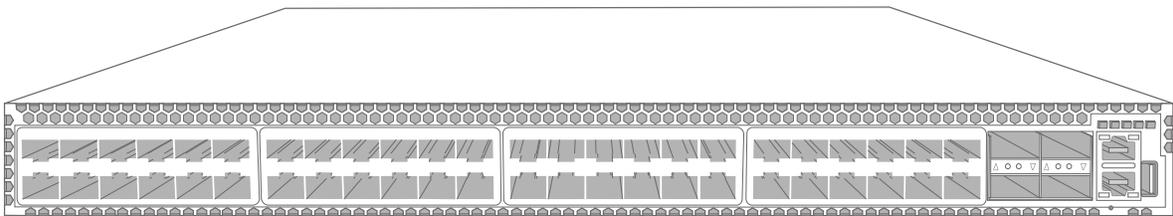


EVPN for controller-less VXLAN

BUILDING WEB-SCALE ARCHITECTURES WITH EVPN ON CUMULUS LINUX



 CUMULUS

WHITE PAPER

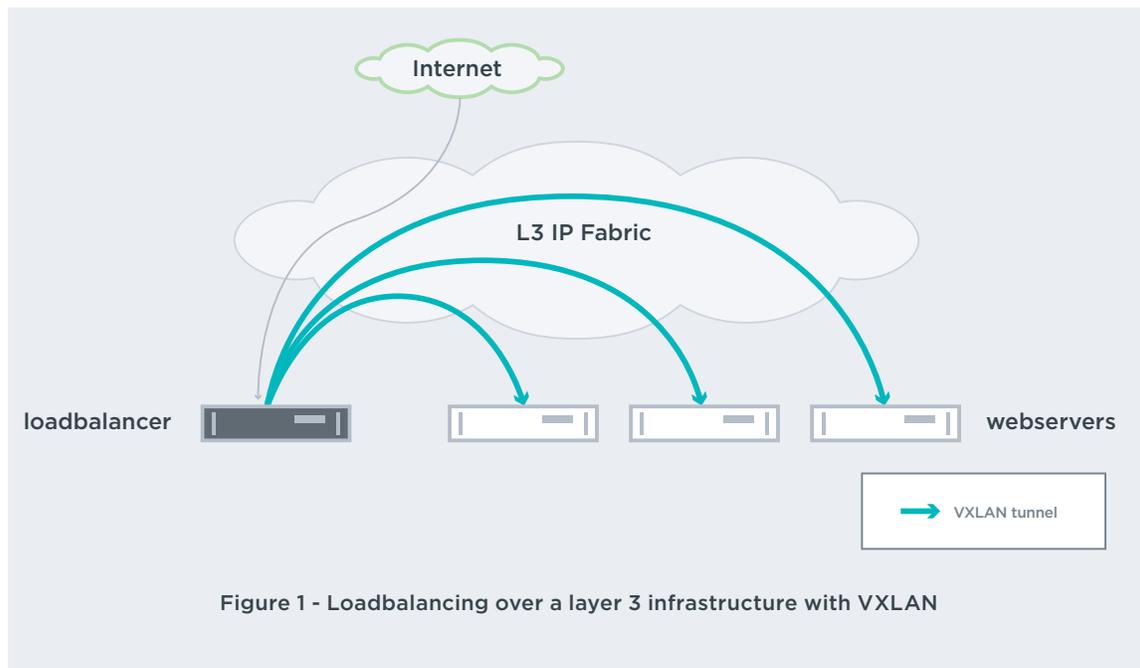
Table of Contents

Introduction	2
Deployment benefits summary	3
EVPN overview and operations	4
EVPN VTEP peer discovery	7
EVPN MAC address learning/exchange and multi-tenant support ..	8
EVPN VXLAN active-active mode	16
EVPN MAC mobility	17
EVPN deployment scenarios and configuration	20
EVPN in an eBGP environment	21
EVPN in an iBGP environment with OSPF underlay	23
EVPN in an iBGP environment with route reflectors	24
Conclusion	25

Introduction

Many data centers today are moving from a legacy layer 2 design to a modern layer 3 [web-scale IT architecture](#). [Layer 3 designs](#) using traditional routing protocols like OSPF and BGP allow simplified troubleshooting, clear upgrade strategies, multi-vendor support, small failure domains and less vendor lock-in. However, many applications, storage appliances and tenant considerations still require layer 2 adjacency.

Virtual Extensible VLAN (VXLAN) is widely deployed in many layer 3 data centers to provide layer 2 connectivity between hosts for specific applications. For example, as seen in Figure 1, the web servers and the load balancer must be on the same layer 2 network. VXLAN provides that layer 2 connectivity over a layer 3 infrastructure.



However, in many circumstances, a controller is often deployed to configure and control VXLAN tunnels. Many customers prefer a standardized scalable solution that does not require an external controller.

Ethernet Virtual Private Network (EVPN) is a feature now offered by Cumulus Networks that provides a controller-free, scalable and interoperable end-to-end control-plane solution for VXLAN tunnels. It supports redundancy, load

sharing and multi-tenant segmentation. EVPN also provides the benefit of fast convergence for host and VM mobility over VXLAN tunnels.

This white paper discusses deployment benefits, how EVPN works, how to operate EVPN, and different deployment scenarios. This paper also includes sample Cumulus Linux configurations to deploy a scalable, controller-free layer 2 virtualization over a layer 3 IP fabric.

Deployment benefits summary

Deploying EVPN provides many advantages to a layer 3 data center:

Controller-less VXLAN tunnels

No controller is needed for VXLAN tunnels, as EVPN provides peer discovery with authentication natively. This also mitigates the chance of rogue VTEPs in a network and dealing with complicated controller redundancy.

Scale and robustness

EVPN uses the BGP routing protocol. BGP is very mature, scalable, flexible and robust. It is the primary routing protocol for the Internet and can hold a very large number of routes. It is also a preferred routing protocol for [data centers](#). It supports routing policy and filtering, which provides granular control over traffic flow.

Fast convergence and host mobility

Cumulus EVPN supports the new BGP MAC mobility extended community, offering fast convergence and reducing discovery traffic after a MAC or VM move.

Support for VXLAN active-active mode

Cumulus EVPN integrates with MLAG, thereby providing host dual homing for redundancy.

Multitenancy

EVPN uses the mature multi-protocol BGP VPN technology to separate tenants within a data center.

Interoperability between vendors

The standardized [multi-protocol BGP \(MP-BGP\)](#) is used for the EVPN control plane. As long as vendor implementations maintain adherence to both the VXLAN and EVPN standards, interoperability is assured.

EVPN is a standardized control plane protocol that offers controller-less VXLAN tunnels. It also offers scale, redundancy, fast convergence and robustness while reducing BUM traffic across a data center core. More details on the operations providing these benefits are discussed below.

EVPN overview and operations

Customers moving from traditional layer 2 data centers to a layer 3 fabric to overcome one or more of these issues:

- **Large broadcast and failure domains:** A broadcast packet is sent throughout the data center, increasing utilization and a failure can impact the entire data center.
- **Limited redundancy:** MLAG is often deployed for redundancy but it supports only 2 switches.
- **Troubleshooting difficulty:** Spanning tree issues can cause a network meltdown and are difficult to troubleshoot.
- **Limited scale for tenant separation:** A maximum of only 4094 VLANs are supported.

While moving to a layer 3 fabric should overcome these issues, some applications still require layer 2 connectivity between servers, so VXLAN tunnels are often deployed. VXLAN tunnels are identified by IETF RFC 7348 “[Virtual eXtensible Local Area Network \(VXLAN\): A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks.](#)”

VXLAN provides a scalable solution for layer 2 virtualization over a layer 3 routed infrastructure. It allows up to 16 million different VXLANs in the same domain by allocating a 24-bit segment ID called either the VXLAN network identifier (VNI) or the VXLAN-ID. The VNI is used to distinguish between VXLAN tunnels.

