# EVPN for VXLAN Tunnels (Layer 2)

## In This Chapter

This section provides information about Layer 2 and EVPN.

Topics in this section include:

# Applicability

This example is applicable to the 7950 XRS, 7750 SR-c4/c12, 7750 SR-7/12 and 7450 ESS-6/6v/7/12, but it is not supported on the 7750 SR-1, 7450 ESS-1 or 7710 SR. Virtual eXtensible Local Area Network (VXLAN) requires IOM3-XP/IMM or higher-based line cards and chassis-mode D. Ethernet Virtual Private Networks (EVPN) is a control plane technology and does not have line card hardware dependencies.

The configuration was tested in release 12.0.R4.

# Overview

SR OS supports the EVPN control plane with Virtual eXtensible Local Area Network (VXLAN) data plane in VPLS services.

EVPN is an IETF technology (draft-ietf-l2vpn-evpn) that uses a new BGP address family which allows VPLS services to be operated in a similar way to IP-VPNs, where the MAC addresses, IP addresses and the information to set up the flooding tree are distributed by BGP. While EVPN can be used as the control plane for different data plane encapsulations, only VXLAN is supported in SR OS in the release tested.

VXLAN (draft-mahalingam-dutt-dcops-vxlan) is an overlay IP tunneling technology used to carry Ethernet traffic over any IP network, and it is becoming the de-facto standard for overlay data centers and networks. Compared to other IP overlay tunneling technologies, such as GRE, VXLAN supports multi-tenancy and multi-pathing:

- A tenant identifier, the VXLAN Network Identifier (VNI), is encoded in the VXLAN header and allows each tenant to have an isolated Layer 2 domain.
- VXLAN supports multi-pathing scalability through ECMP. VXLAN uses the outer source UDP port as an entropy field that can be used by the core IP routers to balance the load across different paths.

In SR OS, EVPN and VXLAN can be enabled in VPLS or R-VPLS services. In this example, EVPN-VXLAN services will refer to VPLS or R-VPLS services with EVPN and VXLAN enabled. These services can terminate/originate VXLAN tunnels and may have SAPs and/or SDP bindings at the same time. Some other SR OS implementation-specific considerations are listed below:

- VXLAN is only supported on network or hybrid ports with null or dot1q encapsulation on Ethernet/LAG/POS/APS interfaces.
- VXLAN packets are originated/terminated with the system IPv4 address, in other words, a system originating VXLAN packets will use the system IP address as source outer IPv4 address and systems will only process VXLAN packets if their destination outer IPv4 address matches its own system IP address.
- Data plane MAC learning is not supported over VXLAN bindings. Only the control plane (EVPN) will be used for populating the FDB with MAC addresses associated to VXLAN bindings.

- EVPN provides support for the following features that are tested in this document:
  - ç The BGP advertisement of the MAC addresses learned on SAPs, SDP-bindings and conditional static MACs to the remote BGP peers. The advertisement of MAC addresses in BGP can optionally be disabled.
  - ç The optional advertisement of an unknown MAC route, that allows the remote EVPN PEs or Network Virtualization Edge devices (NVEs) to suppress the unknown unicast flooding and send any unknown unicast frame to the owner of the unknown MAC route.
  - ç Ingress replication of Broadcast, Unknown unicast and Multicast (BUM) packets over VXLAN.
  - ç A Proxy-ARP table per service populated by the MAC-IP pairs received in BGP MAC advertisements. When an ARP request is received on a SAP or SDP-binding, the system will perform a lookup on this table and will reply to the ARP request if the lookup yields a valid result.
  - ç MAC mobility and static-mac protection as described in draft-ietf-l2vpn-evpn, as well as MAC duplication detection.
- Multi-homing redundancy for SAPs and SDP-bindings in EVPN-VXLAN services is supported through BGP Multi-homing (L2VPN BGP address family). Only one BGP-MH site is supported in an EVPN-VXLAN service.

One of the main applications for EVPN-VXLAN services in SR OS is the Data Center Gateway (DC GW) function. In such an application, EVPN and VXLAN are expected to be used within the Data Center and VPLS SDP-bindings or SAPs are expected to be used for the connectivity to the WAN. When the system is used as a DC GW a VPLS service is configured per Layer 2 domain that has to be extended to the WAN. In those VPLS services, BGP EVPN automatically sets up the VXLAN auto-bindings that connect the DC GW to the Data Center NVEs. The WAN connectivity is based on regular VPLS constructs where SAPs (null, dot1q and QinQ), spoke-SDPs (FEC type 128 and 129, not BGP-VPLS) and mesh-SDPs are supported. B-VPLS or I-VPLS services are not supported.

Although the DC GW application is one of the most common uses for this feature, this example focuses on the configuration and operation of EVPN-VXLAN for Layer 2 services in general, and its integration with regular VPLS services in MPLS networks.

# Configuration

This section describes the configuration of EVPN-VXLAN on the 7x50 as well as the available troubleshooting and show commands. This example focuses on the following configuration aspects:

- Enabling EVPN and VXLAN in a VPLS service, including the use of BGP-EVPN, BGP-AD (BGP Auto-discovery) and BGP-MH (BGP Multi-homing) in the same VPLS instance.
- Scaling BGP-MH resiliency with the use of operational groups (oper-groups).
- Use of proxy-arp in EVPN-VXLAN services
- MAC mobility, MAC duplication and MAC protection in EVPN-VXLAN services.

The configuration will be shown for PE-71, PE-69 and PE-72 only; the PEs in Overlay-Network-2 (Figure 44) have an equivalent configuration.

# Enabling EVPN-VXLAN in a VPLS Service
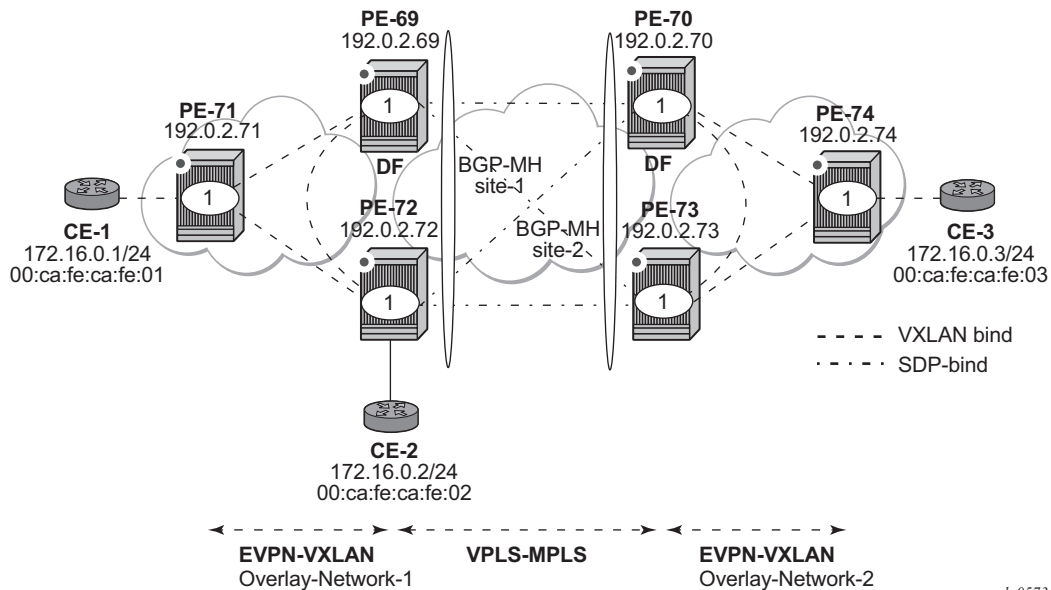
Figure 44 shows the topology used in this example.



