

Architecting a software-defined data center with Big Cloud Fabric and Dell EMC ScaleIO

A Dell EMC Deployment and Best Practices Guide

Dell EMC Networking Solutions Engineering
August 2017

Revisions

Date	Description	Authors
August 2017	Initial release – Version 1.0	Ed Blazek, Curtis Bunch, Dennis Dadey, Jordan Wilson

THIS WHITE PAPER IS FOR INFORMATIONAL PURPOSES ONLY AND MAY CONTAIN TYPOGRAPHICAL ERRORS AND TECHNICAL INACCURACIES. THE CONTENT IS PROVIDED AS IS, WITHOUT EXPRESS OR IMPLIED WARRANTIES OF ANY KIND.

Copyright © 2017 Dell Inc. All rights reserved. Dell and the Dell EMC logo are trademarks of Dell Inc. in the United States and/or other jurisdictions. All other marks and names mentioned herein may be trademarks of their respective companies.

Table of contents

1	Introduction.....	5
1.1	Typographical Conventions.....	6
2	Pod architectures.....	7
2.1	Big Cloud Fabric pod.....	8
2.2	VMware vSphere pod types.....	10
3	Hardware.....	12
4	Management network.....	16
5	Deploy Big Cloud Fabric.....	18
5.1	Big Cloud Fabric controller.....	19
5.2	Big Switch Zero Touch Fabric.....	22
5.3	Tenant and segment configuration.....	26
5.4	VMware integration.....	29
6	Deployment of VMware vSphere.....	35
6.1	vCenter server deployment and design.....	36
6.2	Virtual network design.....	39
7	Deploying Dell EMC ScaleIO.....	44
7.1	Deploy the Dell EMC ScaleIO plug-in.....	46
7.2	Upload Dell EMC ScaleIO OVA Template.....	47
7.3	Deploy Dell EMC ScaleIO.....	48
7.4	Dell EMC ScaleIO GUI.....	51
8	Performance Tuning.....	52
8.1	Maximum Transmission Unit size.....	53
8.2	Quality of Service.....	56
8.3	Network I/O Control.....	59
A	Configuration details.....	61
B	Technical support and resources.....	62
B.1	Dell EMC product manuals and technical guides.....	62
B.2	Dell EMC Solution Briefs.....	62
B.3	Big Switch Networks product manuals and technical guides.....	62
B.4	VMware product manuals and technical guides.....	62
C	Support and feedback.....	63

Executive summary

This document provides best practices and details how to deploy a software-defined data center (SDDC) powered by Dell EMC S-Series switches, Dell EMC PowerEdge servers, Big Switch Networks, VMware vSphere, and Dell EMC ScaleIO. The goal of this document is to:

- Assist administrators in selecting the best hardware and topology for their Big Cloud Fabric™ network
- Deliver detailed instructions and working examples on cabling, configuration, and deploying the BCF network
- Deliver detailed instructions and working examples on deploying and configuring VMware vSphere components
- Deliver step-by-step instructions and working examples on deploying and configuring a sample Dell ScaleIO virtual SAN
- Show conceptual, physical, and logical diagram examples for various networking topologies

1 Introduction

Applications are the engines for modern businesses. They drive innovation, operational efficiency, and revenue generation. They demand an infrastructure that is highly agile and easy to manage while reducing costs. These applications, which include mission critical Enterprise Resource Planning (ERP) systems, multi-tier web applications, and big data, have placed new constraints on the networking infrastructure; support for high east-west traffic bandwidth, virtual machine mobility, and multitenancy.

Infrastructure teams have struggled to respond to these requirements. Unlike the rest of the portfolio they manage, legacy networks remain highly static and require extensive manual intervention and operational overhead. While the speed, scale, and density of equipment offered by traditional networking vendors have increased over the last two decades, the underlying architectures and business operating models have stayed fundamentally unchanged.

Dell EMC is working closely with Big Switch Networks to introduce the industry's first data center leaf-spine IP fabric solution built using Dell EMC open networking switches and Big Cloud Fabric (BCF). This joint solution applies the hardware-software disaggregation enabled by Dell EMC and Big Switch Networks. Software Defined Network (SDN) designs inspired by hyperscale data center architectures, provide significant cost savings and operational efficiencies for enterprise data centers by enabling the VMware centered software-defined data center (SDDC).

With built-in integration for VMware, BCF is an ideal physical network for virtual environments, network virtualization, and hyper-converged Infrastructure (HCI). It is the industry's first SDN-based fabric, using open networking switch hardware that provides intelligent, agile, and flexible networking for the VMware SDDC.

With servers, networking, and virtualization in place, distributed storage remains as the final component of the SDDC model. Dell EMC ScaleIO is a perfect storage foundation for the SDDC. ScaleIO delivers scale-out, block storage using commodity hardware. ScaleIO creates a server-based Storage Area Network (SAN) built from server direct-attached storage to deliver flexible and scalable performance on demand.

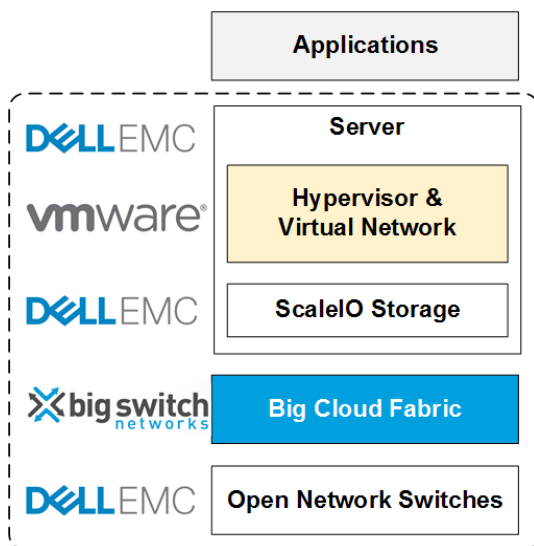


Figure 1 Conceptual view of Big Cloud Fabric (BCF) and Dell EMC solution for VMware environment

1.1 Typographical conventions

This document uses the following typographical conventions:

Monospaced text

Command Line Interface (CLI) examples

Bold monospaced text

Commands entered at the CLI prompt

Italic monospaced text

Variables in CLI examples

Bold text

GUI navigation prompts

2 Pod architectures

A pod is a combination of computing, network, and storage capacity, designed to be deployed as a single unit. As a result, a pod is the largest unit of failure in the software-defined data center (SDDC). Carefully engineered services ensure that each pod has little to no shared vulnerability between pods. While each pod usually spans one rack, it is possible to aggregate multiple pods into a single rack or to span a pod across multiple racks.

There are two different types of pods used in this deployment sample:

- Big Cloud Fabric (BCF) pod – A pair of BCF controllers manages a maximum of 16 racks with redundant leaf switches
- VMware vSphere pods – A collection of the ESXi hosts and virtual machines grouped by function In this document, three VMware pod types are defined: Management, ScaleIO/Compute, and Edge pods

The maximum of 16 racks for the BCF rack is a limit imposed on the deployment covered in this paper. This number is reached due to 32 40GbE QSFP+ interface limit of the Dell EMC S6010-ON used as the spine switch. BCF has a tested maximum of 128 racks. See [Big Cloud Fabric Verified Scale Guide](#) for more information.

Note: Big Switch documentation requires a customer account to access. Contact your Big Switch Networks account representative for assistance.

2.1 Big Cloud Fabric pod

In this example, the Big Cloud Fabric (BCF) pod contains two spine switches and six redundant leaf switches distributed over three racks. Two BCF controller nodes are deployed in an active/standby configuration. Two Dell EMC Networking S6010-ON switches are deployed as spine switches and six Dell EMC Networking S4048-ON switches are servicing three server racks in a redundant configuration. The following image shows a final view of the leaf-spine network architecture used in this example as demonstrated by the BCF Controller graphical user interface (GUI):

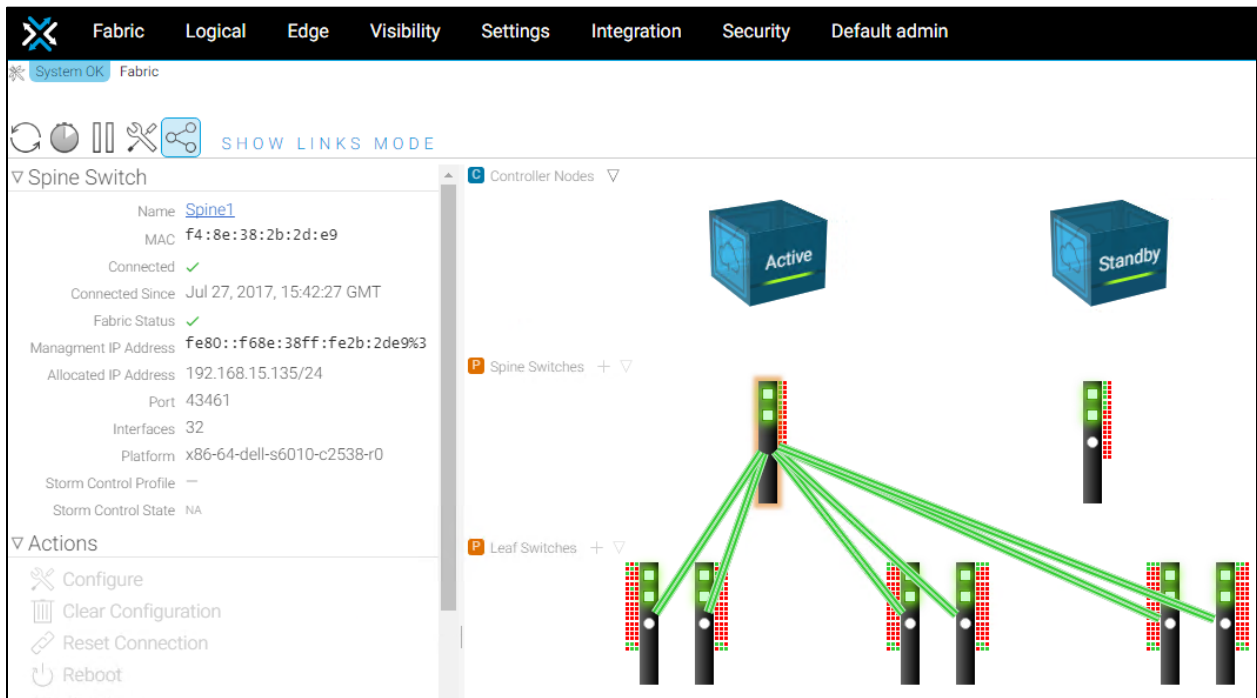


Figure 2 Leaf-Spine switched fabric

The leaf-spine architecture provides a simple and efficient design in response to challenges inherent in the hierarchical data center architecture. The 2-layer leaf-and-spine architecture optimizes bandwidth between switch ports within the data center by creating a high-capacity fabric using multiple spine switches that interconnect the edge ports of each leaf switch. This design provides consistent latency and minimizes the hops between servers in different racks.

The design lends itself well to the creation of an independent, replicable pod that scales without disrupting network traffic. The addition of more leaf switches increases the number of switch edge ports for connecting to servers. Extra spine switches increase the fabric bandwidth and lower oversubscription ratios.

Beyond physical parameters, there are two other considerations made in the BCF pod design:

- Oversubscription ratios
- Fault tolerance

In a leaf-spine network, oversubscription occurs at the leaf layer. Oversubscription is equal to the total amount of bandwidth available to all servers connected to a leaf switch divided by the amount of uplink bandwidth.

$$\text{Oversubscription} = \text{total bandwidth} / \text{uplink bandwidth}$$

The following image shows a typical Dell EMC ScaleIO node uses four 10GbE ports, two links going to each leaf switch. If a rack has 19 nodes, each leaf provides a total bandwidth of 380Gbps. Each leaf switch would have 80Gbps uplink bandwidth resulting in a 4.75:1 oversubscription ratio.

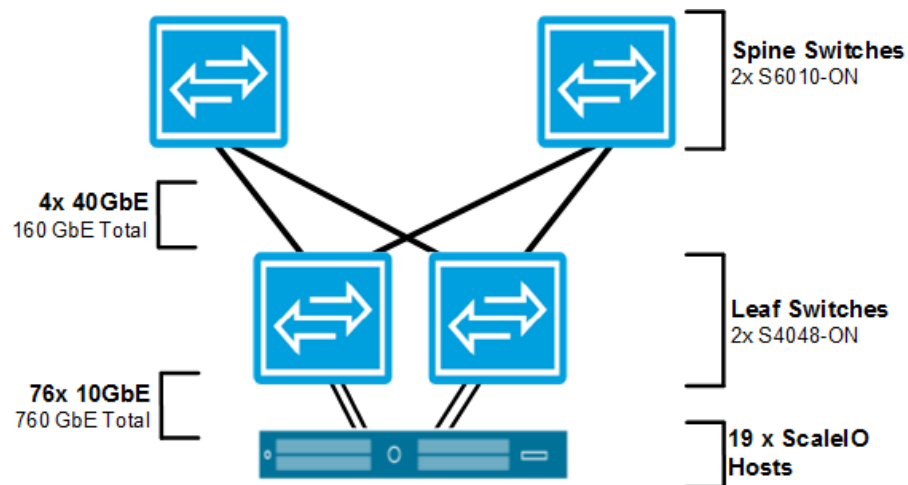


Figure 3 Oversubscription calculations

$$4.75 \text{ (oversubscription)} = 380 \text{ (total bandwidth)} / 80 \text{ (uplink bandwidth)}$$

To decrease oversubscription ratios, more spine switches can be deployed. Using the same leaf/node configuration previously and increasing the number of spine switches to four the total uplink bandwidth doubles to 160Gbps resulting in an oversubscription ratio of 2.4:1.

$$2.4 \text{ (oversubscription)} = 380 \text{ (total bandwidth)} / 160 \text{ (uplink bandwidth)}$$

In addition to decreasing the oversubscription ratio, extra spine switches improve the fabric resiliency. For example, if a spine switch fails, traffic continues across the remaining spine switches. The greater the number of spine switches in the spine, the less additional load the remaining spine switches must take on in the event of a failure in one spine switch. For example, with four spine switches, a failure of a single spine switch only reduces the capacity by 25% vs. a reduction in capacity of 50% with two spine switches.

Tolerable oversubscription ratios vary by each enterprise organization and should be considered part of the design criteria. In this example, an oversubscription ratio of 4.75:1 is considered tolerable.

2.2 VMware vSphere pod types

A VMware pod can be one of three types that attaches to the Big Cloud Fabric (BCF). The [VMware Validated Designs Documentation](#) defines the concept of the VMware pods. The VMware Validated Designs Documentation also contains numerous best practices for VMware vSphere deployment. The virtualized environment in the following example contains three different pods:

- Management pod
- ScaleIO and Compute pod
- Edge Pod

The following image shows a physical layout with the three VMware pod types shown. In this example, the BCF pod is shown in the gray block and contains two BCF controller appliances, six leaf switches, and four spine switches. The spine switches are shown in the right rack at the top with the two leaf switches directly below. The ScaleIO and Compute pod, shown in red, can be expanded to an extra 14 racks while being included in the same BCF pod. The VMware vSphere management pod is shown in blue.

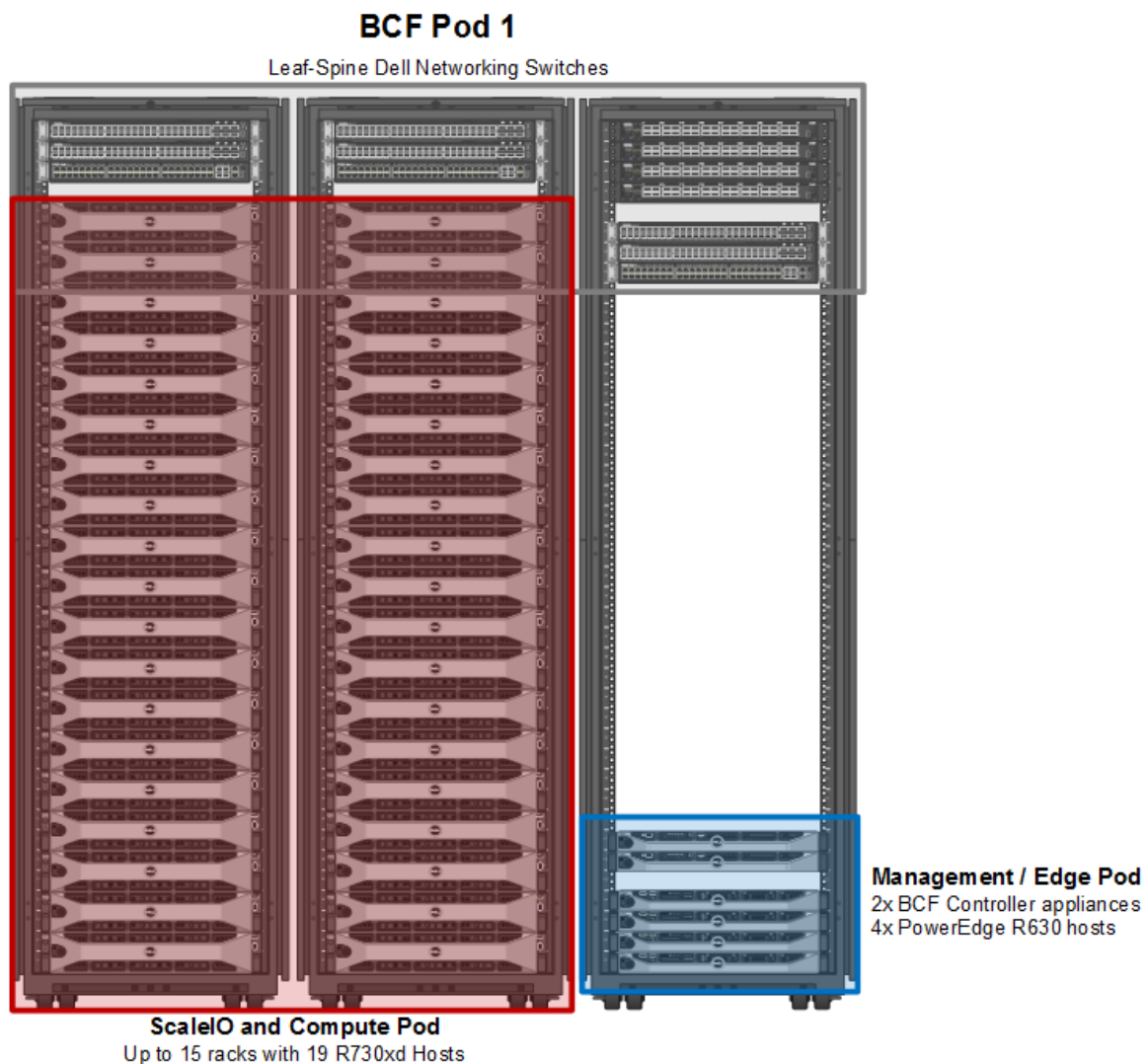


Figure 4 Big Cloud Fabric (BCF) and the VMware pod design

The **Management pod** runs the virtual machines and BCF Controller appliances that manage the entire environment. Management, monitoring, and infrastructure services that are virtual machines, are provisioned to a vSphere cluster, and provide high availability for these services. Permissions on the management cluster limit access to administrators. This limitation protects these virtual machines and provides clear administrative boundaries in the environment.

