

Troubleshooting Hyperflex Storage NFS all paths down(APD) issue

Contents

[Introduction](#)

[How are HX Datastores Mounted on ESXI?](#)

[All Paths Down](#)

[Problem Description](#)

[Troubleshooting Workflow](#)

[Checks in vCenter Server:](#)

[Checks in all StCtlVMs:](#)

[StCtlVM: StCtlVM of an Affected ESXi Host](#)

[Checks in ESXi host:](#)

Introduction

This document gives you quick understanding and troubleshooting steps that can be performed in order to assess the source of the problem if you see "NFS all paths down" error message in vCenter to which Hyperflex cluster is integrated with.

How are HX Datastores Mounted on ESXI?

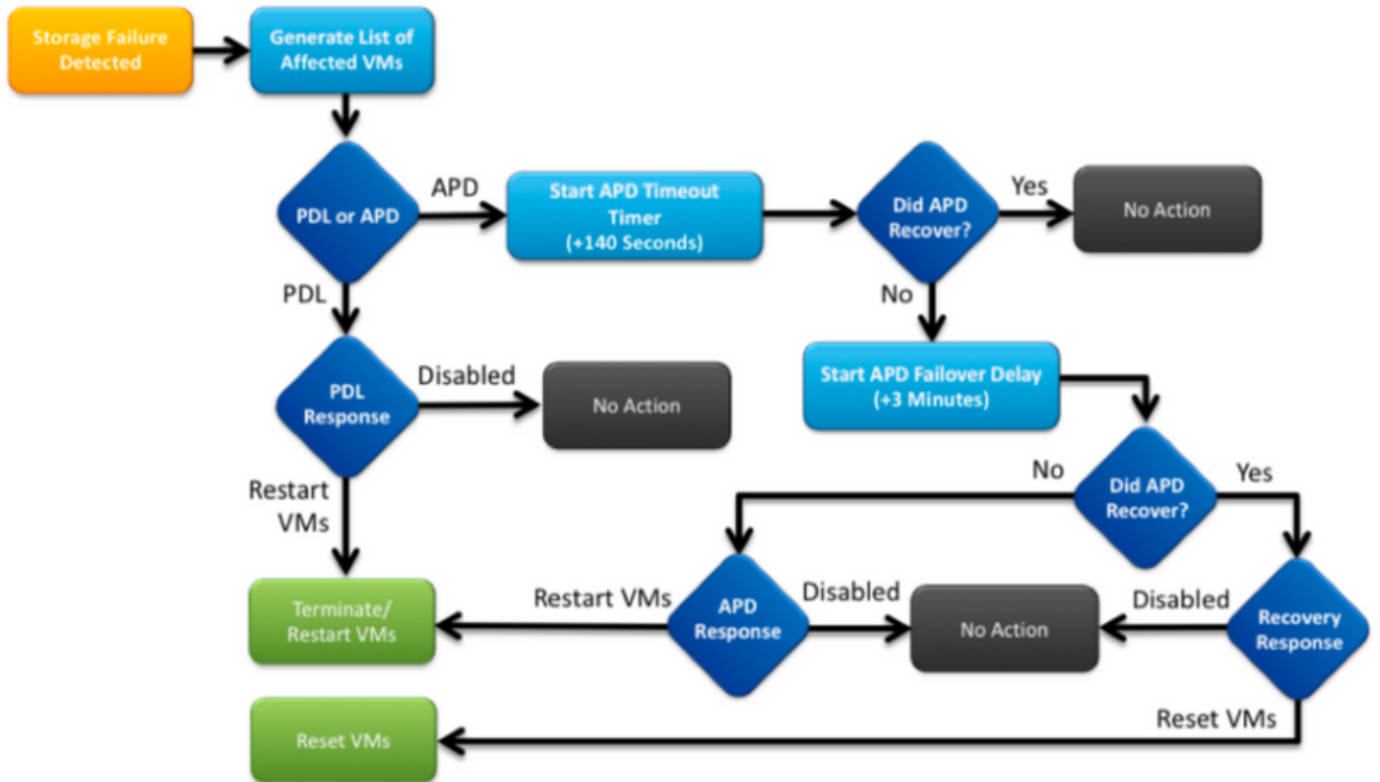
- Hyperflex Datastores are mounted on the ESXi hosts as NFS mounts, in order to mount an NFS datastore we need the NFS Server IP which in our case is the eth1:0 virtual floating interface.
- Hyperflex cluster leverages the use of **virtual floating IP** both for **management** (*eth0:mgmtip*) and **storage data** (*eth1:0*) on which each IP will be assigned to one particular Storage Controller VM (**StCtlVM**). Please note they may end up in different StCtlVMs.
- The importance of this is that the **cluster storage data IP (eth1:0)** is the one used to **mount** the **datastore(s)** created in the **Hyperflex cluster**. Thus it is essential to have it **assigned** and **reachable** from all the **nodes** of the cluster.
- Please note that in case of failure of the StCtlVM that currently owns eth1:0 virtual IP, it should "migrate" to another available StCtlVM working in a similar way as an FHRP (First Hop Redundancy Protocol).

