



Configuring MPLS Layer 3 VPN Load Balancing

This chapter describes how to configure load balancing for Multiprotocol Label Switching (MPLS) Layer 3 virtual private networks (VPNs) on Cisco Nexus 9508 switches.

- [Information About MPLS Layer 3 VPN Load Balancing, on page 1](#)
- [Prerequisites for MPLS Layer 3 VPN Load Balancing, on page 6](#)
- [Guidelines and Limitations for MPLS Layer 3 VPN Load Balancing, on page 6](#)
- [Default Settings for MPLS Layer 3 VPN Load Balancing, on page 7](#)
- [Configuring MPLS Layer 3 VPN Load Balancing, on page 8](#)
- [Configuration Examples for MPLS Layer 3 VPN Load Balancing, on page 11](#)

Information About MPLS Layer 3 VPN Load Balancing

Load balancing distributes traffic so that no individual router is overburdened. In an MPLS Layer 3 network, you can achieve load balancing by using the Border Gateway Protocol (BGP). When multiple iBGP paths are installed in a routing table, a route reflector advertises only one path (next hop). If a router is behind a route reflector, all routes that are connected to multihomed sites are not advertised unless a different route distinguisher is configured for each virtual routing and forwarding instance (VRF). (A route reflector passes learned routes to neighbors so that all iBGP peers do not need to be fully meshed.)

iBGP Load Balancing

When a BGP-speaking router configured with no local policy receives multiple network layer reachability information (NLRI) from the internal BGP (iBGP) for the same destination, the router chooses one iBGP path as the best path and installs the best path in its IP routing table. iBGP load balancing enables the BGP-speaking router to select multiple iBGP paths as the best paths to a destination and to install multiple best paths in its IP routing table.

eBGP Load Balancing

When a router learns two identical eBGP paths for a prefix from a neighboring autonomous system, it chooses the path with the lower route ID as the best path. The router installs this best path in the IP routing table. You can enable eBGP load balancing to install multiple paths in the IP routing table when the eBGP paths are learned from a neighboring autonomous system instead of picking one best path.

During packet switching, depending on the switching mode, the router performs either per-packet or per-destination load balancing among the multiple paths.

Layer 3 VPN Load Balancing

Layer 3 VPN load balancing for both eBGP and iBGP allows you to configure multihomed autonomous systems and provider edge (PE) routers to distribute traffic across both external BGP (eBGP) and iBGP multipaths.

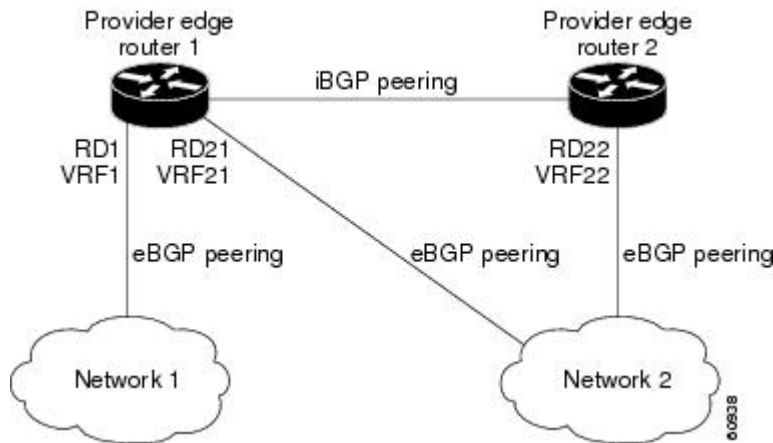
Layer 3 VPN load balancing supports IPv4 and IPv6 for the PE routers and VPNs.

BGP installs up to the maximum number of multipaths allowed. BGP uses the best path algorithm to select one path as the best path, inserts the best path into the routing information base (RIB) and advertises the best path to BGP peers. The router can insert other paths into the RIB but selects only one path as the best path.

Layer 3 VPNs load balance on a per-packet or per-source or destination pair basis. To enable load balancing, configure the router with Layer 3 VPNs that contain VPN routing and forwarding instances (VRFs) that import both eBGP and iBGP paths. You can configure the number of paths separately for each VRF.