**Figure 44: EVPN-VXLAN Topology**

The network topology shows two overlay (VXLAN) networks interconnected by an MPLS network:

- PE-69, PE-71 and PE-72 are part of Overlay-Network-1
- PE-70, PE-73 and PE-74 are part of Overlay-Network-2

CE-1, CE-2 and CE-3 belong to the same IP subnet, hence Layer 2 connectivity must be provided to them.

Note that the above topology can illustrate a DCI (Data Center Interconnect) example, where Overlay-Network-1 and Overlay-Network-2 are two Data Centers interconnected through an MPLS WAN. In this application, CE-1, CE-2 and CE-3 would simulate Virtual Machines or appliances, PE-69/70/72/73 would act as DC GWs and PE-71/74 as NVEs (or virtual PEs running on compute infrastructure).

The following protocols and objects are configured beforehand:

- The ports interconnecting the six PEs in Figure 44 are configured as network ports (or hybrid) and have router network interfaces defined on them. Only the ports connected to the CEs are configured as access ports.
- The six PEs shown in the Figure 44 are running IS-IS for the global routing table with the four core PEs interconnected using IS-IS Level-2 point-to-point interfaces and each overlay network is using IS-IS Level-1 point-to-point interfaces.
- LDP is used as the MPLS protocol to signal transport tunnel labels among PE-69, PE-72, PE-70 and PE-73. There is no LDP running in the two overlay networks.
- Note that the network port MTU (in all the ports sending/receiving VXLAN packets) must be at least 50-bytes (54 if dot1q encapsulation is used) greater than the service-mtu in order to accommodate the size of the VXLAN header.

Once the IGP infrastructure and LDP are enabled in the core, BGP has to be configured. In this example, two BGP families have to be enabled: EVPN within each overlay-network for the exchange of MAC/IP addresses and setting up the flooding domains, and L2-VPN for the use of BGP-MH and BGP-AD in the VPLS-MPLS network.

As an example, the following CLI output shows the relevant BGP configuration of PE-71, which only needs the EVPN family. PE-74 would have a similar BGP configuration. Note that the use of Route-Reflectors (RRs) in these type of scenarios is common. Although this example does not use RRs, an EVPN RR could have been used in Overlay-Network-1 and Overlay-Network-2 and an L2-VPN RR could have been used in the core VPLS-MPLS network.

```
A:PE-71>config>router>bgp# info
--------------------------------------------
            vpn-apply-import
            vpn-apply-export
            min-route-advertisement 1
            enable-peer-tracking
            rapid-withdrawal
```

```
            rapid-update evpn
            group "DC"
                family evpn
                type internal
                neighbor 192.0.2.69
                exit
                neighbor 192.0.2.72
                exit
            exit
            no shutdown
---------------------------------------------
```

The BGP configuration of PE-69 and PE-72 follows (PE-70 and PE-73 have an equivalent configuration).

```
A:PE-69>config>router>bgp# info
---------------------------------------------
            vpn-apply-import
            vpn-apply-export
            min-route-advertisement 1
            enable-peer-tracking
            rapid-withdrawal
            rapid-update l2-vpn evpn
            group "DC"
                family l2-vpn evpn
                type internal
                neighbor 192.0.2.71
                exit
                neighbor 192.0.2.72
                exit
            exit
            group "WAN"
                family l2-vpn
                type internal
                neighbor 192.0.2.70
                exit
                neighbor 192.0.2.73
                exit
            exit
            no shutdown
---------------------------------------------
A:PE-72>config>router>bgp# info
---------------------------------------------
            vpn-apply-import
            vpn-apply-export
            min-route-advertisement 1
            enable-peer-tracking
            rapid-withdrawal
            rapid-update l2-vpn evpn
            group "DC"
                family l2-vpn evpn
                type internal
                neighbor 192.0.2.69
                exit
                neighbor 192.0.2.71
                exit
            exit
```

```
                    group "WAN"
                        family l2-vpn
                        type internal
                        neighbor 192.0.2.70
                        exit
                        neighbor 192.0.2.73
                        exit
                    exit
                    no shutdown
    --------------------------------------------
```

Figure 45 shows the BGP peering sessions among the PEs and the enabled BGP families. Note that, for instance, PE-71 will only establish an EVPN peering session with its peers (only the EVPN family is enabled on PE-71), even though PE-69 and PE-72 have EVPN and L2-VPN families configured.



**Figure 45: BGP Adjacencies and Enabled Families**

Once the network infrastructure is running properly, the actual service configuration can be carried out. The following CLI outputs show the configuration of VPLS 1 in PE-71, PE-69 and PE-72 as per the topology illustrated in Figure 44.

VPLS 1 in those three PEs are interconnected using VXLAN bindings, whereas PE-69 and PE-72 are connected to the remote PEs by means of BGP-AD SDP-bindings. Although BGP-AD SDP-bindings are used in this example for the connectivity of the EVPN-VXLAN PEs to a regular VPLS network, SAPs, manual spoke-SDPs or mesh-SDPs could have been used instead. BGP-VPLS cannot be enabled in an EVPN-VXLAN VPLS service.

VPLS 1 configuration of PE-71 is shown below:

```
A:PE-71>config>service>vpls# info
```

```
--------------------------------------------
          vxlan vni 1 create
          exit
          bgp
              route-distinguisher 192.0.2.71:1
              route-target export target:64500:12 import target:64500:12
          exit
          bgp-evpn
              vxlan
                  no shutdown
              exit
          exit
          stp
              shutdown
          exit
          sap 1/1/1:1 create
          exit
          no shutdown
--------------------------------------------
```

EVPN-VXLAN is enabled by the configuration of a valid VXLAN Network Identifier (VNI) and the **bgp-evpn>vxlan>no shutdown** command. These two commands, along with the required bgp route-distinguisher (RD) and route-target (RT) information, are the minimum mandatory attributes:

- The VNI is a 24-bit identifier with valid values in the [1..16777215] range. This defines the VNI that the 7x50 will use in the EVPN routes generated for the VPLS service, and therefore the VNI that the system expects to see in the VXLAN packets destined to that particular VPLS service. Note that the configured VNI determines the VNI that has to be received in the packets for the VPLS service, but not the VNI that will be sent in VXLAN packets to remote PEs for the service. In other words, in this example, VPLS 1 is configured with VNI=1 in all the PEs, however each PE could have used a different VNI. Note that the VNI is a system-wide significant value and two VPLS services cannot be configured with the same VNI.