VXLAN provides a scalable solution for layer 2 virtualization over a layer 3 routed infrastructure

Virtual Tunnel Endpoints (VTEPs) are used to originate and terminate the VXLAN tunnel and map end devices such as hosts and VMs to VXLAN segments. The VTEP provides the encapsulation of layer 2 frames into User Datagram Protocol (UDP) segments to traverse across a layer 3 fabric. Likewise, the VTEP also de-encapsulates the UDP segments from a VXLAN tunnel to send to a local host. A VTEP requires an IP address (often a loopback address) and uses this address as the source/destination tunnel IP address. The VTEP IP address must be advertised into the routed domain so the VXLAN tunnel endpoints can reach each other as shown in Figure 1. Each switch that hosts a VTEP must have a VXLAN-supported chipset such as Mellanox Spectrum or Broadcom Trident II, Trident II+ or Tomahawk. A list of our compatible hardware can be found in the [Hardware Compatibility List](#). Though it's not depicted in Figure 2, you can have multiple VNIs (VXLANs) using one VTEP IP address.

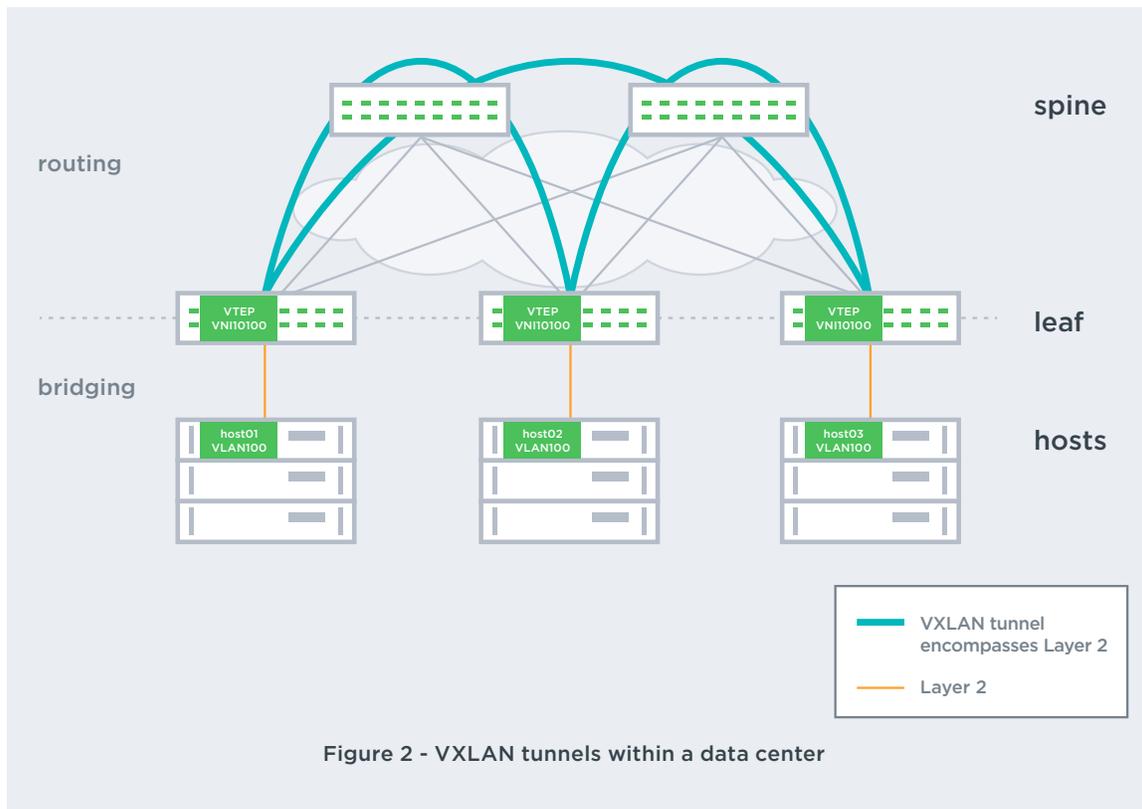
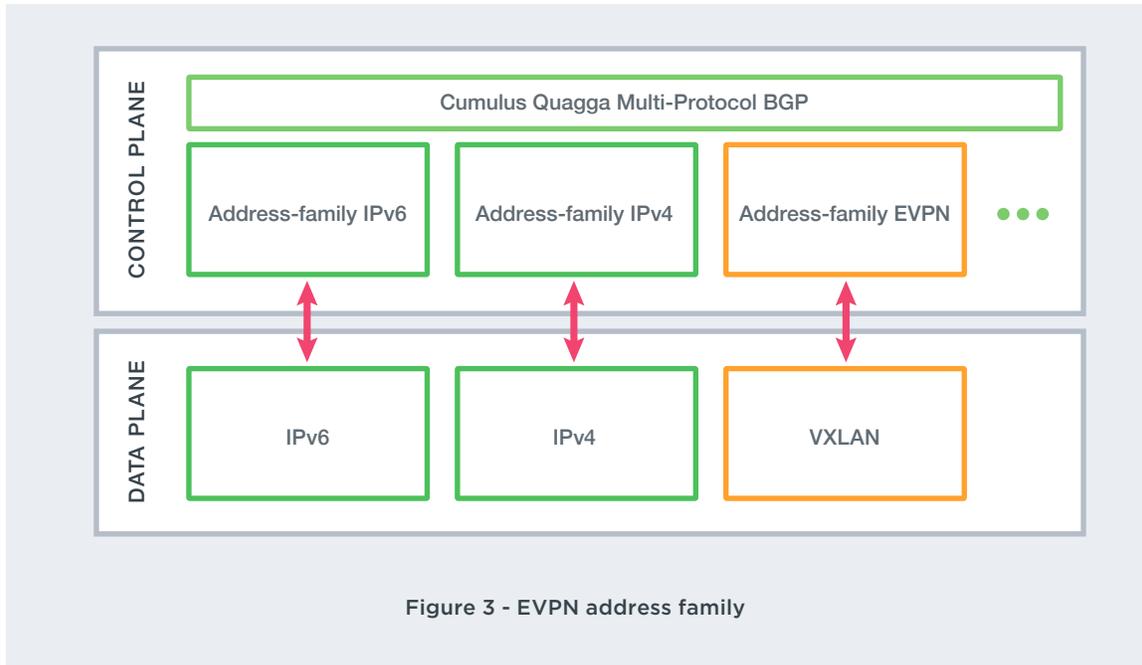


Figure 2 - VXLAN tunnels within a data center

In traditional VXLAN, as seen above in Figure 2, the control plane and the data plane are integrated together — meaning MAC address learning happens over the data plane, often called *flood and learn*. This causes limitations, including limited load balancing and slow convergence times especially for VM and host mobility. Further, broadcast, unknown unicast, and multicast (BUM) traffic, such as ARP, is required to traverse across the tunnel for discovery purposes, thereby increasing overall data center traffic.

The Cumulus Networks EVPN implementation provides a separate control plane for VXLAN tunnels. EVPN provides exchange of MAC/IP addresses between VTEPs through the use of a separate control plane, similar to pure IP routing. Cumulus EVPN is an open and standards based solution that implements IETF RFC 7432 “[BGP MPLS-Based Ethernet VPN](#)” along with IETF draft “[A Network Virtualization Overlay Solution using EVPN](#)” for a VXLAN tunnel control plane. EVPN introduces a new address family to the MP-BGP protocol family, as depicted in Figure 3.



EVPN provides remote VTEP discovery, thus it doesn't require an external controller. Learning control plane information independently of the data plane offers greater redundancy, load sharing and multipathing while also supporting MAC address filtering and traffic engineering, which can provide granular control of traffic flow. EVPN also provides faster convergence for mobility. The greater redundancy and multipathing can be achieved because all the possible paths are exchanged across the control plane, not just from one data plane path.

When EVPN is implemented with the VXLAN data plane, the evpn address family can exchange either just the MAC layer control plane information (that is, MAC addresses) or it can exchange both the MAC address and IP address information in its updates between VTEPs. Exchanging IP information can allow for ARP suppression at the local switch, thereby reducing the broadcast traffic in a data center.

EVPN VTEP PEER DISCOVERY

One large advantage of deploying EVPN is the ability to deploy controller-free VXLAN tunnels. EVPN uses **type 3 EVPN routes** to exchange information about the location of the VTEPs on a per-VNI basis, thereby enabling automatic discovery. It also reduces or eliminates the chance of a rogue VTEP being introduced in the data center.