The **ScaleIO and Compute pod** runs the required ScaleIO software components and hosts the tenant's virtual machine workloads. The pod scales by adding more nodes, which increases computing and storage capacity linearly.

The **Edge pod** uses two leaf switches to route to the core of the data center. BCF is designed to easily connect to an existing Layer 3 device in the data center, such as a core router. Initially, the Edge pod consists of a pair of leaf switches, shared by the Management pod. It can be expanded to include servers if a workload requires it. An example of the expansion would be for the deployment of VMware NSX Edge Service Gateways (ESG), enabling VXLAN overlay networks.

Note: Edge pod deployment is beyond the scope of this document. See [VMware Validated Designs Documentation](#) for instructions on deploying the Edge pod.

3 Hardware

The three switch models that are listed, support the Open Network Install Environment (ONIE) and are on the Big Switch Networks hardware compatibility list. The six Dell EMC Networking S4048-ON switches are deployed in three redundant leaf pairs across three racks, with the seventh deployed as an aggregation switch for the management network. One Dell EMC Networking S3048-ON switch is deployed per rack and serves as an out-of-band access switch. The following table outlines the Dell EMC Networking switches used in this example:

Table 1 Dell EMC Networking ONIE switches

Switch model	Switch port types	Roles	Count
Dell EMC Networking S3048-ON	48 x 1GbE BASE-T and 4 x 10GbE SFP+ ports	BCF switch control plane, BCF management plane, PowerEdge iDRAC access	3
Dell EMC Networking S4048-ON	48 x 10GbE SFP+ and 6 x 40GbE QSFP+ ports	Leaf switch, management network aggregation	7
Dell EMC Networking S6010-ON	32 ports of 40GbE QSFP+ ports	Spine switch	2

The S3048-ON is a 1-Rack Unit (RU) switch with forty-eight 1GbE Base-T ports and four 10GbE SFP+ ports and deploys as an out-of-band switch.

Note: The 24-port Dell EMC Networking S3024-ON is interchangeable with the S3048-ON and can be used based on switch port requirements.

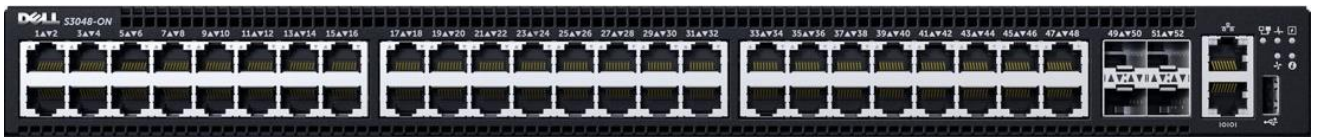


Figure 5 Dell EMC Networking S3048-ON

The S4048-ON is a 1-RU, multilayer switch with forty-eight 10GbE SFP+ ports and six 40GbE QSFP+ ports and deploys as leaf switches in the examples in this guide.

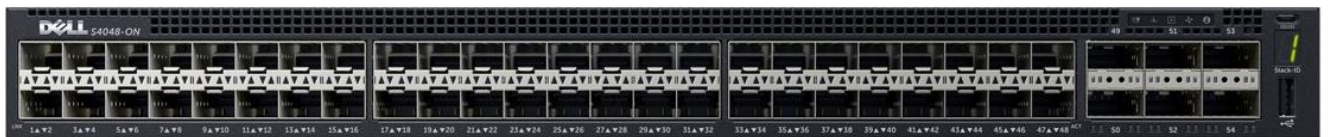


Figure 6 Dell EMC Networking S4048-ON

The S6010-ON switch is a 1-RU, multilayer switch with thirty-two 40GbE QSFP+ ports and deploys as spine switches.



Figure 7 Dell EMC Networking S6010-ON

The Dell PowerEdge R630 hosts the Management and Edge pod workloads while the Dell PowerEdge R730xd hosts the ScaleIO and Compute pod workloads. The BCF controllers were provided by Big Switch Networks as part of the BCF offering and are a part of the infrastructure requirements. The following table shows the servers used, their function, and quantity:

Table 2 Dell EMC PowerEdge Servers

Server model	Function	Count	Operating system
Dell PowerEdge R630	Management pod	4	ESXi 6.0 U3
Dell PowerEdge R730xd	ScaleIO and compute pod	4	ESXi 6.0 U3
Big Cloud Fabric (BCF) controller	BCF management	2	n/a

Shown are two of the four R630 servers, mgmt01esx01 and mgmt01esx04. ESX01 connects to Ethernet 1 on both of the leaf switches, while ESX04 connects to Ethernet port 7. The port highlighted in red is the PowerEdge iDRAC interface and connects to the S3048-ON management switch. The IP fabric automatically defines the BCF MLAG when the BCF controller discovers the connection between the two leaf switches. The following image shows the interfaces provided by the R630 hardware, and the connections made to the leaf switches:

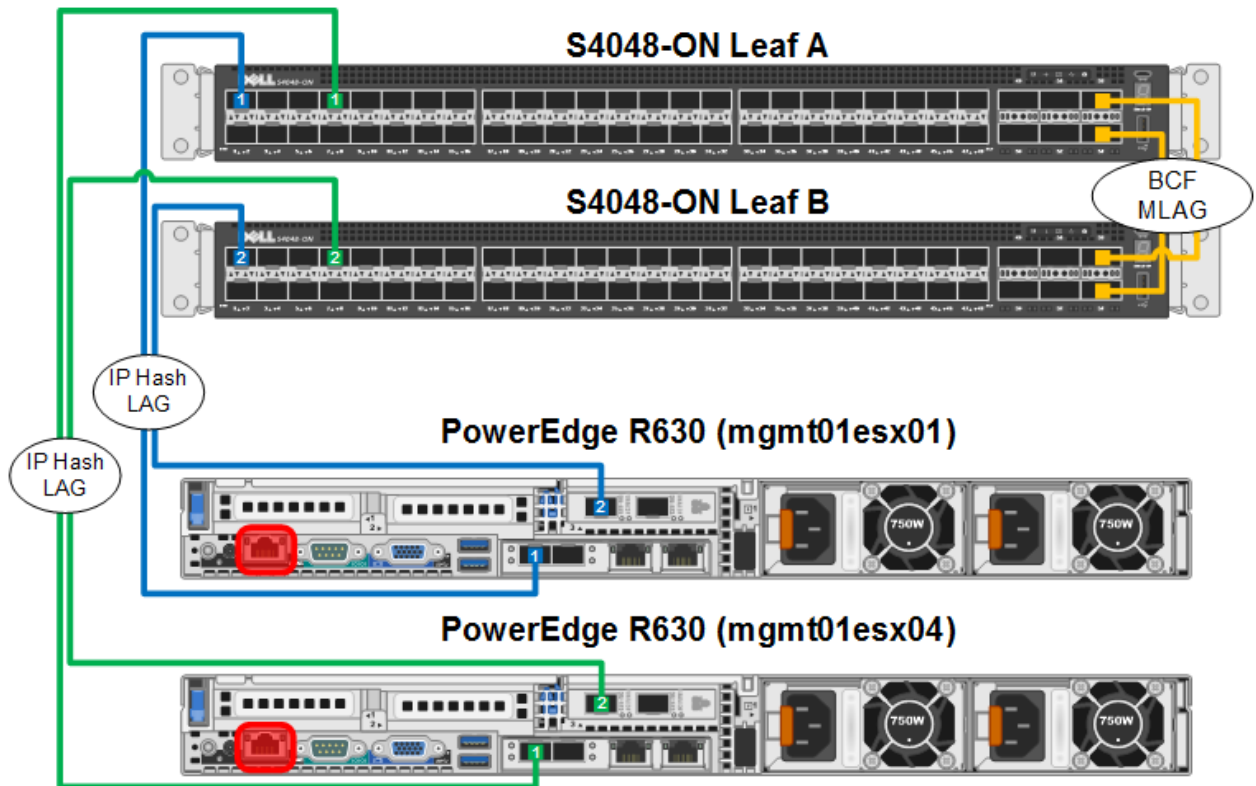


Figure 8 Dell PowerEdge R630 interface connections

The following table provides a list of the interfaces provided by the R630 server:

Table 3 Dell PowerEdge R630 interfaces

Port set	Port count/type	Color	Function	Connection location
Management	2 x 10GbE	Blue/Green	ESXi management traffic, VMware vSphere vMotion traffic	S4048-ON leaf switch pair for the Management pod
iDRAC	1 x 1GbE	Red </td <td>Onboard controller for out-of-band (OOB) management</td> <td>PowerEdge iDRAC port on S3048-ON management switch</td>	Onboard controller for out-of-band (OOB) management	PowerEdge iDRAC port on S3048-ON management switch

Two of the possible 19 R730xd servers are shown, sio01esx01 and sio01esx19. ESX01 connects to Ethernet 1 and 2 on both leaf switches while ESX19 connects to Ethernet 19 and 20. The port highlighted in red is the PowerEdge iDRAC interface which connects to the S3048-ON management switch. The BCF Multi-Chassis Link Aggregation (MLAG) is automatically defined in the IP fabric when the BCF controller discovers the connection between the two leaf switches. The following image shows the interfaces provided by the R730xd hardware and the connections made to the corresponding S4048-ON leaf switch pair:

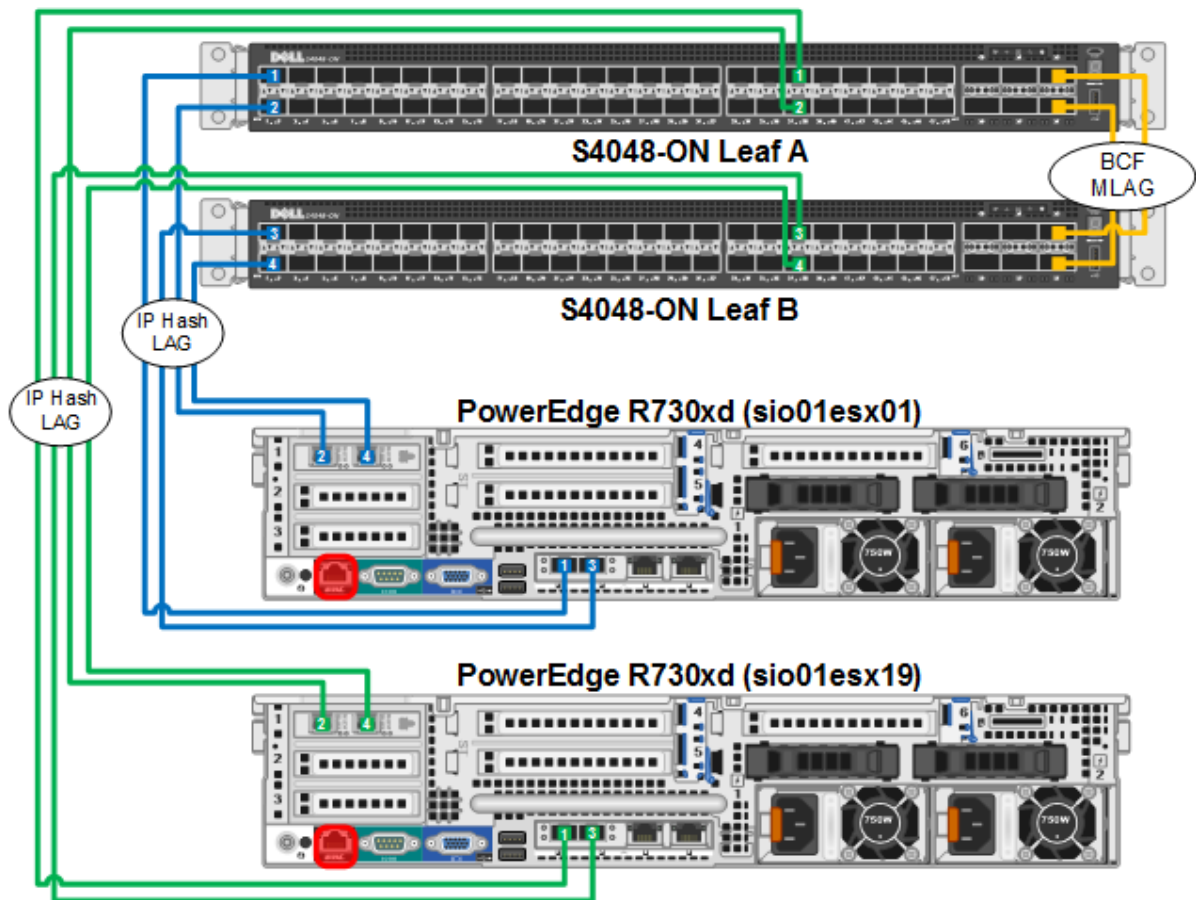


Figure 9 Dell PowerEdge R730xd interface connections

The following table provides a listing of the interfaces provided by the PowerEdge R730xd server:

Table 4 PowerEdge R730xd interfaces

Port set	Port count/type	Color	Function	Connection location
ScaleIO and Compute	4 x 10GbE	Blue/Green	ESXi management, VMware vSphere vMotion, ScaleIO management, tenant traffic, and ScaleIO data network traffic	S4048-ON leaf switch pair for the ScaleIO and Compute pod
iDRAC	1 x 1GbE	Red	Onboard controller for out-of-band (OOB) management	PowerEdge iDRAC port on S3048-ON management switch

All interfaces on the BCF controllers are connected to the management network. The blue links are used to manage the switches in the IP fabric, and the green links are used to access the BCF controllers through the management network. The following image shows the interfaces provided by the BCF controller hardware:

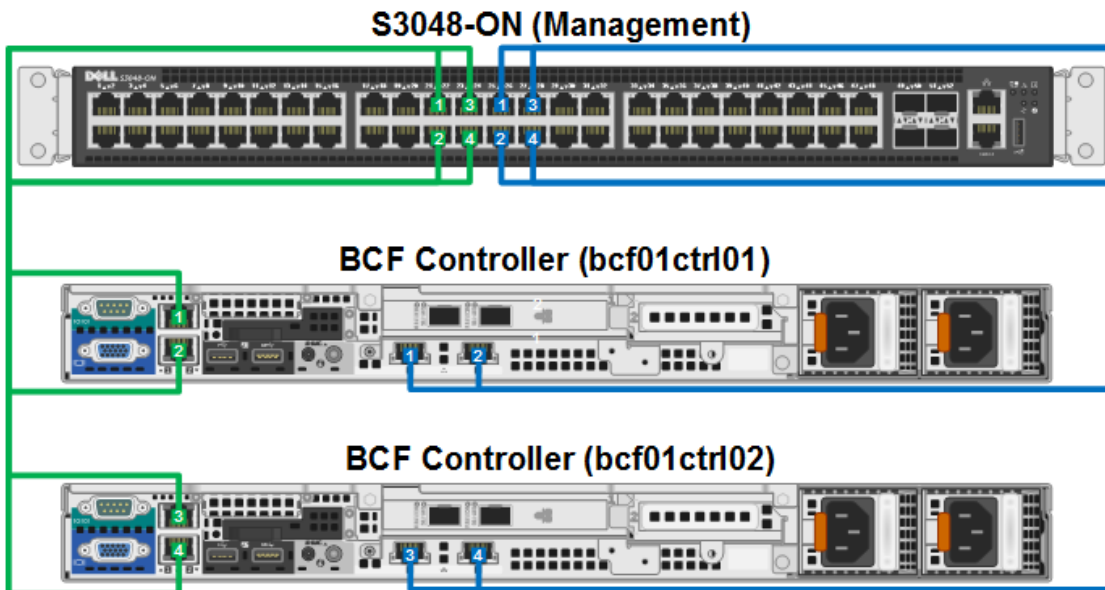


Figure 10 BCF controller hardware appliance interface connections

The following table is a list of the interfaces provided by the BCF controller:

Table 5 BCF controller interfaces

Port set	Port count/type	Color	Function	Connection location
P-switch	2 x 1GbE	Blue	Control traffic between BCF controller and fabric switches	S3048-ON management switch
BCF mgmt.	2 x 1GbE	Green	Administrative access to the BCF controller	S3048-ON management switch

4 Management network

The traffic flows for a typical data center management network tend to be north-south oriented. As a result, the management network used in this deployment example is built using a traditional 3-tier hierarchal model of access, aggregation, and core. BCF requires that all ports connected to the controllers and fabric switches, be in the same broadcast domain and that the network supports both IPv4 and IPv6. The BCF physical switch network does not require any routing to function properly, the aggregation/core is outlined here to provide integration with other management systems found in a typical data center.

