



# Network Fundamentals for Cloud

**BITS Pilani**  
Pilani Campus

Nishit Narang  
WILPD-CSIS



# **CC ZG503: Network Fundamentals for Cloud**

## **Lecture No. 11: DCN Control Plane Protocols**



# RECAP: L2MP



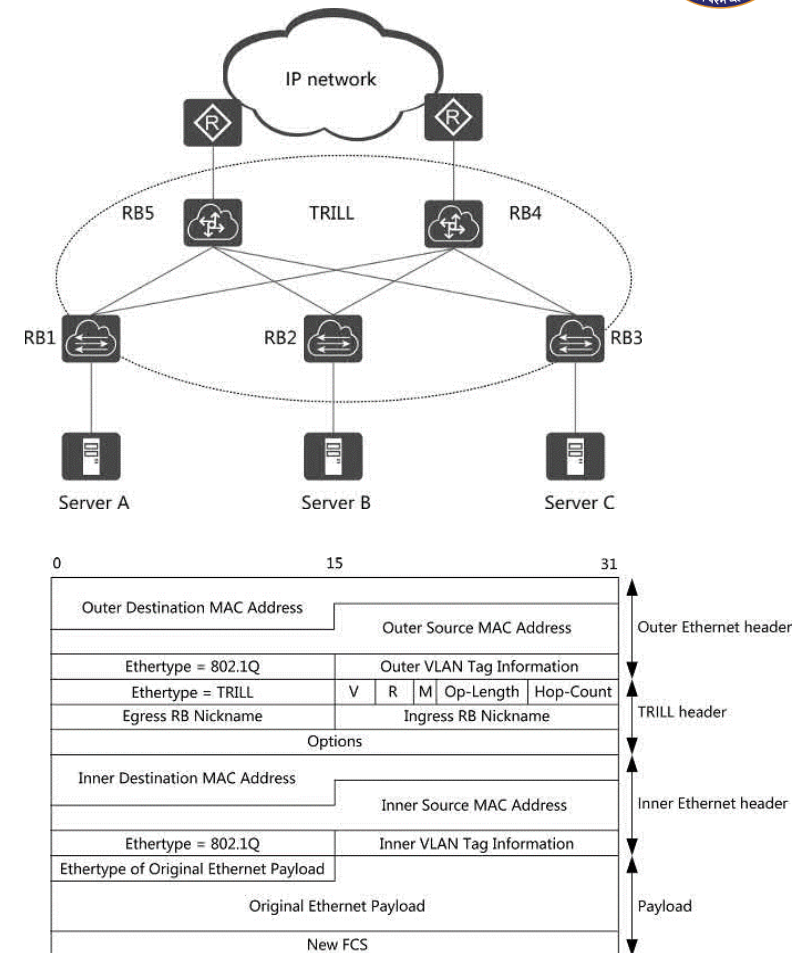
- Layer 2 Multi-Pathing (L2MP) Technologies
  - L2MP technologies attempt to address the below twin challenges of xSTP and Virtual Chassis technologies:
    - They cannot support large DCNs with massive amounts of data.
    - Their link utilization is low.
  - It is recommended that link- state routing protocols widely used on Layer 3 networks be employed
  - These protocols not only support a large number of devices, but they are also loop-free and have high link utilization
    - Example: OSPF and IS-IS → they support ECMP load balancing and use the Shortest Path First (SPF) algorithm
  - The basic principle of L2MP technologies is to introduce mechanisms of routing technologies used on Layer 3 networks to Layer 2 networks



# RECAP: TRILL



- TRILL: Basic Concepts
  - Stands for **Transparent Interconnection of Lots of Links**
  - Is implemented by devices called TRILL switches
  - TRILL combines techniques from bridging and routing, and is the application of link-state routing to L2 networks
  - To apply link-state routing protocols to Ethernet networks, a frame header needs to be added to an Ethernet header for the addressing of the link-state routing protocols
  - TRILL uses MAC in TRILL in MAC encapsulation
    - I.e. In addition to the original Ethernet header, a TRILL header that provides an addressing identifier and an outer Ethernet header used to forward a TRILL packet on an Ethernet network are added



Source: Lei Zhang, Le Chen. Cloud Data Center Network Architectures and Technologies, CRC Press 2021

# RECAP: L2MP Disadvantages

- Disadvantages of L2MP Technologies
  - Limited number of tenants: Similar to xSTP, TRILL uses VLAN IDs to identify tenants. Because the VLAN ID field has only 12 bits, a TRILL network supports only about 4000 tenants.
    - a solution to the tenant problem was considered at the beginning of TRILL design. A field was reserved in the TRILL header for tenant identification, but the problem has not yet been resolved because the protocol has not been continuously evolved.
  - Increased deployment costs: L2MP technologies introduce new forwarding identifiers or add new forwarding processes, which inevitably requires the upgrade of forwarding chips.
  - Mechanism-related challenges: The Operations, Administration, and Maintenance (OAM) mechanism and multicast mechanism of TRILL have not been defined into formal standards, restricting further protocol evolution