- The **bgp-evpn>vxlan>no shutdown** command enables the use of EVPN for VXLAN. It requires the previous configuration of the VNI, RD and RT. As soon as this command is executed, EVPN will advertise an inclusive multicast route to all of the BGP EVPN peers (regardless of the existing SAP/SDP-binding operational status). The exchange of inclusive multicast routes allows the establishment of the VXLAN bindings among the PEs.

Upon the reception of the EVPN inclusive multicast routes from PE-69 and PE-72, PE-71 will automatically setup its VXLAN bindings for VPLS-1. A VXLAN binding is represented by an (egress VTEP, egress VNI) pair, where VTEP is a VXLAN Termination End Point. This can be shown with the following show commands:

```
*A:PE-71# show service id 1 vxlan
===============================================================================
VPLS VXLAN, Ingress VXLAN Network Id: 1
===============================================================================
```

```
Egress VTEP, VNI
===============================================================================
VTEP Address           Egress VNI     Num. MACs     In Mcast List? Oper State
-------------------------------------------------------------------------------
192.0.2.69             1              1             Yes            Up
192.0.2.72             1              1             Yes            Up
-------------------------------------------------------------------------------
Number of Egress VTEP, VNI : 2
-------------------------------------------------------------------------------
===============================================================================

*A:PE-71# show service vxlan
===============================================================================
VXLAN Tunnel Endpoints (VTEPs)
===============================================================================
VTEP Address                 Number of Egress VNIs   Oper State
-------------------------------------------------------------------------------
192.0.2.69                   2                       Up
192.0.2.72                   2                       Up
-------------------------------------------------------------------------------
Number of VTEPs: 2
-------------------------------------------------------------------------------
===============================================================================
```

As can be seen in the CLI output, PE-71 has two VXLAN bindings, one to PE-69 and one to PE-72. Both use egress VNI=1 (the actual VNI used in its egress VXLAN packets) and both are part of the flooding multicast list for VPLS 1 and are UP.

- The **In Mcast List? = Yes** entry is set when the proper inclusive multicast route is received from the remote VTEP. If the entry is **No,** the VXLAN binding will not be use to flood BUM (Broadcast, Unknown unicast, Multicast) packets.
- The **Oper State** is based on the existence of the VTEP in the global routing table.

The VPLS 1 configuration of PE-69 and PE-72 is shown below:

```
A:PE-69>config>service>vpls# info
----------------------------------------------
        vxlan vni 1 create
        exit
        bgp
            route-distinguisher 192.0.2.69:1
            vsi-export "vsi-policy-1"
            vsi-import "vsi-policy-1"
            pw-template-binding 1 split-horizon-group "CORE"
            exit
        exit
        bgp-ad
            vpls-id 64500:1
            no shutdown
        exit
        bgp-evpn
            unknown-mac-route
            vxlan
                no shutdown
            exit
```

```
            exit
            stp
                shutdown
            exit
            site "site-1" create
                site-id 1
                split-horizon-group CORE
                no shutdown
            exit
            no shutdown
---------------------------------------------

A:PE-72>config>service>vpls# info
---------------------------------------------
            vxlan vni 1 create
            exit
            bgp
                route-distinguisher 192.0.2.72:1
                vsi-export "vsi-policy-1"
                vsi-import "vsi-policy-1"
                pw-template-binding 1 split-horizon-group "CORE"
                exit
            exit
            bgp-ad
                vpls-id 64500:1
                no shutdown
            exit
            bgp-evpn
                unknown-mac-route
                vxlan
                    no shutdown
                exit
            exit
            proxy-arp
                no age-time
                no send-refresh
                no shutdown
            exit
            stp
                shutdown
            exit
            site "site-1" create
                site-id 1
                split-horizon-group CORE
                no shutdown
            exit
            sap 1/1/1:1 create
            exit
            no shutdown
---------------------------------------------
```

In addition to the VNI and **bgp-evpn>vxlan>no shutdown** commands for enabling EVPN-VXLAN in VPLS 1, PE-69 and PE-72 require the configuration of BGP-AD for the discovery and establishment of FEC129 spoke SDPs to the remote PEs in the core, as well as BGP-MH for redundancy. As outlined in Figure 44, there are two BGP-MH sites defined in the network: site-1 is used on PE-69/PE-72 and site-2 is used on PE-70/PE-73. Only one of the two gateway PEs in

each Overlay-Network will be the Designated Forwarder (DF) for VPLS 1, and only the DF will send/receive traffic for VPLS 1 in the Overlay-Network. The following considerations must be taken into account when configuring the connectivity of EVPN-VXLAN services to regular VPLS objects:

- As discussed, in this example, BGP-AD spoke-SDPs are used but SAPs, manual spoke-SDPs or mesh-SDPs are also supported.

- In this example, BGP-AD spoke-SDPs are auto-instantiated using **pw-template-binding 1 split-horizon-group "CORE".**

  ç Although not shown above, this requires the creation of the pw-template 1 (**config>service>pw-template 1 create**).

- The split-horizon-group CORE is added to the BGP-MH site "site-1". This statement will ensure that all the spoke SDPs automatically established to the remote PEs are part of the BGP-MH site.

- Although the route-targets for the Overlay-Network and the VPLS-MPLS network can have the same value for the same VPLS service, they are usually different. This example assumes the use of RT-DC-1 in Overlay-Network-1 and RT-WAN-1 in the VPLS-MPLS core for VPLS 1. The **vsi-policy-1** allows the system to export and import the right RTs for VPLS 1:

```
A:PE-69>config>router>policy-options# info
----------------------------------------------
            community "RT-DC-1" members "target:64500:12"
            community "RT-WAN-1" members "target:64500:11"
            policy-statement "vsi-policy-1"
                entry 10  # to import all the evpn routes with RT-DC-1
                    from
                        community "RT-DC-1"
                        family evpn
                    exit
                    action accept
                    exit
                exit
                entry 20 # to import all the bgp-ad/mh routes from the WAN
                    from
                        community "RT-WAN-1"
                        family l2-vpn
                    exit
                    action accept
                    exit
                exit
                entry 30  # to export all the evpn routes with RT-DC-1
                    from
                        family evpn
                    exit
                    action accept
                        community add "RT-DC-1"
                    exit
                exit
                entry 40  # to export all the bgp-ad/mh routes with RT-WAN-1
                    from
```

```
                        family l2-vpn
                   exit
                   action accept
                       community add "RT-WAN-1"
                   exit
               exit
             default-action reject
         exit
```

-----------------------------------------------

Once PE-69 and PE-72 are configured as above, they will setup the spoke SDPs and will run the DF election algorithm to determine the operational status of those spoke SDPs. Refer to **LDP VPLS using BGP-Auto Discovery** and **BGP Multi-Homing for VPLS Networks** for more information about the use of BGP-AD and BGP-MH.