The management network consists of a Dell EMC Networking S3048-ON switch in each rack that acts as an access switch with a 10GbE link to a Dell EMC Networking S4048-ON aggregation switch. The management network is not part of the BCF pod and all management switches run Dell EMC Networking OS9 (DNOS9). The following table shows the three networks that support the deployment:

Table 6 Management network subnets

Network name	Color	Function
iDRAC	Red	PowerEdge onboard controller access
BCF management	Green	Administrative access to the BCF controllers.
BCF p-switch	Blue	Control traffic between BCF controller and fabric switches

The following image shows the logical topology of the management network. Each S3048-ON uses a 10GbE link to carry each VLAN to the S4048-ON. Dell EMC Networking S-Series switches do not block network layer traffic (IPv4/IPv6) by default, which allows for a successful BCF deployment. The S4048-ON serves as the Layer 2/Layer 3 boundary to the management core in the data center.

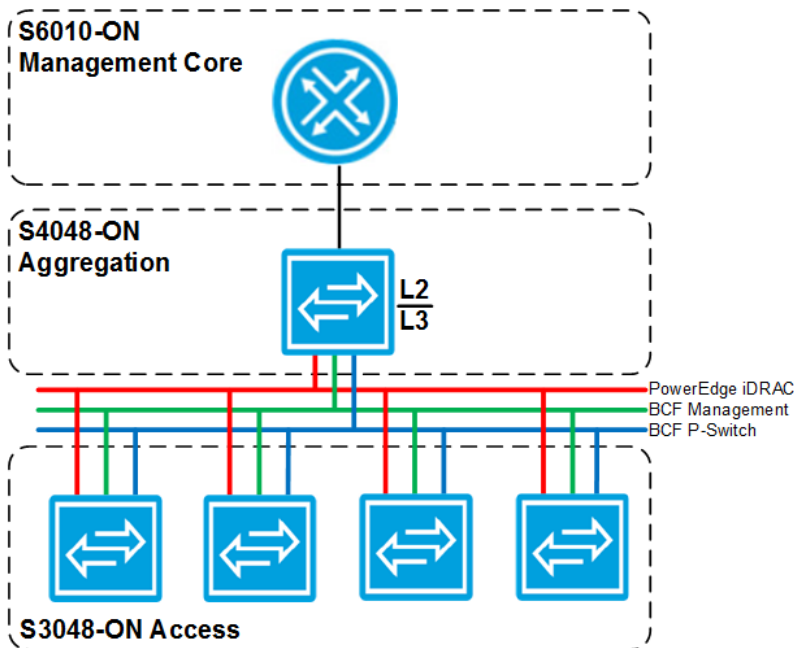


Figure 11 Management network logical topology

The S3048-ON can be preconfigured to divide the interfaces into VLAN/subnet groups. The following figure shows an example of a single S3048-ON switch with groups of interfaces configured as access switch ports that correspond to each VLAN shown in the image preceding. The BCF management VLAN, shown in green, is only required in the Management pod. The entire environment uses the PowerEdge iDRAC and BCF p-switch VLANs.

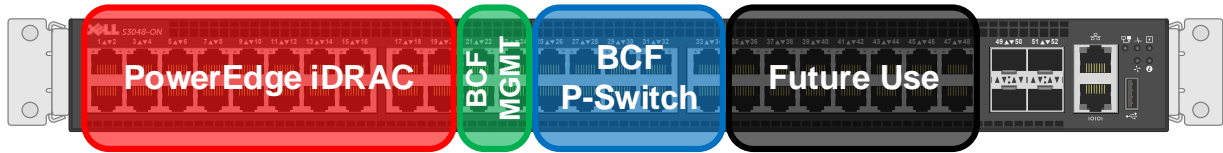


Figure 12 Dell EMC Networking S3048-ON interface division

Note: To configure the S3048-ON, see [Dell Configuration Guide for the S3048-ON System](#).

The following image shows the management ports of the leaf and spine switches in the Management pod, connected to the S3048-ON management switch for the BCF p-switch interfaces.

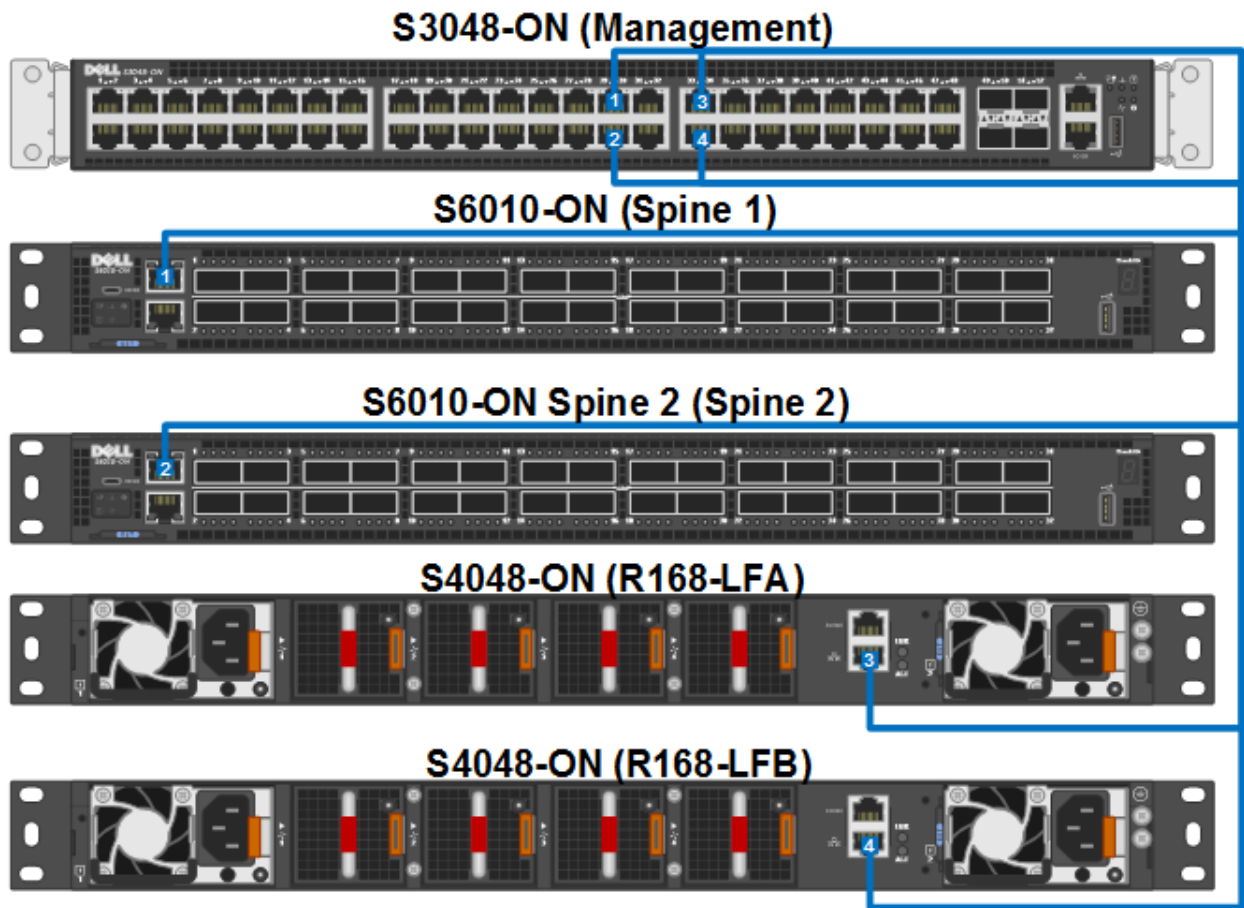


Figure 13 Spine and leaf switch management ports

5 Deploy Big Cloud Fabric

Big Cloud Fabric (BCF) provides a high bisectional bandwidth network. Each fabric device can switch at Layer 2 or route at Layer 3, while the BCF controller centrally provides the intelligence required to make full use of redundant links. Incremental upgrades of the forwarding tables are dynamically pushed to each switch to ensure a stable and dynamic network operation. Spanning tree is not required, and all links are in forwarding mode. The BCF controller prevents loops from forming.

As mentioned in the [Big Cloud Fabric pod](#) section, the networking architecture used by BCF is a leaf-spine design that increases server-to-server bandwidth. The leaf-spine architecture creates a high-performance backplane that can be extended by simply adding more switches. Fabric edge ports are aggregated through static Link Aggregation Groups (LAG) for higher bandwidth to all servers.

A pair of BCF controllers provides functionality including the dual supervisors on a modular chassis. The spine switches provide the functionality of the backplane, while the leaf switches are similar in function to line cards. The following image shows how the entire BCF pod can be thought of as a single larger modular switch:

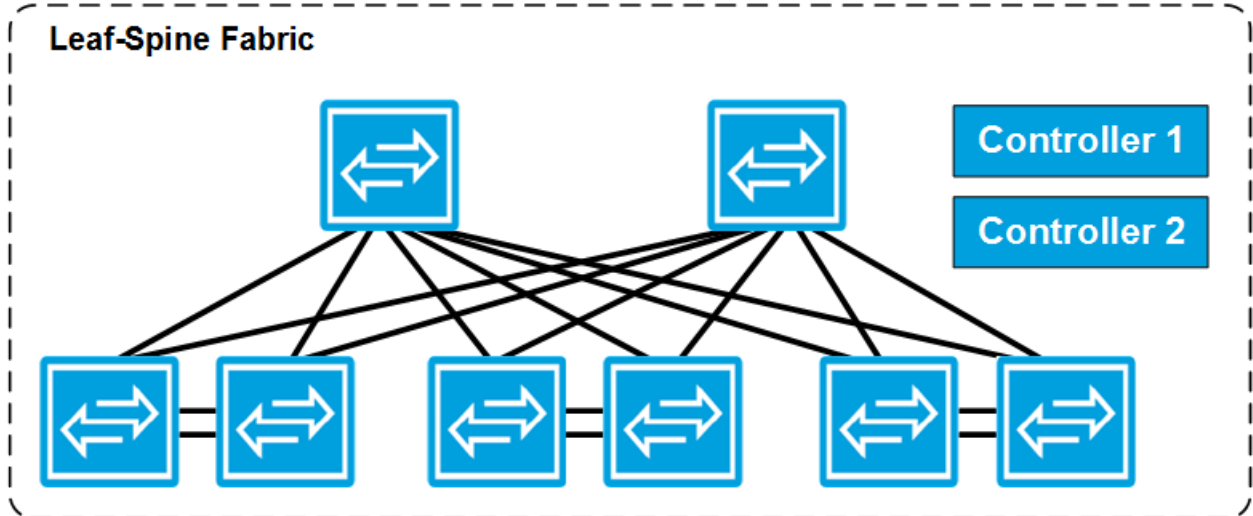


Figure 14 Big Cloud Fabric (BCF) pod

5.1 Big Cloud Fabric controller

The BCF controller provides “single pane of glass” management of all leaf and spine switches. The BCF controller supports a familiar Command Line Interface (CLI), and a web-based graphical user interface (GUI). Any custom orchestration can be executed by using the industry-standard RESTful application programming interface (API).

BCF supports the traditional tools for debugging, including ping, traceroute, show commands, and redirecting packets using port mirroring for fault analysis. Also, the BCF controller supports unique troubleshooting tools, such as Fabric Test Path, and Fabric Analytics, to quickly isolate, identify, and resolve forwarding and application faults.

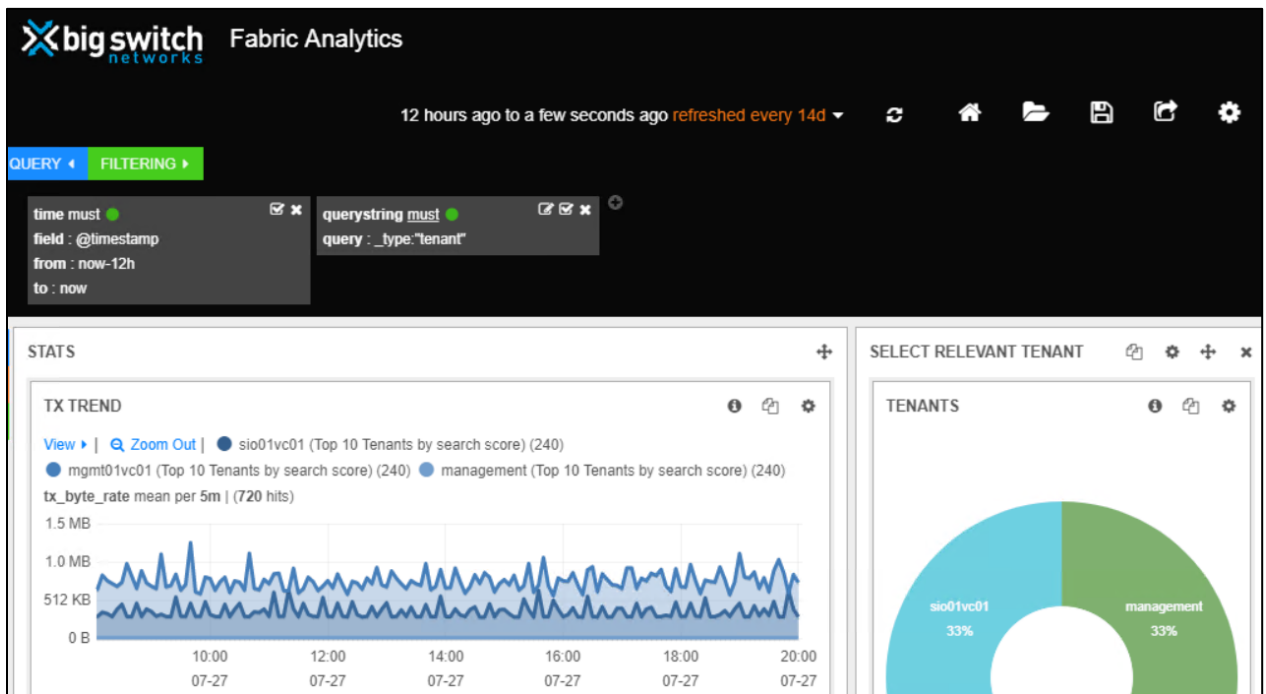


Figure 15 Big Cloud Fabric (BCF) analytics

Note: To deploy the BCF controller, see [Big Cloud Fabric 4.2.0 Deployment Guide](#).

Gateway and DNS settings are optional and are included here for completion. The following table contains the settings used during the initial deployment of both controllers:

Table 7 BCF controller initial configuration settings

Hostname	IP address	IPv4 prefix length	Default Gateway	DNS server address	DNS search domain
bcf01ctrl01	192.168.15.5	24	192.168.15.1	192.168.15.4	dnt.adc.delllabs.net
bcf01ctrl02	192.168.15.6	24	192.168.15.1	192.168.15.4	dnt.adc.delllabs.net

The following table contains the configuration settings for the controllers. Use these values during the deployment of the first controller. The second controller is added to the existing cluster using the IP address of the active controller. At that point, the cluster's name and system time are imported and do not need to be specified.

Table 8 BCF cluster settings

Hostname	Controller clustering	Existing Cluster IP	Cluster name	System time
bcf01ctrl01	Start a new cluster	n/a	BCF-Cluster-01	ntp.dnt.adc.delllabs.net
bcf01ctrl02	Join an existing cluster	192.168.15.5	n/a	n/a

As a best practice, set a Virtual IP (VIP) for the cluster. This allows you to connect to the management port of the active node using an IP address that does not change even if the active controller fails over and the role of the standby controller changes to the active.

