

Pass-Through Technology

This chapter provides best design practices for deploying blade servers using pass-through technology within the Cisco Data Center Networking Architecture, describes blade server architecture, and explores various methods of deployment. It includes the following sections:

- Blade Servers and Pass-Through Technology
- Design Goals
- Design and Implementation Details
- Configuration Details

Blade Servers and Pass-Through Technology

A blade server is an independent server that includes an operating system, memory, one or more processors, network controllers, and optional local storage. Blades servers are designed to reduce the space, power, and cooling requirements within the data center by providing these services within a single chassis. Blade server systems are a key component of data center consolidation that help reduce costs and provide a platform for improving virtualization, automation, and provisioning capabilities.

A primary consideration in any blade server deployment is how the blade server system is connected to the data center network. There are several I/O options for blade server systems, including the following:

- Built-in Ethernet switches (such as the Cisco Ethernet Switch Modules)
- Infiniband switches (such as the Cisco Server Fabric Switch)
- Fibre Channel switches
- Blade Server Chassis Pass-through Modules

Each of these I/O technologies provides a means of network connectivity and consolidation.

This chapter focuses on integrating blade server systems within the Cisco Data Center Architecture using pass-through technology, which allows individual blade servers to communicate directly with resources external to the blade chassis. Both copper and optical pass-through modules are available that provide access to the blade server controllers. It is therefore important to understand the internal connectivity provided by the blade server chassis before discussing the external ramifications of pass-through deployments.

Currently, there is no industry-standard design for blade servers or blade server enclosures. Various blade system architectures are available from various vendors. The following section describes two generic blade server systems, which illustrate many of the design features found in these various architectures:

- System A—Uses octopus cabling to interconnect the blade servers with the data center architecture.
- System B—Passes the blade server signaling to the external network port-to-port.

System A in Figure 3-1 illustrates the front and rear view of a typical blade server chassis.

Figure 3-1 Example Blade Server Architecture – System A



System A is seven rack units (RUs) in height and provides 14 slots to house individual blade servers. The rear of the chassis allows for four individual I/O modules for network connectivity and two management modules to administer the blade system. The blade servers and I/O modules communicate over a midplane, as shown in Figure 3-2.



Figure 3-2 System A Internal Blade Server Connectivity

Each network interface controller (NIC) has a dedicated channel on the midplane connecting it to a specific I/O module bay. Typically, the integrated NICs are Gigabit Ethernet by default, while the I/O expansion card supports host bus adapters (HBA), host channel adapters (HCA), or Gigabit Ethernet NICs.

<u>Note</u>

Note that with this architecture, the I/O expansion card on each blade server must be compatible with the I/O modules installed in Bay 3 and Bay 4.

System B (see Figure 3-3) illustrates another common blade server architecture.



Figure 3-3 Example Blade Server Architecture – System B

This six-RU blade server chassis has ten slots; two for I/O modules, and eight dedicated for server use. The blade server chassis provides four dedicated backplane channels to each server slot. Figure 3-4 shows this design, which supports a total of 32 independent channels for the eight blade server slots.



Figure 3-4 System B Blade Server Connectivity

The architectures represented by System A and System B use different types of pass-through module technology, as shown in Figure 3-5.

Figure 3-5 Pass-Through Module Examples



The pass-through module technology used in System A depends on octopus cables to connect to the external network. This octopus cable allows multiple servers to be supported by a single output cable that connects to the external network with transmit and receive pairs dedicated to each blade server controller.

The architecture represented by System B does not use any physical cabling consolidation. Instead, it simply passes the blade server signaling to the external network port-to-port.

Both systems provide redundant dedicated connections to the I/O modules over the midplane or backplane. By default, each blade server is dual- or multi-homed to the I/O modules deployed on the blade system. This physical design provides an increased level of availability for the services deployed on the blade servers. In addition, the redundant I/O modules can be used for establishing redundant connections to the external network, which provides an even higher level of availability.

Design Goals

This section describes the key design goals when deploying blade servers with pass-through technology in data centers. It includes the following topics:

- High Availability
- Pass-through technology is a flexible solution that provides blade server high availability by supporting all three NIC teaming modes of operation.
- Manageability

High Availability

This section describes key issues to consider when making design choices for the overall data center architecture as well as for the blade servers.

