



CHAPTER 7

Increasing HA in the Data Center

This chapter provides details of Cisco tested high availability solutions in the enterprise data center. It includes the following topics:

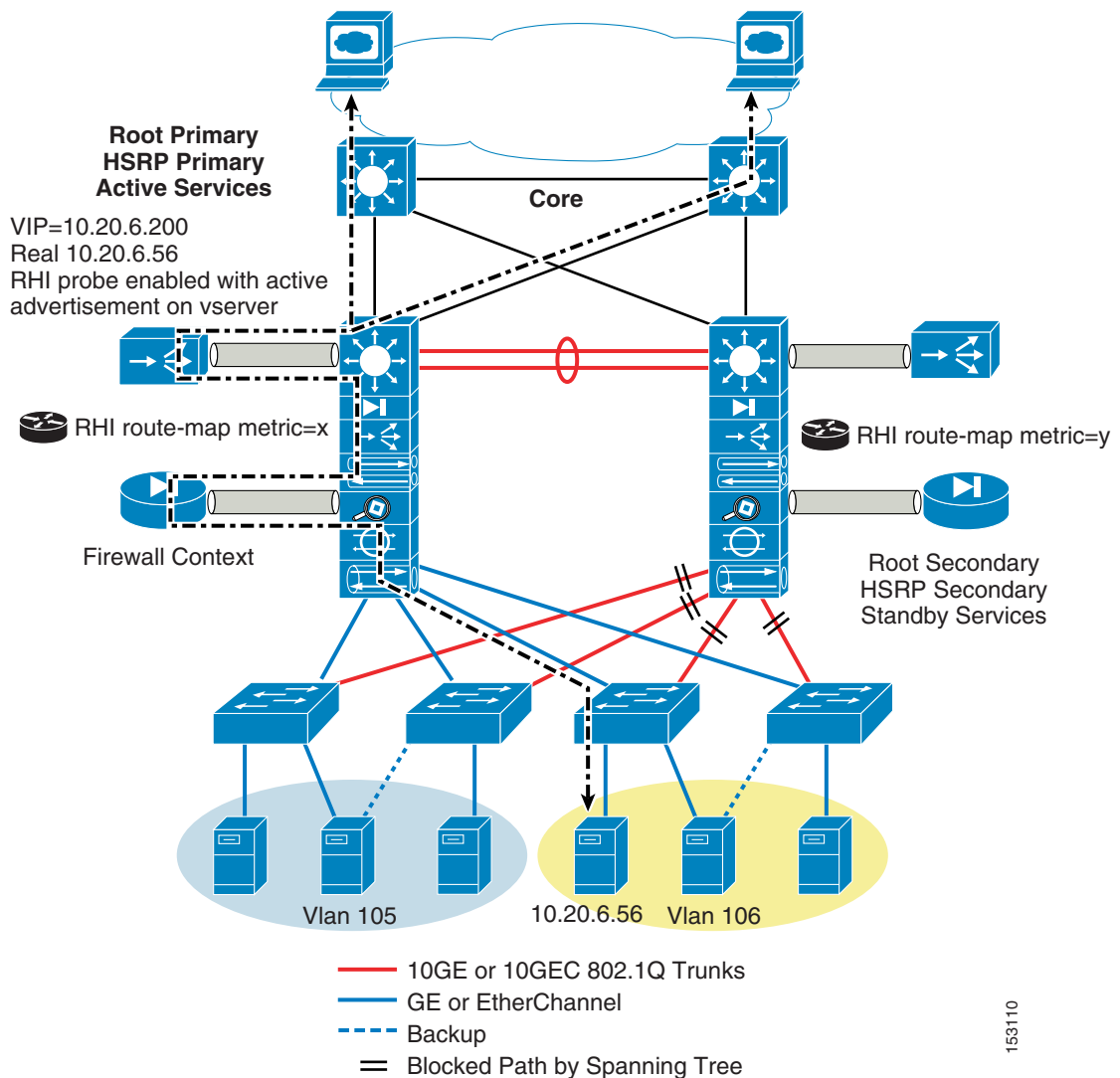
- [Establishing Path Preference with RHI](#)
- [Service Module FT Paths](#)
- [NSF-SSO in the Data Center](#)

Establishing Path Preference with RHI

When active/standby service module pairs are used, it becomes important to align traffic flows such that the active/primary service modules are the preferred path to a particular server application. This is desirable because it creates a design that is more deterministic and easier to troubleshoot but it also becomes particularly important in failure scenarios such as the inter-switch trunk failure described previously.

