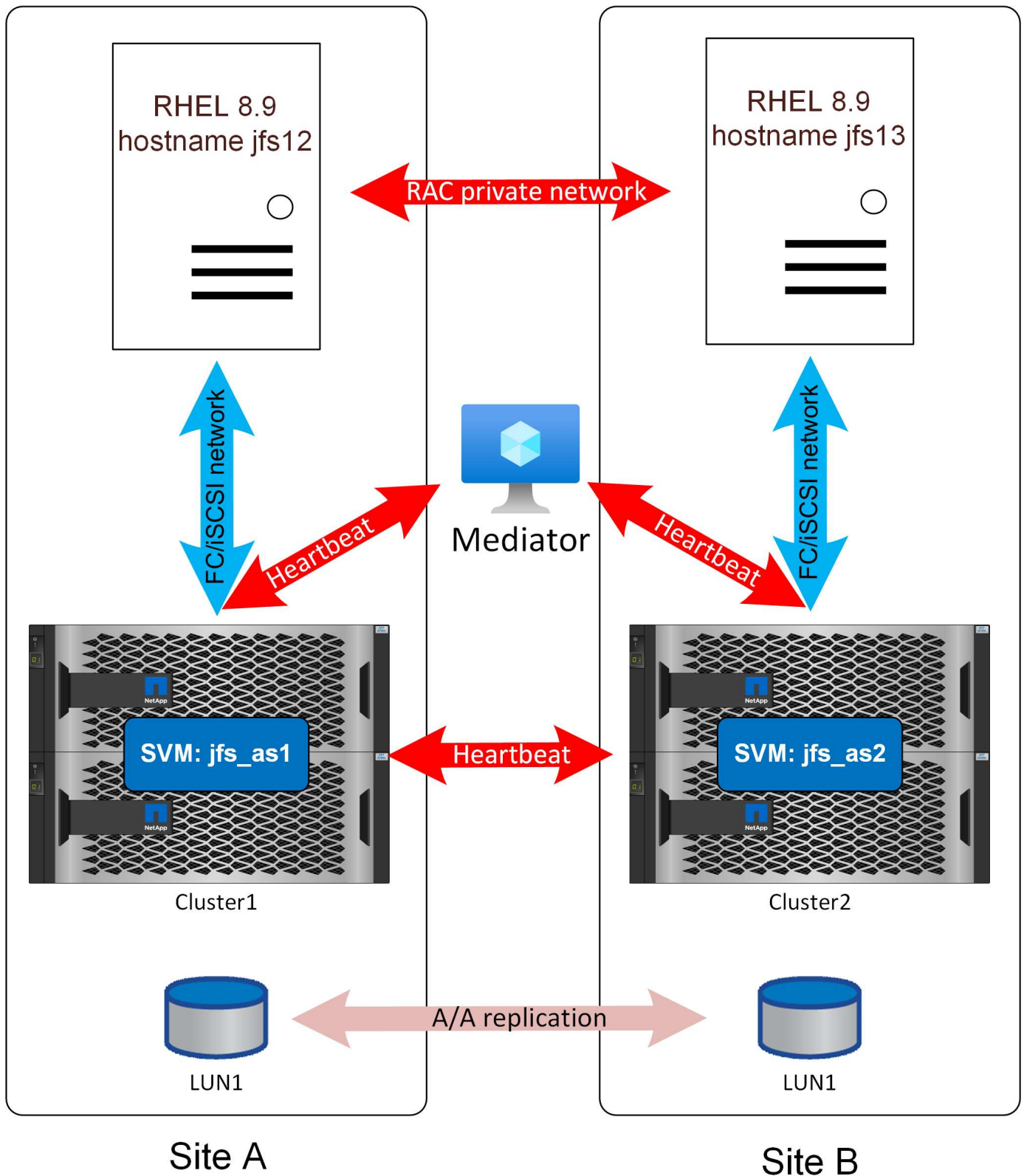
## Sample architecture

The detailed failure examples shown in this sections are based on the architecture shown below.

> ⓘ This is only one of many options for Oracle databases on SnapMirror active sync. This design was chosen because it illustrates some of the more complicated scenarios.

In this design, assume that site A is set at the preferred site.