Achieving Data Center High Availability

Data centers house the critical business applications of the enterprise, which must be accessible for use either at specific times or continuously, and without interruption. The network infrastructure provides the level of availability required by these applications through device and link redundancy and a deterministic topology. Servers are typically configured with multiple NIC cards and dual-homed to the access layer switches to provide backup connectivity to the business applications.

The implementation of blade servers does not change the high availability requirements of the data center. Implementing blade servers with pass-through technology allows non-disruptive deployment. Pass-through deployments do not alter the fast convergence and deterministic traffic patterns provided by Layer 2 and Layer 3 technologies. The connection established between the external network device and the blade server by the pass-through module is neither switched nor blocked. The modules simply expose the blade server NICs to the network without affecting the Layer 2 or Layer 3 network topologies created through spanning tree or routing protocols. When using pass-through modules, the blade servers function as servers that happen to be located inside a blade server chassis.

Achieving Blade Server High Availability

Blade server enclosures provide high availability to local blade servers through a multi-homed architecture. Each blade server achieves an increased level of accessibility by using NIC teaming software. NIC teaming lets you create a virtual adapter consisting of one to eight physical NIC interfaces, which can typically support up to 64 VLANs. NIC teaming is a high availability mechanism that can provide both failover and local load balancing services. There are the following three primary modes of NIC teaming:

- Fault tolerant mode
- Load balancing mode
- Link aggregation mode

Fault tolerant mode, also known as *active/standby*, creates a virtual NIC by grouping one or more network controllers in a team. One adapter is the primary or active NIC, leaving all other network controllers in the team as secondary or standby interfaces. The standby interface becomes active if the primary NIC fails because of probe or physical link problems. Fault tolerance can be achieved in the following two ways:

- Adapter fault tolerance (AFT)—The NIC team is homed to a single switch.
- Network or switch fault tolerance (NFT/SFT)—The NIC team is dual-homed to two different switches.

Pass-through I/O modules support both configurations. Figure 3-6 illustrates AFT and NFT configurations with blade servers using pass-through.

Figure 3-6 Pass-Through Module with Fault Tolerant NIC Teaming



Blade Server with Pass-thru I/O Modules

The solid yellow links represent the active links, and the dotted grey links represent the standby links that are not being used by the server. NFT provides a greater level of network availability. However, neither configuration optimizes the bandwidth available to the blade server.

The second method of NIC teaming (load balancing mode) builds upon the high availability features of NFT by configuring all the NICs in a team to participate in the transmission of data. This feature lets the server utilize more of the available bandwidth.

In load balancing mode, a single primary NIC receives all incoming traffic while egress traffic is load balanced across the team by a Layer 2- or Layer 3-based algorithm. Pass-through technology permits this configuration, as shown in Figure 3-7.



Figure 3-7 Pass-Through Module with Load Balance NIC Teaming

The solid yellow lines indicate the primary interfaces of the blade server that both receive and transmit traffic. The dotted green lines are standby interfaces that only transmit traffic. A hashing algorithm, usually based on the source and destination IP addresses, determines which NIC is responsible for transmitting the traffic for any given transaction. The standby controller becomes responsible for both ingress and egress traffic only if the primary NIC fails.

The third method of NIC teaming, link aggregation mode (channeling), extends the load balancing functionality of NIC teaming by allowing all interfaces to receive and transmit traffic. This requires the switch to load balance traffic across the ports connected to the server NIC team. This mode of NIC teaming provides link redundancy and the greatest bandwidth utilization. However, the access switch in this design represents a single point of failure. Figure 3-8 shows a channel established between the blade servers and the access switch.





NIC teaming software typically supports the configuration of Gigabit EtherChannel (GEC) or IEEE 802.3ad (LACP) channels.

Pass-through technology is a flexible solution that provides blade server high availability by supporting all three NIC teaming modes of operation.

Scalability

Network designs incorporating blade server devices must ensure network scalability to allow for increases in server density. Scalability allows increases in services or servers without requiring fundamental modifications to the data center infrastructure. The choice of I/O module used by the blade servers (integrated switches or pass-through modules) is a critical deployment factor that influences the overall data center design.