[Figure 7-1](#) shows an aggregation layer with route preference established toward the primary service modules in an active/standby pair.

Figure 7-1 Route Preference toward Primary Service Modules in an Active/Standby Pair



By using Route Health Injection (RHI) combined with specific route map attributes, a path preference is established with the core so that all sessions to a particular VIP go to agg1 where the primary service modules are located.

The RHI configuration is accomplished by defining a probe type and related values in the Cisco Content Switching Module (CSM) portion of the Cisco IOS configuration. The following probe types are supported:

```

dns          slb dns probe
ftp          slb ftp probe
http         slb http probe
icmp         slb icmp probe
kal-ap-tcp   KAL-AP TCP probe
kal-ap-udp   KAL-AP UDP probe
name         probe with this name
real         SLB probe real suspects information
script       slb script probe
smtp         slb smtp probe
tcp          slb tcp probe

```

```
telnet      slb telnet probe
udp        slb udp probe
```

In the following configuration, a simple ICMP probe is defined, after which it is attached to the server farm configuration. This initiates the probe packets and the monitoring of each real server defined in the server farm.

The last step in configuring RHI is indicating that you want the VIP address to be advertised based on the probe status being operational. If the probe determines that the servers are active and healthy, it inserts a /32 static route for the VIP address into the MSFC configuration.

Aggregation 1 CSM Configuration

```
module ContentSwitchingModule 3
  ft group 1 vlan 102
  priority 20
  heartbeat-time 1
  failover 3
  preempt
!
vlan 44 server
  ip address 10.20.44.42 255.255.255.0
  gateway 10.20.44.1
  alias 10.20.44.44 255.255.255.0
!
probe RHI icmp
interval 3
failed 10
!
serverfarm SERVER200
  nat server
  no nat client
  real 10.20.6.56
  inservice
probe RHI
!
vserver SERVER200
  virtual 10.20.6.200 any
  vlan 44
  serverfarm SERVER200
advertise active
  sticky 10
  replicate csrp sticky
  replicate csrp connection
  persistent rebalance
  inservice
```

With the static host route installed by the CSM into the MSFC based on the health of the server farm, you can now advertise the host route to the core with a route path preference to the active VIP in aggregation 1. This is accomplished with the **redistribute** command in the router process that points to a specific route map. The route map points to an access list that identifies the VIP addresses to match against, and also permits metric attributes to be set to establish path preference. In the following OSPF configuration, **set metric-type type-1** is used to set the host route advertisement to an OSPF external type-1. Alternative methods can include setting the actual OSPF metric value. By setting the host route to an OSPF external type-1, the route appears in the core with the actual accumulated cost of the path used. This approach could prove to be more desirable than attempting to set specific values because it better reflects actual path costs in the case of link failures.

**Note**

The VIP server subnet itself should not be included in the router network statements. If 10.20.6.0 were to be advertised, it would be the next most exact route if the VIP host route, 10.20.6.200, were to not be advertised in an actual failure situation. This would defeat the purpose of advertising the VIP as healthy, allowing sessions to continue to be directed to the agg1 switch.

Aggregation 1 OSPF and Route Map Configurations

```
router ospf 10
  log-adjacency-changes
  auto-cost reference-bandwidth 10000
  nsf
  area 10 authentication message-digest
  area 10 nssa
  timers throttle spf 1000 1000 1000
  redistribute static subnets route-map rhi
  passive-interface default
  no passive-interface Vlan3
  no passive-interface TenGigabitEthernet7/2
  no passive-interface TenGigabitEthernet7/3
  network 10.10.20.0 0.0.0.255 area 10
  network 10.10.40.0 0.0.0.255 area 10
  network 10.10.110.0 0.0.0.255 area 10
  (note: server subnet 10.20.6.0 is not advertised)
  access-list 44 permit 10.20.6.200 log
  route-map rhi permit 10
  match ip address 44
  set metric-type type-1
  set metric +(value) (on Agg2)
```