Note that in the configuration for VPLS 1, both gateway PEs, PE-69 and PE-72 will attempt to establish two parallel Layer 2 paths between each other (a BGP-AD spoke SDP and a EVPN VXLAN binding). Since that would create a Layer 2 loop, the SR OS implementation gives priority to the EVPN path and only the VXLAN binding will be active. In other words, when an (egress VTEP, VNI) and a spoke SDP are attempted to be set up to the same far-end IP address at the same time, the VXLAN path will prevail and the spoke SDP will be kept down. The spoke SDP will only be brought up if the VXLAN (egress VTEP, VNI) goes down.

This behavior can be easily observed in this setup by using the following show commands. In PE-69, the spoke SDP to far-end PE-72 will be down with a **EvpnRouteConflict** Flag. The (egress VTEP, VNI) = (192.0.2.72, 1) VXLAN bind will be UP.

```
A:PE-69# show service id 1 base
===============================================================================
Service Basic Information
===============================================================================
Service Id        : 1                   Vpn Id           : 0
Service Type      : VPLS
Name              : (Not Specified)
Description       : (Not Specified)
Customer Id       : 1                   Creation Origin  : manual
Last Status Change: 07/17/2014 00:03:52
Last Mgmt Change  : 07/17/2014 19:02:39
Etree Mode        : Disabled
Admin State       : Up                  Oper State       : Up
MTU               : 1514                Def. Mesh VC Id  : 1
...
-------------------------------------------------------------------------------
Service Access & Destination Points
-------------------------------------------------------------------------------
Identifier                         Type       AdmMTU  OprMTU  Adm  Opr
-------------------------------------------------------------------------------
sdp:17405:4294967290 SB(192.0.2.73)  BgpAd      0       8974    Up   Up
sdp:17406:4294967292 SB(192.0.2.72)  BgpAd      0       8974    Up   Down
sdp:17407:4294967294 SB(192.0.2.70)  BgpAd      0       8974    Up   Up
===============================================================================
A:PE-69# show service id 1 all | match Flag
Flags             : None
```

```
Flags               : EvpnRouteConflict
Flags               : None

A:PE-69# show service id 1 vxlan
===============================================================================
VPLS VXLAN, Ingress VXLAN Network Id: 1
===============================================================================
Egress VTEP, VNI
===============================================================================
VTEP Address        Egress VNI   Num. MACs   In Mcast List?  Oper State
-------------------------------------------------------------------------------
192.0.2.71          1            0           Yes             Up
192.0.2.72          1            0           Yes             Up
-------------------------------------------------------------------------------
Number of Egress VTEP, VNI : 2
-------------------------------------------------------------------------------
===============================================================================
```

At the non-DF, PE-72, all the spoke SDPs will be down due to BGP-MH:

```
A:PE-72# show service id 1 base
===============================================================================
Service Basic Information
===============================================================================
Service Id        : 1                   Vpn Id            : 0
Service Type      : VPLS
Name              : (Not Specified)
Description       : (Not Specified)
Customer Id       : 1                   Creation Origin   : manual
Last Status Change: 07/17/2014 00:03:42
Last Mgmt Change  : 07/17/2014 19:02:50
Etree Mode        : Disabled
Admin State       : Up                  Oper State        : Up
MTU               : 1514                Def. Mesh VC Id   : 1
SAP Count         : 1                   SDP Bind Count    : 3
...


-------------------------------------------------------------------------------
Service Access & Destination Points
-------------------------------------------------------------------------------
Identifier                         Type      AdmMTU  OprMTU  Adm  Opr
-------------------------------------------------------------------------------
sap:1/1/1:1                        q-tag     1518    1518    Up   Up
sdp:17405:4294967290 SB(192.0.2.73) BgpAd    0       8974    Up   Down
sdp:17406:4294967292 SB(192.0.2.69) BgpAd    0       8974    Up   Down
sdp:17407:4294967294 SB(192.0.2.70) BgpAd    0       8974    Up   Down
===============================================================================

A:PE-72# show service id 1 all | match Flag
Flags             : StandbyForMHProtocol
Flags             : StandbyForMHProtocol
Flags             : StandbyForMHProtocol
Flags             : None
```

# MAC Learning and unknown-mac-route

Once the VPLS service (VPLS 1) is configured, the network allows the CEs to exchange unicast and BUM traffic over the Overlay and VPLS-MPLS service infrastructure. BUM traffic sent by CE-1 will be ingress-replicated to PE-69 and PE-72 by PE-71, and propagated by PE-69 (the DF) to the remote network. From this point on, MAC addresses will be learned on active SAPs and spoke SDPs and advertised in EVPN MAC routes. No data plane MAC learning is carried out on VXLAN bindings. MACs associated with (egress VTEP, VNI) bindings will always be learned through EVPN.

The following CLI output shows the reception of an EVPN MAC route and how the (CE-2) MAC address appears in the FDB for VPLS 1.

```
33 2014/07/17 22:06:08.48 UTC MINOR: DEBUG #2001 Base Peer 1: 192.0.2.72
"Peer 1: 192.0.2.72: UPDATE
Peer 1: 192.0.2.72 - Received BGP UPDATE:
    Withdrawn Length = 0
    Total Path Attr Length = 88
    Flag: 0x90 Type: 14 Len: 44 Multiprotocol Reachable NLRI:
        Address Family EVPN
        NextHop len 4 NextHop 192.0.2.72
        Type: EVPN-MAC Len: 33 RD: 192.0.2.72:1 ESI: 0:0:0:0:0:0:0:0:0:0, tag: 1
, mac len: 48 mac: 00:ca:fe:ca:fe:02, IP len: 0, IP: NULL, label: 0
    Flag: 0x40 Type: 1 Len: 1 Origin: 0
    Flag: 0x40 Type: 2 Len: 0 AS Path:
    Flag: 0x80 Type: 4 Len: 4 MED: 0
    Flag: 0x40 Type: 5 Len: 4 Local Preference: 100
    Flag: 0xc0 Type: 16 Len: 16 Extended Community:
        target:64500:12
        bgp-tunnel-encap:VXLAN
"

*A:PE-71# show service id 1 fdb detail
===============================================================================
Forwarding Database, Service 1
===============================================================================
ServId    MAC               Source-Identifier       Type     Last Change
                                                    Age
-------------------------------------------------------------------------------
1         00:ca:fe:ca:fe:01 sap:1/1/1:1             L/0      07/17/14 22:06:08
1         00:ca:fe:ca:fe:02 vxlan:                  Evpn     07/17/14 22:06:08
                            192.0.2.72:1
-------------------------------------------------------------------------------
No. of MAC Entries: 2
-------------------------------------------------------------------------------
Legend:  L=Learned O=Oam P=Protected-MAC C=Conditional S=Static
===============================================================================
*A:PE-71#
```

When a frame destined to 00:ca:fe:ca:fe:02 enters SAP 1/1/1:1, it is encapsulated into a VXLAN packet with outer destination IP 192.0.2.72 and VNI 1, and sent on the wire.