EVPN offers peer discovery, thus requiring no external controller

For example, in Figure 4, the VTEPs are automatically discovered via eBGP and do not need to be explicitly configured or controlled as peers. The spine switches do not need to be configured for VLAN or VXLAN at all. All the discovered VTEPs within a VXLAN can easily be

seen from one participating VTEP with a simple show command. The command in Figure 4 below displays all the remote VTEPs associated with a specific VNI that are automatically discovered, including any rogue VTEPs.

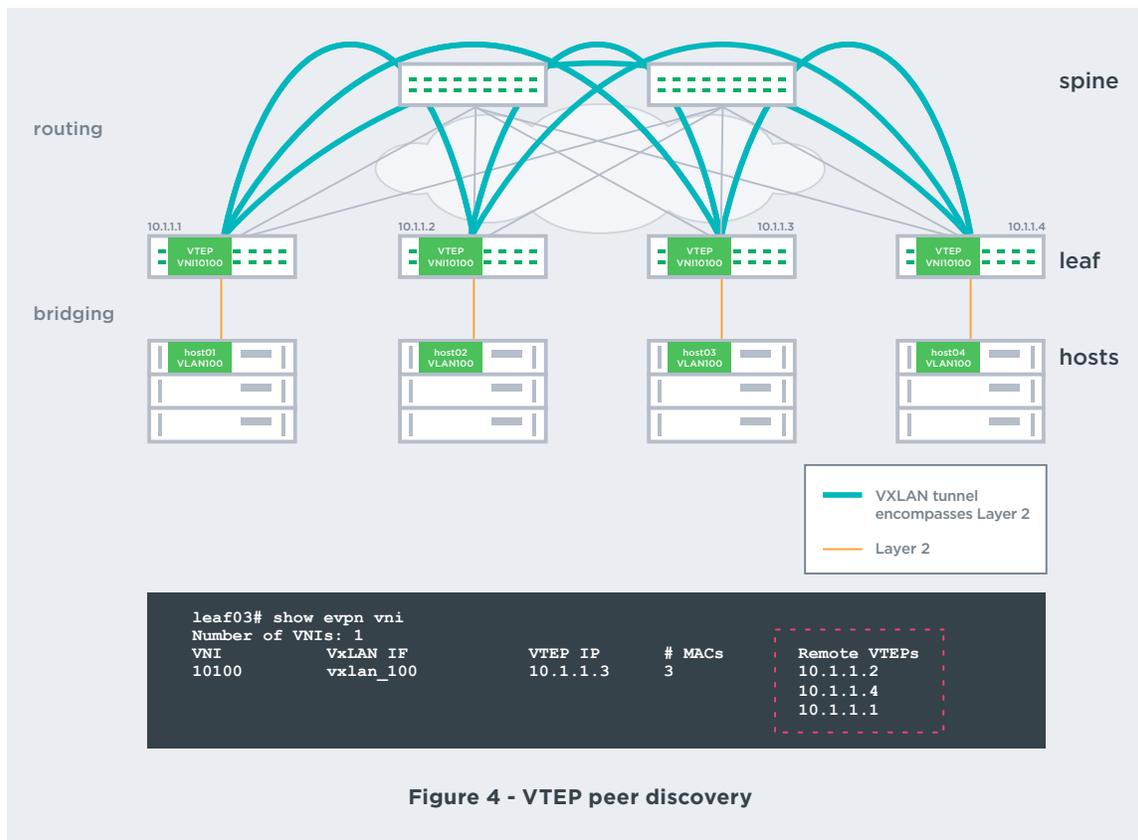
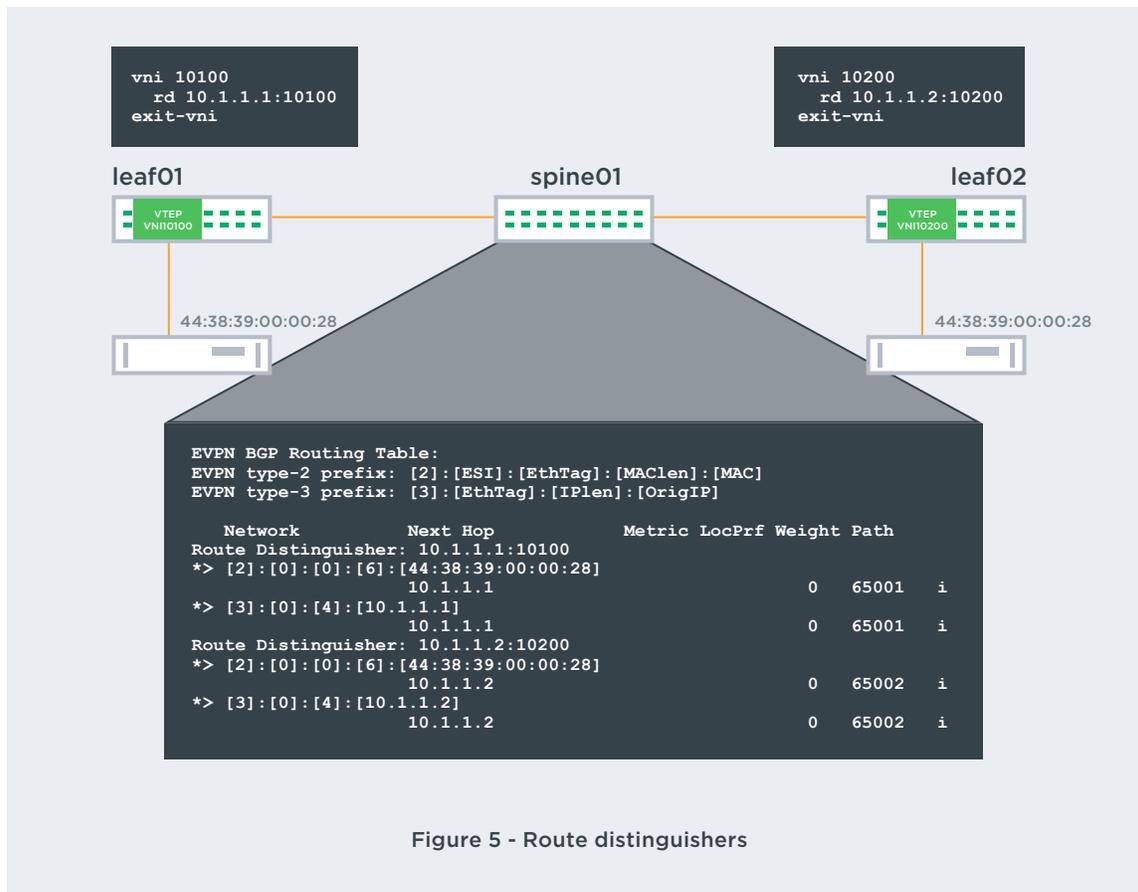


Figure 4 - VTEP peer discovery

EVPN MAC ADDRESS LEARNING/EXCHANGE AND MULTI-TENANT SUPPORT

The new EVPN address family also provides multi-tenant separation and allows for overlapping addresses between tenants. To maintain the separation, it uses mature MP-BGP VPN technology: Route Distinguishers (RDs) and Route-Targets (RTs).

The RD makes overlapping routes from different tenants look unique to the data center spine switches to provide proper routing. A per-VXLAN 8-byte RD is prepended to each advertised route before the route is sent to its BGP EVPN peer. In Figure 5, the same MAC address is advertised from 2 hosts in separate tenants, but the spine router can distinguish between the routes since they have different route distinguishers.



EVPN also makes use of the RT extended community for route filtering and separating tenants. The RT is advertised in the BGP update message along with the EVPN routes. The RT community distinguishes which routes should be exported from and imported into a specific VNI route table on a VTEP. If the export RT in the

received update matches the import RT of a VNI instance on the VTEP receiving the update, the corresponding routes will be imported into that VNI's EVPN Route table. If the RTs do not match, the route will not be imported into that VNI's EVPN route table.

EVPN provides multi-tenant separation with one protocol instance

Figure 6 below depicts leaf01 sending a BGP EVPN route to leaf02 with the attached route-target community. As seen, four MAC addresses are sent in the advertisement, two each originating from different VNIs on leaf01. Since leaf02 only has the one route-target import, 65001:1, it will receive only those routes associated with 65001:1, and the route with route-target 65001:2 will not be installed as there is no matching import route-target within VNI 10100 located on leaf02.