The following figure shows an MPLS provider network that uses BGP. In the figure, two remote networks are connected to PE1 and PE2, which are both configured for VPN unicast iBGP peering. Network 2 is a multihomed network that is connected to PE1 and PE2. Network 2 also has extranet VPN services configured with Network 1. Both Network 1 and Network 2 are configured for eBGP peering with the PE routers.

Figure 1: Provider MPLS Network Using BGP



You can configure PE1 so that it can select both iBGP and eBGP paths as multipaths and import these paths into the VPN routing and forwarding instance (VRF) of Network 1 to perform load balancing.

Traffic is distributed as follows:

- IP traffic that is sent from Network 2 to PE1 and PE2 is sent across the eBGP paths as IP traffic.
- IP traffic that is sent from PE1 to PE2 is sent across the iBGP path as MPLS traffic.
- Traffic that is sent across an eBGP path is sent as IP traffic.

Any prefix that is advertised from Network 2 will be received by PE1 through route distinguisher (RD) 21 and RD22.

- The advertisement through RD21 is carried in IP packets.

- The advertisement through RD22 is carried in MPLS packets.

The router can select both paths as multipaths for VRF1 and insert these paths into the VRF1 RIB.

Layer 3 VPN Load Balancing with Route Reflectors

Route reflectors reduce the number of sessions on PE routers and increase the scalability of Layer 3 VPN networks. Route reflectors hold on to all received VPN routes to peer with PE routers. Different PEs can require different route target-tagged VPNv4 and VPNv6 routes. The route reflector may also need to send a refresh for a specific route target to a PE when the VRF configuration has changed. Storing all routes increases the scalability requirements on a route reflector. You can configure a route reflector to only hold routes that have a defined set of route target communities.

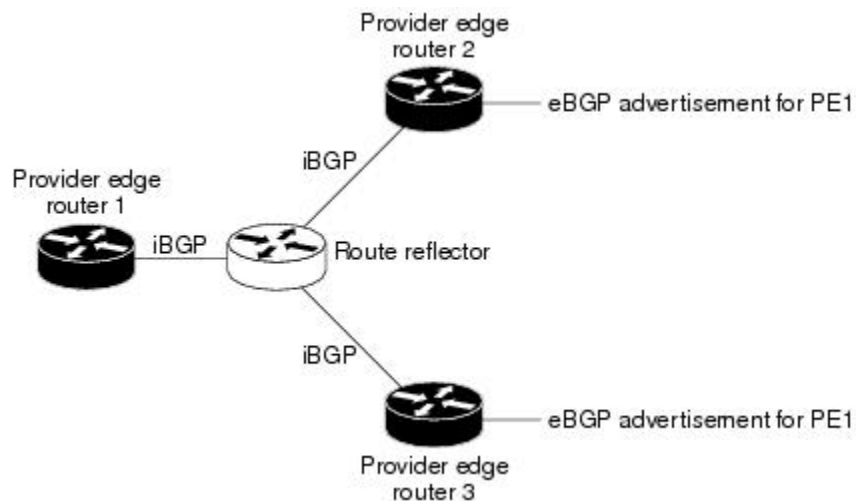
You can configure route reflectors to service a different set of VPNs and configure a PE to peer with all route reflectors that service the VRFs configured on the PE. When you configure a new VRF with a route target that the PE does not already hold routes for, the PE issues route refreshes to the route reflectors and retrieves the relevant VPN routes.

The following figure shows a topology that contains three PE routers and a route reflector, all configured for iBGP peering. PE2 and PE3 each advertise an equal preference eBGP path to PE1. By default, the route reflector chooses only one path and advertises PE1.



Note The route reflectors do not need to be in the forwarding path, but you must configure unique route distinguisher (RDs) for VPN sites that are multihomed.

Figure 2: Topology with a Route Reflector



For all equal preference paths to PE1 to be advertised through the route reflector, you must configure each VRF with a different RD. The prefixes received by the route reflector are recognized differently and advertised to PE1.

Layer 2 Load Balancing Coexistence

The load balance method that is required in the Layer 2 VPN is different from the method that is used for Layer 3 VPN. Layer 3 VPN and Layer 2 VPN forwarding is performed independently using two different

types of adjacencies. The forwarding is not impacted by using a different method of load balancing for the Layer 2 VPN.