In virtualized data center networks where all the MACs are known beforehand (all the virtual machine and appliance MACs are distributed by EVPN before any traffic flows), unknown MAC addresses are always outside the data center. If that is the case, the DC GWs can make use of the **unknown-mac-route** so that the DC NVEs supporting the concept of this route send the unknown unicast traffic only to the DC GW. This minimizes the flooding within the Data Center, as explained in draft-rabadan-l2vpn-dci-evpn-overlay.

In this example the unknown-mac-route is configured in the gateway PEs (PE-69, PE-72 and PE-70, PE-73) in the following way:

```
*A:PE-69>config>service>vpls# bgp-evpn unknown-mac-route
*A:PE-69>config>service>vpls#
27 2014/07/17 22:15:54.94 UTC MINOR: DEBUG #2001 Base Peer 1: 192.0.2.71
"Peer 1: 192.0.2.71: UPDATE
Peer 1: 192.0.2.71 - Send BGP UPDATE:
    Withdrawn Length = 0
    Total Path Attr Length = 88
    Flag: 0x90 Type: 14 Len: 44 Multiprotocol Reachable NLRI:
        Address Family EVPN
        NextHop len 4 NextHop 192.0.2.69
        Type: EVPN-MAC Len: 33 RD: 192.0.2.69:1 ESI: 0:0:0:0:0:0:0:0:0:0, tag: 1
, mac len: 48 mac: 00:00:00:00:00:00, IP len: 0, IP: NULL, label: 0
    Flag: 0x40 Type: 1 Len: 1 Origin: 0
    Flag: 0x40 Type: 2 Len: 0 AS Path:
    Flag: 0x80 Type: 4 Len: 4 MED: 0
    Flag: 0x40 Type: 5 Len: 4 Local Preference: 100
    Flag: 0xc0 Type: 16 Len: 16 Extended Community:
        target:64500:12
        bgp-tunnel-encap:VXLAN
"

...
```

Note that:

- Although the 7x50 can generate the unknown-mac-route, it will never honor it and normal flooding applies when an unknown unicast packet arrives at an ingress sap/sdp-binding.

- When unknown-mac-route is configured, it will ONLY be generated when: a) no BGP-MH site is configured within the same VPLS service or b) a site is configured AND the site is DF (Designated Forwarder) in the PE. If the site becomes a non-DF site, the unknown-mac-route will be withdrawn.

- If the unknown-mac-route is used in the DC GW and all the NVEs in the DC understand it, the advertisement of MAC addresses can be disabled with the [**no**] **mac-advertisement** command. If so, the 7x50 will only advertise the unknown-mac-route.

```
*A:PE-72>config>service>vpls>bgp-evpn# info
----------------------------------------------
    unknown-mac-route
    no mac-advertisement
    vxlan
        no shutdown
    exit
----------------------------------------------
```

# Scaling BGP-MH Resiliency with the Use of Operational Groups

In Figure 44, VPLS 1 in PE-69/PE-72 is configured with a BGP-MH site that controls which of the two PEs forwards the traffic to the remote PEs (in this case PE-69 is the DF and the gateway responsible for forwarding packets to the remote PEs).

When new VPLS services are required in PE-69/PE-72 the same BGP-MH configuration can be used. However, if the number of VPLS services grows significantly, the use of individual BGP-MH sites per service will not scale. Since all the services in these two PEs share the same physical topology, the use of oper-groups can provide a simple and scalable way of providing resiliency to as many services as the user needs (up to the maximum number of VPLS services per system).

The way oper-groups can be used to scale this type of deployments is the following (using the network topology in Figure 1 and focusing on Overlay-Network-1):

- A control-VPLS service is defined in PE-69 and PE-72. For instance, VPLS 1.
    - ç This service is configured with a BGP-MH site in both PEs.
    - ç An oper-group **control-vpls-1** is created and associated to the pw-template-binding 1 in VPLS 1.
- Data VPLS services are defined in both PEs. For instance: VPLS 2, VPLS 3,... VPLS 999.
    - ç In all these services, the pw-template-binding is configured with **monitor-oper-group "control-vpls-1".**
    - ç The status of the spoke SDPs in the data VPLS services depends on the status of the oper-group. If there is a DF switchover in VPLS 1 and VPLS 1 spoke SDPs go down on PE-69, all the spoke SDPs in all the data VPLS services controlled by **control-vpls-1** in PE-69 will go down too. In the same way, the spoke SDPs in PE-72 will come up.
- To allow per-service load balancing a second control-VPLS service with a different BGP-MH site should be configured.
    - ç For instance, VPLS 1 might have PE-69 as the DF and VPLS 1000 might be a second control-VPLS service with PE-72 as the DF.
    - ç Each control-VPLS would control a group of data VPLS services based on the definition and association of a second oper-group.

The following example shows the configuration of VPLS 1 as the control-VPLS and VPLS 2 as a data-VPLS. VPLS 1 controls the VPLS 2 spoke SDP status.

```
*A:PE-69>config>service# info
---------------------------------------------
        customer 1 create
            description "Default customer"
        exit
        pw-template 1 create
        exit
        oper-group "control-vpls-1" create
        exit
        vpls 1 customer 1 create
            description "control-VPLS"
            vxlan vni 1 create
            exit
            bgp
                route-distinguisher 192.0.2.69:1
                vsi-export "vsi-policy-1"
                vsi-import "vsi-policy-1"
                pw-template-binding 1 split-horizon-group "CORE"
                    oper-group "control-vpls-1"
                exit
            exit
            bgp-ad
                vpls-id 64500:1
                no shutdown
            exit
            bgp-evpn
                unknown-mac-route
                vxlan
                    no shutdown
                exit
            exit
            stp
                shutdown
            exit
            site "site-1" create
                site-id 1
                split-horizon-group CORE
                no shutdown
            exit
            no shutdown
        exit
        vpls 2 customer 1 create
            description "data-VPLS"
            vxlan vni 2 create
            exit
            bgp
                route-distinguisher 192.0.2.69:2
                vsi-export "vsi-policy-2"
                vsi-import "vsi-policy-2"
                pw-template-binding 1
                    monitor-oper-group "control-vpls-1"
                exit
            exit
            bgp-ad
```

```
                          vpls-id 64500:2
                          no shutdown
                   exit
                   bgp-evpn
                       unknown-mac-route
                       vxlan
                             no shutdown
                       exit
                   exit
                   stp
                       shutdown
                   exit
                   no shutdown
           exit
       ----------------------------------------------
```

## Use of Proxy-ARP in EVPN-VXLAN Services

EVPN-VXLAN services support proxy-ARP functionality that is enabled by the **proxy-arp [no] shutdown** command. The default value is shutdown. When proxy-arp is enabled:

- MAC and IP addresses contained in the received valid EVPN MAC routes are populated in the proxy-ARP table.