Aggregation Inter-switch Link Configuration

The design in [Figure 7-1](#) uses VLAN 3 between the aggregation switches to establish a Layer 3 OSPF peering between them. This VLAN traverses a 10GE-802.1Q trunk. With VLANs that cross a 10GE trunk, OSPF sets the bandwidth value equal to a GE interface, not a 10GE interface as one might expect. If the bandwidth on the VLAN configuration is not adjusted to reflect an equal value to the aggregation-core 10GE links, the route from agg2 to the active VIP appears better via the core instead of via VLAN 3 on the inter-switch trunk. At first, this might not appear to be a real problem because the core is going to use the preferred paths directly to agg1 anyway. However, certain link failures can create a scenario where sessions come through the agg2 switch, which would then need to be routed to agg1, so it makes sense to keep the optimal path via the inter-switch link rather than hopping around unnecessarily back to the core. The following configuration reflects the bandwidth changes to show this:

```
interface Vlan3
  description AGG1_to_AGG2_L3-RP
  bandwidth 10000000
  ip address 10.10.110.1 255.255.255.0
  no ip redirects
  no ip proxy-arp
  ip pim sparse-dense-mode
  ip ospf authentication message-digest
  ip ospf message-digest-key 1 md5 C1sC0!
  ip ospf network point-to-point
  ip ospf hello-interval 1
  ip ospf dead-interval 3
  logging event link-status
```

Aggregation 2 Route Map Configuration

The aggregation 2 switch requires the same configurations as those outlined for aggregation 1. The only difference is with the route map configuration for RHI. Because you want path preference to be toward the active VIP in the Aggregation 1 switch, you need to adjust the metric for the advertised RHI route to be less favorable in agg2. The reason this is necessary is to make sure that in an active-active service module scenario, there is symmetry in the connections to prevent asymmetrical flows that would break through the CSM and Cisco Firewall Service Module (FWSM). The following route map adds cost to the OSPF route advertised by using the `set metric +` command. Note that the `type-1` is also set for the same reasons mentioned previously.

```
route-map rhi permit 10
  match ip address 44
  set metric +30
  set metric-type type-1
```

The following command can be used to view the status of an RHI probe configuration:

```
Aggregation-1#sh module contentswitchingModule 3 probe detail
```

probe	type	port	interval	retries	failed	open	receive
RHI	icmp	3	3	3	10		10
real		vserver		serverfarm		policy	status
10.20.6.56:0		SERVER200		SERVER200		(default)	OPERABLE
10.20.6.25:0		SERVER201		SERVER201		(default)	OPERABLE

```
Aggregation-1#
```

Service Module FT Paths

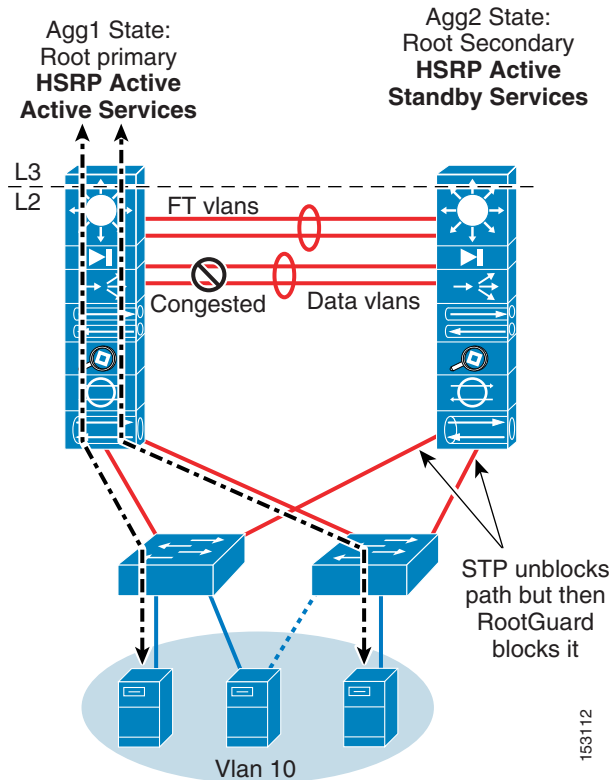
Service module redundant pairs monitor each other to ensure availability as well as to maintain state for all sessions that are currently active. The availability is provided by a hello protocol that is exchanged between the modules across a VLAN. Session state is provided by the active service module replicating the packet headers to the standby across the same VLAN as the hellos, as in the case of the CSM, or on a separate VLAN as in the case of the FWSM.

If the availability path between the redundant pairs becomes heavily congested or misconfigured, it is likely that the service modules believe the other has failed. This can create a split-brain scenario where both service modules move into an active state, which creates undesirable conditions including asymmetric connection attempts.

The CSM exchanges hello and session state on a common VLAN. The FWSM uses separate VLANs for both hello exchange and state replication. The standby FWSM assumes the active role when it no longer sees its redundant peer on at least two VLAN interfaces. This VLAN could be a combination of the context, failover, or state VLANs.

If a second inter-switch link were added with only service module FT vlans provisioned across it as shown in [Figure 7-2](#), the chance of a split-brain scenario because of these conditions is reduced.