Note Load balancing is not supported at the ingress PE for Layer 2 VPNs

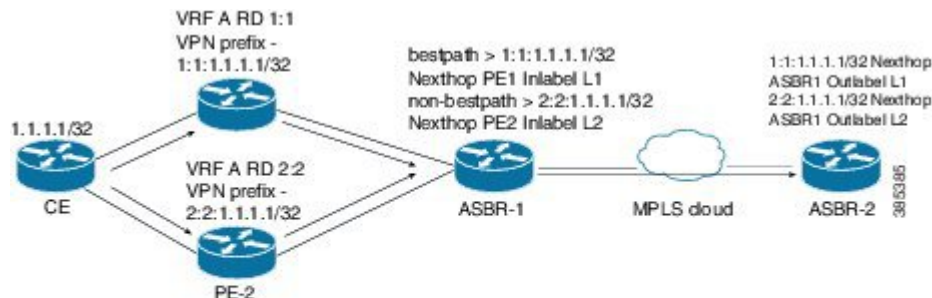
BGP VPNv4 Multipath

BGP VPNv4 Multipath feature helps to achieve Equal Cost Multi-Path (ECMP) for traffic flowing from an Autonomous System Border Router (ASBR) towards the Provider Edge (PE) device in an Multi-Protocol Label Switching (MPLS) cloud network by using a lower number of prefixes and MPLS labels. This feature configures the maximum number of multipaths for both eBGP and iBGP paths. This feature can be configured on PE devices and Route Reflectors in an MPLS topology.

Consider a scenario in which a dual homed Customer Edge (CE) device is connected to 2 PE devices and you have to utilize both the PE devices for traffic flow from ASBR-2 to the CE device.

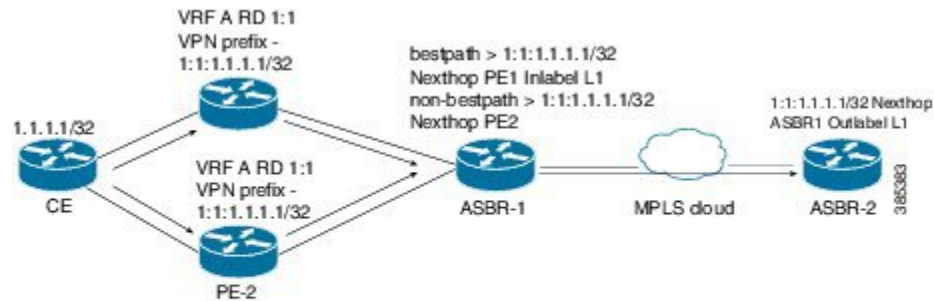
Currently, as shown in following figure, Virtual Routing and Forwarding (VRF) on each PE is configured using separate Route Distinguishers (RD). The CE device generates a BGP IPv4 prefix. The PE devices are configured with 2 separate RDs and generate two different VPN-IPv4 prefixes for the BGP IPv4 prefix sent by the CE device. ASBR-1 receives both the VPN-IPv4 prefixes and adds them to the routing table. ASBR-1 allocates Inter-AS option-B labels, Inlabel L1 and Inlabel L2, to both the VPN routes and then advertises both VPN routes to ASBR-2. To use both PE devices to maintain traffic flow, ASBR-1 has to utilize two Inter-AS option-B labels and two prefixes which limits the scale that can be supported.

Figure 3: Virtual Routing and Forwarding (VRF) on each PE configured using separate Route Distinguishers



Using the BGP VPN Multipath feature, as shown in Figure 22-4, you can enable the VRF on both PE devices to use the same RD. In such a scenario, ASBR-1 receives the same prefix from both the PE devices. ASBR-1 allocates only one Inter-AS option-B label, Inlabel L1, to the received prefix and advertises the VPN route to ASBR-2. In this case, the scale is enhanced as traffic flow using both PE devices is established with only one prefix and label on ASBR-1.