Introduction of NVO3 technologies (like VXLAN) eventually led to the downfall of L2MP technology

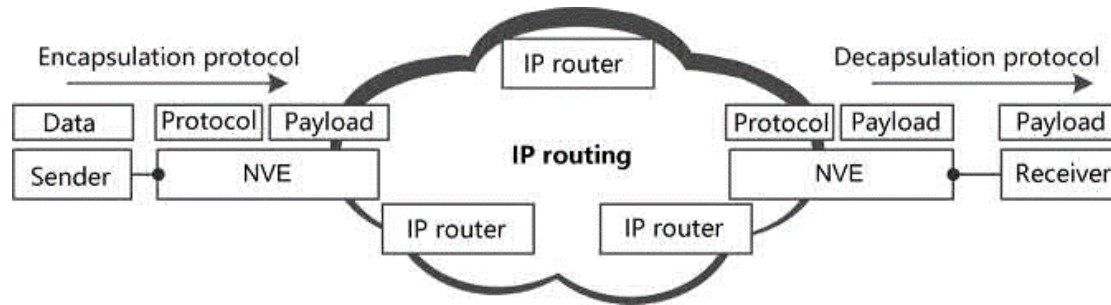


# RECAP: DCN continues to evolve....

- As DCNs scale and traffic patterns change to more east-west traffic....
- ...DCN topologies evolve towards Clos Topology and its variants....
- ....DCN protocols evolve towards use of L3 technologies (IP & associated networking/control plane protocols)

# RECAP: NVO3

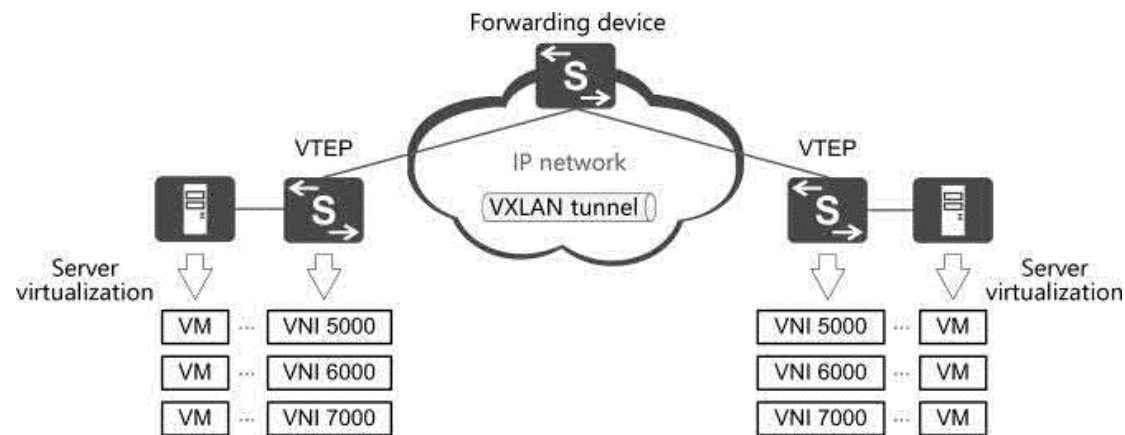
- Network Virtualization Overlays (NVO3) Technologies



- The sender in the figure is an endpoint, which may be a VM or physical server
- An NVE may be a physical switch or a virtual switch on a hypervisor
- The sender can be connected to an NVE directly or through a switching network
- NVEs are connected through a tunnel

# RECAP: VXLAN

- VXLAN is an NVO3 technology that enables Layer 2 forwarding over a Layer 3 network by using L2 over L4 (MAC-in-UDP) encapsulation
- Defined by the IETF, it allows VMs to migrate over a large Layer 2 network and isolates tenants in a DC



Source: Lei Zhang, Le Chen. Cloud Data Center Network Architectures and Technologies, CRC Press 2021



# RECAP: VXLAN Control Plane

- Two widely adopted control planes are used with VXLAN:
  - the VXLAN Flood and Learn Multicast-Based Control Plane and
  - the VXLAN MPBGP EVPN Control Plane.