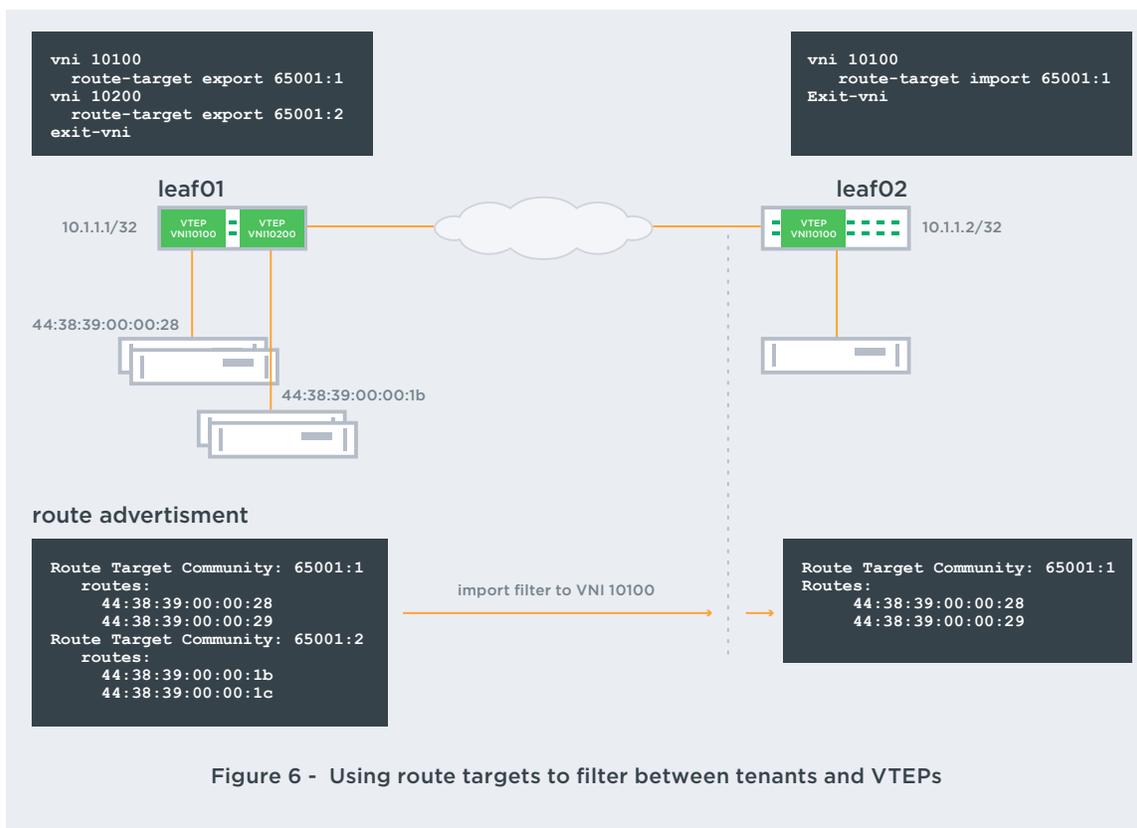


Figure 6 - Using route targets to filter between tenants and VTEPs

Cumulus Linux supports either a default RD and/or RT for configuration ease, or configuring explicit RD and/or RT within BGP for each VNI to allow flexibility. By default, the switch automatically derives the RD and RT from the VNI. In the default case, the RD would be *Router ID: VNI*, and the export RT is set at *AS:VNI*. The import RT community is set to *<any>:VNI*, which allows all routes from that same VXLAN to be imported. If more granular control of importing routes, compatibility with other vendors, and/or if globally unique VNIs are not configured, the RD and RT community can be manually configured as well. Manually configuring the RT and RD overrides the default RD and RT values.

Each local VLAN is mapped to a VNI. When a local switch learns a new MAC address on a particular VLAN, either via Gratuitous ARP (GARP) or via the first data packet, which is typically an ARP request, it is placed into the local switch's bridge forwarding table. The local

MP-BGP process learns every new local MAC address from the local forwarding table and advertises them to the remote VTEPS via [Type 2 EVPN routes](#).

On the remote end, the MAC addresses that BGP learns are placed into the BGP table. From there, if the route target community sent with the route matches a local VNI route-target import, the route will be placed into that switch's MAC forwarding table with the appropriate VXLAN tunnel as its destination. This process separates the data and control planes, allowing dynamic learning of MAC addresses, allows overlapping MAC addresses between tenants, and allows granular filtering of MAC addresses, all without requiring a data plane packet to traverse each switch first.

To walk through an example of the MAC address being propagated through the network, consider the example network in Figure 7 where there are two leaf switches, each participating in two independent VXLAN tunnels.

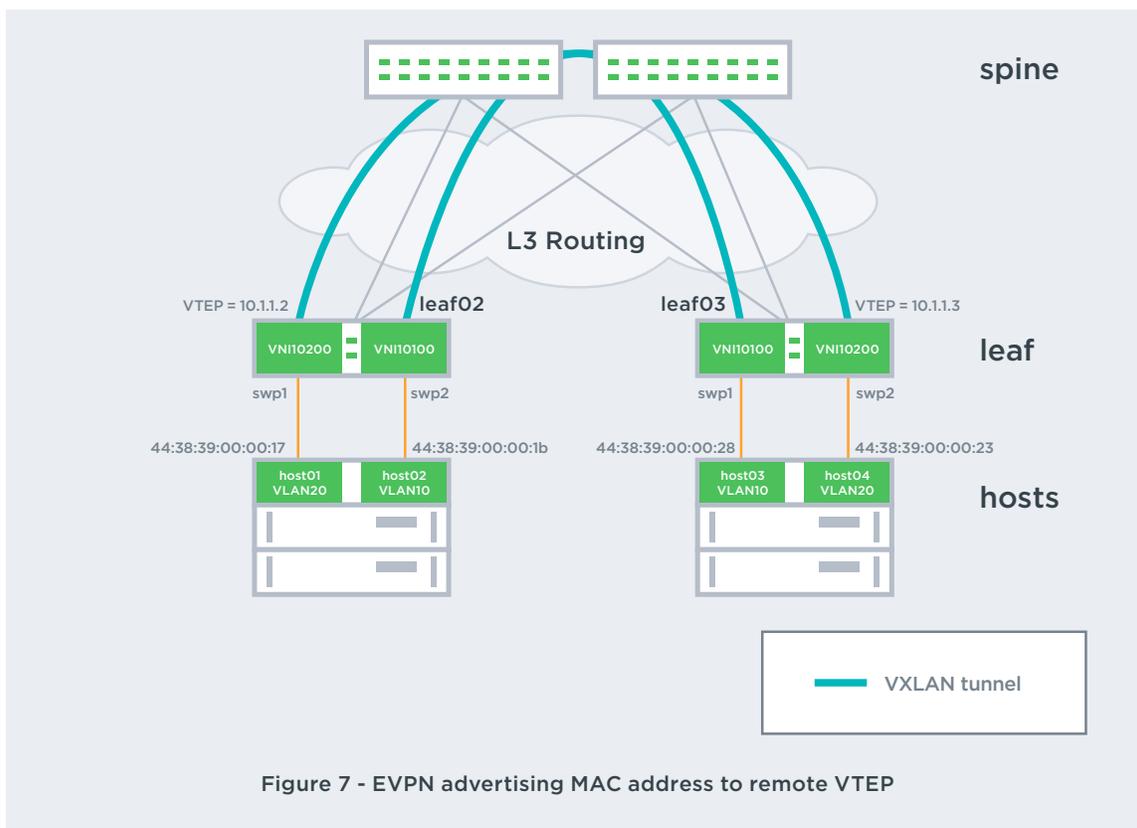


Figure 7 - EVPN advertising MAC address to remote VTEP

Following the route through the network, leaf02 learns host01's MAC address (44:38:39:00:00:17) and host02's MAC address (44:38:39:00:00:1b) into its bridge forwarding database. It can be seen here:

```
cumulus@leaf02:~$ net show bridge mac
```