RHEL 8.9
hostname jfs12

RHEL 8.9
hostname jfs13

RAC private network

FC/iSCSI network

FC/iSCSI network

Mediator

Heartbeat

Heartbeat

SVM: jfs_as1

SVM: jfs_as2

Heartbeat

Cluster1

Cluster2

A/A replication

LUN1

LUN1

Site A

Site B

**RAC interconnect failure**

Loss of the Oracle RAC replication link will produce a similar result to loss of SnapMirror connectivity, except the timeouts will be shorter by default. Under default settings, an Oracle RAC node will wait 200 seconds after loss of storage connectivity before evicting,

but it will only wait 30 seconds after loss of the RAC network heartbeat.

The CRS messages are similar to those shown below. You can see the 30 second timeout lapse. Since css_critical was set on jfs12, located on site A, that will be the site to survive and jfs13 on site B will be evicted.

```
2024-09-12 10:56:44.047 [ONMD(3528)]CRS-1611: Network communication with
node jfs13 (2) has been missing for 75% of the timeout interval.  If this
persists, removal of this node from cluster will occur in 6.980 seconds
2024-09-12 10:56:48.048 [ONMD(3528)]CRS-1610: Network communication with
node jfs13 (2) has been missing for 90% of the timeout interval.  If this
persists, removal of this node from cluster will occur in 2.980 seconds
2024-09-12 10:56:51.031 [ONMD(3528)]CRS-1607: Node jfs13 is being evicted
in cluster incarnation 621599354; details at (:CSSNM00007:) in
/gridbase/diag/crs/jfs12/crs/trace/onmd.trc.
2024-09-12 10:56:52.390 [CRSD(6668)]CRS-7503: The Oracle Grid
Infrastructure process 'crsd' observed communication issues between node
'jfs12' and node 'jfs13', interface list of local node 'jfs12' is
'192.168.30.1:33194;', interface list of remote node 'jfs13' is
'192.168.30.2:33621;'.
2024-09-12 10:56:55.683 [ONMD(3528)]CRS-1601: CSSD Reconfiguration
complete. Active nodes are jfs12 .
2024-09-12 10:56:55.722 [CRSD(6668)]CRS-5504: Node down event reported for
node 'jfs13'.
2024-09-12 10:56:57.222 [CRSD(6668)]CRS-2773: Server 'jfs13' has been
removed from pool 'Generic'.
2024-09-12 10:56:57.224 [CRSD(6668)]CRS-2773: Server 'jfs13' has been
removed from pool 'ora.NTAP'.
```

## SnapMirror communication failure

If the SnapMirror active sync replication link, write IO cannot be completed because it would be impossible for a cluster to replicate changes to the opposite site.

### Site A

The result on site A of a replication link failure will be an approximately 15 second pause in write IO processing as ONTAP attempts to replicate writes before it determines that the replication link is genuinely inoperable. After the 15 seconds elapses, the ONTAP cluster on site A resumes read and write IO processing. The SAN paths will not change, and the LUNs will remain online.

### Site B

Since site B is not the SnapMirror active sync preferred site, its LUN paths will become unavailable after about 15 seconds.

The replication link was cut at the timestamp 15:19:44. The first warning from Oracle RAC arrives 100 seconds later as the 200 second timeout (controlled by the Oracle RAC parameter disktimeout) approaches.

```
2024-09-10 15:21:24.702 [ONMD(2792)]CRS-1615: No I/O has completed after
50% of the maximum interval. If this persists, voting file
/dev/mapper/grid2 will be considered not functional in 99340 milliseconds.
2024-09-10 15:22:14.706 [ONMD(2792)]CRS-1614: No I/O has completed after
75% of the maximum interval. If this persists, voting file
/dev/mapper/grid2 will be considered not functional in 49330 milliseconds.
2024-09-10 15:22:44.708 [ONMD(2792)]CRS-1613: No I/O has completed after
90% of the maximum interval. If this persists, voting file
/dev/mapper/grid2 will be considered not functional in 19330 milliseconds.
2024-09-10 15:23:04.710 [ONMD(2792)]CRS-1604: CSSD voting file is offline:
/dev/mapper/grid2; details at (:CSSNM00058:) in
/gridbase/diag/crs/jfs13/crs/trace/onmd.trc.
2024-09-10 15:23:04.710 [ONMD(2792)]CRS-1606: The number of voting files
available, 0, is less than the minimum number of voting files required, 1,
resulting in CSSD termination to ensure data integrity; details at
(:CSSNM00018:) in /gridbase/diag/crs/jfs13/crs/trace/onmd.trc
2024-09-10 15:23:04.716 [ONMD(2792)]CRS-1699: The CSS daemon is
terminating due to a fatal error from thread:
clssnmvDiskPingMonitorThread; Details at (:CSSSC00012:) in
/gridbase/diag/crs/jfs13/crs/trace/onmd.trc
2024-09-10 15:23:04.731 [OCSSD(2794)]CRS-1652: Starting clean up of CRSD
resources.
```

Once the 200 second voting disk timeout has been reached, this Oracle RAC node will evict itself from the cluster and reboot.