- ARP-request messages received on SAPs and SDP-binds are intercepted and the target IP address is looked up. If the IP address is found, an ARP reply will be issued based on the information found in the proxy-ARP table, otherwise the ARP request would be flooded in the VPLS service (except for the source SAP/SDP binding).

- ARP-reply messages received on SAPs and SDP-bindings are also intercepted and sent to the CPM. These ARP-reply messages are re-injected in the data plane and forwarded based on the FDB information to the destination MAC address. If the destination MAC address is not in the FDB, the ARP-reply message will be flooded in the VPLS service (except for the source SAP/SDP binding).

The following CLI output shows the proxy-ARP configuration in PE-72 and a received valid MAC route that includes the MAC and IP of CE-1. This MAC-IP pair is installed in the proxy-ARP table for VPLS 1.

```
*A:PE-72>config>service>vpls# info
----------------------------------------------
           vxlan vni 1 create
           exit
           bgp
               route-distinguisher 192.0.2.72:1
               vsi-export "vsi-policy-1"
               vsi-import "vsi-policy-1"
               pw-template-binding 1 split-horizon-group "CORE"
                   oper-group "control-vpls-1"
               exit
           exit
```

```
                bgp-ad
                    vpls-id 64500:1
                    no shutdown
                exit
                bgp-evpn
                    unknown-mac-route
                    vxlan
                        no shutdown
                    exit
                exit
                proxy-arp
                    no shutdown
                exit
...

27 2014/07/17 23:15:54.85 UTC MINOR: DEBUG #2001 Base Peer 1: 192.0.2.71
"Peer 1: 192.0.2.71: UPDATE
Peer 1: 192.0.2.71 - Received BGP UPDATE:
    Withdrawn Length = 0
    Total Path Attr Length = 92
    Flag: 0x90 Type: 14 Len: 48 Multiprotocol Reachable NLRI:
        Address Family EVPN
        NextHop len 4 NextHop 192.0.2.71
        Type: EVPN-MAC Len: 37 RD: 192.0.2.71:1 ESI: 0:0:0:0:0:0:0:0:0:0, tag: 1
, mac len: 48 mac: 00:ca:fe:ca:fe:01, IP len: 4, IP: 172.16.0.1, label: 0
    Flag: 0x40 Type: 1 Len: 1 Origin: 0
    Flag: 0x40 Type: 2 Len: 0 AS Path:
    Flag: 0x80 Type: 4 Len: 4 MED: 0
    Flag: 0x40 Type: 5 Len: 4 Local Preference: 100
    Flag: 0xc0 Type: 16 Len: 16 Extended Community:
        target:64500:12
        bgp-tunnel-encap:VXLAN
"

*A:PE-72# show service id 1 proxy-arp
-------------------------------------------------------------------------------
VPLS Proxy Arp Table
-------------------------------------------------------------------------------
IP Address          Mac Address
-------------------------------------------------------------------------------
172.16.0.1          00:ca:fe:ca:fe:01
-------------------------------------------------------------------------------
Number of entries : 1
-------------------------------------------------------------------------------
===============================================================================
*A:PE-72#
```

Note that in the tested release, the 7x50 does not include a host IP address in any EVPN MAC advertisement for a MAC learned on a SAP or SDP-bind. Host IP addresses are only included in the EVPN MAC advertisements corresponding to R-VPLS IP interfaces. When deployed as DC GW in a Nuage architecture, the Nuage Networks VSC (Virtual Services Controller) or VSG (Virtual Services Gateway) will send virtual machine and host MAC/IP pairs in EVPN MAC routes. Please refer to the Alcatel-Lucent Nuage documentation for more information about the Nuage DC architecture. The 7x50 DC GW will populate the proxy-ARP tables with those MAC/IP pairs. In the CLI excerpt above, assume that PE-71 is replaced by a Nuage VSC that sends the

pair <172.16.0.1, 00:ca:fe:ca:fe:01> in an EVPN MAC route. PE-72 receives the advertisement and adds the entry to its proxy-ARP table for VPLS 1.

# MAC Mobility, MAC Duplication and MAC Protection in EVPN

MAC mobility, duplication and protection are fully supported as specified in draft-ietf-l2vpn-evpn. Figure 46 illustrates the concept of mobility (Virtual Machine VM-1 moves from PE-71 to PE-72).
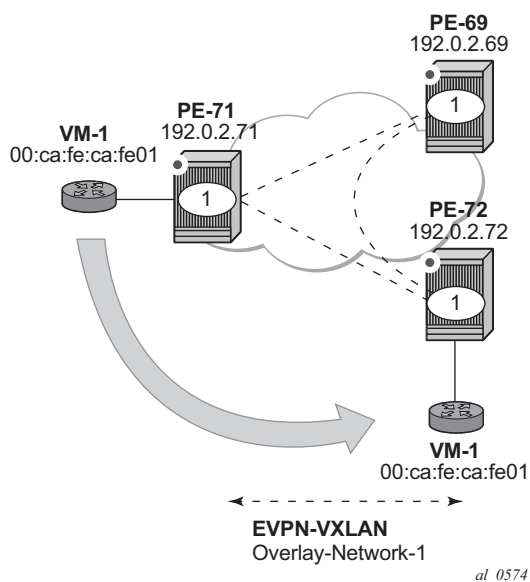


**Figure 46: EVPN MAC Mobility**

MAC mobility is handled in EVPN by the use of sequence numbers in the MAC routes. When 00:ca:fe:ca:fe:01 moves from PE-71 to PE-72, SR OS will gracefully handle it in this way:

- 00:ca:fe:ca:fe:01 moves to PE-72 SAP 1/1/1:1

- PE-72 advertises 00:ca:fe:ca:fe:01 using a higher sequence number (the first time a MAC is advertised, EVPN uses sequence number 0).

- PE-69 at this point has two valid MAC routes for 00:ca:fe:ca:fe:01. It picks up the one coming from PE-72 since the sequence number is higher.

- PE-71 receives the MAC route, and since the sequence number is higher than the one for its own route, it updates the FDB and withdraws its own MAC route.

However, if MAC 00:ca:fe:ca:fe:01 is constantly learned on the PE-71 and PE-72 SAPs, the process above causes an endless exchange of MAC route advertisements and withdraws that has a negative impact on all the PEs in the EVPN network. This issue is known as "MAC duplication" and is originated by a loop at the access or a duplicated MAC address in two hosts of the same service. SR OS solves this issue through the use of the mac-duplication detection feature. Note that mac-duplication is always enabled with the following default settings:

```
*A:PE-71>config>service>vpls>bgp-evpn# info detail
---------------------------------------------
              no unknown-mac-route
              mac-advertisement
              no ip-route-advertisement
              mac-duplication
                  detect num-moves 5 window 3
                  retry 9
              exit
              vxlan
                  no shutdown
              exit
---------------------------------------------
```

Where:

- **num-moves** — Identifies the number of MAC moves in a VPLS service. The counter is incremented when a given MAC is locally relearned in the FDB or flushed from the FDB due to the reception of a better remote EVPN route for that MAC. When the threshold is reached for a given MAC, this MAC is put in hold-down state (this 'hold-down' state is described below). Range: <3..10>. Default value: 5.