Figure 4: Enabling the VRF on both PE devices to use the same RD



BGP Cost Community

The BGP cost community is a nontransitive extended community attribute that is passed to iBGP and confederation peers but not to eBGP peers. (A confederation is a group of iBGP peers that use the same autonomous system number to communicate to external networks.) The BGP cost community attributes includes a cost community ID and a cost value. You can customize the BGP best path selection process for a local autonomous system or confederation by configuring the BGP cost community attribute. You configure the cost community attribute in a route map with a community ID and cost value. BGP prefers the path with the lowest community ID, or for identical community IDs, BGP prefers the path with the lowest cost value in the BGP cost community attribute.

BGP uses the best path selection process to determine which path is the best where multiple paths to the same destination are available. You can assign a preference to a specific path when multiple equal cost paths are available.

Since the administrative distance of iBGP is worse than the distance of most Interior Gateway Protocols (IGPs), the unicast Routing Information Base (RIB) may apply the same BGP cost community compare algorithm before using the normal distance or metric comparisons of the protocol or route. VPN routes that are learned through iBGP can be preferred over locally learned IGP routes.

The cost extended community attribute is propagated to iBGP peers when an extended community exchange is enabled.

How the BGP Cost Community Influences the Best Path Selection Process

The cost community attribute influences the BGP best path selection process at the point of insertion (POI). The POI follows the IGP metric comparison. When BGP receives multiple paths to the same destination, it uses the best path selection process to determine which path is the best path. BGP automatically makes the decision and installs the best path into the routing table. The POI allows you to assign a preference to a specific path when multiple equal cost paths are available. If the POI is not valid for local best path selection, the cost community attribute is silently ignored.

You can configure multiple paths with the cost community attribute for the same POI. The path with the lowest cost community ID is considered first. All of the cost community paths for a specific POI are considered, starting with the one with the lowest cost community ID. Paths that do not contain the cost community (for the POI and community ID being evaluated) are assigned with the default community cost value.

Applying the cost community attribute at the POI allows you to assign a value to a path originated or learned by a peer in any part of the local autonomous system or confederation. The router can use the cost community as a tie breaker during the best path selection process. You can configure multiple instances of the cost community for separate equal cost paths within the same autonomous system or confederation. For example, you can apply a lower cost community value to a specific exit path in a network with multiple equal cost exits points, and the BGP best path selection process prefers that specific exit path.

Cost Community and EIGRP PE-CE with Back-Door Links

BGP prefers back-door links in an Enhanced Interior Gateway Protocol (EIGRP) Layer 3 VPN topology if the back-door link is learned first. A back-door link, or a route, is a connection that is configured outside of the Layer 3 VPN between a remote and main site.

The pre-best path point of insertion (POI) in the BGP cost community supports mixed EIGRP Layer 3 VPN network topologies that contain VPN and back-door links. This POI is applied automatically to EIGRP routes that are redistributed into BGP. The pre-best path POI carries the EIGRP route type and metric. This POI influences the best-path calculation process by influencing BGP to consider this POI before any other comparison step.

Prerequisites for MPLS Layer 3 VPN Load Balancing

MPLS Layer 3 VPN load balancing has the following prerequisites:

- You must enable the MPLS and L3VPN features.
- You must install the correct license for MPLS.

Guidelines and Limitations for MPLS Layer 3 VPN Load Balancing

MPLS Layer 3 VPN load balancing has the following configuration guidelines and limitations:

- You can configure MPLS Layer 3 VPN load balancing for Cisco Nexus 9508 platform switches with the N9K-X9636C-R, N9K-X9636C-RX, and N9K-X9636Q-R line cards.
- Beginning with Cisco NX-OS Release 9.3(3), you can configure MPLS Layer 3 VPN load balancing on Cisco Nexus 9364C-GX, Cisco Nexus 9316D-GX, and Cisco Nexus 93600CD-GX switches.
- From Cisco NX-OS Release 10.4(1)F, you can configure mpls load-balancing for switch under port-channel load balance. This feature is supported on Cisco Nexus 9300 -EX/FX/FX2/ FX3/GX/GX2 TOR and EOR platform switches. For more information on configuration, refer to *Cisco Nexus 9000 Series NX-OS Interfaces Configuration Guide*.
- Cisco Nexus 9348GC-FX3PH switch has feature limitations due to full-duplex for ports 41-48.
- Cisco Nexus C93108TC-FX3 switch has feature limitations due to half-duplex for ports 41-48.
- If you place a router behind a route reflector and it is connected to multihomed sites, the router will not be advertised unless separate VRFs with different RDs are configured for each VRF.