On the active controller, set the VIP by using the `virtual-ip` command from the `config-controller` submode:

```
bcf01ctrl01> enable
bcf01ctrl01# config
bcf01ctrl01 (config)# controller
bcf01ctrl01 (config-controller)# virtual-ip 192.168.15.7
```

To verify the cluster settings, enter the `show controller` command from the active controller. Verify that the VIP is reporting correctly and that the cluster status is in a redundant state.

```
bcf01ctrl01> show controller
Cluster Name           : BCF-Cluster-01
Cluster Virtual IP    : 192.168.15.7
Redundancy Status     : redundant
Last Role Change Time : 2017-07-11 15:18:43.016000 UTC
Failover Reason       : Changed connection state: cluster$
Cluster Uptime        : 3 weeks, 5 days
# IP                  @ State    Uptime
-|-----|-----|-----|
1 192.168.15.5 * active 1 day, 20 hours
2 192.168.15.6  standby 1 day, 3 hours
```

At this point, the BCF GUI can be accessed using the VIP address of the cluster. This is the address that is used through the rest of the document when referring to management and the configuration of IP fabric. In the sample deployment, the hostname of the cluster was used to access the control cluster VIP as shown in the following figure:

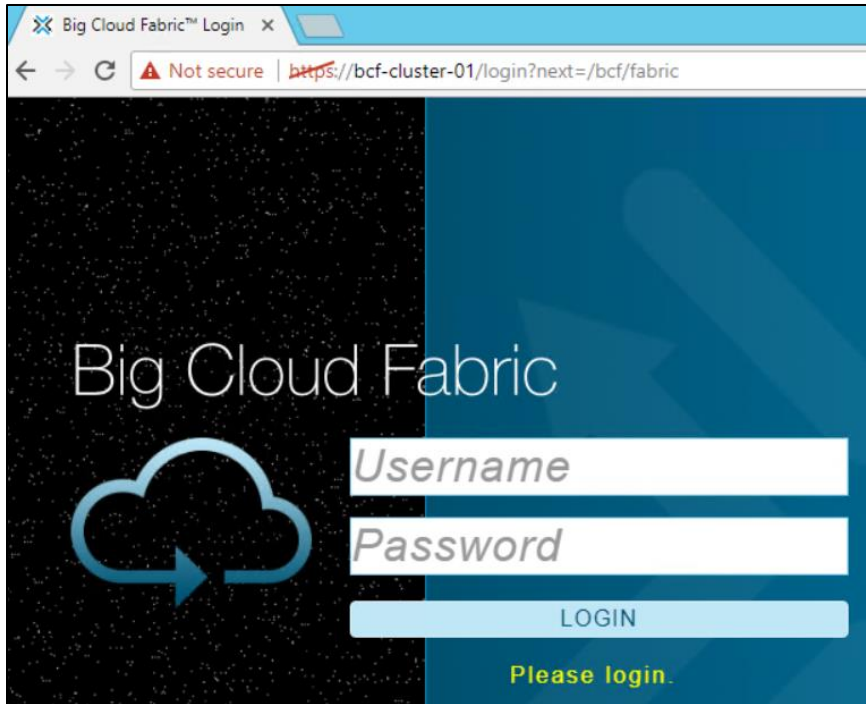


Figure 16 Connecting to the Big Cloud Fabric (BCF) graphical user interface (GUI)

Note: The BCF controller uses a self-signed certificate, replaceable through the BCF CLI. See [Big Cloud Fabric 4.2.0 CLI Reference Guide](#) to properly secure the web connection.

5.2 Big Switch Zero Touch Fabric

Big Switch Zero Touch Fabric (ZTF) uses the Open Networking Install Environment (ONIE) boot loader to automate switch installation and configuration. ONIE makes deploying many switches in a data center easier and less prone to errors. The ZTF process uses ONIE to automatically install the correct version of Switch Light OS on each switch when the switch is powered on and connected to the Big Cloud Fabric (BCF) controller. The Dell EMC Networking switches used in this example do not have BCF Switch Light OS installed initially. In the following steps, the BCF controller deploys the OS to the fabric switches.

Switch Light OS is a complete SDN operating system based on Open Network Linux (ONL) and is bundled with the BCF software distribution. This ensures that the software running on the switch is compatible with the version of the controller software. The following image shows the OS deployment steps, followed by the summarization of the steps which are provided in the following table:

Note: For more information about this process, see Chapter 4 of the [BCF 4.2 Users Guide](#).

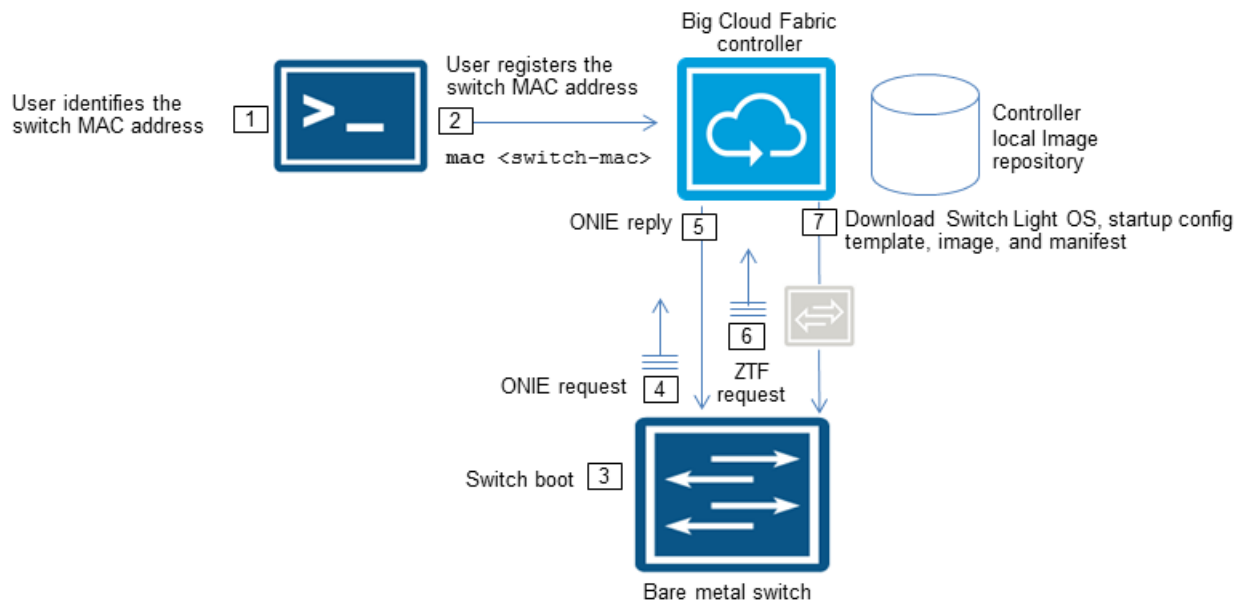


Figure 17 Big Cloud Fabric (BCF) switch registration workflow

Table 9 Big Cloud Fabric (BCF) switch steps summary

Step number	Description
1	Collect switch MAC addresses from the Dell EMC S-Series Inventory tag on the switch
2	Register switch MAC addresses using the BCF GUI or CLI
3	Power cycle switch by briefly pulling the power cables out from each PSU
4	ONIE request to BCF controller
5	ONIE loader generates an IPv6 neighbor discovery message on the local network segment
6	The controller responds to the ONIE request from the switch and instructs the switch to download the Switch Light OS loader and begins the installation
7	After Switch Light loader reboots, it broadcasts a ZTF request
8	The ZTF server sends the Switch Light OS image, manifest, and startup-config to the switch

The switch downloads the startup-config from the controller, which includes the following configuration information:

- Hostname
- Switch MAC address
- Controller IP addresses
- NTP, logging, and SNMP configuration

The BCF GUI is used to configure the Dell EMC Networking S4048-ON leaf switches and the Dell EMC Networking S6010-ON spine switches that comprise the IP fabric. To bring up these fabric switches, navigate to **Fabric > Switches** and add a switch. Before deploying the IP fabric, collect all the switch MAC addresses. Also, place all fabric switches in ONIE boot mode. The following table lists the switch model, fabric role, and MAC addresses from this sample deployment.

Table 10 Switch fabric configuration details

Switch name	Switch model	MAC address	Fabric role
Spine1	S6010-ON	f4:8e:38:2b:0b:69	Spine
Spine2	S6010-ON	f4:8e:38:2b:2d:e9	Spine
R168-LFA	S4048-ON	f4:8e:38:45:ad:22	Leaf
R168-LFB	S4048-ON	f4:8e:38:45:ae:22	Leaf
R170-LFA	S4048-ON	64:00:6a:e6:b6:14	Leaf
R170-LFB	S4048-ON	f4:8e:38:20:44:29	Leaf
R171-LFA	S4048-ON	14:18:77:e0:69:31	Leaf
R171-LFB	S4048-ON	f4:8e:38:45:b5:22	Leaf

To create an IP fabric, each switch in the topology must have the management Ethernet interface connected to the Dell EMC Networking S3048-ON management switch in each rack. See the [Management network](#) section for more information.

The remaining steps of the installation and configuration process happen automatically after you turn-on the registered switch. To verify successful connectivity, navigate to **Fabric > Switches** to view all fabric switches, MAC addresses, name connection, and fabric status and fabric role. The following image shows this view:

The screenshot shows the 'Switches' page with a 'Summary of Firmware Versions' section and an 'IP Address Allocation' section. Below these is a table with the following columns: MAC, Name, Description, Connected, Fabric Status, Fabric Role, Spine, Leaf, Virtual, and Leaf Group. The table lists eight switches, including two spine switches (Spine1 and Spine2) and six leaf switches (R168-LFA, R168-LFB, R170-LFA, R170-LFB, R171-LFA, R171-LFB).

MAC	Name	Description	Connected	Fabric Status	Fabric Role	Spine	Leaf	Virtual	Leaf Group
f4:8e:38:45:ad:22	R168-LFA	R168U29	✓	✓	Leaf	–	✓	–	Rack168
f4:8e:38:45:ae:22	R168-LFB	R168U28	✓	✓	Leaf	–	✓	–	Rack168
64:00:6a:e6:b6:14	R170-LFA	R170U35	✓	✓	Leaf	–	✓	–	Rack170
f4:8e:38:20:44:29	R170-LFB	R170U34	✓	✓	Leaf	–	✓	–	Rack170
14:18:77:e0:69:31	R171-LFA	R171U35	✓	✓	Leaf	–	✓	–	Rack171
f4:8e:38:45:b5:22	R171-LFB	R171U34	✓	✓	Leaf	–	✓	–	Rack171
f4:8e:38:2b:2d:e9	Spine1	R168U36	✓	✓	Spine	✓	–	–	NA
f4:8e:38:2b:0b:69	Spine2	R168U35	✓	✓	Spine	✓	–	–	NA

Figure 18 Big Cloud Fabric (BCF) switches

The IP fabric topology is auto-discovered through Link Layer Discovery Protocol (LLDP), and the BCF controller creates Link Aggregation Groups (LAG) automatically from the links coming from the PowerEdge servers. These redundant links span multiple leaf switch pairs in the rack, creating an MLAG similar in function Dell EMC Networking OS9 Virtual Link Trunking (VLT). To view these links, enter the `show link` command on the BCF CLI:

```
bfc01ctrl01> show link
#  Switch Name IF Name      Switch Name IF Name      Link Type
--|-----|-----|-----|-----|-----|
 1 R168-LFA      ethernet53 R168-LFB      ethernet53 peer
 2 R168-LFA      ethernet54 R168-LFB      ethernet54 peer
 3 R170-LFA      ethernet49 Spine1         ethernet1 leaf-spine
 4 R170-LFA      ethernet50 Spine2         ethernet1 leaf-spine
 5 R170-LFA      ethernet53 R170-LFB      ethernet53 peer
 6 R170-LFA      ethernet54 R170-LFB      ethernet54 peer
 7 R170-LFB      ethernet49 Spine1         ethernet3 leaf-spine
 8 R170-LFB      ethernet50 Spine2         ethernet3 leaf-spine
 9 R171-LFA      ethernet49 Spine1         ethernet5 leaf-spine
10 R171-LFA      ethernet50 Spine2         ethernet5 leaf-spine
11 R171-LFA      ethernet53 R171-LFB      ethernet53 peer
12 R171-LFA      ethernet54 R171-LFB      ethernet54 peer
13 Spine1        ethernet11 R168-LFB      ethernet49 leaf-spine
14 Spine1        ethernet7  R171-LFB      ethernet49 leaf-spine
15 Spine1        ethernet9  R168-LFA      ethernet49 leaf-spine
16 Spine2        ethernet11 R168-LFB      ethernet50 leaf-spine
17 Spine2        ethernet7  R171-LFB      ethernet50 leaf-spine
18 Spine2        ethernet9  R168-LFA      ethernet50 leaf-spine
```

Note: The VMware vSphere virtual standard switches and distributed switches use the Cisco Discovery Protocol (CDP). BCF supports both CDP and LLDP.

ZTF initially configures the switches using link local IPv6 addressing. An IPv4 address is also required to reach external services such as SNMP, NTP, and the Syslog servers found on the management network. The BCF controller automatically assigns the IPv4 addresses, Default Gateway, and DNS server addresses from a defined IPv4 pool. Configuration of this pool is completed under **Fabric > Switches > IP Address Allocation > Configuration**. The following table shows the information used.

Table 11 BCF switch IPv4 allocation

Status	DNS server address	Gateway address	IP range	CIDR prefix length
Enabled	192.168.15.4	192.168.15.1	192.168.15.129 - .254	24

5.3 Tenant and segment configuration

In Big Cloud Fabric (BCF), a tenant is similar in function to a Virtual Routing and Forwarding (VRF) entity. Each tenant establishes a Layer 3 boundary that separates traffic from other tenants through a logical router. Within each tenant, separate segments establish a Layer 2 boundary for each tier. Logical ports are assigned membership based on VLAN IDs. The devices connected to a logical port on a leaf switch are an end-point. For example, the VMware vSphere VMkernels and VMnics are end-points. The following image shows two tenants from the sample deployment and their respective segments.

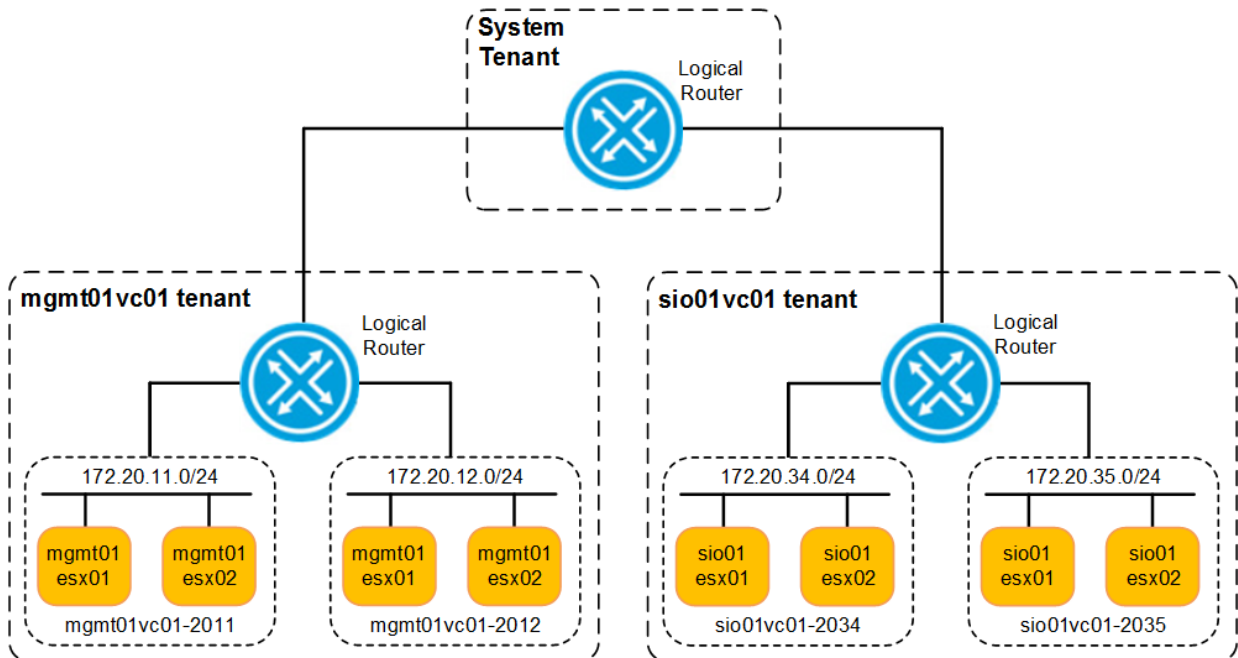


Figure 19 Big Cloud Fabric (BCF) tenant and segments

Separate tenants are used for each vCenter Server. Each tenant, in turn, has separately defined segments. The following table shows these tenants and segments and the mapping between a VLAN ID and the segment name:

Table 12 Big Cloud Fabric (BCF) tenant and segment configuration

Tenant	Tenant function	Logical segment name	VLAN ID	Subnet
mgmt01vc01	Management pod	mgmt01vc01-2011	2011	172.20.11.1/24
		mgmt01vc01-2012	2012	172.20.12.1/24
sio01vc01	ScaleIO and Compute pod	sio01vc01-2031	2031	172.20.31.1/24
		sio01vc01-2032	2032	172.20.32.1/24
		sio01vc01-2033	2033	172.20.33.1/24
		sio01vc01-2034	2034	172.20.34.1/24
		sio01vc01-2035	2035	172.20.35.1/24
management	Management network access	management	2010	172.20.10.1/24

Note: The management tenant is used to extend the existing management network into the BCF pod.

A logical router is automatically assigned to each tenant when it is defined. A tenant has two types of interfaces: tenant interfaces, and a segment interface. The segment interface exists on the tenant logical router and acts as the gateway for the subnet enabling forwarding between segments within a tenant and routing traffic to other tenants through the system tenant. Also, a system-reserved tenant, called the system tenant, has a logical router that supports a tenant interface for each tenant configured in the fabric.

In the following Big Cloud Fabric (BCF) logical router configuration and the System tenant interfaces tables, the tenant interfaces are outlined from the perspective of the tenant as well as the system router. Navigate to **Logical Tenants > Tenant > Segment Interfaces** to define tenant interfaces.

Table 13 Big Cloud Fabric (BCF) logical router configuration

Tenant	Next Hop Tenant	Next Hop Group	Default Route	Next Hop Address
mgmt01vc01	system	Tenant iface system	0.0.0.0/0	n/a
sio01vc01	system	Tenant iface system	0.0.0.0/0	n/a
system	Management	MgmtCoreRouters	0.0.0.0/0	172.20.10.1

Table 14 System tenant interfaces

Tenant	Export Routes
mgmt01vc01	No
sio01vc01	No
management	Yes

Under segment interfaces, the three segment names are shown as well as the respective gateway IP address and subnets (i.e. 172.20.11.1/24). Each of these segments represents a Layer 3 boundary to the tenant. The

logical segments show the Layer 2 segment names and the number of interface groups associated with that tenant. The following image shows the tenant mgmt01vc01:

The screenshot displays two sections: 'Segment Interfaces' and 'Logical Segments'.