VLAN	Master	Interface	MAC	TunnelDest	State	Flags	LastSeen
10	bridge	swp2	44:38:39:00:00:1b				00:00:06
10	bridge	vxlan_10	44:38:39:00:00:28				00:01:00
20	bridge	swp1	44:38:39:00:00:17				00:00:25
20	bridge	vxlan_20	44:38:39:00:00:23				00:01:00
untagged		vxlan_10	00:00:00:00:00:00	10.1.1.3	permanent	self	00:01:00
untagged		vxlan_10	44:38:39:00:00:28	10.1.1.3		self	00:01:00
untagged		vxlan_20	00:00:00:00:00:00	10.1.1.3	permanent	self	00:01:00
untagged		vxlan_20	44:38:39:00:00:23	10.1.1.3		self	00:01:00
untagged	bridge	swp1	44:38:39:00:00:18		permanent		00:01:03
untagged	bridge	swp2	44:38:39:00:00:1c		permanent		00:01:03
untagged	bridge	vxlan_10	c6:95:d1:01:6e:7c		permanent		00:01:03
untagged	bridge	vxlan_20	42:38:5c:aa:b8:39		permanent		00:01:03

In this case, we can see the local MAC address 44:38:39:00:00:1b is located in VLAN 10, and the local MAC address 44:38:39:00:00:17 is located in VLAN 20. The remote MAC addresses can also be seen across the tunnels. For example, host04's MAC address 44:38:39:00:00:23 in VLAN 20, is reachable through interface vxlan_20 and is behind VTEP 10.1.1.3. Host03's MAC address 44:38:39:00:00:28 resides in VLAN

10, reachable through the vxlan_10 interface and is also behind VTEP 10.1.1.3.

The 00:00:00:00:00:00 MAC address associated with vxlan_20 and the 00:00:00:00:00:00 MAC address associated with vxlan_10 entries are added by EVPN when the VTEP is discovered. These entries are the head end replication entries and should never age out as long as a remote VTEP is active.

To propagate the local MAC addresses to the remote VTEP, the local MAC addresses will be learned by BGP, as seen in the following. The type 2 routes are advertising the MAC addresses, and the type 3 routes are advertising the location of the VTEPs in the network.

```
leaf02# show bgp evpn route
BGP table version is 0, local router ID is 10.1.1.2
Status codes: s suppressed, d damped, h history, * valid, > best, i - internal
Origin codes: i - IGP, e - EGP, ? - incomplete
EVPN type-2 prefix: [2]:[ESI]:[EthTag]:[MAClen]:[MAC]
EVPN type-3 prefix: [3]:[EthTag]:[IPlen]:[OrigIP]

      Network          Next Hop          Metric LocPrf Weight Path
Route Distinguisher: 10.1.1.2:10100
*> [2]:[0]:[0]:[6]:[44:38:39:00:00:1b]
      10.1.1.2
      32768 i
*> [3]:[0]:[4]:[10.1.1.2]
      10.1.1.2
      32768 i
Route Distinguisher: 10.1.1.2:10200
*> [2]:[0]:[0]:[6]:[44:38:39:00:00:17]
      10.1.1.2
      32768 i
*> [3]:[0]:[4]:[10.1.1.2]
      10.1.1.2
      32768 i
Route Distinguisher: 10.1.1.3:10100
* [2]:[0]:[0]:[6]:[44:38:39:00:00:28]
      10.1.1.3
      0 65000 65003 i
*> [2]:[0]:[0]:[6]:[44:38:39:00:00:28]
      10.1.1.3
      0 65000 65003 i
* [3]:[0]:[4]:[10.1.1.3]
      10.1.1.3
      0 65000 65003 i
*> [3]:[0]:[4]:[10.1.1.3]
      10.1.1.3
      0 65000 65003 i
Route Distinguisher: 10.1.1.3:10200
* [2]:[0]:[0]:[6]:[44:38:39:00:00:23]
      10.1.1.3
      0 65000 65003 i
*> [2]:[0]:[0]:[6]:[44:38:39:00:00:23]
      10.1.1.3
      0 65000 65003 i
* [3]:[0]:[4]:[10.1.1.3]
      10.1.1.3
      0 65000 65003 i
*> [3]:[0]:[4]:[10.1.1.3]
      10.1.1.3
      0 65000 65003 i

Displayed 8 prefixes (12 paths)
```

As seen above, the routes are separated per tenant, and is identified by the route distinguishers. The local routes naturally have no path, whereas the remote ones do show the AS path to the MAC address and VTEP IP addresses.

From here, the local routes are advertised to the remote BGP neighbor (usually a spine in the case of eBGP) and then propagated to the remote leaf. The eBGP EVPN output from a remote leaf looks like the following:

```
leaf03# show bgp evpn route
BGP table version is 0, local router ID is 10.1.1.3
Status codes: s suppressed, d damped, h history, * valid, > best, i - internal
Origin codes: i - IGP, e - EGP, ? - incomplete
EVPN type-2 prefix: [2]:[ESI]:[EthTag]:[MACLen]:[MAC]
EVPN type-3 prefix: [3]:[EthTag]:[IPLen]:[OrigIP]
```

Network	Next Hop	Metric	LocPrf	Weight	Path
Route Distinguisher: 10.1.1.2:10100					
* [2]:[0]:[0]:[6]:[44:38:39:00:00:1b]	10.1.1.2			0 65000	65002 i
*> [2]:[0]:[0]:[6]:[44:38:39:00:00:1b]	10.1.1.2			0 65000	65002 i
* [3]:[0]:[4]:[10.1.1.2]	10.1.1.2			0 65000	65002 i
*> [3]:[0]:[4]:[10.1.1.2]	10.1.1.2			0 65000	65002 i
Route Distinguisher: 10.1.1.2:10200					
* [2]:[0]:[0]:[6]:[44:38:39:00:00:17]	10.1.1.2			0 65000	65002 i
*> [2]:[0]:[0]:[6]:[44:38:39:00:00:17]	10.1.1.2			0 65000	65002 i
* [3]:[0]:[4]:[10.1.1.2]	10.1.1.2			0 65000	65002 i
*> [3]:[0]:[4]:[10.1.1.2]	10.1.1.2			0 65000	65002 i
Route Distinguisher: 10.1.1.3:10100					
*> [2]:[0]:[0]:[6]:[44:38:39:00:00:28]	10.1.1.3			32768	i
*> [3]:[0]:[4]:[10.1.1.3]	10.1.1.3			32768	i
Route Distinguisher: 10.1.1.3:10200					
*> [2]:[0]:[0]:[6]:[44:38:39:00:00:23]	10.1.1.3			32768	i
*> [3]:[0]:[4]:[10.1.1.3]	10.1.1.3			32768	i

Displayed 8 prefixes (12 paths)

As seen above, 44:38:39:00:00:1b and 44:38:39:00:00:17 are now remote MAC addresses, as expected.

Based upon the configured import route targets, BGP then places certain routes within specific VNIs. For example, in this case, we have an import route target of <any>:10200 to be imported

into VNI 10200, and an import route-target of <any>:10100 to be imported into VNI 10100, so all the MAC addresses with the same route target will be imported into the respective VNI.

```
leaf03# show bgp evpn import-rt
Route-target: 0:10200
List of VNIs importing routes with this route-target:
  10200
Route-target: 0:10100
List of VNIs importing routes with this route-target:
  10100
```

Finally, looking at leaf03’s forwarding database, those MAC addresses are now reachable through the VXLAN tunnels. They are identified per VLAN and/or VXLAN:

```
cumulus@leaf03:~$ net show bridge mac
```

VLAN	Master	Interface	MAC	TunnelDest	State	Flags	LastSeen
10	bridge	swp1	44:38:39:00:00:28				00:00:21
10	bridge	vxlan_10	44:38:39:00:00:1b				never
20	bridge	swp2	44:38:39:00:00:23				00:00:16
20	bridge	vxlan_20	44:38:39:00:00:17				00:04:46
untagged		vxlan_10	00:00:00:00:00:00	10.1.1.2	permanent	self	00:10:53
untagged		vxlan_10	44:38:39:00:00:1b	10.1.1.2		self	00:00:49
untagged		vxlan_20	00:00:00:00:00:00	10.1.1.2	permanent	self	00:10:53
untagged		vxlan_20	44:38:39:00:00:17	10.1.1.2		self	00:00:47
untagged	bridge	swp1	44:38:39:00:00:29		permanent		01:16:49
untagged	bridge	swp2	44:38:39:00:00:24		permanent		01:16:49
untagged	bridge	vxlan_10	de:de:e3:a7:dc:8d		permanent		01:16:49
untagged	bridge	vxlan_20	c2:a0:61:d2:77:b1		permanent		01:16:49

We can also look at the MAC addresses per VNI:

```
leaf03# show evpn mac vni 10200
Number of MACs (local and remote) known for this VNI: 2
MAC                Type  Intf/Remote VTEP      VLAN
44:38:39:00:00:17  remote 10.1.1.2
44:38:39:00:00:23  local  swp2                20

leaf03# sh evpn mac vni 10100
Number of MACs (local and remote) known for this VNI: 2
MAC                Type  Intf/Remote VTEP      VLAN
44:38:39:00:00:1b  remote 10.1.1.2
44:38:39:00:00:28  local  swp1                10
```

As clearly seen above, EVPN is able to learn and exchange MAC addresses via the MP-BGP routing protocol while keeping tenant separation.