## Total network interconnectivity failure

If the replication link between sites is completely lost, both SnapMirror active sync and Oracle RAC connectivity will be interrupted.

Oracle RAC split-brain detection has a dependency on the Oracle RAC storage heartbeat. If loss of site-to-site connectivity results in simultaneous loss of both the RAC network heartbeat and storage replication services, the result is the RAC sites will not be able to communicate cross-site via either the RAC interconnect or the RAC voting disks. The result in an even-numbed set of nodes may be eviction of both sites under default settings. The exact behavior will depend on the sequence of events and the timing of the RAC network and disk heartbeat polls.

The risk of a 2-site outage can be addressed in two ways. First, a tiebreaker configuration can be used.

If a 3rd site is not available, this risk can be addressed by adjusting the misscount parameter on the RAC cluster. Under the defaults, the RAC network heartbeat timeout is 30 seconds. This normally is used by RAC to identify failed RAC nodes and remove them from the cluster. It also has a connection to the voting disk heartbeat.

If, for example, the conduit carrying intersite traffic for both Oracle RAC and storage replication services is cut by a backhoe, the 30 second misscount countdown will begin. If the RAC preferred site node cannot reestablish contact with the opposite site within 30 seconds, and it also cannot use the voting disks to confirm

the opposite site is down within that same 30 second window, then the preferred site nodes will also evict. The result is a full database outage.

Depending on when the misscount polling occurs, 30 seconds may not be enough time for SnapMirror active sync to time out and allow storage on the preferred site to resume services before the 30 second window expires. This 30 second window can be increased.

```
[root@jfs12 ~]# /grid/bin/crsctl set css misscount 100
CRS-4684: Successful set of parameter misscount to 100 for Cluster
Synchronization Services.
```

This value allows the storage system on the preferred site to resume operations before the miscount timeout expires. The result will then be eviction only of the nodes at the site where the LUN paths were removed. Example below:

```
2024-09-12 09:50:59.352 [ONMD(681360)]CRS-1612: Network communication with
node jfs13 (2) has been missing for 50% of the timeout interval.  If this
persists, removal of this node from cluster will occur in 49.570 seconds
2024-09-12 09:51:10.082 [CRSD(682669)]CRS-7503: The Oracle Grid
Infrastructure process 'crsd' observed communication issues between node
'jfs12' and node 'jfs13', interface list of local node 'jfs12' is
'192.168.30.1:46039;', interface list of remote node 'jfs13' is
'192.168.30.2:42037;'.
2024-09-12 09:51:24.356 [ONMD(681360)]CRS-1611: Network communication with
node jfs13 (2) has been missing for 75% of the timeout interval.  If this
persists, removal of this node from cluster will occur in 24.560 seconds
2024-09-12 09:51:39.359 [ONMD(681360)]CRS-1610: Network communication with
node jfs13 (2) has been missing for 90% of the timeout interval.  If this
persists, removal of this node from cluster will occur in 9.560 seconds
2024-09-12 09:51:47.527 [OHASD(680884)]CRS-8011: reboot advisory message
from host: jfs13, component: cssagent, with time stamp: L-2024-09-12-
09:51:47.451
2024-09-12 09:51:47.527 [OHASD(680884)]CRS-8013: reboot advisory message
text: oracssdagent is about to reboot this node due to unknown reason as
it did not receive local heartbeats for 10470 ms amount of time
2024-09-12 09:51:48.925 [ONMD(681360)]CRS-1632: Node jfs13 is being
removed from the cluster in cluster incarnation 621596607
```

Oracle Support strongly discourages altering with the misscount or disktimeout parameters to solve configuration problems. Changing these parameters can, however, be warranted and unavoidable in many cases, including SAN booting, virtualized, and storage replication configurations. If, for example, you had stability problems with a SAN or IP network that was resulting in RAC evictions you should fix the underlying problem and not charge the values of the misscount or disktimeout. Changing timeouts to address configuration errors is masking a problem, not solving a problem. Changing these parameters to properly configure a RAC environment based on design aspects of the underlying infrastructure is different and is consistent with Oracle support statements. With SAN booting, it is common to adjust misscount all the way up to 200 to match disktimeout. See this link for additional information.

## Site failure

The result of a storage system or site failure is nearly identical to the result of losing the replication link. The surviving site should experience a roughly 15 second IO pause on writes. Once that 15 second period elapses, IO will resume on that site as usual.