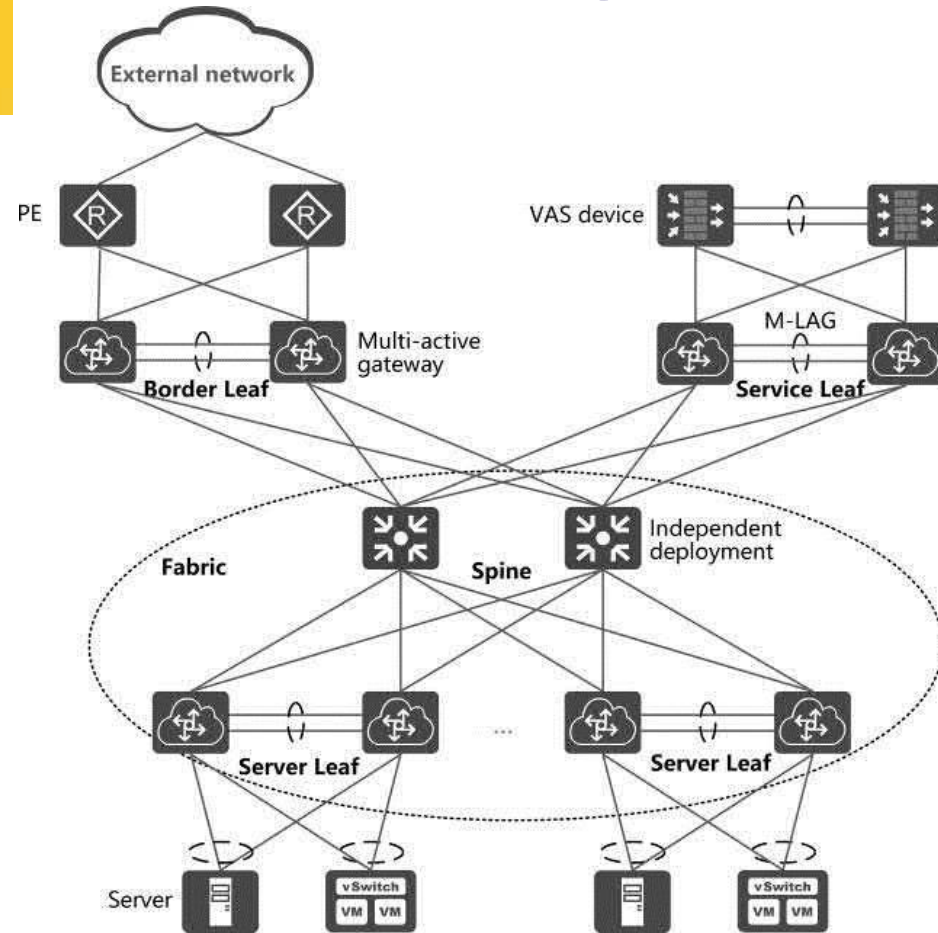
More on this later.....



# Control Plane Protocols in Data Center Networks

---

# Constructing DCN Underlay Network



(a) Role	Function
(b) Fabric	Network failure domain that is managed by an SDN controller. It contains one or more spine-leaf architectures
Spine	Core node on a VXLAN fabric network. It provides high-speed IP forwarding and connects to leaf nodes through high-speed interfaces
Leaf	Access node on a VXLAN fabric network. It connects various network devices to the VXLAN fabric network
Service leaf	Functional node that connects VAS devices, such as firewalls and LBs, to a VXLAN fabric network
Server leaf	Functional node that connects virtual and physical servers to a VXLAN fabric network
Border leaf	Functional node that connects to routers or transmission devices outside a DC to forward traffic from an external network to a VXLAN fabric network in a DC

# Constructing DCN Underlay Network (contd..)



- Leaf nodes and spine nodes are connected through Layer 3 routed interfaces
- They communicate at Layer 3 by configuring a dynamic routing protocol. OSPF or BGP is recommended.
- ECMP is recommended for implementing load balancing and link backup
  - In this case, leaf nodes forward data traffic to spine nodes through multiple ECMP paths, guaranteeing reliability while ensuring that network bandwidth improves.

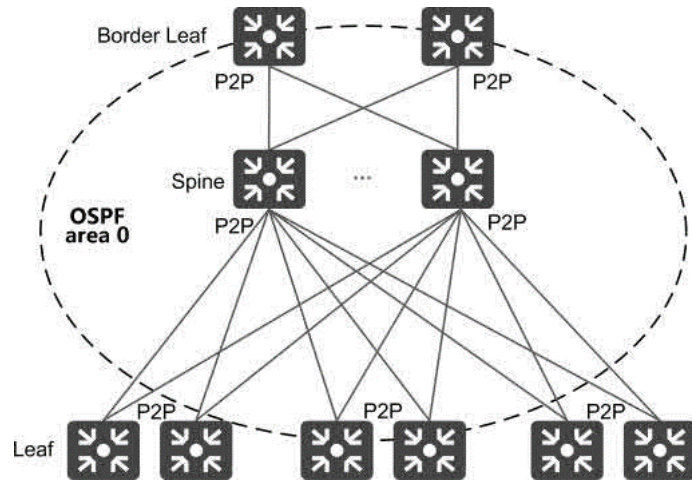


Source: Lei Zhang, Le Chen. Cloud Data Center Network Architectures and Technologies, CRC Press 2021