Figure 7-2 Congestion on FT Path



The bandwidth required for the FT link must be considered. The maximum possible required can be equal to the CSM bus interface (4G) or the FWSM bus interface (6G) as a worst case. The required bandwidth is based on the amount and type of sessions that are being replicated.

**Note**

More detail on access layer design is covered in [Chapter 6, “Data Center Access Layer Design.”](#)

NSF-SSO in the Data Center

**Note**

The testing performed in support of this section included the use of CSM one-arm mode and FWSM transparent mode design, which influences the behavior of NSF/SSOs failover characteristics.

The data center solutions that are covered in this and other data center guides are designed for a high level of resiliency. For example, the core and aggregation layer switches are always in groups of two, and are interconnected such that no individual module or full system failure can bring the network down. The service modules and other software and hardware components are configured in redundant pairs that are located in each of the aggregation switches to further remove any single point of failure. The access layer also has many options that permit dual homing to the aggregation layer and leverage spanning tree, FlexLinks, and NIC teaming to achieve high availability.

The main objective in building a highly available data center network design is to avoid TCP session breakage while providing convergence that is unnoticeable, or as fast as possible. Each of the TCP/IP stacks that are built into the various operating systems have a different level of tolerance for determining when TCP will break a session. The least tolerant are the Windows Server and Windows XP client stacks, which have been determined to have a ~9 second tolerance. Other TCP/IP stacks such as those found in Linux, HP, and IBM are more tolerant and have a longer window before tearing down a TCP session.

This does not necessarily mean that the data center network should be designed to converge in less than 9 seconds, but it could serve as a guideline to a worst case. The most optimal acceptable failure convergence time is zero, of course. Although each network has its own particular convergence time requirements, many network designers usually break acceptable convergence times into the following categories:

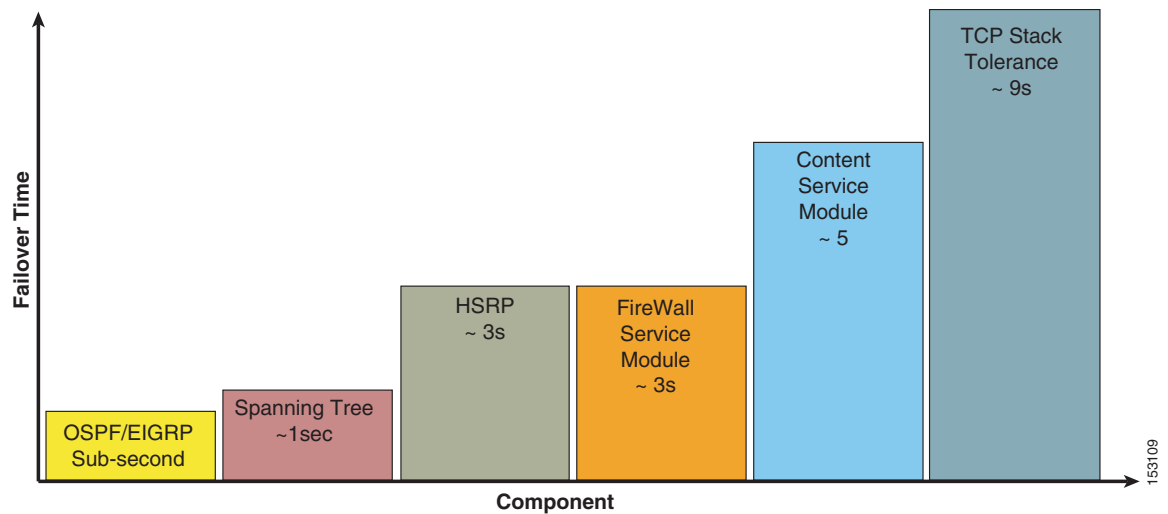
- **Minor failures**—Failures that would be expected to happen because of more common events, such as configuration errors or link outages because of cable pulls or GBIC/Xenpak failure, might be considered to be a minor failure. A minor failure convergence time is usually expected to be in the sub-second to 1 second range.
- **Major failures**—Any failure that can affect a large number of users or applications is considered a major failure. This could be because of a power loss, supervisor, or module failure. This type of failure usually has a longer convergence time and is usually expected to be under 3–5 seconds.