Segment Interfaces Table:

Status	State	Segment Name	Segment Group	Description	Private	Subnets	IPv6 Addresses
Up	Active	mgmt01vc01-2011	-	vSphere management segment interface	-	172.20.11.1/24	SLAAC
Up	Active	mgmt01vc01-2012	-	vSphere vMotion segment interface	-	172.20.12.1/24	SLAAC
Up	Active	mgmt01vc01-2013	-	VSAN segment interface	-	172.20.13.1/24	SLAAC

Jul 12, 2017, 20:10:10 GMT

Logical Segments Table:

Name	Tenant	Description	Member VNI	Interface Group Membership Rules	Switch
mgmt01vc01-2011	mgmt01vc01	1 vSphere portgroups: vDS-Mgmt-Management	-	4	
mgmt01vc01-2012	mgmt01vc01	1 vSphere portgroups: vDS-Mgmt-vMotion	-	4	
mgmt01vc01-2013	mgmt01vc01	1 vSphere portgroups: vDS-Mgmt-VSAN	-	4	

Jul 12, 2017, 20:10:11 GMT

Figure 20 Tenant mgmt01vc01 segment interfaces and logical segments

With the configuration in place, there are two primary tenants, one for the Management Pod and one for the ScaleIO and Compute pod and their respective segments. Each segment can reach adjacent segments through the shared tenant logical router. This router can in turn forward traffic to the other tenant router through the system logical router.

5.4 VMware integration

Big Cloud Fabric (BCF) provides an integrated solution for VMware vSphere environments using BCF as the underlying physical network. The following objectives are necessary to complete the VMware integration:

- BCF controller integration
- Connecting VMware vCenter to the Big Cloud Fabric

Completion of the VMware integration provides the following benefits:

- vCenter instance monitoring from the BCF controller
- Ability to install the vCenter BCF controller plug-in

In this example, vCenter management traffic uses BCF to temporarily manage the manual tenant and segment configuration. Initially, each ESXi host requires a single link for the ESXi management. This unique link connects to Leaf-A. See the [Hardware](#) section for the ESXi host wiring diagrams. After established, perform the remaining configuration steps:

1. Create the tenant that hosts the management pod traffic.
2. Create a temporary segment that hosts the management pod traffic.
3. Create interface switch port memberships to connect to the ESXi hosts.
4. Connect the tenant through the system router to the general management tenant that has connectivity to the BCF controller network.
5. Complete vSphere integration with BCF.

In the following Big Cloud Fabric (BCF) tenants and temporary segments table, two vCenter servers used in this sample deployment are listed. Their tenant names, temporary segment names used, and the segment interface IP are listed. The temporary switch port memberships are provided in the following BCF segment switch port memberships table:

Table 15 Big Cloud Fabric (BCF) tenants and temporary segments

vCenter instance	Tenant name	Temp segment name	Segment interface
mgmt01vc01	mgmt01vc01	tempESXiMgmt	172.20.11.1/24
sio01vc01	mgmt01vc01	tempSIOMgmt	172.20.31.1/24

Table 16 Big Cloud Fabric (BCF) segment switch port memberships

Temp segment name	VLAN	Switch	Interface name
tempESXiMgmt	2011	R168-LFA	Ethernet 1
	2011	R168-LFA	Ethernet 3
	2011	R168-LFA	Ethernet 5
	2011	R168-LFA	Ethernet 7
tempSIOMgmt	2031	R170-LFA	Ethernet 2
	2031	R170-LFA	Ethernet 4
	2031	R171-LFA	Ethernet 2
	2031	R171-LFA	Ethernet 4

Note: In ESXi, the system default virtual switch, vSwitch0, is tagged with the corresponding VLAN listed above (2011 or 2031). The Direct Console User Interface (DCUI) is used for VLAN tagging on vSwitch0 during the ESXi installation.

The dashed line between the BCF controller and the first leaf switch is the management tenant that serves as a bridge between the BCF virtual environments and the physical management network. The following image illustrates how the active BCF controller communicates with the temporary segments through the created tenants:

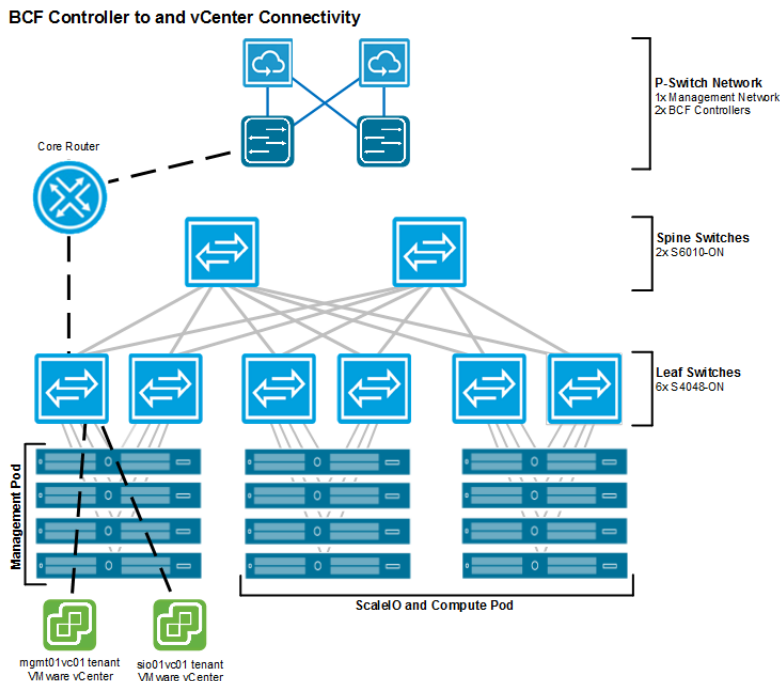


Figure 21 Connecting VMware vCenter to Big Cloud Fabric (BCF) controller through BCF

After initial connectivity is established, all ESXi hosts can now communicate with the BCF controller through the IP fabric. This can be tested by pinging from the BCF controller CLI to the IP address of any ESXi host:

```
bfc01ctrl01> ping mgmt01esx01.dnt.adc.delllabs.net controller-management
PING mgmt01esx01.dnt.adc.delllabs.net (172.20.11.11) 56(84) bytes of data.
64 bytes from mgmt01esx01 (172.20.11.11): icmp_seq=1 ttl=60 time=0.530 ms
64 bytes from mgmt01esx01 (172.20.11.11): icmp_seq=2 ttl=60 time=0.579 ms
64 bytes from mgmt01esx01 (172.20.11.11): icmp_seq=3 ttl=60 time=0.618 ms
64 bytes from mgmt01esx01 (172.20.11.11): icmp_seq=4 ttl=60 time=0.640 ms
64 bytes from mgmt01esx01 (172.20.11.11): icmp_seq=5 ttl=60 time=0.671 ms
```

Or ping a ScaleIO and Compute pod host:

```
bfc01ctrl01> ping sio01esx04 controller-management
PING sio01esx04 (172.20.31.14) 56(84) bytes of data.
64 bytes from sio01esx04 (172.20.31.14): icmp_seq=1 ttl=60 time=0.643 ms
64 bytes from sio01esx04 (172.20.31.14): icmp_seq=2 ttl=60 time=0.626 ms
64 bytes from sio01esx04 (172.20.31.14): icmp_seq=3 ttl=60 time=0.864 ms
64 bytes from sio01esx04 (172.20.31.14): icmp_seq=4 ttl=60 time=0.695 ms
64 bytes from sio01esx04 (172.20.31.14): icmp_seq=5 ttl=60 time=0.661 ms
```

Note: At this point, deploy and configure VMware vCenter before continuing with BCF vSphere integration. See vCenter Server deployment and design for deployment instructions for both of the vCenter servers used in this example.

With both the VMware vCenter and the ESXi management traffic now reachable across the IP fabric, the VMware integration can continue. To start integration, navigate to **Integration > Orchestration > VMware vCenters**.

The following table shows the information used to connect to both vCenter servers. With automation set to **Full**, this allows the BCF configuration to be automatically updated in response to changes on vCenter. The vCenter Plugin Access sets the permission level for the vCenter BCF plug-in. The **Read-Write** option allows the plug-in to be used similarly to the BCF GUI. With the sio01vc01 setting configured to **Read-Only**, security best practices limit administrative access.

Table 17 VMware vCenter connection details

Name	Tenant	Hostname	vCenter plug-in access right	BCF config automation level
mgmt01vc01	mgmt01vc01	mgmt01vc01.dnt.adc.delllabs.net	Read-Write	Full
sio01vc01	sio01vc01	sio01vc01.dnt.adc.delllabs.net	Read-Only	Full

The following figure shows both vCenter Servers connected to the BCF controller. The Operating Mode shows the tenant to have a **Normal** function with no problems shown on this screen.

Name	Operating Mode	Description	Status	Status Detail	Hostname
mgmt01vc01	✓ Normal	—	✓ Connected and authenticated	—	mgmt01vc01.dnt.adc.delllabs.net
sio01vc01	✓ Normal	—	✓ Connected and authenticated	—	sio01vc01.dnt.adc.delllabs.net

Figure 22 Big Switch Fabric vCenter integration status

At this point, the BCF controller configures any remaining segment automatically. For instance, any port groups and VLANs defined by any host attached to vCenter. After, any changes made in either vCenter Server are propagated by the BCF controller to the IP fabric to reflect these changes.

Clicking a vCenter instance displays a summary of the configuration of the current vCenter instance. The following image shows that the sio01vc01 vCenter instance has four hosts, eight distributed switches, 50 endpoints, and six networks.

Info

- Graphic
- Hosts
- Virtual Switches
- Physical Connections
- Endpoints
- Network Host Connection Details
- Networks

Options & Shortcuts

Selected Content Operation

Replace Append

Selected Content Placement

Top of content area

Bottom of content area

In order listed above

Graphic

Summary	Configuration
4 Hosts	Name sio01vc01
8 Virtual Switches	Operational Mode Maintenance <input checked="" type="checkbox"/> Normal
50 Endpoints	Host Name sio01vc01.dnt.adc.delllabs.net
6 Networks	User Name administrator@vsphere.local
	Tenant sio01vc01
	Last Updated Today, 15:23:32 GMT
	Status ✓ Connected and authenticated
	Status Detail —
	Version 6.0.0

Figure 23 sio01vc01 vCenter information

The Virtual Switches screen presents a graphical view of the distributed switch used in the ScaleIO and Compute pod.

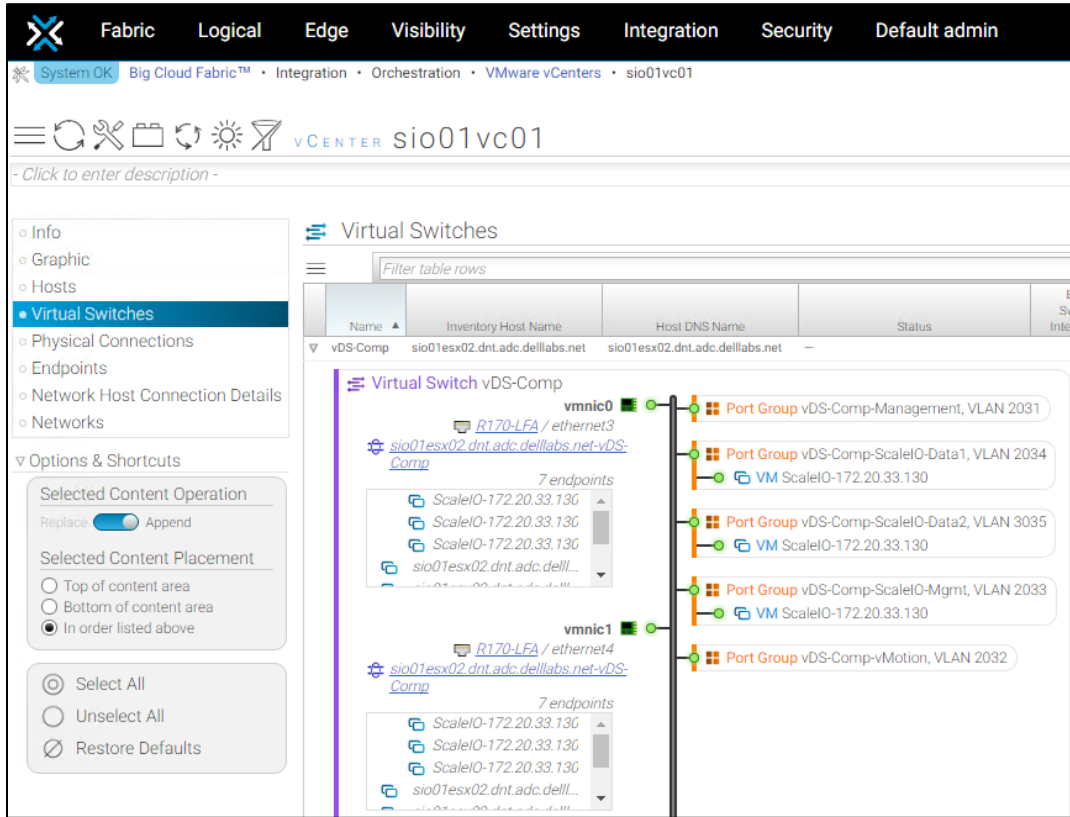


Figure 24 sio01vc01 virtual switch information

Note: See the [Virtual network design](#) section for the VMware vSphere virtual distributed switch configuration.

The vCenter GUI plug-in for VMware vCenter lets you monitor the IP fabric from the vCenter instance. The GUI installation wizard is accessed by going to **Integration > Orchestration > VMware vCenters > Deploy vCenter GUI Plugin**. After installation, the vSphere web client shows the BCF plug-in from the **Home** screen. For more information about the BCF plug-in, see [Big Cloud Fabric User Guide](#).

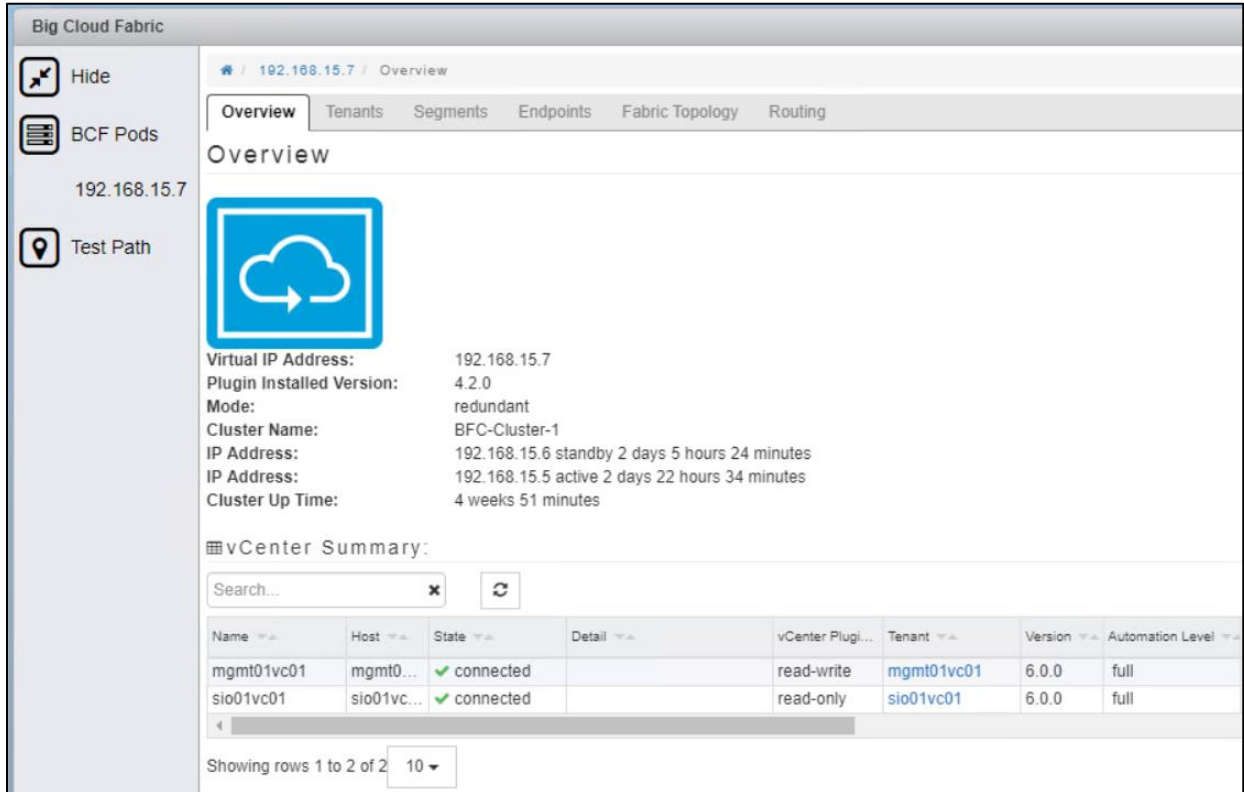


Figure 25 BCF plug-in for vSphere web client

6 Deployment of VMware vSphere

The following image shows a logical representation of the VMware vSphere deployment and how it integrates into the Big Cloud Fabric (BCF) implementation. The mgmt01vc01 and sio01vc01 tenant objects created in the previous section represent the two vCenter servers deployed in the corresponding Management pod, ScaleIO, and the Compute pod. Also shown are the five port groups created on the ScaleIO and Compute virtual distribute switch along with the associated tenant and segments from BCF.