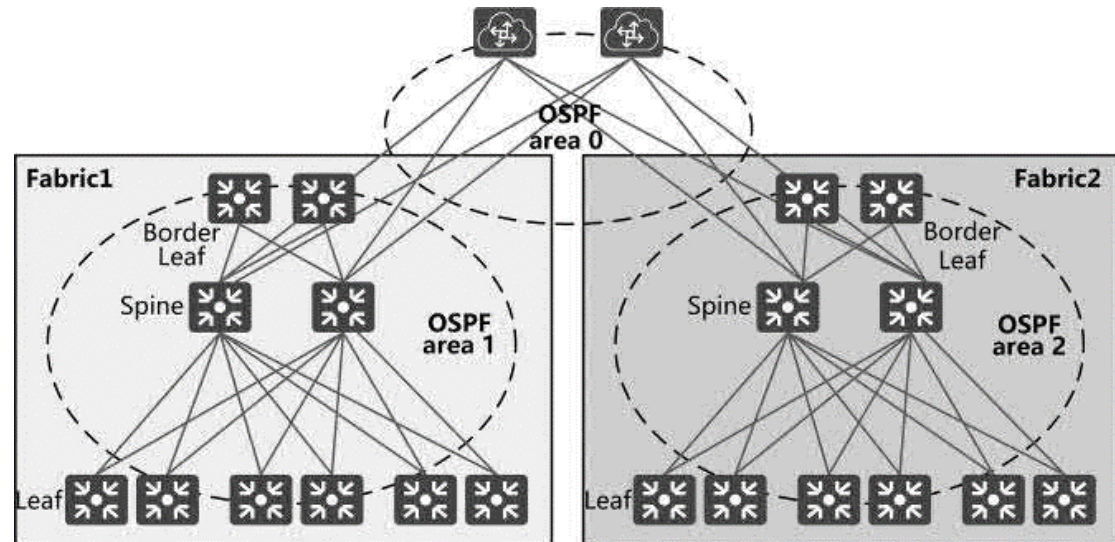
# Routing Protocol Selection

- In most cases, either OSPF or EBGp can be used on an underlay network
- OSPF is preferred for most small to medium sized networks
  - When the number of leaf nodes is less than 100, OSPF is recommended on the underlay network
- If the scale of the network is large, EBGp is recommended as the underlay network needs to be partitioned into areas, and flexible control of BGP is required
  - As the size of the network increases, and/or the number of prefixes to be advertised increases, BGP becomes the go-to routing protocol
  - For advertising fewer than 32,000 prefixes, OSPF is a fine choice as a routing protocol. For larger prefixes, BGP should be preferred.

# OSPF Deployment Scenarios



Recommended OSPF planning for a single fabric network

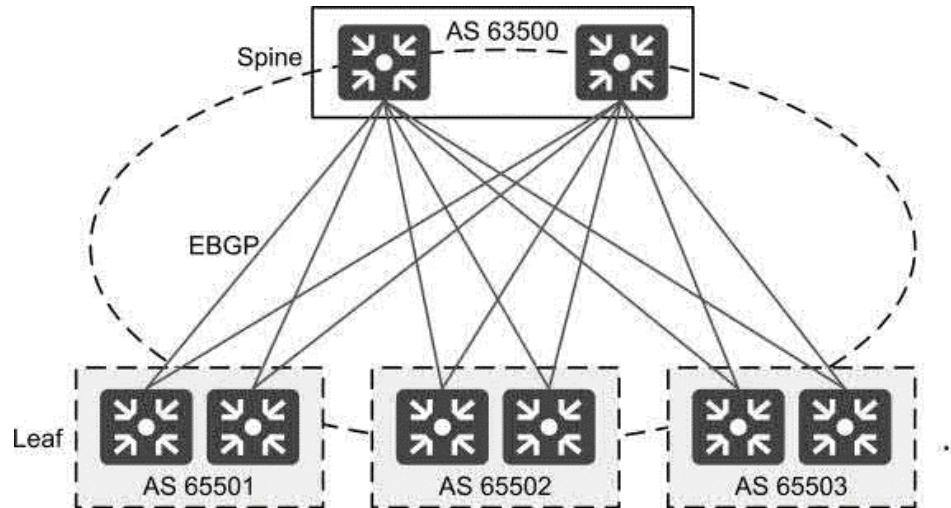


Recommended OSPF planning for multi-fabric deployment (single VXLAN domain).

If multiple fabric networks form two VXLAN domains on the overlay network (that is, two DCNs are managed through two sets of management interfaces, but they need to be interconnected), then:

- OSPF be deployed on each fabric network, and
- interconnected devices between fabric networks can exchange routes through BGP

# EBGP Deployment in DCN Underlay Networks



Recommended EBGP planning for a single fabric network.