- **window** — Identifies the timer within which a MAC is considered duplicate if it reaches the configured num-moves. Range: <1..15> minutes. Default value: 3 minutes.

- **Retry** — The timer after which the MAC in hold-down state is automatically flushed and the mac-duplication process starts again. This value is expected to be equal to two times or more than the window. If no retry is configured, this implies that, once mac-duplication is detected, MAC updates for that MAC will be held down until the user intervenes or a network event (that flushes the MAC) occurs. Range: <2..60> minutes. Default value: 9 minutes.

When a MAC is considered a duplicate or in the 'hold-down' state, no further BGP advertisements are issued for this MAC and an alarm is triggered (by the first MAC in hold-down state). The following CLI output shows how PE-72 detects that MAC 00:ca:fe:ca:fe:01 is a duplicate (after reaching the **num-moves** in **window**) and the corresponding alarm. The **show service id bgp-evpn** command shows the mac-duplication settings and the list of duplicate MACs on hold-down.

```
41 2014/07/17 23:50:45.83 UTC MINOR: DEBUG #2001 Base Peer 1: 192.0.2.71
"Peer 1: 192.0.2.71: UPDATE
Peer 1: 192.0.2.71 - Send BGP UPDATE:
    Withdrawn Length = 0
    Total Path Attr Length = 96
    Flag: 0x90 Type: 14 Len: 44 Multiprotocol Reachable NLRI:
        Address Family EVPN
        NextHop len 4 NextHop 192.0.2.72
        Type: EVPN-MAC Len: 33 RD: 192.0.2.72:1 ESI: 0:0:0:0:0:0:0:0:0:0, tag: 1
, mac len: 48 mac: 00:ca:fe:ca:fe:01, IP len: 0, IP: NULL, label: 0
    Flag: 0x40 Type: 1 Len: 1 Origin: 0
    Flag: 0x40 Type: 2 Len: 0 AS Path:
    Flag: 0x80 Type: 4 Len: 4 MED: 0
    Flag: 0x40 Type: 5 Len: 4 Local Preference: 100
    Flag: 0xc0 Type: 16 Len: 24 Extended Community:
        target:64500:12
        bgp-tunnel-encap:VXLAN
        mac-mobility:Seq:4
"


2 2014/07/17 23:50:46.62 UTC MINOR: SVCMGR #2331 Base
"VPLS Service 1 has MAC(s) detected as duplicates by EVPN mac-duplication detect
ion."

*A:PE-72# show service id 1 bgp-evpn

===============================================================================
BGP EVPN Table
===============================================================================
MAC Advertisement  : Enabled        Unknown MAC Route  : Enabled
VXLAN Admin Status : Enabled        Creation Origin    : manual
MAC Dup Detn Moves : 5              MAC Dup Detn Window: 3
MAC Dup Detn Retry : 9              Number of Dup MACs : 1
IP Route Advertise*: Disabled



-------------------------------------------------------------------------------
Detected Duplicate MAC Addresses          Time Detected
-------------------------------------------------------------------------------
00:ca:fe:ca:fe:01                         07/17/2014 23:50:47
-------------------------------------------------------------------------------
===============================================================================
* indicates that the corresponding row element may have been truncated.
```

The 7x50 stops sending and processing any BGP MAC Advertisement routes for that MAC address until:

- The MAC is flushed due to a local event (SAP/SDP-binding associated to the MAC fails) or the reception of a remote withdraw for the MAC (due to a MAC flush at the remote 7x50) or

- The **retry <in_minutes>** timer expires, which flushes the MAC and restart the process.

When the last duplicate MAC address is removed from the duplicate list, the system will show the following message:

```
*A:PE-72#
3 2014/07/17 23:56:21.71 UTC MINOR: SVCMGR #2332 Base
"VPLS Service 1 no longer has MAC(s) detected as duplicates by EVPN mac-duplicat
ion detection."
```

EVPN also provides a mechanism to protect certain MACs that do not move for which connectivity must be guaranteed. These addresses must be protected in case there is an attempt to dynamically learn them in a different place in the EVPN-VXLAN VPLS service (on the same or different PE).

The protected MACs are configured in SR OS as conditional static MACs. A conditional static MAC defined in an EVPN-VXLAN VPLS service is advertised by BGP-EVPN as a static address. An example of the configuration of a conditional static MAC is shown below:

```
*A:PE-71>config>service>vpls# info
----------------------------------------------
            vxlan vni 1 create
            exit
            bgp
                route-distinguisher 192.0.2.71:1
                route-target export target:64500:12 import target:64500:12
            exit
            bgp-evpn
                vxlan
                    no shutdown
                exit
            exit
            proxy-arp
                no shutdown
            exit
            sap 1/1/1:1 create
            exit
            static-mac
                mac 00:ca:fe:ca:fe:05 create sap 1/1/1:1 monitor fwd-status
            exit
            no shutdown
----------------------------------------------
```