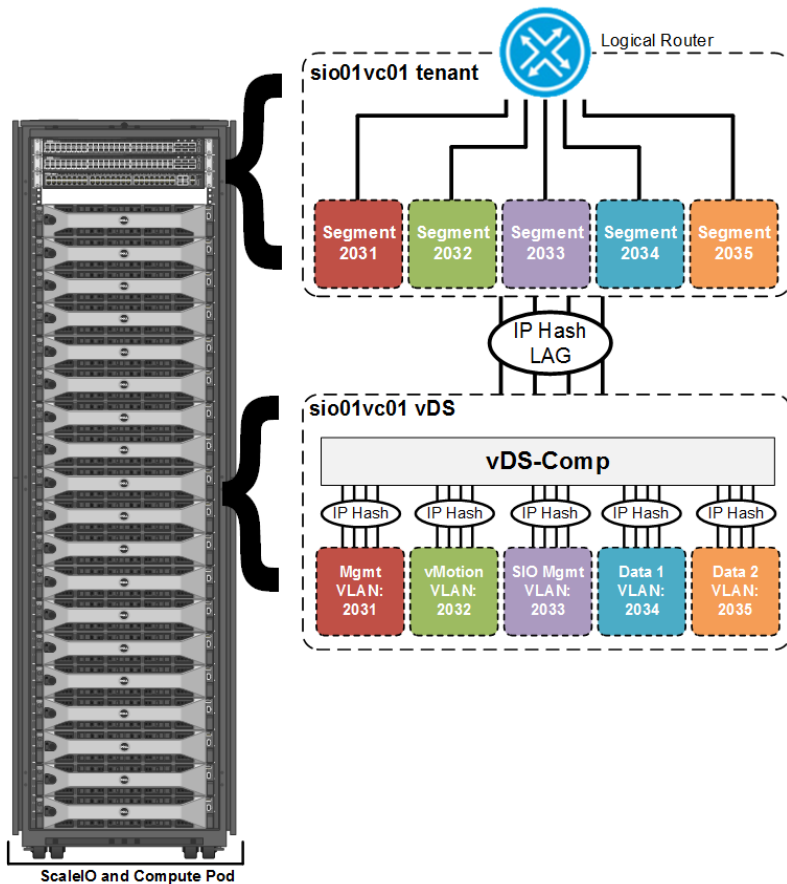


Figure 26 SDD logical design

vSphere is a critical component of the deployment of the software-defined data center (SDDC). Best practices and design decisions in this section follow the guidance outlined in the VMware Validated Designs (VDD) guide. See [VMware Validated Design 4.0 Documentation Center](#) for more information.

This section covers the following topics:

- vCenter Server deployment and design
- Virtual network design

6.1 vCenter server deployment and design

In this deployment example, two vCenter Server appliances are deployed:

- mgmtvc01.dnt.adc.delllabs.net – supports the ESXi hosts that compose the Management pod
- siovc01.dnt.adc.delllabs.net – supports the ESXi hosts that compose the ScaleIO and Compute pod

The deployment of two VMware vCenter servers provides administrative and failure isolation benefits. By dividing along an administrative boundary, a separate security policy can be applied to either of the vCenter servers to reflect administrative functions that would typically be completed by separate organizations. As a secondary benefit, capacity planning for ScaleIO compute workloads is simplified by removing the management workloads from consideration. With this configuration, maintenance becomes easier and the ScaleIO workloads remain available during management workload maintenance windows.

Each vCenter Server is deployed using the Linux-based vCenter Server Appliance (VCSA). A VCSA allows for rapid deployment, enables scalability, and reduces the Microsoft licensing requirements.

The individual vCenter servers deploy with an external Platform Services Controller (PSC) which can be replicated when configured in external mode. With each PSC joined to a single vCenter Single Sign-On domain, the controllers function as a cluster and provide authentication to all components.

vCenter servers are assigned static IP addresses and hostnames during installation and include a valid DNS registration with reverse name resolution. The following table shows the configuration information for the two vCenter Server components and the PSC:

Table 18 vCenter Server component DNS

vCenter Server component	FQDN	IP address
Management PSC	mgmt01psc01.dnt.adc.delllabs.net	172.20.11.61
Management vCenter	mgmt01vc01.dnt.adc.delllabs.net	172.20.11.62
ScaleIO PSC	sio01psc01.dnt.adc.delllabs.net	172.20.11.63
ScaleIO vCenter	sio01vc01.dnt.adc.delllabs.net	172.20.11.64

A vCenter Server has multiple sizing options available for selection during the deployment process. In this example, mgmt01vc01 is built using the small appliance size, while sio01vc01 is built using the large appliance size. The two appliance size characteristics are provided in the following Small vCenter Server Appliance specifications and in the large vCenter Server Appliance specifications tables:

Table 19 Small vCenter Server Appliance specifications

Attribute	Specification
vCenter Server Version	6.0
Appliance Size	Small (up to 100 hosts / 1,000 VMs)
Platform Services Controller	External
Number of CPUs	4
Memory	16 GB
Disk Space	106 GB

Table 20 Large vCenter Server Appliance specifications

Attribute	Specification
vCenter Server Version	6.0
Appliance Size	Large (up to 1,000 hosts / 10,000 VMs)
Platform Services Controller	External
Number of CPUs	16
Memory	32 GB
Disk Space	295 GB

After both vCenter servers are deployed, create a vSphere cluster, one under each vCenter Server. One cluster for management components, and the other for all ScaleIO and Compute pod components. The following table outlines these two clusters, initial host membership, and count:

Table 21 Initial cluster specifications

vCenter Server	Data center	Cluster name	Initial hosts
Management vCenter	mgmt01	mgmt01	mgmt01esx01, mgmt01esx02, mgmt01esx03, mgmt01esx04
ScaleIO vCenter	sio01	sio01	sio01esx01, sio01esx02, sio01esx03, sio01esx04

Note: Initial hosts are truncated and do not include the dnt.adc.delllabs.net domain name.

The following image shows the two vSphere vCenter servers and their associated data centers, clusters, and hosts for each, as well as the virtual machines in the Management pod.

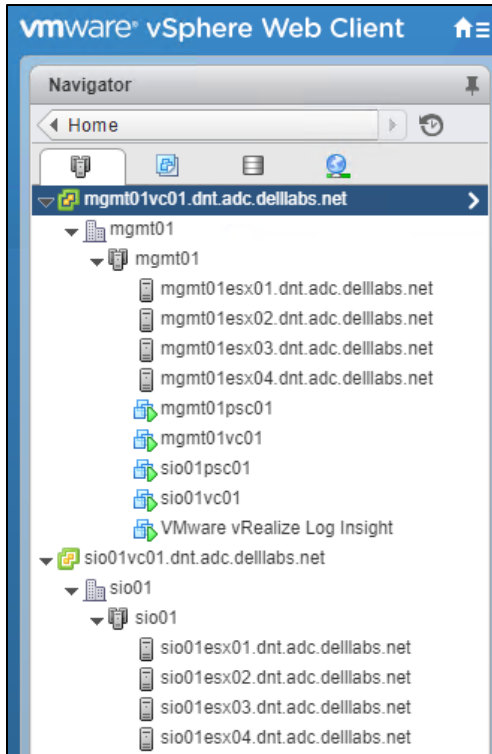


Figure 27 vSphere Cluster design

Protecting vCenter Server systems is important as they are the central point of administration and monitoring. In this example, vSphere High Availability (HA) is enabled on the management cluster to protect both vCenter servers and the PSC.

With vSphere HA enabled on the sio01 cluster, the environment has a robust level of protection for both hosts and virtual machines. Sufficient resources on the remaining hosts are required so that virtual machines can be migrated to those hosts in the event of a host outage. To enable vSphere HA, right-click on the appropriate cluster then select **Settings > Edit vSphere HA properties > vSphere HA**, then click to place a check in the **Turn on vSphere HA** box.

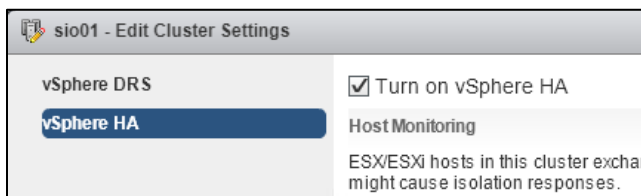


Figure 28 Enable vSphere HA

6.2 Virtual network design

When building the virtual network counterpart to Big Cloud Fabric (BCF), a few principles must be followed to ensure that the design meets a diverse set of requirements while keeping operational complexity to a minimum:

- The separation of network services to achieve greater security and performance by attaching each service to port groups with a different VLAN ID
- The use of network I/O control and traffic shaping that provides bandwidth to critical workloads
- That all virtual machines use the VMXNET3 virtual NIC drivers

Network services require the definition of a well-defined IP Address Management (IPAM) scheme. The following table shows the VLAN IDs and IP subnets for the various traffic types.

Note: Within BCF, each VMware vCenter instance has a separate tenant.

Table 22 VLAN and IP subnet configuration

Pod	VLAN function	VLAN ID	Subnet	Gateway
Management pod	ESXi management	2011	172.20.11.0/24	172.20.11.1
	vSphere vMotion	2012	172.20.12.0/24	172.20.12.1
ScaleIO and Compute pod	ESXi management	2031	172.20.31.0/24	172.20.31.1
	vSphere vMotion	2032	172.20.32.0/24	172.20.32.1
	ScaleIO management	2033	172.20.33.0/24	172.20.33.1
	Data network 1	2034	172.20.34.0/24	n/a
	Data network 2	2035	172.20.35.0/24	n/a

Note: ScaleIO data networks share a TCP/IP stack with the ScaleIO management network and cannot be assigned a default gateway.

Subnet-to-VLAN mapping uses the [RFC1918](#) defined private network 172.20.0.0/12 as the base for all subnets. The second and third octets represent the VLAN ID. For instance, 172.20.33.0/24 would have an associated VLAN ID of 2033. This algorithm ensures that each subnet and VLAN pairing is unique.

For each defined subnet, the first ten (1-10) host addresses are reserved for subnet-specific services. For example, 172.20.11.1 is used as the gateway address for the mgmt01vc01 tenant interface. 172.20.11.11 is the first address assigned to the ESXi host in the subnet for the mgmt01esx01 tenant interface.

The next section provides details regarding VMware vSphere Distributed Switch (VDS). Following best practices, each VMware pod (Management, ScaleIO, and Compute) has a single VDS to keep operational complexity to a minimum. For the Management pod, the VDS is named vDS-Mgmt. For the ScaleIO and Compute pod, the VDS is named vDS-Comp.

In this example, the load balancing algorithm used for all port groups, regardless of the VDS, is IP Hash. BCF automatically creates a Link Aggregation Group (LAG) on leaf switches with all physical links to the host. All uplinks must be active, otherwise, hashing may occur.

Note: See [VMware vCenter Server™ 6.0 Deployment Guide](#) for more information.

6.2.1 VDS-Mgmt configuration details

The following tables contain the pre-installation and post-installation configuration details for the VMware vSphere Distributed Switch (VDS) used for the Management pod:

Table 23 Virtual switch for the Management cluster

vDS name	Function	Network I/O control	Physical NIC port count	MTU
vDS-Mgmt	<ul style="list-style-type: none"> • ESXi management • vSphere vMotion 	Enabled	2	9000

Table 24 vDS-Mgmt port group configuration settings

Parameter	Setting
Failover detection	Link status only
Notify switches	Enabled
Failback	Yes
Failover Order	Active Uplinks: Uplink1, Uplink2

Table 25 vDS-Mgmt port groups and VLANs

VDS	Port group name	Teaming policy	Active uplinks	VLAN ID
vDS-Mgmt	vDS-Mgmt-Management	Route based on IP hash	1, 2	2011
vDS-Mgmt	vDS-Mgmt-vMotion	Route based on IP hash	1, 2	2012

Table 26 vDS-Mgmt by Physical/Virtual NIC

VDS	Physical NIC	Virtual NIC	Uplink
vDS-Mgmt	Intel 82599 10GbE	4	Uplink 1
vDS-Mgmt	Intel X520 10GbE	1	Uplink 2

6.2.2 vDS-Comp configuration details

The following tables contain the pre and post-installation configuration details for the VDS used for the ScaleIO and Compute pod:

Table 27 Virtual switch for the Management cluster

VDS switch name	Function	Network I/O control	Physical NIC port count	MTU
vDS-Comp	<ul style="list-style-type: none"> • ESXi management • vSphere vMotion • ScaleIO management • SIO data network 1 • SIO data network 2 	Enabled	4	9000

Table 28 vDS-Comp port group configuration settings

Parameter	Setting
Failover detection	Link status only
Notify switches	Enabled
Failback	Yes
Failover Order	Active Uplinks: Uplink1, Uplink2, Uplink3, Uplink4

Table 29 vDS-Comp port groups and VLANs

VDS	Port group name	Teaming policy	Active uplinks	VLAN ID
vDS-Comp	vDS-Comp-Management	Route based on IP hash	1, 2, 3, 4	2031
vDS-Comp	vDS-Comp-vMotion	Route based on IP hash	1, 2, 3, 4	2032
vDS-Comp	vDS-Comp ScaleIO-Mgmt	Route based on IP hash	1, 2, 3, 4	2033
vDS-Comp	vDS-Comp-ScaleIO-Data1	Route based on IP hash	1, 2, 3, 4	2034
vDS-Comp	vDS-Comp-ScaleIO-Data2	Route based on IP hash	1, 2, 3, 4	2035

Table 30 vDS-Comp by Physical/Virtual NIC

VDS	Physical NIC	Virtual NIC	Uplink
vDS-Comp	Intel 82599 10 GbE	1	Uplink 1
vDS-Comp	Intel X520 10 GbE	4	Uplink 2
vDS-Comp	Intel 82599 10 GbE	0	Uplink 3
vDS-Comp	Intel X520 10 GbE	3	Uplink 4

6.2.3 VMware vSphere VMkernel configuration

In this section, VMkernels are created and associated with VMware vSphere Distributed Switch (VDS) port groups. A VMkernel provides connectivity to hosts and handles the standard system traffic of VMware vSphere vMotion, the ESXi management, and the ScaleIO data storage traffic.

The following table shows the VMkernel configuration details for vDS-Mgmt. Initially, two VMkernels are defined. During the ESXi installation, a default VMkernel is created called management. The vMotion VMkernel is set up to provide virtual machine mobility across the pod.

Table 31 vDS-Mgmt VMkernel adapters

VDS	Network label	Connected port group	Enabled services	TCP/IP stack	MTU
vDS-Mgmt	Management	vDS-Mgmt-Management	Management traffic	Default	1500
vDS-Mgmt	vMotion	vDS-Mgmt-vMotion	vMotion traffic	vMotion	9000

As a best practice, the VMware vMotion TCP/IP stack is used to isolate traffic for vMotion and allows a dedicated default gateway for vMotion traffic. By using a separate TCP/IP stack, vMotion traffic can be routed using a different default gateway than other ESXi services. Using a separate stack also allows the use of a separate set of buffers and sockets and avoids route table conflicts that may occur otherwise.

The following table contains the configuration details for the ScaleIO and Compute pod. Four VMkernel adapters are assigned. Following best practices, the vMotion TCP/IP stack handles vSphere vMotion traffic. Both ScaleIO Data VMkernels work with the ScaleIO Data Client (SDC) driver installed on ESXi hosts to access the ScaleIO Virtual SAN.

Note: See the Deploy ScaleIO section for more information information about SDC configuration.

Table 32 vDS-Comp VMkernel adapters

VDS	Network label	Connected port group	Enabled services	TCP/IP stack	MTU
vDS-Comp	Management	vDS-Comp-Management	Management traffic	Default	1500
vDS-Comp	vMotion	vDS-Comp-vMotion	vMotion traffic	vMotion	9000
vDS-Comp	ScaleIO-Data1	vDS-Comp-ScaleIO-Data1	None	Default	9000
vDS-Comp	ScaleIO-Data2	vDS-Comp-ScaleIO-Data2	None	Default	9000

The following image shows the completed topology of vDS-Comp for the ScaleIO and Compute pod showing port groups and VLAN assignments, VMkernels and IP addresses, and physical NIC uplinks.

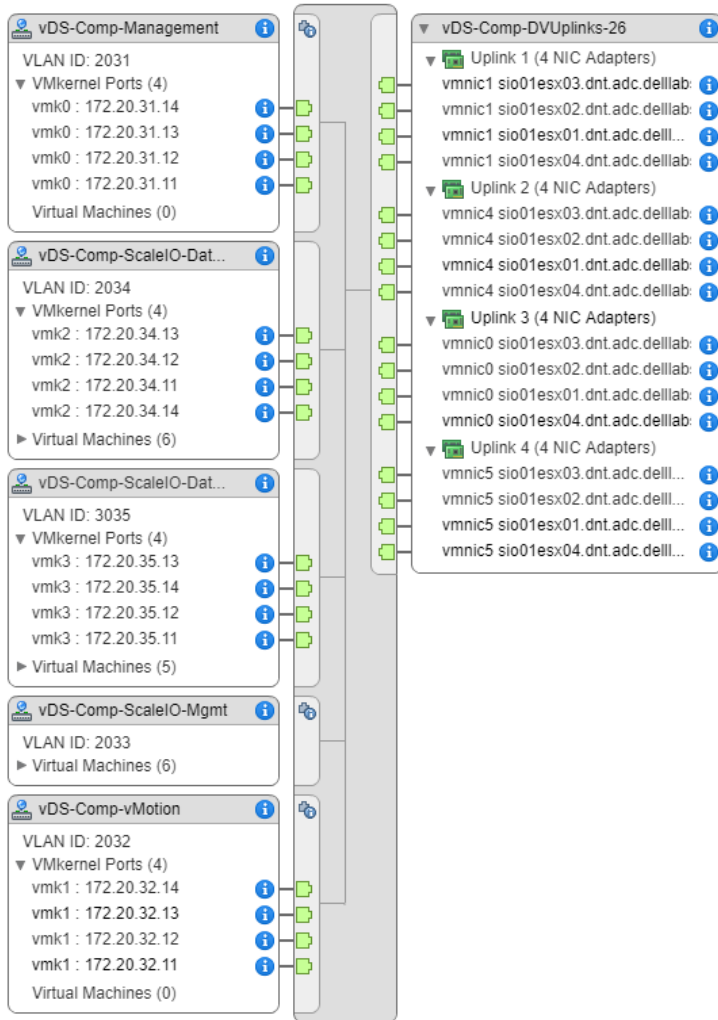


Figure 29 vDS-Comp distributed switch topology

7 Deploying Dell EMC ScaleIO

Dell EMC ScaleIO is a software-only solution that uses the local disks and LAN of an existing server to create a virtual SAN with all the benefits of external storage at a fraction of the cost and complexity. ScaleIO uses the existing local storage devices and converts them into shared block storage. For many workloads, the ScaleIO storage is comparable to, or better than, external shared block storage.