EVPN VXLAN ACTIVE-ACTIVE MODE

EVPN can also exchange control plane information in a VXLAN active-active mode environment, as depicted in Figure 8. **Multi-chassis Link Aggregation Group (MLAG)** is configured between the two leaf switches and a logical VTEP is configured using a shared, anycast IP address representing the VTEP. EVPN interacts with MLAG transitions and advertises and withdraws routes appropriately.

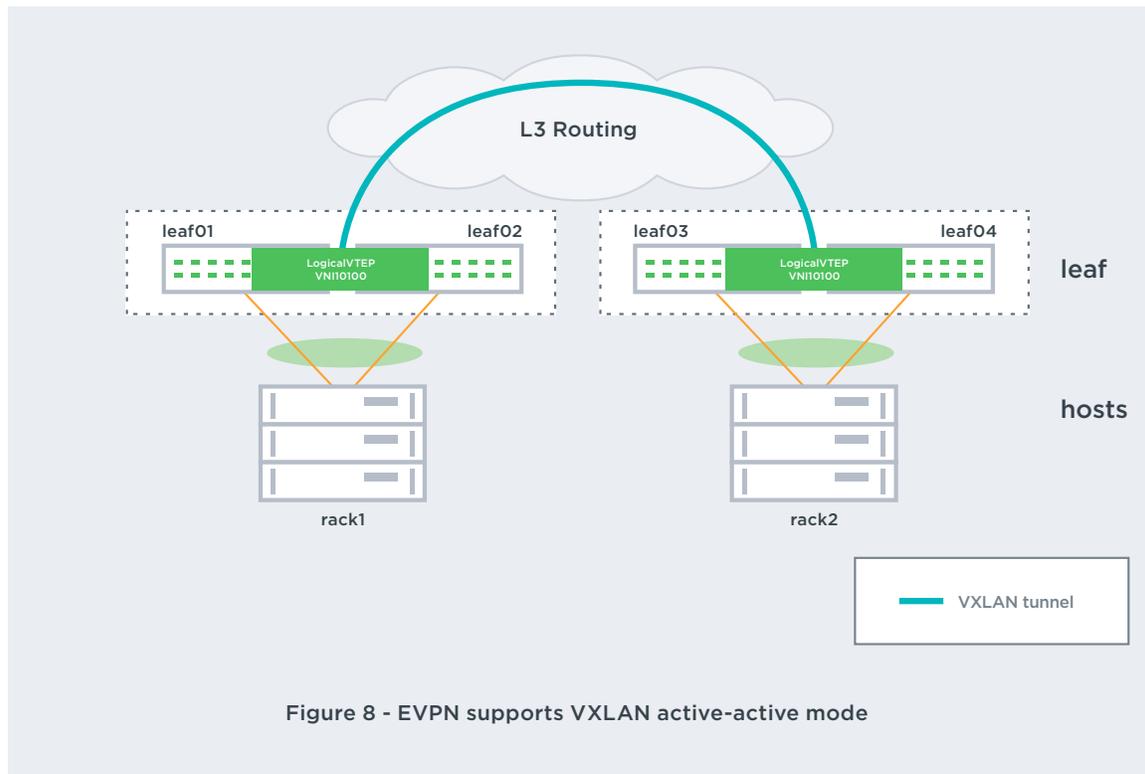


Figure 8 - EVPN supports VXLAN active-active mode

EVPN MAC MOBILITY

Cumulus EVPN supports a new “MAC Mobility” extended BGP community that enables quick sub-second convergence during host or VM moves within the data center. This community conveys sequence numbers along with the MAC address and is advertised in the Type 2 routes.

EVPN offers fast convergence during host or VM moves in a data center

As a local switch learns a new MAC address, MP-BGP sends the EVPN route to its peers without the MAC mobility community. However, upon a first move, an initial sequence number is sent in the MAC mobility community along with the update. The BGP table stores the sequence number along with the route. See Figure 9 below:

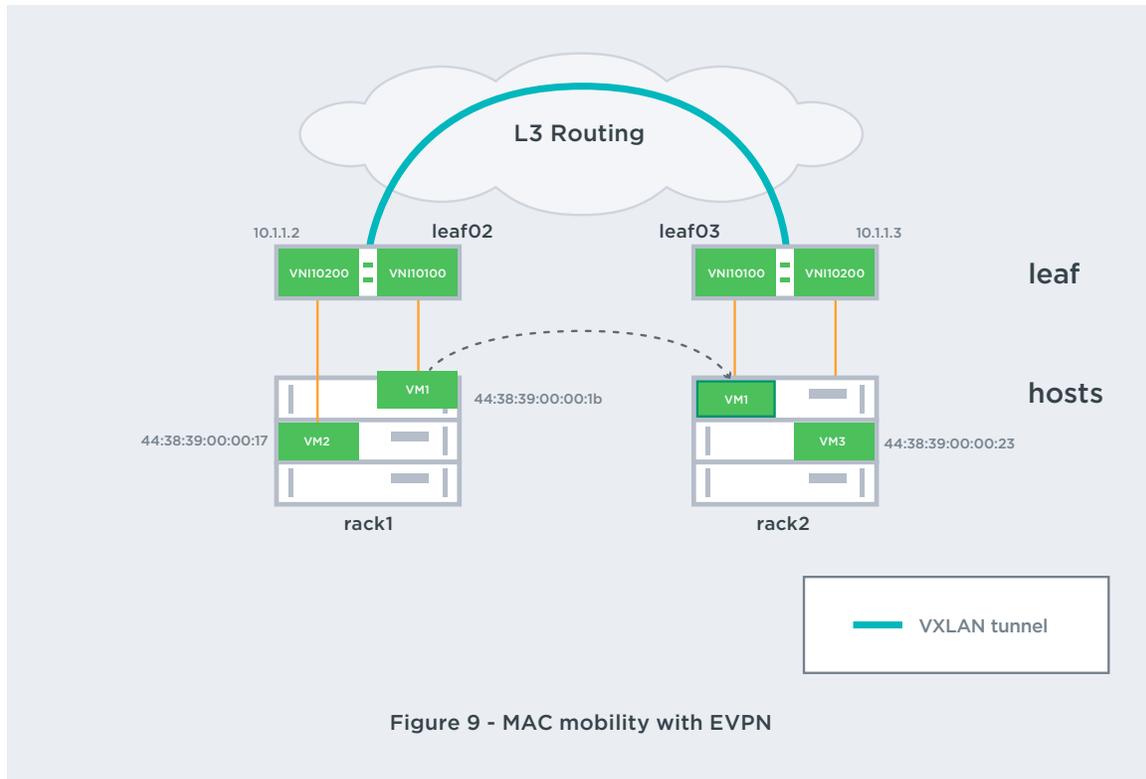


Figure 9 - MAC mobility with EVPN

In this case, VM1 (MAC address 44:38:39:00:00:1b) is moved to Rack 2. Before it moves, we can see on leaf03 there are two routes to this VM, one through spine01 and one through spine02.

```
leaf03# show bgp evpn route rd 10.1.1.2:10100 mac 44:38:39:00:00:1b
BGP routing table entry for 10.1.1.2:10100:[2]:[0]:[0]:[6]:[44:38:39:00:00:1b]
Paths: (2 available, best #2)
  Advertised to non peer-group peers:
    spine01(swp51) spine02(swp52)
Route [2]:[0]:[0]:[6]:[44:38:39:00:00:1b] VNI 10100
  65000 65002
    10.1.1.2 from spine02(swp52) (10.10.2.2)
      Origin IGP, localpref 100, valid, external
      Extended Community: RT:65002:10100 ET:8      AddPath ID: RX 0, TX 35
      Last update: Fri Feb 3 21:50:17 2017

Route [2]:[0]:[0]:[6]:[44:38:39:00:00:1b] VNI 10100
  65000 65002
    10.1.1.2 from spine01(swp51) (10.10.2.1)
      Origin IGP, localpref 100, valid, external, bestpath-from-AS 65000, best
      Extended Community: RT:65002:10100 ET:8
      AddPath ID: RX 0, TX 22
      Last update: Fri Feb 3 21:11:58 2017
```