- Each IP routing table entry for a BGP prefix that has multiple iBGP paths uses additional memory. We recommend that you do not use this feature on a router with a low amount of available memory or when it is carrying a full Internet routing table.
- You should not ignore the BGP cost community when a back-door link is present and EIGRP is the PE-CE routing protocol.
- A maximum of 16K VPN prefixes is supported on Cisco Nexus 9508 platform switches with N9K-X9636Q-R and N9K-X9636C-R line cards, and a maximum of 470K VPN prefixes is supported on Cisco Nexus 9508 platform switches with N9K-X9636C-RX line cards.
- 4K VRFs are supported.
- Beginning with Cisco NX-OS Release 10.1(1), on Cisco Nexus 9300-FX2, 9300-GX, 9300-GX2 platform switches, addition or deletion of dot1q tag is not supported when packet is received on an interface enabled with mpls ip forwarding. For previous releases, addition or deletion of dot1q tag is not supported when the CLI **feature mpls segment-routing** is enabled or **mpls load-sharing [label-only | [label-ip]** is configured.
- On Cisco Nexus 9300-EX, 9300-FX, 9300-EX-LC, 9300-FX-LC, and also N9K-C9364C, N9K-C9508-FM-E2, N9K-C9516-FM-E2, and N9K-C9332C platform switches, addition or deletion of dot1q tag is not supported when the CLI **feature mpls segment-routing** is enabled or **mpls load-sharing [label-only | [label-ip]** is configured.
- On Cisco Nexus 9300-EX and 9300-EX-LC platform switches, port-channel and ecmp load-sharing based on mpls label or SRC/DST-IP does not work even when the CLI **mpls load-sharing label-ip** is configured; however, **label-only** works.
- VXLAN BUM traffic should not traverse through a Pure L2 switch with mpls load-balancing enabled (**mpls load-sharing [label-only | [label-ip]**).

Default Settings for MPLS Layer 3 VPN Load Balancing

The following table lists the default settings for MPLS Layer 3 VPN load balancing parameters.

Table 1: Default MPLS Layer 3 VPN Load Balancing Parameters

Parameters	Default
Layer 3 VPN feature	Disabled
BGP cost community ID	128
BGP cost community cost	2147483647
maximum multipaths	1
BGP VPNv4 Multipath	Disabled

Configuring MPLS Layer 3 VPN Load Balancing

Configuring BGP Load Balancing for eBGP and iBGP

You can configure a Layer 3 VPN load balancing for an eBGP or iBGP network.

Procedure

	Command or Action	Purpose
Step 1	configure terminal Example: switch# configure terminal switch(config)#	Enters global configuration mode.
Step 2	feature-set mpls Example: switch(config)# feature-set mpls	Enables the MPLS feature-set.
Step 3	feature mpls l3vpn Example: switch(config)# feature mpls l3vpn	Enables the MPLS Layer 3 VPN feature.
Step 4	feature bgp Example: switch(config)# feature bgp switch(config)#	Enables the BGP feature.
Step 5	router bgp <i>as - number</i> Example: switch(config)# router bgp 1.1 switch(config-router)#	Configures a BGP routing process and enters router configuration mode. The <i>as-number</i> argument indicates the number of an autonomous system that identifies the router to other BGP routers and tags the routing information passed along. The AS number can be a 16-bit integer or a 32-bit integer in the form of a higher 16-bit decimal number and a lower 16-bit decimal number in xx.xx format.
Step 6	bestpath cost-community ignore remote-as <i>as-number</i> Example: switch(config-router)# bestpath cost-community ignore#	(Optional) Ignores the cost community for BGP bestpath calculations.

	Command or Action	Purpose
Step 7	address-family { ipv4 ipv6 } unicast Example: <pre>switch(config-router)# address-family ipv4 unicast switch(config-router-af)#</pre>	Enters address family configuration mode for configuring IP routing sessions.
Step 8	maximum-paths [bgp] number-of-paths Example: <pre>switch(config-router-af)# maximum-paths 4</pre>	Configures the maximum number of multipaths allowed. Use the ibgp keyword to configure iBGP load balancing. The range is from 1 to 16.
Step 9	show running-config bgp Example: <pre>switch(config-router-vrf-neighbor-af)# show running-config bgp</pre>	(Optional) Displays the running configuration for BGP.
Step 10	copy running-config startup-config Example: <pre>switch(config-router-vrf)# copy running-config startup-config</pre>	(Optional) Copies the running configuration to the startup configuration.