The lightweight ScaleIO software components are installed on the ScaleIO and Compute pod nodes and communicate via Big Cloud Fabric (BCF) to handle the application I/O requests sent to the ScaleIO block volumes. An efficient decentralized block I/O flow combined with a distributed, sliced volume layout, results in a massively parallel I/O system that can scale up to thousands of nodes.

Because ScaleIO is hardware agnostic, the software works efficiently with various types of disks, including:

- Hard Disk Drives (HDD)
- Solid-state Disks (SSD)
- Flash PCI Express (PCIe) cards

Three key software components comprise ScaleIO:

- Meta Data Manager (MDM)
- ScaleIO Data Server (SDS)
- ScaleIO Data Client (SDC)

The MDM configures and monitors the ScaleIO system. To support high availability, three or more instances of MDM run on different nodes. In a multiple MDM environment, one MDM is given the Master role, and the others act as Slave or Tiebreaker MDMs.

The MDM cluster has the following requirements and options:

- A minimum of three nodes are installed with the MDM package
- The MDM package installs with the SDS package in the same VM

The following image shows a 5-node MDM deployment: The MDM cluster has three copies of the repository and can withstand two MDM cluster member failures.



Figure 30 5-node MDM cluster

The SDS manages the capacity of a single server and acts as a backend for data access. The SDS is installed on all nodes contributing storage devices to the ScaleIO system. The SDS is a virtual machine running on top of ESXi with all local storage disks mapped to the virtual machine.

The SDC is a lightweight device driver that exposes the ScaleIO volumes as block devices to the application and resides on the same server as the SDS.

In addition to the software components, the hardware that ScaleIO ultimately runs on top of is also important. While ScaleIO is hardware agnostic, Dell EMC has developed a ScaleIO Ready Node that is a converged hardware and software architecture. It uses the VMware Hypervisor with the ScaleIO solution, supplying a converged virtualization option.

Automated Management Services (AMS) manages Ready Nodes and is a dedicated Management Server that resides outside of the ScaleIO Ready Node converged system and serves GUI, CLI, and REST clients. The AMS enables the full range of the ScaleIO Ready Node configuration: Protection Domain and Storage Pool configuration, volume and snapshot management, addition and removal of servers, hardware monitoring, alerts, and access to a read/write cache model using SanDisk DAS Cache.

AMS enables the following:

- Graphical user interface (GUI) client
- Command Line Interface (CLI) client - SNMP Command Line Interface (SCLI) and Automated Management Services (AMS) Command Line Interface (CLI)
- Representational State Transfer (REST) application programming interface (API)

AMS enables reporting via the EMC proprietary EMC Secure Remote Support (ESRS), SNMP, and Syslog. The ScaleIO Ready Node gives administrators a fully converged solution that is easy to deploy, configure, and upgrade. The following image shows a ScaleIO Ready Node:



Figure 31 Dell EMC ScaleIO Ready Node

Note: Dell EMC recommends the use of ScaleIO Ready Node for the ScaleIO implementation. The ScaleIO Ready Node provides a consistent hardware configuration, reduces operational complexities, and creates a stable environment. For more information, see [Dell EMC ScaleIO Ready Node](#).

ScaleIO is installed in an existing infrastructure and in greenfield configurations. Deploying ScaleIO in this environment consists of the following topics:

- Register the ScaleIO plug-in
- Upload the ScaleIO Open Virtual Appliance (OVA) template
- Deploy ScaleIO

7.1 Deploy Dell EMC ScaleIO plug-in

At this point, all physical connectivity and configuration of the VMware vSphere distributed switches are complete. Also, all the VLAN and port groups for the Dell EMC ScaleIO have automatically propagated into the IP fabric. The ScaleIO plug-in for vSphere simplifies the installation and management of the ScaleIO system in a vSphere environment.

Note: For more information about registering the ScaleIO plug-in, see the ScaleIO Software Only: Documentation Library and access the [ScaleIO v2.0.x Deployment guide](#).

In this sample deployment, the following software versions were used to deploy the ScaleIO plug-in:

Table 33 Dell EMC ScaleIO setup software versions

Product	Version
VMware vSphere PowerCLI	6.3.0 R1 Patch 1
ScaleIO vSphere Plug-in Installer	2.0.1.3
ScaleIO virtual machine OVA	2.0.13000.211.ova

Use the following parameters during the installation:

Table 34 Dell EMC ScaleIO VMware vSphere plug-in parameters

Parameter	Setting
vCenter Server	sio01vc01.dnt.adc.delllabs.net
Registration mode	Standard

The ScaleIO vCenter instance, sio01vc01, provides a separate administrative domain. Standard registration mode deploys an embedded Tomcat web server on the workstation for the target vCenter Server to pull the plug-in from on next login. The following image shows the available ScaleIO vCenter plug-in after login:

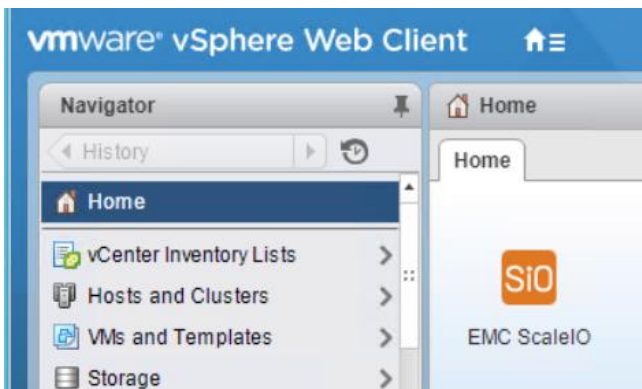


Figure 32 VMware vSphere Dell EMC ScaleIO plug-in

7.2 Upload Dell EMC ScaleIO OVA template

Once the VMware vCenter Dell EMC ScaleIO plug-in installation is completed, upload the ScaleIO virtual machine Open Virtual Appliance (OVA). The OVA serves as a virtual machine template for deploying all the software components of ScaleIO. From the VMware PowerCLI program, select **Create SVM Template**. During the installation process, use the parameters shown in the following table:

Table 35 Dell EMC ScaleIO VMware vSphere plug-in parameters

Parameter	Setting
vCenter Server	sio01vc01.dnt.adc.delllabs.net
Data center name	sio01
Path to OVA	<<local OVA path including file name>>
Datastore Names	LDS-SIO01ESX01

Note: While it is possible to specify a datastore that corresponds to each ScaleIO host, this sample deployment uses a single datastore. VMware vSphere vMotion is used to copy the template to each remaining host during the ScaleIO implementation.

7.3 Deploy Dell EMC ScaleIO

The information in this section describes how the deployment wizard is used in the deployment example provided in this guide. For detailed information on each deployment step and how to apply changes based on a specific topology or hardware configuration, see [EMC ScaleIO 2.0.X Deployment guide](#).

ScaleIO deployment has four separate steps:

- ScaleIO Data Client (SDC) deployment and configuration
- ScaleIO advanced configuration settings
- Deploy the ScaleIO environment
- Install the ScaleIO GUI (optional)

Before an ESXi host can consume the Virtual SAN, ScaleIO provides, a kernel driver must be installed on each ESXi host, regardless of the role that host is playing. The process installs the SDC driver on the target host and which time the host is rebooted.

To start the installation wizard, perform the following steps:

1. From the **Basic tasks** section of the **EMC ScaleIO** screen, click **Install SDC on ESX**.
2. Select all hosts under the **sio01** data center as targets for the installation.
3. Once complete, reboot all hosts before continuing with the deployment.

Before using the deployment wizard, use the **Advanced Settings** link to **Enable RDMs on nonparallel SCSI controllers** (check the box). The Remote Device Mapping (RDM) setting enables non-SCSI controller devices as RDM devices.

Note: Do not select this option if the device does not support SCSI Inquiry Vital Data Product (VPD) page code 0x83.

To deploy ScaleIO, perform the following steps:

1. From the **Basic tasks** section of the **EMC ScaleIO** screen, click **Deploy ScaleIO environment**.
2. Using the following table, assign the settings listed to the parameters provided.

Note: The parameters and settings provided in the table address the selections necessary through step 4 of the installation wizard. A setting that is not listed indicates that the default setting has been applied.

Table 36 ScaleIO deployment settings, part 1

Parameter	Setting
Select installation	Create a new system
System name	SIO01
Admin password	ScaleIO Admin Password
vCenter server	sio01vc01.dnt.adc.delllabs.net
Host selection	sio01esx01, sio01esx02, sio01esx03, sio01esx04
ScaleIO components	3-node mode
Initial Master MDM	sio01esx01
Manager MDM	sio01esx02
TieBreaker MDM	sio01esx03
DNS Server 1	172.20.10.4
DNS Server 2	172.20.10.5

- Using the following table, select the ScaleIO wizard parameter settings for steps 5 through 7.

Table 37 ScaleIO deployment settings, part 2

Parameter	Setting
Protection domain name	PD1
RAM read cache size per SDS	1,024 MB
Storage pools	HDD01, SSD01
Enable zero padding	True
SDS host selection	sio01esx01, sio01esx02, sio01esx03, sio01esx04
Selected devices	All empty device categorized into the appropriate storage pool.
SDC host selection	sio01esx01, sio01esx02, sio01esx03, sio01esx04
Enable/Disable SCSI LUN	Enable

- Using the following table, select the ScaleIO wizard parameter settings to complete the wizard setup:

Table 38 ScaleIO deployment settings, part 3

Parameter	Setting
Host for ScaleIO gateway	sio01esx04
Gateway admin password	ScaleIO Admin Password
Gateway LIA password	ScaleIO Admin Password
Select OVA template	EMC ScaleIO SVM Template (v2.0.130000.211) 1
OVA root password	ScaleIO Admin Password
OVA LIA password	ScaleIO Admin Password
Management network label	vDS-Comp-ScaleIO-Mgmt
Data network label	vDS-Comp-ScaleIO-Data1
2nd data network label	vDS-Comp-ScaleIO-Data2

Table 39 ScaleIO networking addressing

ESX name	Management IP	Default gateway	Data 1 IP	Data 2 IP
sio01esx04 (ScaleIO Gateway)	172.20.33.11/24	172.20.33.1	172.20.34.11/24	172.20.35.11/24
sio01esx01 (Master MDM)	172.20.33.12/24	172.20.33.1	172.20.34.12/24	172.20.35.12/24
sio01esx02 (Slave 1 MDM)	172.20.33.13/24	172.20.33.1	172.20.34.13/24	172.20.35.13/24
sio01esx03 (TieBreaker 1)	172.20.33.14/24	172.20.33.1	172.20.34.14/24	172.20.35.14/24
sio01esx04	172.20.33.15/24	172.20.33.1	172.20.34.15/24	172.20.35.15/24

Table 40 Scale IO networking virtual IP addresses

Parameter	Setting
Data (vDS-Comp-ScaleIO-Data1)	172.20.34.4
2nd Data (vDS-Comp-ScaleIO-Data2)	172.20.35.4

Once the summary screen displays, the deployment begins. During deployment, the SVM template that was previously uploaded, creates each ScaleIO virtual machine.

7.4 Dell EMC ScaleIO GUI

The Dell EMC ScaleIO graphical user interface (GUI) can be installed on the management workstation to provide an easy way to monitor and configure the ScaleIO system. Once installed, an IP address or hostname of one MDM-enabled SVM can be used to access the host. The installation file is part of the ScaleIO for Windows download. The following image shows the ScaleIO GUI during the testing phases of a sample deployment:

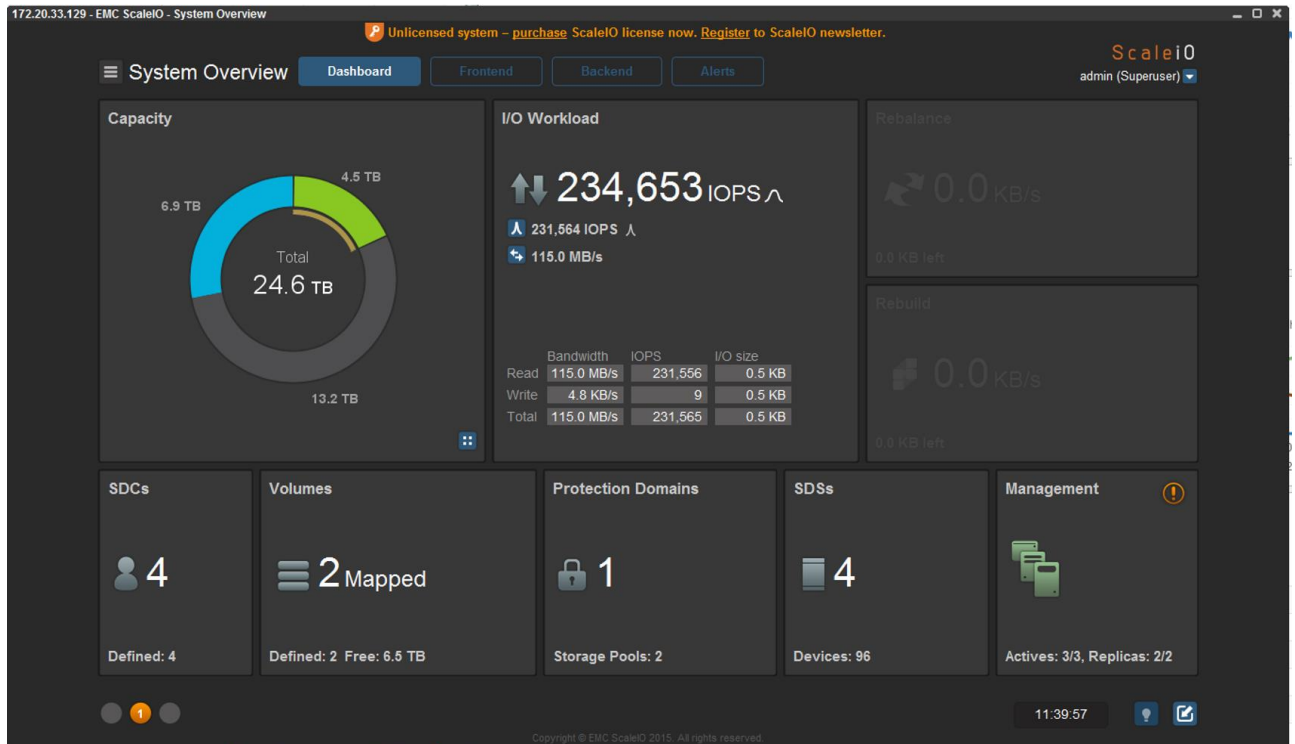


Figure 33 ScaleIO 4 node cluster under load

Note: For more information about using the Dell EMC ScaleIO GUI, see the ScaleIO Software Only: Documentation Library and access the [ScaleIO v2.0.x User Guide](#).

8 Performance tuning

The post-installation information provided in this section consists of the following:

- Increase the Maximum Transmission Unit (MTU) for Big Cloud Fabric (BCF), VMware vSphere, and Dell EMC ScaleIO
- Implementation of Quality of Service (QoS) for the IP fabric
- Configure Network I/O Control (NIOC) for VMware vSphere

8.1 Maximum Transmission Unit size

Maximum Transmission Unit (MTU) refers to the maximum packet size allowed over a network. The default Ethernet size is approximately 1,500 bytes with a maximum of approximately 9,000 bytes. In an IP storage network like Dell EMC ScaleIO, increasing the MTU to the maximum allowed size decreases CPU utilization due to decreased number of frames needed to complete a similarly sized workload. The following image illustrates this concept.

Big Cloud Fabric (BCF) supports an MTU value up to 9,216 based on the LLDP/CDP values received from the VMware vSphere virtual switches. No additional configuration is required in BCF to support the increased MTU value.

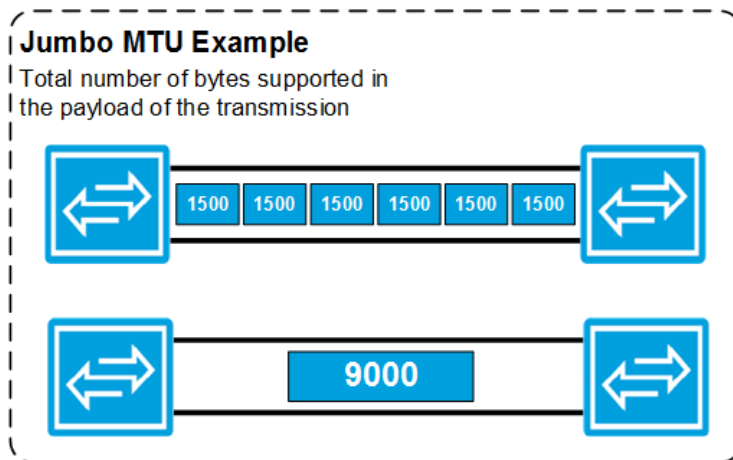


Figure 34 Example of standard frames versus jumbo frames

In this environment, both distributed switches, vDS-Mgmt and vDS-Comp, are assigned an MTU value of 9000. Also, any storage-related port group has an MTU value of 9000. The following table summarizes the port groups that have an MTU value of 9000:

Table 41 VDS port groups with modified MTU value

VDS	Network label	Connected port group	MTU
vDS-Mgmt	vMotion	vDS-Mgmt-vMotion	9000
vDS-Comp	vMotion	vDS-Comp-vMotion	9000
vDS-Comp	ScaleIO-Data1	vDS-Comp-ScaleIO-Data1	9000
vDS-Comp	ScaleIO-Data2	vDS-Comp-ScaleIO-Data2	9000

To verify that jumbo frames are working in the environment, the ESXi CLI tool `vmkping` is used. After establishing an SSH connection with `sio01esx01.dnt`, a non-defragment capable ping with an MTU value of 8972 is sent from the host using the ScaleIO-Data1 VMkernel adapter to `sio01esx02`.