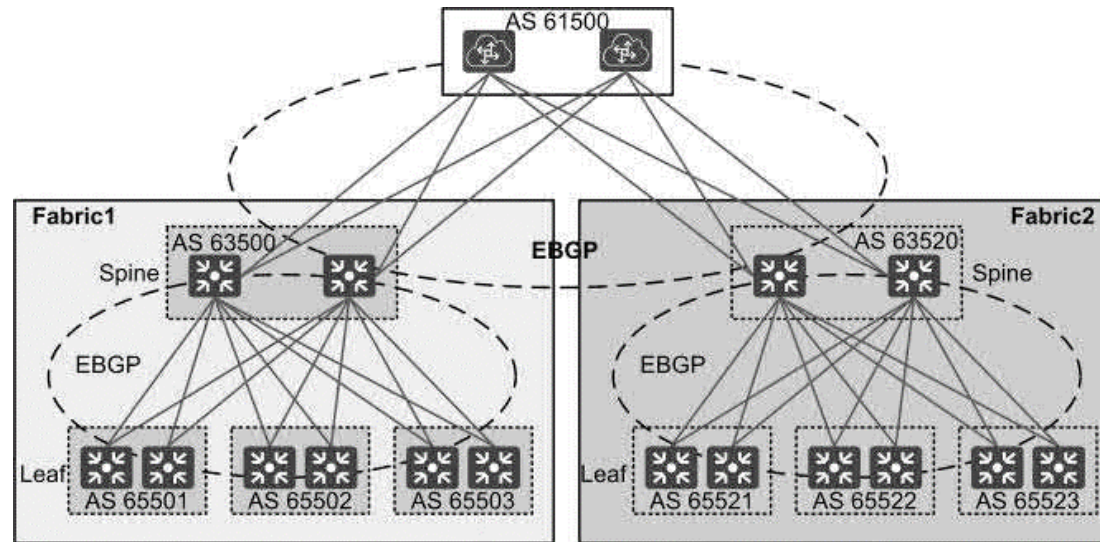
On a single fabric network:

- spine nodes are added to the same autonomous system (AS),
- each group of leaf nodes is added to an AS
- EBGP peer relationships are established between leaf nodes and all spine nodes.





# EBGP Deployment in DCN Underlay Networks



Recommended EBGP planning for multiple fabric networks

Similar concept as Single Fabric Network



Source: Lei Zhang, Le Chen. Cloud Data Center Network Architectures and Technologies, CRC Press 2021



# Back to.....VXLAN Control Plane

- Two widely adopted control planes are used with VXLAN:
  - the VXLAN Flood and Learn Multicast-Based Control Plane and
  - the VXLAN MPBGP EVPN Control Plane.