Pass-through modules allow blade servers to connect directly to a traditional access layer. The access layer should provide the port density to support the data center servers and the flexibility to adapt to increased demands for bandwidth or server capacity. For more information on data center scalability, see Design Goals, page 3-5, or the *Cisco Data Center Infrastructure SRND* at the following URL: http://www.cisco.com/en/US/docs/solutions/Enterprise/Data_Center/DC_Infra2_5/DCI_SRND_2_5_b ook.html.

Manageability

Management of the data center, including administration of software versions and hardware configurations, is also a key design consideration. From this point of view, blade server pass-through modules reduce the administrative complexity of the data center. Pass-through modules do not require configuration, which eliminates configuration errors on the devices and reduces the need for configuration backups. The blade server chassis may also provide limited diagnostic information and the ability to enable or disable external ports on the module. The availability of these features and the level of diagnostic information depends on the manufacturer of the blade system.



Serial Over LAN (SOL) is a blade server management feature available on the IBM BladeCenter chassis. However, SOL requires the use of an integrated switch and is not currently available with the IBM pass-through modules. SOL leverages the trunk that exists between the management module and the Ethernet blade switch to allow console access to the individual blade servers.

Design and Implementation Details

The following section discusses the following two access switch attachment options available when using blade server pass-through I/O modules:

- Modular Access Switches
- One Rack Unit Access Switches

These network designs emphasize high availability in the data center by eliminating any single point of failure, by providing deterministic traffic patterns, and through predictable network convergence behavior. The configuration examples use Cisco Catalyst 6500s as the aggregation layer platform. This Layer 2/Layer 3 switching platform supports high slot density and integrated network services, which are important features for data centers deploying blade systems.

The introduction of pass-through modules into an existing Layer 2/Layer 3 design does not require much modification to the existing data center architecture, which allows blade server systems to be easily inserted into the network. However, the disadvantage is that cable consolidation and use of shared interconnects, which are important benefits that can be provided by blade systems, are not fully realized.

Modular Access Switches

Figure 3-9 illustrates the use of a modular access switch with pass-through I/O modules in a blade server system. The blade servers are dual-homed to the access layer switches. The modular access provides port density and 10 Gigabit Ethernet uplinks to the aggregation layer where intelligent services such as security, load balancing, and network analysis reside.

The advantages of this design include the following:

- Proven high availability design
- Scalable approach to blade server deployments
- Improved operational support

Figure 3-9



Modular Access Design with Pass-Through I/O Modules

is no single point of failure in this network topology. The access layer switches are dual-homed to the aggregation layer switches, which provide redundant network paths. Spanning tree manages the physical loops created by the uplinks between the aggregation and access switches, assuring a predictable and fast converging topology.

The design in Figure 3-9 uses a classic triangular topology with the modular access layer switches. There

In Figure 3-9, the solid black lines represent uplinks in a spanning tree forwarding state, while the dotted red lines represent uplinks in blocking state. Rapid Per VLAN Spanning Tree Plus (RPVST+) is recommended for this design because of its high availability features. RPVST+ provides fast convergence (less than 1 second) in device or uplink failure scenarios. In addition, RPVST+ offers enhanced Layer 2 features for the access layer with integrated capabilities equivalent to the UplinkFast and BackboneFast features in the previous version of Spanning Tree Protocol (STP).

<u>Note</u>

For more information on Layer 2 design and RPVST+, see the *Data Center Infrastructure SRND* at the following URL:

http://www.cisco.com/en/US/docs/solutions/Enterprise/Data_Center/DC_Infra2_5/DCI_SRND_2_5_b ook.html.

In addition to the high availability attributes of this proven access layer design, the modular access switch provides a very high level of port density. This allows the server farm to scale as it addresses future data center needs without an exponential increase in administrative overhead (see Figure 3-10).



Figure 3-10 Scalability with a Modular Access Switch Design

Access Layer Modules

In Figure 3-10, ten modular access switches are connected to a pair of aggregation layer switches. Each pair of access switches supports 12 fully populated blade server systems housed in a set of three racks. In this example, each blade system requires 32 ports on the access switch for the pass-through links from the blade servers.

With four blade systems per rack, 128 access layer ports are required to support a single rack. Dual-homing the blade servers to each modular access switch means that each access layer switch must provide 64 ports per rack or 192 ports for three racks. Four 48-port line-cards (192/48 = 4) are required on each modular access switch to support this configuration.