Note: The maximum frame size that vmkping can send is 8972 because of the Ethernet encapsulation.

```
[root@sio01esx01:~] vmkping -d -s 8972 -I vmk2 172.20.34.12
PING 172.20.34.12 (172.20.34.12): 8972 data bytes
8980 bytes from 172.20.34.12: icmp_seq=0 ttl=64 time=0.360 ms
8980 bytes from 172.20.34.12: icmp_seq=1 ttl=64 time=0.373 ms
8980 bytes from 172.20.34.12: icmp_seq=2 ttl=64 time=0.451 ms
```

When the ScaleIO domain installation is successful, each SDS virtual machine is modified to enable Jumbo Frames. To enable Jumbo Frames for the SDS virtual machines, perform the following steps:

1. Run the `ifconfig` command to get the NIC information. The following is an example from a ScaleIO SDS deployed in this solution, ScaleIO-172.20.33.129:

```
ScaleIO-172-20-33-129:~ # ifconfig eth1 && ifconfig eth2
eth1      Link encap:Ethernet  HWaddr 00:50:56:80:6F:5A
          inet addr:172.20.34.129  Bcast:172.20.34.255  Mask:255.255.255.0
          inet6 addr: fe80::250:56ff:fe80:6f5a/64 Scope:Link
          UP BROADCAST RUNNING MULTICAST  MTU:1500  Metric:1
          RX packets:237074680 errors:0 dropped:6 overruns:0 frame:0
          TX packets:234992149 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:1000
          RX bytes:80461295825 (76733.8 Mb)  TX bytes:79612175717 (75924.0 Mb)

eth2      Link encap:Ethernet  HWaddr 00:50:56:80:D2:2C
          inet addr:172.20.35.129  Bcast:172.20.35.255  Mask:255.255.255.0
          inet6 addr: fe80::250:56ff:fe80:d22c/64 Scope:Link
          UP BROADCAST RUNNING MULTICAST  MTU:1500  Metric:1
          RX packets:237529290 errors:0 dropped:7 overruns:0 frame:0
          TX packets:235409078 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:1000
          RX bytes:83165001180 (79312.3 Mb)  TX bytes:78017648972 (74403.4 Mb)
```

In this example, `eth1` and `eth2` correspond to ScaleIO Data Network 1 (Subnet 172.20.34.0/24, VLAN 2034) and ScaleIO Data Network 2 (Subnet 172.20.35.0/24, VLAN 2035).

2. Using the interface name, edit the appropriate network configuration files, and append MTU value of 9000 to the end of the configuration. The following is an example for interface `eth1`:

```
ScaleIO-172-20-33-129:~ # vi /etc/sysconfig/network/ifcfg-eth1
DEVICE=eth1
STARTMODE=onboot
USERCONTROL=no
BOOTPROTO=static
NETMASK=255.255.255.0
IPADDR=172.20.34.129
MTU=9000
```

3. Save the file (:qw [ENTER]) then enter the following command to restart the network services for the virtual machine:

```
ScaleIO-172-20-33-129:~ # service network restart
Shutting down network interfaces:
  eth0      device: VMware VMXNET3 Ethernet Controller      done
  eth1      device: VMware VMXNET3 Ethernet Controller      done
  eth2      device: VMware VMXNET3 Ethernet Controller      done
Shutting down service network . . . . . done
Hint: you may set mandatory devices in /etc/sysconfig/network/config
Setting up network interfaces:
  eth0      device: VMware VMXNET3 Ethernet Controller
  eth0      IP address: 172.20.33.129/24                      done
  eth1      device: VMware VMXNET3 Ethernet Controller
  eth1      IP address: 172.20.34.129/24                      done
  eth2      device: VMware VMXNET3 Ethernet Controller
  eth2      IP address: 172.20.35.129/24                      done
Setting up service network . . . . . done
```

4. Use the ping command to validate jumbo frames connectivity to another, already-configured, SDS virtual machine:

Note: The maximum frame size that vmkping can send is 8972 because of Ethernet frame overhead.

```
ScaleIO-172-20-33-129:~ # ping -M do -s 8972 172.20.34.130
PING 172.20.34.130 (172.20.34.130) 8972(9000) bytes of data.
8980 bytes from 172.20.34.130: icmp_seq=1 ttl=64 time=0.393 ms
8980 bytes from 172.20.34.130: icmp_seq=2 ttl=64 time=0.398 ms
8980 bytes from 172.20.34.130: icmp_seq=3 ttl=64 time=0.366 ms
```

Note: For further information on performance tuning ScaleIO, see [ScaleIO v2.0.x Performance Fine-Tuning Technical Notes](#).

8.2 Quality of Service

Big Cloud Fabric (BCF) provides three ways to prioritize traffic through Quality of Service (QoS):

- Segment-based QoS
- Differentiated Service Code Point (DSCP)-based QoS
- Priority Flow Control (PFC)

Note: For information on Big Cloud Fabric (BFC) and Quality of Service (QoS), see [BCF 4.2.0 User Guide](#). Big Switch documentation requires a customer account to access. Contact your Big Switch Networks account representative for assistance.

In this design, segment-based QoS is used to assign queue values due to the extensive use of VLANs used. After enabling QoS and assigning a traffic class to a segment, all traffic received on that segment is allocated to the associated queue. Segment-based allocation provides an end-to-end QoS solution for the fabric that ensures different segments are guaranteed the appropriate bandwidth during contention.

When QoS is disabled by default on the fabric, 95% of the available bandwidth is allocated to Queue 0, and Queue 5 is assigned 5% for span-fabric traffic. Queue 7 and Queue 8 are for in-band and controller management traffic and receive strict priority.

When QoS is enabled, four additional queues are enabled, with the default weight assigned to each queue as follows:

- Queue 0 (Traffic Class 0): 10
- Queue 1 (Traffic Class 1): 20
- Queue 2 (Traffic Class 2): 30
- Queue 3 (Traffic Class 3): 30
- Queue 4 (PFC Traffic Class 4): 5

To enable QoS, go to **Settings > QoS** and then change the **Enabled** slider to **Y** (Yes). By default, segment-based QoS is used.

The system displays the dialog as shown in the following image:

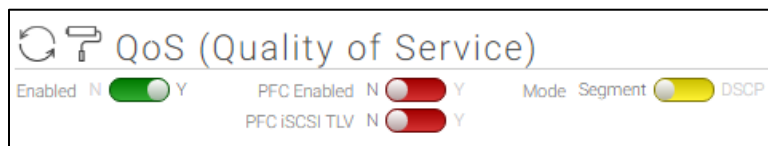


Figure 35 Enabling QoS

The Traffic Classes table at the bottom of the page displays the number of segments assigned to a traffic class. The following figure shows Traffic Class 0 as the only traffic class with segments:

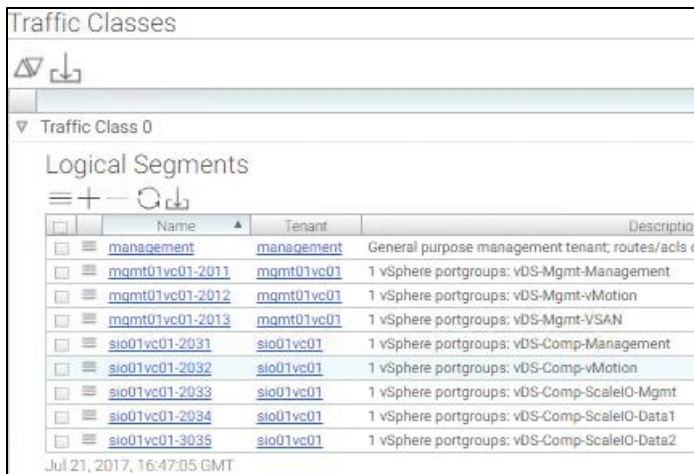


Figure 36 Traffic classes

Note: VMware VSAN was used for the Management Pod and is beyond the scope of this document. See [Dell EMC TechCenter Networking Guides](#) for more information.

In this sample deployment, define a custom queueing profile. The QoS management page selects the add symbol above **Queue Profiles**. The following table provides Traffic Class weights:

Table 42 Create QoS Queue Profile

Name	Traffic Class 0 weight	Traffic Class 1 weight	Traffic Class 2 weight	Traffic Class 3 weight	Traffic Class PFC weight	Traffic Class span fabric weight
Queue Profile	25	40	5	20	5	5

In this example, the ScaleIO Management traffic is assigned to Traffic Class 1. The following figure shows how the segment's traffic class is modified:

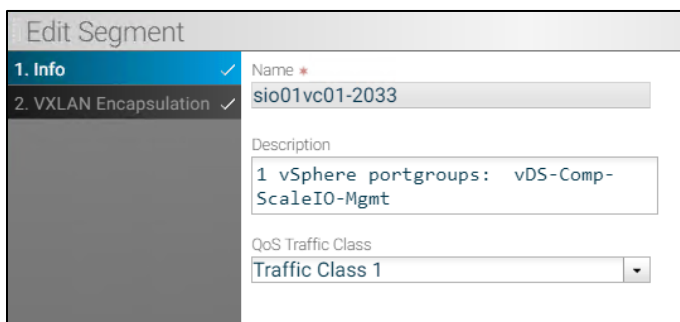


Figure 37 Assigning ScaleIO Management traffic to Traffic Class 1

A new segment is automatically assigned to Traffic Class 0 by default. As a best practice, reserve a minimum of 25% available bandwidth for the default traffic class. Multi-tenant traffic is initially reserved and used in

Traffic Class 3. The following table shows the Traffic Class assignments by tenant/segment in this sample deployment:

Table 43 Tenant and Segment Traffic Class mapping

Tenant	Segment	Purpose	Traffic Class
management	Management	Management	2
mgmt01vc01	mgmt01vc01-2011	ESXi Management	2
	mgmt01vc01-2012	Management vSphere vMotion	0
sio01vc01	sio01vc01-2031	ESXi Management	2
	sio01vc01-2032	ScaleIO vSphere vMotion	0
	sio01vc01-2033	ScaleIO Management	1
	sio01vc01-2034	ScaleIO Data 1	1
	sio01vc01-2035	ScaleIO Data 2	1

The following image shows the nine total segments assigned to the appropriate Traffic Classes:

The screenshot shows a web interface titled "Traffic Classes". It contains a table with two columns: "Name" and "Segments". The table lists several traffic classes and their corresponding segment counts:

Name	Segments
▶ Traffic Class 0	3
▶ Traffic Class 1	3
Traffic Class 2	0
▶ Traffic Class 3	3
Priority-Based Flow Control (PFC) Traffic Class	0
Span Fabric Traffic Class	0

At the bottom of the interface, there is a timestamp "Jul 21, 2017, 17:49:02 GMT" and a "Show:" dropdown menu with options "10", "25", "100", and "All", followed by "(1 - 6 / 6)".

Figure 38 Assigned Traffic Classes

8.3 Network I/O Control

In VMware vSphere, Network I/O Control (NIOC) enforces the share value specified for the different traffic types only when there is a network contention event. When contention occurs, NIOC applies the share values set to each traffic type. As a result, less relevant traffic, as defined by the share percentage, is throttled allowing more important traffic types to gain access to more network resources.

NIOC allows either shares or limits for bandwidth allocation restriction. It is a best practice to use shares instead of limits. Limits impose hard restrictions on the amount of bandwidth traffic flows utilizes, even when network bandwidth is available. To locate the configuration, navigate to the VDS and select **Manage > Settings > Resource Allocation**. The following figure shows the default service enabled when a new distributed switch is created:

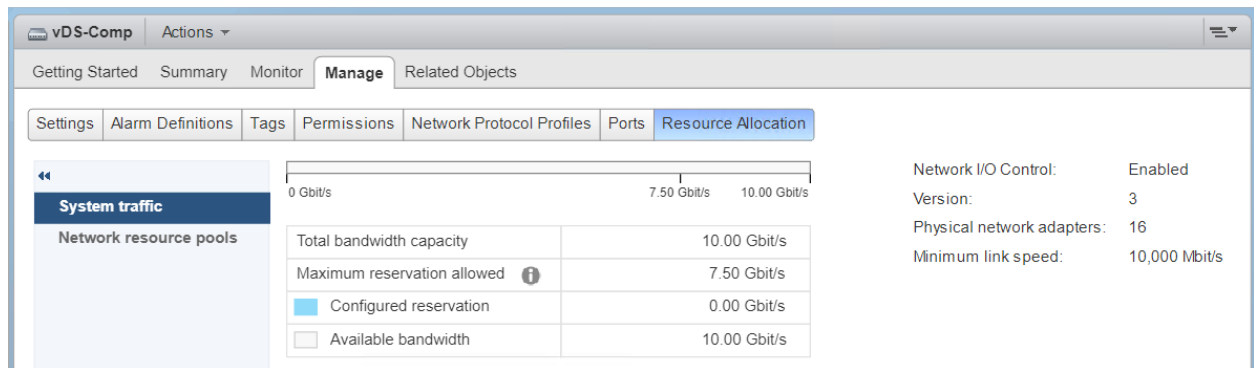


Figure 39 vDS-Comp resource allocation

In this sample deployment, NIOC and Quality of Service (QoS) are both deployed. The result is an environment that increases resiliency and performance of the network. The following Network I/O control traffic type settings table shows the values configured for both vDS-Mgmt and vDS-Comp distributed switches. Management traffic is set to a **Normal** value while **Virtual Machine Traffic** is given a **High** value due to storage traffic being handled by ScaleIO virtual machines.

Table 44 Network I/O control traffic type settings

Traffic type	Value
vSAN traffic	Low
NFS traffic	Low
vMotion traffic	Low
vSphere replication traffic	Low
Management traffic	Normal
vSphere Data Protection Backup traffic	Low
Virtual Machine traffic	High
Fault Tolerance traffic	Low
iSCSI traffic	Normal

A Configuration details

Table 1 Component table example

Product group	Product name	Component	Product version
Dell EMC Networking	S3048-ON	CPLD	8
		DNOS	9.11.2.1
	S4048-ON	CPLD	15.12.5
		ONIE	3.21.1.2
	S6010-ON	CPLD	12.12.5
		ONIE	3.26.1.0
	Z9100-ON	CPLD	5.4.4.4
		ONIE	3.23.1.3
Dell EMC PowerEdge	R630	BIOS	2.2.5
		PERC	13.17.03.00
		iDRAC	2.40.40.40
	R730	BIOS	2.2.5
		PERC	13.17.03.00
		iDRAC	2.40.40.40
Big Switch Networks	Big Switch Fabric		4.2.0
	Big Switch Light		4.2.0
Dell EMC ScaleIO	Advanced Management Service		2.0.1.2
	EMC ScaleIO		2.0.1.2
VMware vSphere Enterprise Plus	ESXi		6.0 U3
	vCenter Server Appliance		6.0
VMware vRealize Log Insight	vRealize Log Insight		4.3.0
	vRealize Log Insight Content Pack for BSN		1.0

Note: The Dell EMC Z9100-ON was tested briefly during the deployment of this solution as a spine switch using BCF 4.2. The Z9100-ON could be substituted for the S6010-ON to provide 25/100GbE.

B Technical support and resources

[Dell.com/support](https://dell.com/support) focuses on meeting customer needs with proven services and support.

[Dell TechCenter](https://delltechcenter.com) is an online technical community where IT professionals have access to numerous resources for Dell EMC software, hardware, and services.

B.1 Dell EMC product manuals and technical guides

[Manuals and documentation for Dell Networking S3048-ON](#)

[Manuals and documentation for Dell Networking S4048-ON](#)

[Manuals and documentation for Dell Networking Z9100-ON](#)

[Manuals and Documentation for PowerEdge R730xd](#)

[Manuals and documentation for PowerEdge R630](#)

[Dell EMC ScaleIO Software Only: Documentation Library](#)

[Dell EMC ScaleIO Ready Node \(Dell\) with AMS: Documentation Library](#)

B.2 Dell EMC Solution Briefs

[Enabling Modern Data Centers with Hyperscale Networking](#)

B.3 Big Switch Networks product manuals and technical guides

[Big Cloud Fabric: A Next-Generation Data Center Switching Platform](#)

[Big Switch Networks support portal](#)

[Big Switch Networks + Dell: Ideal SDN Fabric for VMware SDDC](#)

B.4 VMware product manuals and technical guides

[VMware vSphere 6.0 Documentation Center](#)

[VMware vCenter Server 6.0 Deployment Guide](#)

[VMware Validated Design 4.0](#)

C Support and feedback

Contacting Technical Support

Support Contact Information

Web: <http://Support.Dell.com/>

Telephone: USA: 1-800-945-3355

Feedback for this document

We encourage readers to provide feedback on the quality and usefulness of this publication by sending an email to Dell_Networking_Solutions@Dell.com.