The protected MACs advertised in EVPN are shown in the receiving BGP RIB as Static (MAC mobility extended community with Sequence 0 and sticky bit set) and **EvpnS** (Evpn Static) in the FDB. The advertising PE shows the protected MAC as **CStatic** (Conditional Static) in the FDB:

```
# advertising PE
 *A:PE-71>config>service>vpls# show service id 1 fdb detail
 ===============================================================================
 Forwarding Database, Service 1
 ===============================================================================
 ServId    MAC                 Source-Identifier       Type     Last Change
                                                        Age
 -------------------------------------------------------------------------------
 1         00:ca:fe:ca:fe:05 sap:1/1/1:1               CStatic  07/18/14 00:29:34
 -------------------------------------------------------------------------------
 No. of MAC Entries: 1
 -------------------------------------------------------------------------------
 Legend:  L=Learned O=Oam P=Protected-MAC C=Conditional S=Static
 ===============================================================================


# receiving PE

*A:PE-72# show service id 1 fdb detail
 ===============================================================================
 Forwarding Database, Service 1
 ===============================================================================
 ServId    MAC                 Source-Identifier       Type     Last Change
                                                        Age
 -------------------------------------------------------------------------------
 1         00:ca:fe:ca:fe:05 vxlan:                    EvpnS    07/18/14 00:29:35
                             192.0.2.71:1
 -------------------------------------------------------------------------------
 No. of MAC Entries: 1
 -------------------------------------------------------------------------------
 Legend:  L=Learned O=Oam P=Protected-MAC C=Conditional S=Static
 ===============================================================================


*A:PE-72# show router bgp routes evpn mac mac-address 00:ca:fe:ca:fe:05 hunt
 ===============================================================================
  BGP Router ID:192.0.2.72      AS:64500      Local AS:64500
 ===============================================================================
  Legend -
  Status codes  : u - used, s - suppressed, h - history, d - decayed, * - valid
  Origin codes  : i - IGP, e - EGP, ? - incomplete, > - best, b - backup

 ===============================================================================
 BGP EVPN Mac Routes
 ===============================================================================
 -------------------------------------------------------------------------------
 RIB In Entries
 -------------------------------------------------------------------------------
 Network      : N/A
 Nexthop      : 192.0.2.71
 From         : 192.0.2.71
 Res. Nexthop : N/A
 Local Pref.  : 100                    Interface Name : NotAvailable
 Aggregator AS : None                  Aggregator     : None
 Atomic Aggr.  : Not Atomic            MED            : 0
 AIGP Metric   : None
 Connector     : None
 Community     : target:64500:12 bgp-tunnel-encap:VXLAN
                   mac-mobility:Seq:0/Static
 Cluster       : No Cluster Members
 Originator Id : None                  Peer Router Id : 192.0.2.71
```

```
Flags          : Used  Valid  Best  IGP
Route Source   : Internal
AS-Path        : No As-Path
EVPN type      : MAC
ESI            : 0:0:0:0:0:0:0:0:0   Tag            : 1
IP Address     : N/A                 Route Dist.    : 192.0.2.71:1
Mac Address    : 00:ca:fe:ca:fe:05
MPLS Label1    : 0                   MPLS Label2    : N/A
Route Tag      : 0
Neighbor-AS    : N/A
Orig Validation: N/A
Source Class   : 0                   Dest Class     : 0


-------------------------------------------------------------------------------
RIB Out Entries
-------------------------------------------------------------------------------
-------------------------------------------------------------------------------
Routes : 1
===============================================================================
*A:PE-72#
```

The following procedures are supported in order to protect the configured static MAC addresses:

- All the SAP/SDP-bindings are internally configured as MAC protect restrict-protected-src as soon as bgp-evpn is enabled in the VPLS service.

- Local static MACs or remote EVPN Static MACs are considered as protected.

- If a frame with a source MAC address matching one of the protected MACs is received on a different SAP/SDP-binding than the owner of the protected MAC, the frame is discarded and an alarm triggered. Note that this MAC protection is not performed for frames received on VXLAN bindings.

- The same throttled alarm mechanism used in MAC protect for restrict-protected-src with discard-frame is used here: the offending frames are captured to a list to be polled by the CPM every ~10min.

In this example, PE-72 has 00:ca:fe:ca:fe:05 in its FDB as EvpnS. If SAP 1/1/1:1 receives a frame with source MAC address 00:ca:fe:ca:fe:05, the frame is discarded and an alarm triggered:

```
*A:PE-72#
4 2014/07/18 00:33:49.05 UTC MINOR: SVCMGR #2208 Base Slot 1
"Protected MAC 00:ca:fe:ca:fe:05 received on SAP 1/1/1:1 in service 1. "
```

# Debug and Show Commands

In addition to the previously mentioned **show service id vxlan, show service id bgp-evpn and show service id fdb detail** commands, the following commands provide valuable information when troubleshooting an EVPN-VXLAN VPLS service.

The **show router bgp routes evpn** command supports filtering by route type as well as many other route fields.

```
*A:PE-72# show router bgp routes evpn
  - evpn <evpn-type>

      inclusive-mcast - Display BGP EVPN Inclusive-Mcast Routes
      ip-prefix       - Display BGP EVPN IP-Prefix Routes
      mac             - Display BGP EVPN Mac Routes

*A:PE-72# show router bgp routes evpn mac
{hunt|detail}
 hunt    detail
rd <rd>
next-hop <ip-address>
mac-address <mac-address>
community <comm-id>
tag <vni-id>

*A:PE-72# show router bgp routes evpn mac tag 1
===============================================================================
 BGP Router ID:192.0.2.72       AS:64500       Local AS:64500
===============================================================================
 Legend -
 Status codes  : u - used, s - suppressed, h - history, d - decayed, * - valid
 Origin codes  : i - IGP, e - EGP, ? - incomplete, > - best, b - backup


===============================================================================
BGP EVPN Mac Routes
===============================================================================
Flag  Route Dist.         ESI                   Tag        MacAddr
                          NextHop                          IpAddr
                                                           Mac Mobility
-------------------------------------------------------------------------------
u*>i  192.0.2.69:1        0:0:0:0:0:0:0:0:0:0    1          00:00:00:00:00:00
                          192.0.2.69                       N/A
                                                           Seq:0

u*>i  192.0.2.71:1        0:0:0:0:0:0:0:0:0:0    1          00:ca:fe:ca:fe:05
                          192.0.2.71                       N/A
                                                           Static


-------------------------------------------------------------------------------
Routes : 2
===============================================================================
```

The **tools dump service id vxlan** displays the number of times a service could not add a VXLAN binding or <VTEP, Egress VNI> due to the following limits:

- The per System VTEP limit has been reached
- The per System (egress VTEP, egress VNI) limit has been reached
- The per Service (egress VTEP, egress VNI) limit has been reached
- The per System Bind limit: Total bind limit or VXLAN bind limit has been reached.

**Tools dump service vxlan usage** displays the consumed VXLAN resources in the system, whereas **tools dump service vxlan dup-vtep-egrvni** displays the (egress VTEP, egress VNI) bindings that have been detected as duplicate attempts, in other words, an attempt to add the same binding to more than one service:

```
*A:PE-72# tools dump service id 1 vxlan

VTEP, Egress VNI Failure statistics at 001 01:06:55.950:

statistics last cleared at 000 00:00:00.000:

Failures: None

*A:PE-72# tools dump service vxlan usage

VXLAN usage statistics at 001 01:07:59.790:

VTEP                      :      2/8191
VTEP, Egress VNI          :      4/131071
Sdp Bind + VTEP, Egress VNI :    10/196607
RVPLS Egress VNI          :      0/40959

*A:PE-72# tools dump service vxlan dup-vtep-egrvni

Duplicate VTEP, Egress VNI usage attempts at 001 01:08:04.080:

1. 192.0.2.71:100
```

# Conclusion

SR OS supports the EVPN control plane for VXLAN tunnels terminated in VPLS services. VXLAN is an overlay IP tunneling mechanism that is being used in data center, data center interconnect and other applications. EVPN is a scalable and flexible control plane that provides control over the MACs being learned and advertised, as well as other mechanisms to optimize Layer 2 services such as proxy-ARP, MAC mobility, MAC duplication detection and MAC protection. SR OS provides a resilient and scalable EVPN-VXLAN solution for Layer 2 services, including interoperability to existing VPLS networks. This example showed all of those functions and how they are configured and operated.