If only the storage system was affected, the Oracle RAC node on the failed site will lose storage services and enter the same 200 second disktimeout countdown before eviction and subsequent reboot.

```
2024-09-11 13:44:38.613 [ONMD(3629)]CRS-1615: No I/O has completed after
50% of the maximum interval. If this persists, voting file
/dev/mapper/grid2 will be considered not functional in 99750 milliseconds.
2024-09-11 13:44:51.202 [ORAAGENT(5437)]CRS-5011: Check of resource "NTAP"
failed: details at "(:CLSN00007:)" in
"/gridbase/diag/crs/jfs13/crs/trace/crsd_oraagent_oracle.trc"
2024-09-11 13:44:51.798 [ORAAGENT(75914)]CRS-8500: Oracle Clusterware
ORAAGENT process is starting with operating system process ID 75914
2024-09-11 13:45:28.626 [ONMD(3629)]CRS-1614: No I/O has completed after
75% of the maximum interval. If this persists, voting file
/dev/mapper/grid2 will be considered not functional in 49730 milliseconds.
2024-09-11 13:45:33.339 [ORAAGENT(76328)]CRS-8500: Oracle Clusterware
ORAAGENT process is starting with operating system process ID 76328
2024-09-11 13:45:58.629 [ONMD(3629)]CRS-1613: No I/O has completed after
90% of the maximum interval. If this persists, voting file
/dev/mapper/grid2 will be considered not functional in 19730 milliseconds.
2024-09-11 13:46:18.630 [ONMD(3629)]CRS-1604: CSSD voting file is offline:
/dev/mapper/grid2; details at (:CSSNM00058:) in
/gridbase/diag/crs/jfs13/crs/trace/onmd.trc.
2024-09-11 13:46:18.631 [ONMD(3629)]CRS-1606: The number of voting files
available, 0, is less than the minimum number of voting files required, 1,
resulting in CSSD termination to ensure data integrity; details at
(:CSSNM00018:) in /gridbase/diag/crs/jfs13/crs/trace/onmd.trc
2024-09-11 13:46:18.638 [ONMD(3629)]CRS-1699: The CSS daemon is
terminating due to a fatal error from thread:
clssnmvDiskPingMonitorThread; Details at (:CSSSC00012:) in
/gridbase/diag/crs/jfs13/crs/trace/onmd.trc
2024-09-11 13:46:18.651 [OCSSD(3631)]CRS-1652: Starting clean up of CRSD
resources.
```

The SAN path state on the RAC node that has lost storage services looks like this:

```
oradata7 (3600a0980383041334a3f55676c697347) dm-20 NETAPP,LUN C-Mode
size=128G features='3 queue_if_no_path pg_init_retries 50' hwhandler='1
alua' wp=rw
|-+- policy='service-time 0' prio=0 status=enabled
| `- 34:0:0:18 sdam 66:96  failed faulty running
`-+- policy='service-time 0' prio=0 status=enabled
  `- 33:0:0:18 sdaj 66:48  failed faulty running
```

The linux host detected the loss of the paths much quicker than 200 seconds, but from a database perspective the client connections to the host on the failed site will still be frozen for 200 seconds under the default Oracle RAC settings. Full database operations will only resume after the eviction is completed.

Meanwhile, the Oracle RAC node on the opposite site will record the loss of the other RAC node. It otherwise continues to operate as usual.

```
2024-09-11 13:46:34.152 [ONMD(3547)]CRS-1612: Network communication with
node jfs13 (2) has been missing for 50% of the timeout interval.  If this
persists, removal of this node from cluster will occur in 14.020 seconds
2024-09-11 13:46:41.154 [ONMD(3547)]CRS-1611: Network communication with
node jfs13 (2) has been missing for 75% of the timeout interval.  If this
persists, removal of this node from cluster will occur in 7.010 seconds
2024-09-11 13:46:46.155 [ONMD(3547)]CRS-1610: Network communication with
node jfs13 (2) has been missing for 90% of the timeout interval.  If this
persists, removal of this node from cluster will occur in 2.010 seconds
2024-09-11 13:46:46.470 [OHASD(1705)]CRS-8011: reboot advisory message
from host: jfs13, component: cssmonit, with time stamp: L-2024-09-11-
13:46:46.404
2024-09-11 13:46:46.471 [OHASD(1705)]CRS-8013: reboot advisory message
text: At this point node has lost voting file majority access and
oracssdmonitor is rebooting the node due to unknown reason as it did not
receive local heartbeats for 28180 ms amount of time
2024-09-11 13:46:48.173 [ONMD(3547)]CRS-1632: Node jfs13 is being removed
from the cluster in cluster incarnation 621516934
```

## Mediator failure

The mediator service does not directly control storage operations. It functions as an alternate control path between clusters. It exists primarily to automate failover without the risk of a split-brain scenario.