Configuring BGPv4 Multipath

Procedure

	Command or Action	Purpose
Step 1	configure terminal Example: <pre>switch# configure terminal switch(config)#</pre>	Enters global configuration mode.
Step 2	feature bgp Example: <pre>switch(config)# feature bgp</pre>	Enables the BGP feature.
Step 3	router bgp as - number Example: <pre>switch(config)# router bgp 2 switch(config-router)#</pre>	Assigns an autonomous system (AS) number to a router and enter the router BGP configuration mode.

	Command or Action	Purpose
Step 4	address-family vpnv4 unicast Example: <pre>switch(config-router)# address-family vpnv4 unicast switch(config-router-af)#</pre>	Enters address family configuration mode for configuring routing sessions, such as BGP, that use standard VPNv4 address prefixes.
Step 5	maximum-paths eibgp parallel-paths Example: <pre>switch(config-router-af)# maximum-paths eibgp 3</pre>	Specifies the maximum number of BGP VPNv4 multipaths for both eBGP and iBGP paths. The range is from 1 to 32.

Configuring MPLS ECMP Load Sharing

Beginning Cisco NX-OS Release 9.3(1), you can configure MPLS ECMP load sharing based on labels. This feature is supported on Cisco Nexus 9200, Cisco Nexus 9300-EX, Cisco Nexus 9300-FX, and Cisco Nexus 9500 platform switches with Cisco Nexus N9K-X9700-EX and N9K-X9700-FX line cards.

Beginning with Cisco NX-OS Release 9.3(3), this feature is supported on Cisco Nexus 9364C-GX, Cisco Nexus 9316D-GX, and Cisco Nexus 93600CD-GX switches.

Procedure

	Command or Action	Purpose
Step 1	configure terminal Example: <pre>switch# configure terminal switch(config)#</pre>	Enters global configuration mode.
Step 2	feature-set mpls Example: <pre>switch(config)# feature-set mpls</pre>	Enables the MPLS feature-set.
Step 3	mpls load-sharing [label-only [label-ip] Example: <pre>switch(config)# mpls load-sharing label-only switch(config)# mpls load-sharing label-ip</pre>	Configures the load sharing based on the mpls labels. The label-only option configures the load sharing based on the labels, while the label-ip option configures it based on the label and the IP address.
Step 4	copy running-config startup-config Example: <pre>switch(config)# copy running-config startup-config</pre>	(Optional) Copies the running configuration to the startup configuration.

Verifying MPLS ECMP Load Sharing

To display the mpls ECMP load sharing configuration, perform one of the following tasks:

Table 2: Verifying MPLS ECMP Load Sharing Configuration

Command	Purpose
<code>show mpls load-sharing</code>	Displays the number of labels that are used for the mpls hashing and the IP fields that are used for the hashing.

Configuration Examples for MPLS Layer 3 VPN Load Balancing

Example: MPLS Layer 3 VPN Load Balancing

The following example shows how to configure iBGP load balancing:

```
configure terminal
feature-set mpls
feature mpls l3vpn
feature bgp
router bgp 1.1
bestpath cost-community ignore
address-family ipv6 unicast
maximum-paths ibgp 4
```

Example: BGP VPNv4 Multipath

The following example shows how to configure a maximum of 3 BGP VPNv4 multipaths:

```
configure terminal
router bgp 100
address-family vpnv4 unicast
maximum-paths eibgp 3
```

Example: MPLS Layer 3 VPN Cost Community

The following example shows how to configure the BGP cost community:

```
configure terminal
feature-set mpls
feature mpls l3vpn
feature bgp
route-map CostMap permit
set extcommunity cost 1 100
router bgp 1.1
router-id 192.0.2.255
neighbor 192.0.2.1 remote-as 1.1
address-family vpnv4 unicast
send-community extended
route-map CostMap in
```