When the host or VM moves, the new local switch (in this case leaf03) learns of the change via GARP or the data plane. The local switch then installs the MAC address into its bridge forwarding database, and BGP reads it. BGP then compares the MAC address with its current BGP table. If the MAC address is already there, BGP then increments the sequence number in this community (it is assumed to be 0 if the community is not there) before advertising the address to remote peers. The remote peers similarly compare the BGP extended MAC mobility community's

sequence numbers between two identical routes (the new one versus any that are already in the table), and install the route with the highest sequence number into the EVPN table, which then gets installed to the local bridge forwarding database. Use of the MAC mobility community with the sequence numbers ensure all applicable VTEPs converge quickly on the latest route to the MAC address. Below shows the output on the new local leaf (leaf03) after the move. The MAC Mobility community (MM) is now shown and the MAC address has moved once.

```

leaf03# show bgp evpn route vni 10100 mac 44:38:39:00:00:1b
BGP routing table entry for [2]:[0]:[0]:[6]:[44:38:39:00:00:1b]
Paths: (1 available, best #1)
  Not advertised to any peer
  Route [2]:[0]:[0]:[6]:[44:38:39:00:00:1b] VNI 10100
  Local
    10.1.1.3 from 0.0.0.0 (10.1.1.3)
      Origin IGP, localpref 100, weight 32768, valid, sourced, local, bestpath-
from-AS Local, best
      Extended Community: RT:65003:10100 ET:8 MM:1
      AddPath ID: RX 0, TX 18
      Last update: Sat Feb  4 02:26:56 2017

Displayed 1 paths for requested prefix

```

The new remote leaf (leaf02) shows the following:

```

leaf02# show bgp evpn route vni 10100 mac 44:38:39:00:00:1b
BGP routing table entry for [2]:[0]:[0]:[6]:[44:38:39:00:00:1b]
Paths: (2 available, best #2)
  Not advertised to any peer
  Route [2]:[0]:[0]:[6]:[44:38:39:00:00:1b] VNI 10100
  Imported from 10.1.1.3:10100:[2]:[0]:[0]:[6]:[44:38:39:00:00:1b]
65000 65003
    10.1.1.3 from spine01(swp51) (10.10.2.1)
      Origin IGP, localpref 100, valid, external
      Extended Community: RT:65003:10100 ET:8 MM:1
      AddPath ID: RX 0, TX 68
      Last update: Sun Feb  5 18:35:37 2017

  Route [2]:[0]:[0]:[6]:[44:38:39:00:00:1b] VNI 10100
  Imported from 10.1.1.3:10100:[2]:[0]:[0]:[6]:[44:38:39:00:00:1b]
65000 65003
    10.1.1.3 from spine02(swp52) (10.10.2.2)
      Origin IGP, localpref 100, valid, external, bestpath-from-AS 65000, best
      Extended Community: RT:65003:10100 ET:8 MM:1
      AddPath ID: RX 0, TX 66
      Last update: Sun Feb  5 18:35:37 2017

```

EVPN deployment scenarios and configuration

EVPN is used as the control plane solution for extending layer 2 connectivity across a data center using layer 3 fabric or it can be used to provide layer 2 connectivity between data centers.

Naturally, VTEPs must be configured on the leaf switches for the data plane traffic. Below is a snippet of a sample VXLAN configuration in a leaf with VXLAN active-active mode. The MLAG and layer 3 configurations are left off for brevity.

```
interface lo
  address 10.0.0.11/32
  clagd-vxlan-anycast-ip 10.0.0.20

interface swp2
  alias host facing interface
  bridge-access 100

interface swp51
  alias spine facing interface bgp unnumbered

interface vxlan_1
  bridge-access 100
  vxlan-id 10100
  bridge-learning off
  vxlan-local-tunnelip 10.0.0.11

interface bridge
  bridge-ports vxlan_1 swp2
  bridge-vids 100
  bridge-vlan-aware yes
```

As seen on previous page, the active-active mode VXLAN VNI 10100 is configured with the anycast address of 10.0.0.20. There is only one VXLAN tunnel to the remote leaf switch. To prevent data plane learning, bridge-learning is turned off. The locally connected bridge is associated with both the host facing interface (swp2) as well as the VXLAN interface (vxlan_1). A VXLAN interface and bridge interface must be configured on every switch with a desired VTEP. For the active-active scenario, the routing protocol must advertise the anycast VTEP IP address (10.0.0.20) to remote VTEPs. More information about configuring VXLAN in active-active mode with EVPN can be found in the [Cumulus Linux user guide](#).

The MP-BGP EVPN control plane running Cumulus Linux can be deployed in three layer 3 routed environments:

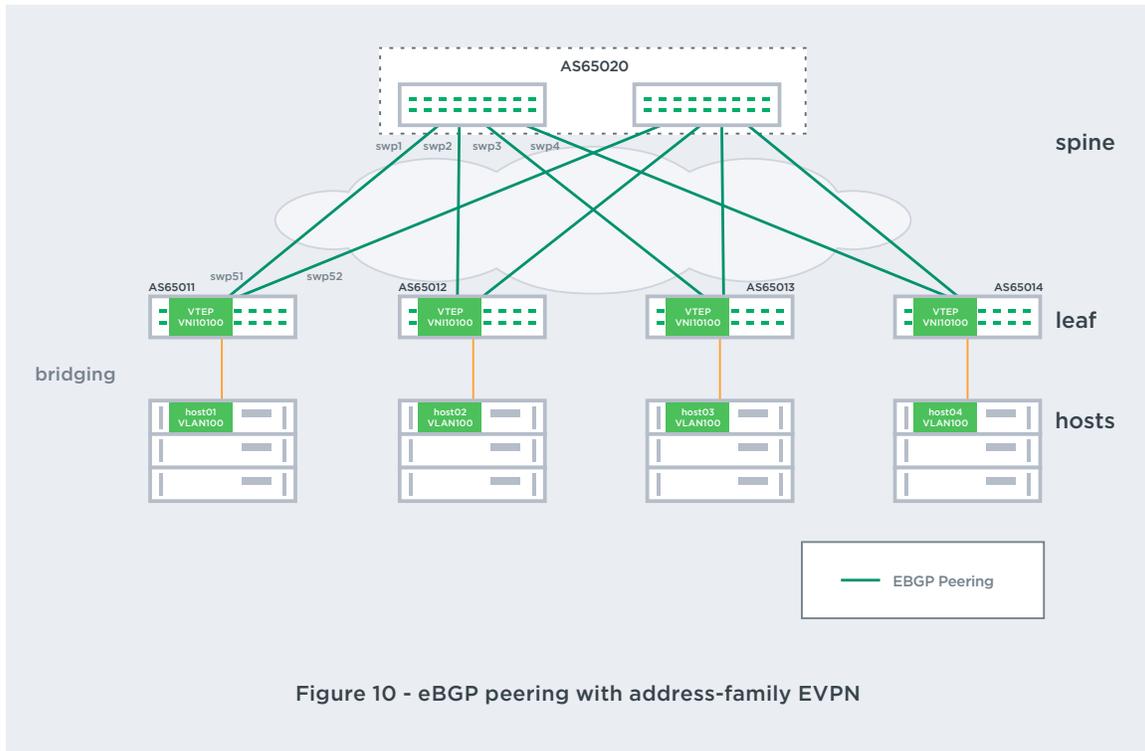
- **eBGP between the VTEPs (leafs) and spines**
- **iBGP between the VTEPs (leafs) with OSPF underlay**
- **iBGP between the VTEPs (leafs) and route reflectors (spines)**

Although Cumulus Linux supports all options mentioned above, Cumulus Networks recommends deploying eBGP for greenfield deployments. eBGP is already the most preferred data center routing protocol for the underlay network and the same session can carry the overlay EVPN routes also.