All Paths Down

- APD means that the host cannot reach the storage and there is no Permanent Device Lost (PDL) SCSI code returned from the storage array.
- As it does not know if the loss is temporary or not, it will keep trying to establish communication for more 140s by default (timeout) + 3min (Delay for failover) the ESXi Host begins to fail any non-virtual machine I/O traffic that is being sent to the storage device.
- For more details with regards to APD in vSphere 6.x environment, refer to VMware KB

article [Permanent Device Loss \(PDL\) and All-Paths-Down \(APD\) in vSphere 5.x and 6.x](#)

• Below image explains each intermittent stage:



A typical error message in vCenter will be as follows.

	Status	Name	Defined In
:02.corplex....	Alert	NFS all paths down	SV-VC

Problem Description

Once you see APD alerts on your hosts, obtain the below information to better understand the problem description:

- If one/several/all hosts impacted, and if some which particular hosts impacted
- If any changes were performed previously (configuration/upgrade/etc)
- The timestamp on when the problem first observed and if the issue is recurrent

Troubleshooting Workflow

In order to Troubleshoot APD we need to look into 3 components - vCenter, SCVM, and ESXi host.

These steps are a suggested workflow in order to pinpoint or narrow down the source of the All Paths Down symptom observed. Please note this order does not have to be meticulously followed and you may adequate it as per the particular symptoms observed on the customer environment.

Checks in vCenter Server:

Connect to vCenter Server (VCS) and navigate to an **affected host**

1. **Related Objects -> Virtual Machines** and confirm the StCtlVM is up and running
2. **Related Objects -> Datastores** and confirm if NFS datastores show "**inaccessible**". If datastores seem to be **accessible** and status you may try on **Summary** tab to "Reset to Green" the APD event and later verify if the alert pops back later
3. **Monitor -> Issues** and **Monitor -> Events** should provide information on when the APD was first spotted.

Checks in all StCtlVMs:

Connect to **all** the **StCtlVMs** and verify the below pointers, you may use [MobaXterm](#) software.

1. Verify if all StCtlVMs have the same time using **date** or **ntpq -p**. Time skew on StCtlVM may lead to issues with zookeeper database sync and thus it is paramount to have it in sync among all StCtlVMs. The astrick sign in front of the ntp server denotes that the NTP of your SCVM is synced.

```
root@SpringpathControllerPZTMTRSH7K:~# date
Tue May 28 12:47:27 PDT 2019
```

```
root@SpringpathControllerPZTMTRSH7K:~# ntpq -p -4
remote refid st t when poll reach delay offset jitter
=====
*abcdefghijkl.GNSS. 1 u 429 1024 377 225.813 -1.436 0.176
```

2. If APD occurred during an **upgrade** you might consider to verify which **StCtlVMs** have **not been completely upgraded** and particularly identify the one that last failed. It is possible that it was the one holding the eth1:0 previously Use **dpkg -l | grep -i springpath** to identify the StCtlVMs not completely upgraded as they will have mixed version springpath packages.

```
root@SpringpathControllerPZTMTRSH7K:~# dpkg -l | grep -i springpath
ii storfs-appliance 4.0.1a-33028 amd64 Springpath Appliance
ii storfs-asup 4.0.1a-33028 amd64 Springpath ASUP and SCH
ii storfs-core 4.0.1a-33028 amd64 Springpath Distributed Filesystem
ii storfs-fw 4.0.1a-33028 amd64 Springpath Appliance
ii storfs-mgmt 4.0.1a-33028 amd64 Springpath Management Software
ii storfs-mgmt-cli 4.0.1a-33028 amd64 Springpath Management Software
ii storfs-mgmt-hypervcli 4.0.1a-33028 amd64 Springpath Management Software
ii storfs-mgmt-ui 4.0.1a-33028 amd64 Springpath Management UI Module
ii storfs-mgmt-vcplugin 4.0.1a-33028 amd64 Springpath Management UI and vCenter Plugin
ii storfs-misc 4.0.1a-33028 amd64 Springpath Configuration
ii storfs-pam 4.0.1a-33028 amd64 Springpath PAM related modules
ii storfs-replication-services 4.0.1a-33028 amd64 Springpath Replication Services
ii storfs-restapi 4.0.1a-33028 amd64 Springpath REST Api's
ii storfs-robo 4.0.1a-33028 amd64 Springpath Appliance
ii storfs-support 4.0.1a-33028 amd64 Springpath Support
ii storfs-translations 4.0.1a-33028 amd64 Springpath Translations
```