# VXLAN Flood and Learn Multicast-Based Control Plane

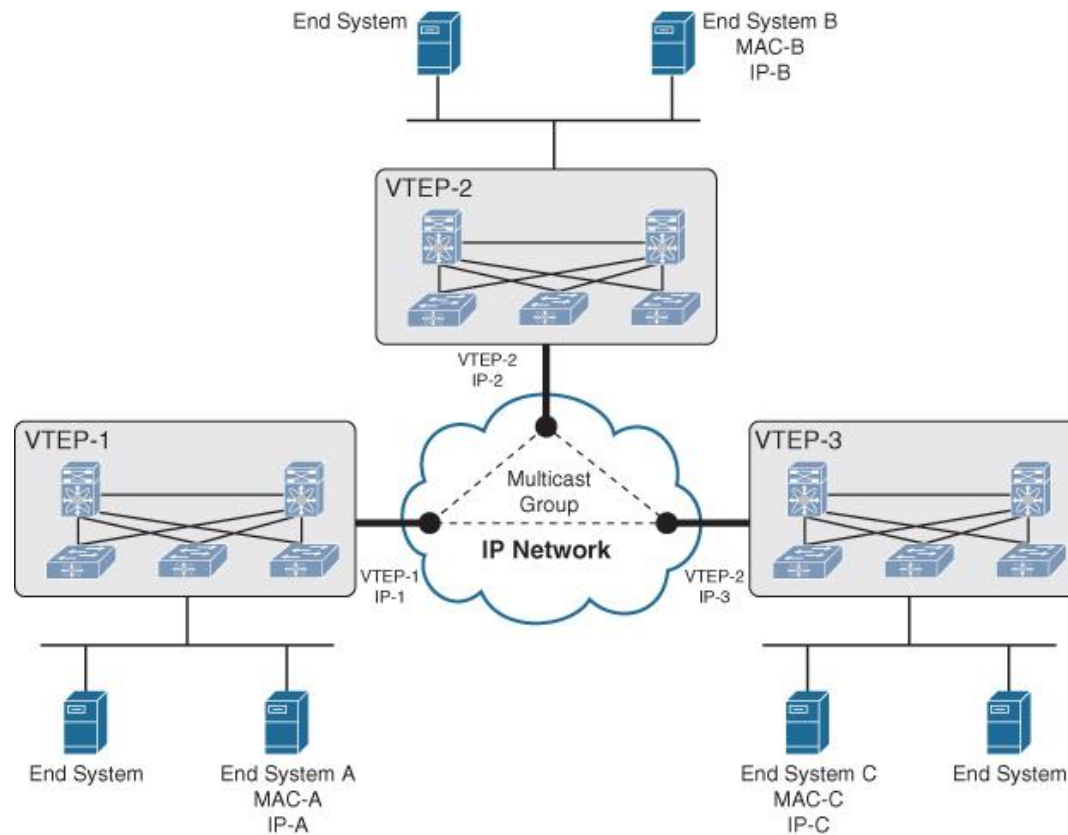


- Switches utilize existing Layer 2 flooding mechanisms and dynamic MAC address learning to
  - Transport broadcast, unknown unicast, and multicast (BUM) traffic
  - Discover remote VTEPs
  - Learn remote-host MAC addresses and MAC-to-VTEP mappings for each VXLAN segment
- IP multicast is used to reduce the flooding scope of the set of hosts that are participating in the VXLAN segment
  - Each VXLAN segment, or VNID, is mapped to an IP multicast group in the transport IP network
  - Each VTEP device is independently configured and joins this multicast group as an IP host through the Internet Group Management Protocol (IGMP).



Source: [Virtual Extensible LAN \(VXLAN\) Overview > Implementing Data Center Overlay Protocols | Cisco Press](#)

# VXLAN Flood and Learn Multicast-Based Control Plane (Contd.)



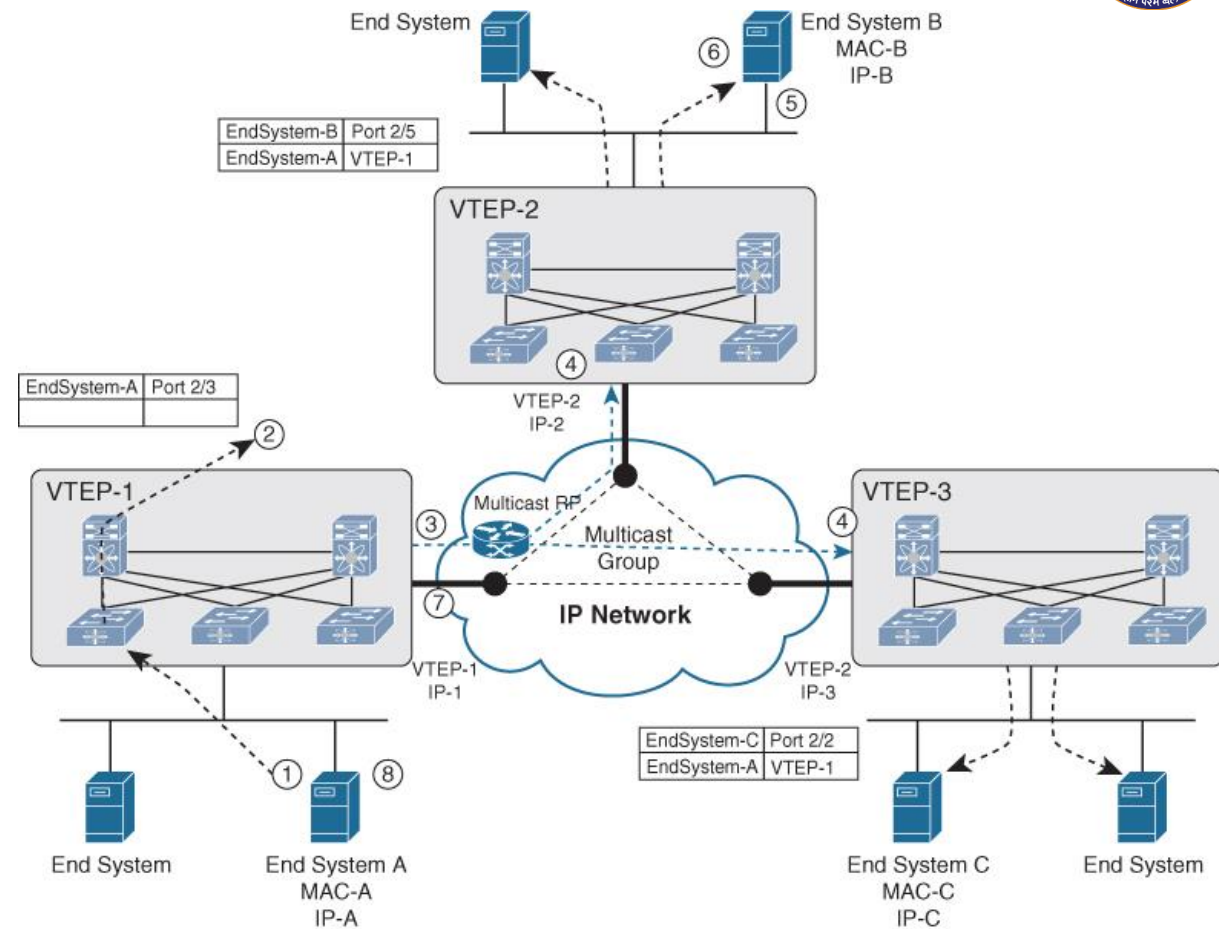
Source: [Virtual Extensible LAN \(VXLAN\) Overview > Implementing Data Center Overlay Protocols | Cisco Press](#)