This example does not consider the scalability of the spanning tree, which depends on the number of active logical interfaces and virtual ports supported per line card. In addition, the acceptable oversubscription ratio for the applications must be taken into account. For more information on scaling the Layer 2 and Layer 3 topologies in the data center, see the *Data Center Infrastructure SRND* at the following URL:

http://www.cisco.com/en/US/docs/solutions/Enterprise/Data_Center/DC_Infra2_5/DCI_SRND_2_5_b ook.html.

A modular access switch design reduces the total number of switches in the network. In the previous example, 120 integrated blade switches would be required, assuming two blade switches per chassis, to support an equivalent number of blade servers. The ramifications of introducing this number of devices into the data center network are obvious. Specifically, the Layer 2 and Layer 3 topologies expand and become more difficult to manage. In addition, there is an increase in the network administration required to support the integrated switch. Using a modular access layer switch reduces these operational and logical complexities.

However, pass-through technology does not provide the benefits of cable consolidation and the use of shared interconnects provided by integrated blade switches. Pass-through modules do not reduce the I/O cabling volume within the rack or lessen the cable bulk from the rack to the modular access switches. Cabling represents an obstruction that restricts the airflow within the data center and may adversely affect the temperature of the room. When using pass-through technology and blade servers, the design and use of an effective cable management system within the facility is necessary to mitigate these issues.

One Rack Unit Access Switches

Figure 3-11 illustrates the use of a 1-RU access switch with pass-through I/O modules in a blade server system. The blade servers are dual-homed to the 1-RU access layer switches. The 1-RU access layer switches provide the port density and uplink connectivity to the aggregation layer required by the blade servers.



Figure 3-11 1-RU Access Layer Switch with Pass-Through Technology

Typically, 1-RU access switches are located at the top of the racks housing the blade server units. This design allows the cable density created with the pass-through modules to remain within the rack, which helps contain the potential problems. This is a distinct advantage compared to the modular access layer

model discussed previously. The uplinks from the 1-RU access layer switches provide a common, consolidated connection to the aggregation layer and its services. This essentially reduces the number of cable runs required between the rack and the aggregation layer switches.

This design also provides a highly available and predictable server farm environment. In Figure 3-11, redundant paths for network traffic are managed with a Layer 2 protocol such as RPVST+ that provides sub-second convergence. In Figure 3-11, Aggregation-1 is the primary root switch with all links forwarding. Aggregation-2 is the secondary root and provides an alternative traffic path for application traffic in the event of a failure. The solid black lines represent links that are forwarding and the dotted red lines represent links in a spanning tree blocking state.

The blade servers provide another level of protection by using the high availability features of NIC teaming, in addition to the high availability functions of the network. The sold yellow lines represent each link to the active NIC while the dotted green lines show the links to each interface in a standby state. In this configuration, the NIC teaming software offers sub-second convergence at the server.



It is also important to consider the high availability features available on the 1-RU switch platform when this is deployed in the data center. Redundant power and hot-swappable fans are recommended to improve Layer 1 availability.

The scalability of 1-RU switches is limited compared to the use of modular switches. A 1-RU switch cannot provide the same level of port density for server connectivity as a modular switch. As a result, the number of switching devices in the data center is increased, compared to the solution using modular switches. This in turn increases the spanning tree domain as well as the administrative overhead required to implement and maintain the solution.

For example, Figure 3-12 demonstrates the potential scale of a 1-RU access switch deployment.

As stated previously, the access switches must provide the local port density required by the blade servers. In this instance, each data center rack houses three blade systems providing 32 individual pass-through connections to the internal blade servers. The blade servers are dual-homed over these pass-through connections to a pair of 1-RU access switches located in the rack. Three blade systems with 16 dual-homed servers per chassis require 96 ports. To provide for network fault tolerance, each 1-RU access layer rack switch should supply 48 ports for server connectivity.

In addition, the access layer switches provide connectivity to the aggregation layer services. The modular aggregation layer switch must furnish the uplink port density for the 1-RU access layer. A Catalyst 6509 would suffice in this scenario. A pair of aggregation layer switches can support 12 1-RU access switches that are dual-homed over ten Gigabit Ethernet uplinks. In this example, 288 servers are supported in the 1-RU access switch model with modular aggregation layer support.