3. Verify if all relevant services are running **service_status.sh**: Some of the main services are Springpath File System (*storfs*), SCVM Client (*scvmclient*), System Management Service (*stMgr*) or Cluster IP Monitor (*cip-monitor*).

```
root@SpringpathController5L0GTCR8SA:~# service_status.sh
Springpath File System ... Running
SCVM Client ... Running
System Management Service ... Running
```

```

HyperFlex Connect Server ... Running
HyperFlex Platform Agnostic Service ... Running
HyperFlex HyperV Service ... Not Running
HyperFlex Connect WebSocket Server ... Running
Platform Service ... Running
Replication Services ... Running
Data Service ... Running
Cluster IP Monitor ... Running
Replication Cluster IP Monitor ... Running
Single Sign On Manager ... Running
Stats Cache Service ... Running
Stats Aggregator Service ... Running
Stats Listener Service ... Running
Cluster Manager Service ... Running
Self Encrypting Drives Service ... Not Running
Event Listener Service ... Running
HX Device Connector ... Running
Web Server ... Running
Reverse Proxy Server ... Running
Job Scheduler ... Running
DNS and Name Server Service ... Running
Stats Web Server ... Running

```

4. If any of these or other relevant service is not up, start it using **start <serviceName>** eg: **start storfs** You may refer to the service_status.sh script to get the service names . Do a **head -n25 /bin/service_status.sh** and identify the service real name.

```

root@SpringpathController5L0GTCR8SA:~# head -n25 /bin/service_status.sh
#!/bin/bash
declare -a upstart_services=("Springpath File System:storfs"\
"SCVM Client:scvmclient"\
"System Management Service:stMgr"\
"HyperFlex Connect Server:hxmanager"\
"HyperFlex Platform Agnostic Service:hxSvcMgr"\
"HyperFlex HyperV Service:hxHyperVSvcMgr"\
"HyperFlex Connect WebSocket Server:zkupdates"\
"Platform Service:stNodeMgr"\
"Replication Services:replsvc"\
"Data Service:stDataSvcMgr"\
"Cluster IP Monitor:cip-monitor"\
"Replication Cluster IP Monitor:repl-cip-monitor"\
"Single Sign On Manager:stSSOMgr"\
"Stats Cache Service:carbon-cache"\
"Stats Aggregator Service:carbon-aggregator"\
"Stats Listener Service:statsd"\
"Cluster Manager Service:exhibitor"\
"Self Encrypting Drives Service:sedsvc"\
"Event Listener Service:storfsevents"\
"HX Device Connector:hx_device_connector");
declare -a other_services=("Web Server:tomcat8"\
"Reverse Proxy Server:nginx"\
"Job Scheduler:cron"\
"DNS and Name Server Service:resolvconf");

```

5. Identify which **StCtIVM** contains the **storage cluster IP** (eth1:0) using **ifconfig -a** If no **StCtIVM** contains that IP possibly the storfs is not running on one or more nodes.

```

root@help:~# ifconfig
eth0:mgmtip Link encap:Ethernet HWaddr 00:50:56:8b:4c:90
inet addr:10.197.252.83 Bcast:10.197.252.95 Mask:255.255.255.224
UP BROADCAST RUNNING MULTICAST MTU:1500 Metric:1

```

6. Verify if **StCtIVM** is in **contact** with **CRMMaster** and if **zookeeper service** is **up** and **running** **echo srvr | nc localhost 2181** and check if mode is Leader, Follower or Standalone and if connections > 0

```

root@help:~# echo srvr | nc localhost 2181

```

```
Zookeeper version: 3.4.12-d708c3f034468a4da767791110332281e04cf6af, built on 11/19/2018
21:16 GMT
Latency min/avg/max: 0/0/137
Received: 229740587
Sent: 229758548
Connections: 13
Outstanding: 0
Zxid: 0x140000526c
Mode: leader
Node count: 3577
```

service exhibitor status or ps -ef | grep -i exhibitor

```
root@help:~# service exhibitor status
exhibitor start/running, process 12519
root@help:~# ps -ef | grep -i exhibitor
root 9765 9458 0 13:19 pts/14 00:00:00 grep --color=auto -i exhibitor
root 12519 1 0 May19 ? 00:05:49 exhibitor
```