The following describes the various recovery times of the components in the Cisco data center design:

- **HSRP**—With recommended timers of Hello=1/Holddown=3, convergence occurs in under 3 seconds. This can be adjusted down to sub-second values, but CPU load must be considered.
- **Routing protocols**—OSPF and EIGRP can both achieve sub-second convergence time with recommended timer configurations.
- **Rapid PVST+ Spanning Tree**—802.1W permits sub-second convergence time for minor failures when logical ports are under watermarks, and usually 1–2 seconds for major failure conditions.
- **CSM**—Convergence time is ~5 seconds with recommended timers.
- **FWSM**—Convergence time is ~3 seconds with recommended timers.

Figure 7-3 shows the various failover times of the components of the Cisco data center design.

Figure 7-3 Failover Times



The worst case convergence time for an individual component failure condition is the CSM at ~5 seconds. In the event of a Supervisor720 failure on the Aggregation 1 switch, all of these components must converge to the Aggregation 2 switch, resulting in a minimum of ~5 second convergence time. This convergence time will most likely be more because of the tables that have to be rebuilt (ARP, CAM tables, and so on), so maximum convergence time can approach the 9 second limit of the Windows TCP/IP stack. This convergence time and possible lost sessions can be avoided by using dual Sup720s with NSF-SSO on the primary aggregation switch of the data center.

Supervisor 720 supports a feature called Non-Stop Forwarding with Stateful Switch-Over (NSF-SSO) that can dramatically improve the convergence time in a Sup720 failure condition. NSF with SSO is a supervisor redundancy mechanism on the Supervisor Engine 720 in Cisco IOS Release 12.2(18)SXD that provides intra-chassis stateful switchover. This technology demonstrates extremely fast supervisor switchover with lab tests resulting in approximately 1.6–2 seconds of packet loss.

The recommended data center design that uses service modules has a minimum convergence time of ~6–7 seconds primarily because of service modules. With NSF/SSO, the service modules do not converge. This alone represents a large reduction in convergence time, making dual supervisors with NSF-SSO a tool for achieving increased high availability in the data center network.

Possible Implications

HSRP

With the current 12.2.(18) release train, HSRP state is not maintained by NSF-SSO. The HSRP static MAC address is statefully maintained between supervisors but the state of HSRP between aggregation nodes is not. This means that if the primary Sup720 fails, the SSO-enabled secondary Sup720 takes over and continues to forward traffic that is sent to the HSRP MAC address, but the HSRP hellos that are normally communicated to the standby HSRP instance on agg2 are not communicated. This means that during a switchover on aggregation 1, the aggregation 2 switch HSRP instances take over as primary during the SSO control plane recovery.

Because the HSRP MAC address was statefully maintained on the agg1 standby Sup720 module, the sessions continue to flow through agg1, regardless of the active state that appears on agg2. After the control plane comes up on the agg1 switch, the HSRP hello messages begin to flow, and preemptively move the active state back to the agg1 switch. The control plane recovery time is ~2 minutes.

With looped access layer topologies (triangle and square) that align HSRP, STP primary root, and active service modules on the agg1 switch, this does not create an issue because the access layer to aggregation layer traffic flow continues to be directed to the agg1 switch. If a loop-free access layer design is used, the active HSRP default gateway instance on agg2 responds to ARP requests.

**Note**

A square looped access also has active-active uplinks, but the FWSM transparent mode active context on agg1 prevents packets from reaching the active HSRP default gateway on agg2. The VLAN on the south side of the context follows the spanning tree path from agg2, across the inter-switch link to the active FWSM on agg1, then to the HSRP default gateway instance on agg1.

IGP Timers

It is possible that IGP timers can be tuned low enough such that NSF/SSO is defeated because the failure is detected by adjacent nodes before it is determined to be an SSO stateful switchover.

Slot Usage versus Improved HA

Placing two Sup720s in a data center aggregation node also has its drawbacks in terms of available slot density. Particularly with 10GE port density challenges, using an available slot that could be available for other purposes might not be a necessary trade-off. Consideration of using dual Sup720s should be based on actual customer requirements. The recommendation is to use the dual supervisor NSF-SSO solution in the primary aggregation node with service modules when slot density is not an issue or when HA is critical at this level. If a Network Analysis Module is used, it can be placed in the agg2 switch while the dual Sup720s are in the agg1 switch, balancing the slot usage across the two aggregation nodes.

Recommendations

NSF-SSO can provide a very robust solution to data centers that require a very high level of resiliency and are willing to use available slots to achieve it. The processes and functions that are involved for NSF/SSO to work correctly are very complex and have many inter-related dependencies. These dependencies involve the access layer design and service module modes used, for example. This guide provides a solution that permits NSF/SSO to provide a lower convergence time than the base recommendation design can provide based on service module failover times. Cisco recommends testing NSF/SSO to ensure that it works as expected in a specific customer design.