This example does not consider the scalability of spanning tree, which depends on the number of active logical interfaces and virtual ports supported per line card. In addition, the acceptable oversubscription ratio for the applications must be taken into account. For more information on scaling the Layer 2 and Layer 3 topologies in the data center, see the *Data Center Infrastructure SRND 2.0* at the following URL: http://www.cisco.com/en/US/docs/solutions/Enterprise/Data_Center/DC_Infra2_5/DCI_SRND_2_5_b ook.html.



Figure 3-12 Scalability with a 1-RU Access Switch Design



The use of three or four 1-RU switches per rack may be necessary depending on the number of ports available on the switching platform and the number of blade server interfaces that must be supported. This affects the scalability of this design by requiring greater uplink port density to connect to the aggregation layer switches.

The 1-RU access switch design reduces the total number of switches in the network. The example shown in Figure 3-12 would require 36 integrated blade switches, assuming two blade switches per chassis, to support an equivalent number of blade servers. The 1-RU access layer switch also reduces the operational and logical complexities (L2/L3 topologies) of the data center when compared to the integrated switch solution. In addition, the design reduces the number of cables required to provide external connectivity to the rack. As previously discussed, the modular access layer switch design requires fewer network devices and topology changes but uses more cabling.

Configuration Details

This section describes the switch configurations necessary to integrate pass-through technology into the Cisco Data Center Architecture. The following configurations are described:

- VLAN Configuration
- RPVST+ Configuration
- Inter-Switch Link Configuration

- Port Channel Configuration
- Trunking Configuration
- Server Port Configuration
- Server Default Gateway Configuration

VLAN Configuration

To configure the VLANs on the switches, complete the following steps:

Step 1 Set the VLAN trunking-protocol administrative domain name and mode as follows:

(config) # vtp domain domain name (config) # vtp mode transparent

Step 2 Configure the server farm VLANs as follows:

(config)# vlan VLAN ID (config-vlan)# name VLAN name (config-vlan)# state active

RPVST+ Configuration

Configure STP to manage the physical loops in the topology. RPVST+ is recommended for its fast convergence characteristics.

Step 1 To set the STP mode an each aggregation switch, enter the following command:

(config)# spanning-tree mode rapid-pvst

The port path cost value represents the media speed of the link and is configurable on a per interface basis, including EtherChannels. To allow for more granular STP calculations, enable the use of a 32-bit value instead of the default 16-bit value. The longer path cost better reflects changes in the speed of channels and allows STP to optimize the network in the presence of loops.

- **Step 2** To configure STP to use 32 bits in port part cost calculations, enter the following command: (config)# spanning-tree pathcost method long
- Step 3 To configure the root switch, enter the following command: (config)# spanning-tree vlan vlan range root primary
 Step 4 To configure the secondary root switch, enter the following command:

(config)# **spanning-tree vlan** vlan range **root secondary**

Inter-Switch Link Configuration

The topologies discussed in this guide require connectivity between the switches. The following two types of inter-switch connections can be used to provide this connectivity:

- Aggregation Switch to Aggregation Switch
- Aggregation Switch to Access Layer Switch

Each of these connections are Layer 2 EtherChannels consisting of multiple physical interfaces bound together as a channel group or port channel. Each of these point-to-point links between the switches is a trunk because they typically carry more than one VLAN.

Port Channel Configuration

Link Aggregate Control Protocol (LACP) is the IEEE standard for creating and managing port channels between switches. Each aggregate switch uses this feature to create a port channel across the line cards. The use of multiple line cards within a single switch reduces the possibility of the point-to-point port channel becoming a single point of failure in the network.