In case of any error or service not running you may verify the below logs and try to start zookeeper service `/var/log/springpath/exhibitor.log` and `/var/log/springpath/stMgr.log`**service exhibitor start** to start zookeeper service

7. Verify if **VC reachable** from all **StCtlVMs** `stcli cluster info | grep -i "url"` to show the URL used containing either FQDN or IP of VC. Verify connectivity to VC using **ping <VC>**

```
root@help:~# stcli cluster info | grep -i "url"
vCenterUrl: https://10.197.252.101
vCenterURL: 10.197.252.101
root@help:~# ping 10.197.252.101
PING 10.197.252.101 (10.197.252.101) 56(84) bytes of data.
64 bytes from 10.197.252.101: icmp_seq=1 ttl=64 time=0.435 ms
```

8. Confirm if **DNS is reachable** in case of cluster using **FQDN stcli services dns show** to list the DNS configured servers on StCtlVM. Test **connectivity** and **resolution** to **DNS** servers using **ping <DNS_IP>** and **host <FQDN> <DNS_IP>**

```
root@help:~# stcli services dns show
1.1.128.140
root@help:~# ping 1.1.128.140
PING 1.1.128.140 (1.1.128.140) 56(84) bytes of data.
64 bytes from 1.1.128.140: icmp_seq=1 ttl=244 time=1.82 ms
```

9. Confirm if all **StCtlVMs** have the same amount of **iptables** entries: **iptables -L | wc -l**. In case they mismatch, please open a TAC case.

```
root@SpringpathControllerI51U7U6QZX:~# iptables -L | wc -l
48
```

10. What are the current cluster status and health **stcli cluster info | less** or **stcli cluster info | grep -i "active\|state\|unavailable"** if trying to find what particular nodes appear with storage unavailable.

```
root@SpringpathControllerI51U7U6QZX:~# stcli cluster info | grep -i
"active\|state\|unavailable"
locale: English (United States)
state: online
upgradeState: ok
healthState: healthy
state: online
state: 1
activeNodes: 3
state: online
```

stcli cluster storage-summary --detail to obtain the storage cluster details

```
root@SpringpathControllerI51U7U6QZX:~# stcli cluster storage-summary --detail
address: 10.197.252.106
name: HX-Demo
state: online
uptime: 185 days 12 hours 48 minutes 42 seconds
activeNodes: 3 of 3
```

```

compressionSavings: 85.45%
deduplicationSavings: 0.0%
freeCapacity: 4.9T
healingInfo:
inProgress: False
resiliencyDetails:
current ensemble size:3
# of caching failures before cluster shuts down:3
minimum cache copies remaining:3
minimum data copies available for some user data:3
minimum metadata copies available for cluster metadata:3
# of unavailable nodes:0
# of nodes failure tolerable for cluster to be available:1
health state reason:storage cluster is healthy.
# of node failures before cluster shuts down:3
# of node failures before cluster goes into readonly:3
# of persistent devices failures tolerable for cluster to be available:2
# of node failures before cluster goes to enospace warn trying to move the existing
data:na
# of persistent devices failures before cluster shuts down:3
# of persistent devices failures before cluster goes into readonly:3
# of caching failures before cluster goes into readonly:na
# of caching devices failures tolerable for cluster to be available:2
resiliencyInfo:
messages:
Storage cluster is healthy.
state: 1
nodeFailuresTolerable: 1
cachingDeviceFailuresTolerable: 2
persistentDeviceFailuresTolerable: 2
zoneResInfoList: None
spaceStatus: normal
totalCapacity: 5.0T
totalSavings: 85.45%
usedCapacity: 85.3G
zkHealth: online
clusterAccessPolicy: lenient
dataReplicationCompliance: compliant
dataReplicationFactor: 3

```

11. What datastores are mounted and available

```

root@bsv-hxaf220m5-sc-4-3:~# stcli datastore list
-----
virtDatastore:
  status:
    EntityRef(idtype=None, confignum=None, type=6, id='235ea35f-6c85-9448-bec7-
06f03b5adf16', name='bsv-hxaf220m5-hv-4-3.cisco.com'):
      accessible: True
      mounted: True
    EntityRef(idtype=None, confignum=None, type=6, id='d124203c-3d9a-ba40-a229-
4dffbe96ae13', name='bsv-hxaf220m5-hv-4-2.cisco.com'):
      accessible: True
      mounted: True
    EntityRef(idtype=None, confignum=None, type=6, id='e85f1980-b3c7-a440-9f1e-
20d7a1110ae6', name='bsv-hxaf220m5-hv-4-1.cisco.com'):
      accessible: True
      mounted: True