In normal operation, each cluster is replicating changes to its partner, and each cluster therefore can verify that the partner cluster is online and serving data. If the replication link failed, replication would cease.

The reason a mediator is required for safe automated operations is because it would otherwise be impossible for a storage clusters to be able to determine whether loss of bidirectional communication was the result of a network outage or actual storage failure.

The mediator provides an alternate path for each cluster to verify the health of its partner. The scenarios are as follows:

- If a cluster can contact its partner directly, replication services are operational. No action required.
- If a preferred site cannot contact its partner directly or via the mediator, it will assume the partner is either actually unavailable or was isolated and has taken its LUN paths offline. The preferred site will then proceed to release the RPO=0 state and continue processing both read and write IO.
- If a non-preferred site cannot contact its partner directly, but can contact it via the mediator, it will take its paths offline and await the return of the replication connection.
- If a non-preferred site cannot contact its partner directly or via an operational mediator, it will assume the partner is either actually unavailable or was isolated and has taken its LUN paths offline. The non-preferred site will then proceed to release the RPO=0 state and continue processing both read and write IO. It will assume the role of the replication source and will become the new preferred site.

If the mediator is wholly unavailable:

- Failure of replication services for any reason will result in the preferred site releasing the RPO=0 state and resuming read and write IO processing. The non-preferred site will take its paths offline.
- Failure of the preferred site will result in an outage because the non-preferred site will be unable to verify that the opposite site is truly offline and therefore it would not be safe for the nonpreferred site to resume services.

## Service restoration

SnapMirror is self-healing. SnapMirror active sync will automatically detect the presence of a faulty replication relationship and bring it back to an RPO=0 state. Once synchronous replication is reestablished, the paths will come online again.

In many cases, clustered applications will automatically detect the return of failed paths, and those applications will also come back online. In other cases, a host-level SAN scan may be required, or applications may need to be brought back online manually.

It depends on the application and how it's configured, and in general such tasks can be easily automated. SnapMirror active sync itself is self-fixing and should not require any user intervention to resume RPO=0 storage operations once power and connectivity is restored.

## Manual failover

The term "failover" does not refer to the direction of replication with SnapMirror active sync because it is a bidirectional replication technology. Instead, 'failover' refers to which storage system will be the preferred site in the event of failure.

For example, you may want to perform a failover to change the preferred site before you shut down a site for maintenance, or before performing a DR test.

Changing the preferred site requires a simple operation. IO will pause for a second or two as authority over replication behavior switches between clusters, but IO is otherwise unaffected.

GUI example:

# Relationships

Local destinations    **Local sources**

Q Search     ↓ Download     👁 Show/hide ⌄     ≡ Filter

| | Source | | Destination | Policy type |
|---|---|---|---|---|
| ⌄ | jfs_as1:/cg/jfsAA | ⋮ | jfs_as2:/cg/jfsAA | Synchronous |
| | Edit | | | |
| | Update | | | |
| | Delete | | | |
| | Failover | | | |

Example of changing it back via the CLI:

```
Cluster2::> snapmirror failover start -destination-path jfs_as2:/cg/jfsAA
[Job 9575] Job is queued: SnapMirror failover for destination
"jfs_as2:/cg/jfsAA                   ".

Cluster2::> snapmirror failover show

Source      Destination                                         Error
Path        Path          Type      Status    start-time end-time Reason
--------    -----------   --------  --------- ---------- ---------- ----------
jfs_as1:/cg/jfsAA
        jfs_as2:/cg/jfsAA
                        planned  completed 9/11/2024  9/11/2024
                                           09:29:22   09:29:32


The new destination path can be verified as follows:


Cluster1::> snapmirror show -destination-path jfs_as1:/cg/jfsAA

                        Source Path: jfs_as2:/cg/jfsAA
                   Destination Path: jfs_as1:/cg/jfsAA
                  Relationship Type: XDP
            Relationship Group Type: consistencygroup
              SnapMirror Policy Type: automated-failover-duplex
                   SnapMirror Policy: AutomatedFailOverDuplex
                         Tries Limit: -
                        Mirror State: Snapmirrored
                 Relationship Status: InSync
```