Step 1 Configure the active LACP members on the aggregation switches that connect to each access layer switch as follows:

```
(config)# interface GigabitEthernet12/1
(config-if)# description Connection to Access Layer Switch
(config-if)# channel-protocol lacp
(config-if)# channel-group 1 mode active
(config)# interface GigabitEthernet11/1
(config-if)# description Connected to Access Layer Switch
(config-if)# channel-protocol lacp
(config-if)# channel-group 1 mode active
```

Step 2 Configure the passive LACP members on the access layer switch as follows:

```
(config) # interface GigabitEthernet0/19
(config-if)# description Connected to Aggregation Layer Switch
(config-if)# channel-group 1 mode on
(config) # interface GigabitEthernet0/20
(config-if)# description Connected to Aggregation Layer Switch
(config-if)# channel-group 1 mode on
```

Trunking Configuration

Use the following guidelines when configuring trunks:

- Allow only those VLANs that are necessary on the trunk
- Use 802.1q trunking
- Tag all VLANs over a trunk from the aggregation switches
- **Step 1** Configure trunks using the standard encapsulation method 802.1q by entering the following command: (config-if)# switchport trunk encapsulation dot1q
- **Step 2** Define the VLANs permitted on a trunk by entering the following command:

(config-if)# switchport trunk allowed vlan VLAN IDs

Step 3 Modify the VLANs allowed on a trunk by using the following commands:

(config-if) # switchport trunk allowed vlan add VLAN IDs (config-if) # switchport trunk allowed vlan remove VLAN IDs

Step 4 Define a port as a trunk port by entering the following command:

(config-if) # switchport mode trunk



• The auto-negotiation of a trunk requires that the ports be in the same VTP domain and be able to exchange DTP frames.

Step 5 To secure and enforce a spanning tree topology, configure the Root Guard feature on the aggregate switch interfaces that connect to the access layer switches. The following is an example of the interface configuration between the aggregate and access layer switch with Root Guard enabled:

```
(config)# interface GigabitEthernet12/13
config-if)# description text
config-if)# no ip address
config-if)# switchport
config-if)# switchport trunk encapsulation dot1q
config-if)# switchport trunk native vlan <vlan id>
config-if)# switchport trunk allowed vlan <vlan id>
config-if)# switchport trunk allowed vlan <vlan id>
config-if)# switchport mode trunk
config-if)# spanning-tree guard root
config-if)# channel-protocol lacp
config-if)# channel-group group id mode active
```

Server Port Configuration

The server ports on the access layer switch may support single VLAN access and/or trunk configuration modes. The operational mode chosen should support the server NIC configuration. In other words, a trunking NIC should be attached to a trunking switch port. Enable PortFast for the edge devices in the spanning tree domain to expedite convergence times.

BPDU Guard disables a port that receives a BPDU. This feature protects the STP topology by preventing the blade server from receiving BPDUs. A port disabled using BPDU Guard must be recovered by an administrator manually. Enable BPDU Guard on all server ports that should not receive BPDUs. The commands required to enable this feature are as follows:

```
interface GigabitEthernet6/1
description blade server port
speed 1000
duplex full
switchport
switchport access vlan VLAN ID
switchport mode access
spanning-tree portfast
end
```

Port Security limits the number of MAC addresses permitted to access the blade switch port. To configure Port Security, configure the maximum number of MAC addresses expected on the port. The NIC teaming driver configuration (the use of a virtual MAC address) must be considered when configuring Port Security.

To enable Port Security, enter the following command:

(config) # switchport port-security maximum maximum addresses

Server Default Gateway Configuration

The default gateway for a server is a Layer 3 device located in the data center aggregation layer. This device may be a firewall, a load balancer, or a router. Using a redundancy protocol, such as HSRP, protects the gateway from becoming a single point of failure and improves data center network availability. HSRP allows the two aggregate switches to act as a single virtual router by sharing a common MAC and IP address between them. To enable HSRP, define a Switched VLAN Interfaces (SVI) on each aggregate switch and use the HSRP address as the default gateway of the server farm.

Step 1 Configure one aggregation switch as the active HSRP router. The **priority** command helps to select this router as the active router by assigning it a higher value.

```
interface Vlan10
description Primary Default Gateway
ip address IP address subnet mask
no ip redirects
no ip proxy-arp
arp timeout 200
standby 1 ip IP address
standby 1 timers 1 3
standby 1 priority 51
standby 1 preempt delay minimum 60
standby 1 authentication password
end
```

Step 2 Configure the second aggregation switch as the standby HSRP router as follows:

```
interface Vlan10
description Standby Default Gateway
ip address 10.10.10.3 255.255.255.0
no ip redirects
no ip proxy-arp
arp timeout 200
standby 1 ip 10.10.10.1
standby 1 timers 1 3
standby 1 priority 50
standby 1 preempt delay minimum 60
standby 1 authentication password
end
```