Cumulus recommends deploying EVPN with eBGP for simplicity

EVPN IN AN EBGP ENVIRONMENT

In this scenario, you peer the leafs and the spines together as in a typical eBGP data center, activating the neighbors in the `evpn` address family. Cumulus Linux also supports **eBGP unnumbered** to further simplify configuration. See Figure 10.



Using the Figure 10 scenario with eBGP unnumbered, an example simple leaf EVPN configuration is shown below using automatic RD/RT assignment.

```

router bgp 65001
  bgp router-id 10.0.0.11
  neighbor swp51 interface remote-as external
  neighbor swp52 interface remote-as external

  !
  address-family ipv4 unicast
    network 10.0.0.11/32
  exit-address-family
  !
  address-family evpn
    neighbor swp51 activate
    neighbor swp52 activate
    advertise-all-vni
  exit-address-family
    
```

A sample spine configuration is shown below.

```
router bgp 65020
  bgp router-id 10.0.0.21
  neighbor swp1 interface remote-as external
  neighbor swp2 interface remote-as external
  neighbor swp3 interface remote-as external
  neighbor swp4 interface remote-as external
  !
  address-family ipv4 unicast
    network 10.0.0.21/32
  exit-address-family
  !
  address-family evpn
    neighbor swp1 activate
    neighbor swp2 activate
    neighbor swp3 activate
    neighbor swp4 activate
  exit-address-family
```

Note the EVPN address family is needed on the spines to forward the EVPN routes, but the command **advertise-all-vni** is not needed on the spines unless VTEPs are also located on the spines.

More information on configuring EVPN with eBGP can be found in the [Cumulus Linux user guide](#).

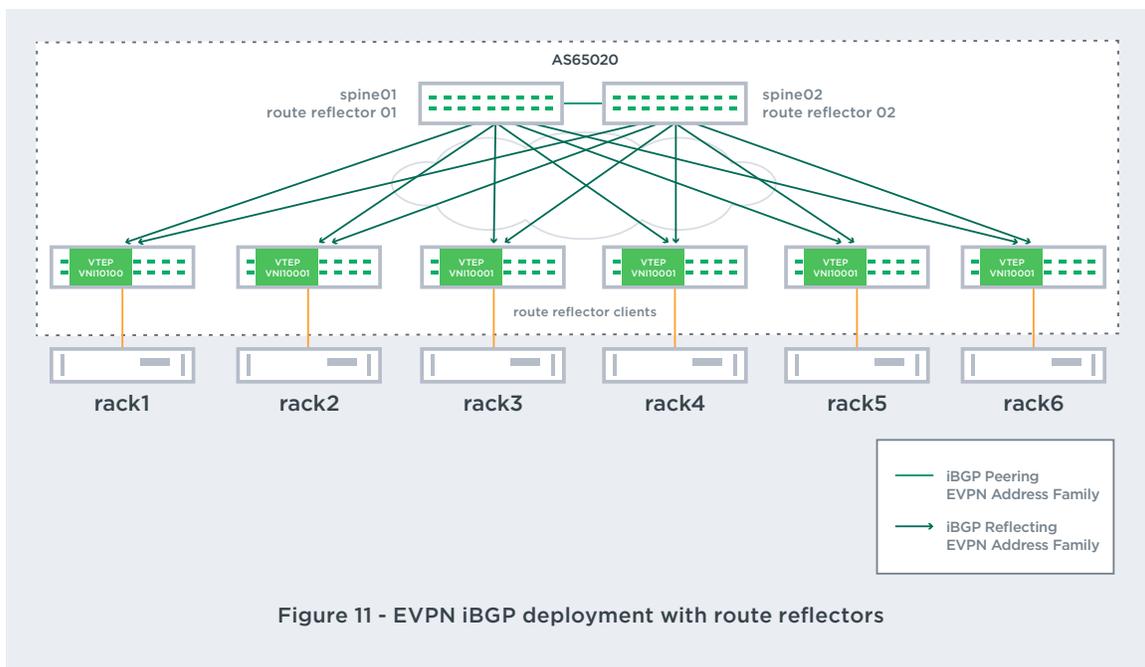
EVPN IN AN IBGP ENVIRONMENT WITH OSPF UNDERLAY

EVPN can also be deployed with an OSPF or static route underlay if needed, but is more complex than the eBGP solution. In this case, iBGP advertises EVPN routes directly between VTEPs and the spines are unaware of EVPN or BGP. The leaf switches peer with each other in a full mesh within the EVPN address family, and generally peer to the leaf loopback addresses which is advertised in OSPF. The receiving VTEP imports routes into a specific VNI with a matching route target community.

EVPN IN AN IBGP ENVIRONMENT WITH ROUTE REFLECTORS

With this scenario, the spines are route reflectors (RR) and reflect EVPN routes between the leaves. This scenario may be necessary for scale, and/or if iBGP is desired with no OSPF underlay. The EVPN address family must be run on the spines (RRs), but the command “advertise-all-vni” is not needed. Although the RRs receive all the MAC address routes associated with the VXLANs, they are not put into hardware on the RRs allowing for greater scale.

To provide redundancy, two spine switches should be configured as RRs within the EVPN address family. It is recommended to use the same cluster-ID on the redundant route reflectors to reduce the total number of stored routes. More information on configuring RRs can be found in [the Cumulus Linux user guide](#). See Figure 11.



If a three tier Clos network is desired without an OSPF underlay, tiers of route reflectors must be deployed.

If more than one pod is needed or the data center expands with all iBGP, the use of additional clusters is recommended. A cluster

consists of one or more route reflectors and their clients. Each route reflector in each cluster peers with each other as well as any other cluster’s route reflectors. Assigning different cluster IDs (a BGP attribute) to each cluster prevents looping of routes between different clusters.

Conclusion

Data centers are moving towards a layer 3 fabric in order to scale, provide ease of troubleshooting, and provide redundancy with multi-vendor interoperability. However, some applications still require layer 2 connectivity. For these applications, VXLAN tunnels are being widely deployed to provide a scalable layer 2 overlay solution over a layer 3 fabric.

Cumulus EVPN is the ideal control plane solution for VXLAN tunnels. It provides controllerless VXLAN tunnels that also scale, provide redundancy and enable fast convergence. It reduces unnecessary BUM traffic, thereby reducing overall data center traffic as well as providing multi-tenant segmentation.

Try it out for yourself on this [ready-to-go demo](#) using [Cumulus VX](#) and Vagrant.

ABOUT CUMULUS NETWORKS*

Cumulus Networks is leading the transformation of bringing web-scale networking to enterprise cloud. Its network switch, Cumulus Linux, is the only solution that allows you to affordably build and efficiently operate your network like the world's largest data center operators, unlocking vertical network stacks. By allowing operators to use standard hardware components, Cumulus Linux offers unprecedented operational speed and agility, at the industry's most competitive cost. Cumulus Networks has received venture funding from Andreessen Horowitz, Battery Ventures, Capital, Peter Wagner and four of the original VMware founders.

For more information visit cumulusnetworks.com or follow [@cumulusnetworks](https://twitter.com/cumulusnetworks).

©2017 Cumulus Networks. All rights reserved. CUMULUS, the Cumulus Logo, CUMULUS NETWORKS, and the Rocket Turtle Logo (the "Marks") are trademarks and service marks of Cumulus Networks, Inc. in the U.S. and other countries. You are not permitted to use the Marks without the prior written consent of Cumulus Networks. The registered trademark Linux[®] is used pursuant to a sublicense from LMI, the exclusive licensee of Linus Torvalds, owner of the mark on a worldwide basis. All other marks are used under fair use or license from their respective owners.

02222017