```

12. In case stcli commands take too long or fail you may try the following sysmtool commands(Don't use if stcli works) ***sysmtool --ns cluster --cmd info sysmtool --ns cluster --cmd healthdetail sysmtool --ns datastore --cmd list***

StCtIVM: StCtIVM of an Affected ESXi Host

Connect to the **StCtlVM** of the affected **ESXi** host

1. Verify **connectivity** to the **storage cluster IP** (eth1:0) and to **other servers** on the **storage network** (eth1 on **StCtlVMs**)

Run **stcli cluster info | grep -i -B 1 "stctl\hypervisor"** to identify all the ESXi Management IP, StCtlVM eth0 (Mgmt) and StCtlVM eth1 (storage data) respectively participating on the cluster. Test the connectivity **ping -I eth1 [-M do -s 8972] <target IP address>, Jumbo frames test between ESXi VMK1 and SCVM eth1.**

2. If problem still not pinpointed you may have a look into following logs
/var/log/springpath/debug-storfs.log Check if any panics, seg fault or critical events **grep -ai "segmentation\critical\panic" debug-storfs.log/var/log/springpath/stmgr.log** Verify if out of memory problem present **grep -i "oom\|out of mem" /var/log/kern.log**
3. Ultimately you may try to reboot the **StCtlVM** of the **node still experiencing the issue** and verify if the problem persists.

Checks in ESXi host:

Connect to an affected **ESXi** host via **SSH** and perform the following actions:

1. **esxcli storage nfs list** or **Esxcfg-nas -l** to list the currently mounted NFS datastores and if they are accessible

```
[root@bsv-hx220m5-hv-4-3:~] esxcli storage nfs list
```

Volume Name	Host	Share	Accessible
Mounted	Read-Only	isPE	Hardware Acceleration
-----	-----	-----	-----
test	8352040391320713352-8294044827248719091	192.168.4.1:test	true
true	false	false	Supported
sradzevi	8352040391320713352-8294044827248719091	192.168.4.1:sradzevi	true
true	false	false	Supported

```
[root@bsv-hx220m5-hv-4-3:~] esxcfg-nas -l
```

```
test is 192.168.4.1:test from 8352040391320713352-8294044827248719091 mounted available
sradzevi is 192.168.4.1:sradzevi from 8352040391320713352-8294044827248719091 mounted available
```

You may also confirm from **/etc/vmware/esx.conf** to verify the consistency in ESXi configuration on the NFS mounted datastores, using command **cat /etc/vmware/esx.conf | grep -l nas**

2. Verify **/var/log/vmkernel.log** and look for example failed state, mount problems or error around the timestamp identified in previous steps
3. Verify the status of IOVisor/NFS Proxy/**SCVMClient** Check if **service** is running on ESXi using command **/etc/init.d/scvmclient status [Optional]** You may verify if any open connections using **esxcli network ip connection list | grep -i "proto\scvmclient"** Confirm if SCVMClient VIB is the same version as your HX cluster, **Esxcli software vib list | grep -i spring**

```
[root@bsv-hx220m5-hv-4-3:~] esxcli software vib list | grep -i spring
scvmclient          3.5.1a-31118          Springpath
VMwareAccepted     2018-12-13
stHypervisorSvc    3.5.1a-31118          Springpath
VMwareAccepted     2018-12-06
vmware-esx-STFSNasPlugin 1.0.1-21              Springpath
VMwareAccepted     2018-11-16
```

Check **/var/log/scvmclient.log** to see if any errors present namely "unable to obtain clustermap" You may restart SCVMClient service if necessary

through ***etc/init.d/scvmclientrestart***

4. Verify network connectivity with other ESXi hosts on vmk1 network, particularly to storage cluster IP eth1:0***esxcfg-vmknic -l*** to obtain information on the vmk nic details, eg IP, mask and MTU***vmkping -l vmk1 [-v -s 8972] -d <target IP address>*** to test connectivity [optionally with jumbo frames] between ESXi hosts on controller data network
5. **esxcli hardware platform get** to obtain server SN which is used on the name of the StCtlVm and may help you to quickly identify on which host a specific StCtlVM is running.