# VXLAN Flood and Learn Multicast-Based Control Plane (Contd.)



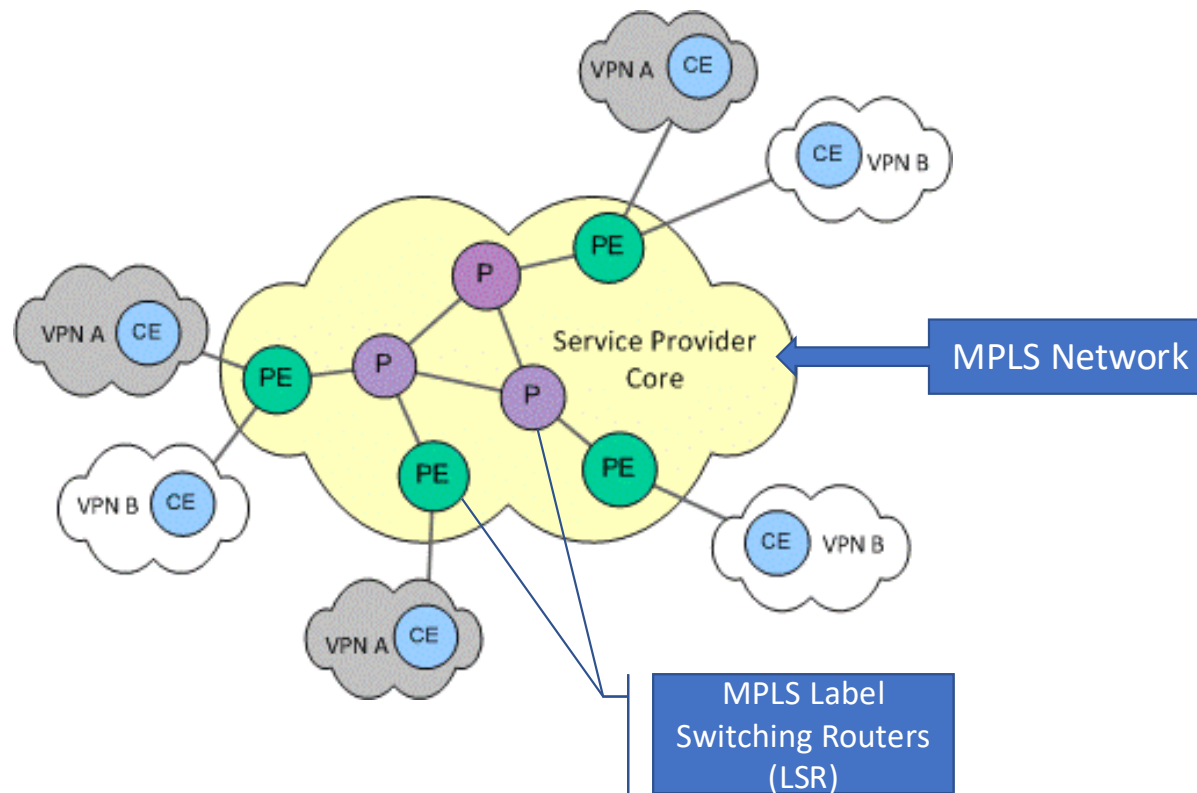
As an example, if End System A wants to talk to End System B, it does the following:

1. End System A generates an ARP request trying to discover the End System B MAC address.
2. When the ARP request arrives at SW1, it will look up its local table, and if an entry is not found, it will encapsulate the ARP request over VXLAN and send it over the multicast group configured for the specific VNI.
3. The multicast RP receives the packet, and it forwards a copy to every VTEP that has joined the multicast group.
4. Each VTEP receives and de-encapsulates the VXLAN packet and learns the System A MAC address pointing to the remote VTEP address.
5. Each VTEP forwards the ARP request to its local destinations.
6. End System B generates the ARP reply. When SW2 VTEP2 receives it, it looks up its local table and finds an entry with the information that traffic destined to End System A must be sent to VTEP1 address. VTEP2 encapsulates the ARP reply with a VXLAN header and unicasts it to VTEP1.
7. VTEP1 receives and de-encapsulates the packet and delivers it to End System A.
8. When the MAC address information is learned, additional packets are fed to the corresponding VTEP address.



Source: [Virtual Extensible LAN \(VXLAN\) Overview > Implementing Data Center Overlay Protocols | Cisco Press](#)

# VXLAN MPBGP EVPN Control Plane: Reference to Layer 3 MPLS VPNs



The various routers within the VPN communicate using BGP, an IP routing protocol that defines how routes can be distributed. BGP transports information about CE routers only to members of the same VPN, ensuring security.

Source: [https://docs.oracle.com/cd/E35292\\_01/doc.72/e47715/con\\_vpn.htm#autold0](https://docs.oracle.com/cd/E35292_01/doc.72/e47715/con_vpn.htm#autold0)

# VXLAN MPBGP EVPN Control Plane

- The Flood-and-Learn method creates significant flooding traffic resulting in network expansion difficulties
- As a solution to these problems, EVPN is introduced on the VXLAN control plane
  - EVPN = Ethernet Virtual Private Network OR Ethernet VPN
  - By referring to the BGP/MPLS IP VPN mechanism, EVPN defines several types of BGP EVPN routes by extending BGP.
  - The PE node role described in BGP MPLS EVPN is equivalent to the VTEP/network virtualization edge (NVE) device
  - It advertises BGP routes on the network to implement automatic VTEP discovery and host address learning.

# EVPN Advantages

- Using EVPN on the control plane offers the following advantages:
  - VTEPs can be automatically discovered and VXLAN tunnels can be automatically established, overall simplifying network deployment and expansion.
  - EVPN can advertise Layer 2 MAC addresses and Layer 3 routing information simultaneously.
  - Flooding traffic is reduced on the network.

# MP-BGP



- Traditional BGP-4 uses Update packets to exchange routing information between peers
  - An Update packet can advertise a type of reachable routes with the same path attributes, placed in Network Layer Reachability Information (NLRI) fields.
- BGP-4 can manage only IPv4 unicast routing information
- As a solution, Multiprotocol Extensions for BGP (MP-BGP) was developed as a means to support multiple network layer protocols, including IPv6 and multicast.
  - MP-BGP extends NLRI fields based on BGP-4.
  - After extension, the description of the address family is added to the NLRI fields to differentiate network layer protocols.
  - These include the IPv6 unicast address family and VPN instance address family.
- The EVPN NLRI defines different types of BGP EVPN routes





# BGP EVPN Route Types

- The BGP EVPN route types are as follows:
  - Type 2 route — MAC/IP route: is used to advertise the MAC addresses, ARP entries, and routing information of hosts
  - Type 3 route — inclusive multicast route: is used for automatic discovery of VTEPs and dynamic establishment of VXLAN tunnels

Route Distinguisher	RD of an EVPN instance.
Ethernet Segment Identifier	Unique ID for defining the connection between local and remote devices.
Ethernet Tag ID	VLAN ID.
MAC Address Length	Length of the host MAC address carried in the route.
MAC Address	Host MAC address carried in the route.
IP Address Length	Mask length of the host IP address carried in the route.
IP Address	Host IP address carried in the route.
MPLS Label1	Layer 2 VNI carried in the route.
MPLS Label2	Layer 3 VNI carried in the route.

## Prefix

Route Distinguisher	RD of an EVPN instance.
Ethernet Tag ID	VLAN ID. Here, the value is 0.
IP Address Length	Mask length of the local VTEP's IP address carried in the route.
Originating Router's IP Address	Local VTEP's IP address carried in the route.

## PMSI attribute

Flags	Flag bit. This field is inapplicable in VXLAN scenarios.
Tunnel Type	Tunnel type carried in the route. The value can only be 6.
MPLS Label	Layer 2 VNI carried in the route.
Tunnel Identifier	Tunnel identifier carried in the route.

## NLRI format of Type 2 and Type 3 Routes

Source: Lei Zhang, Le Chen. Cloud Data Center Network Architectures and Technologies, CRC Press 2021

# BGP EVPN Route Types (Contd.)

- The BGP EVPN route types are as follows:
  - Type 5 route — IP prefix route: is used to advertise imported external routes or advertise routing information of hosts

Route Distinguisher	RD of an EVPN instance.
Ethernet Segment Identifier	Unique ID for defining the connection between local and remote devices.
Ethernet Tag ID	VLAN ID.
IP Prefix Length	Length of the IP prefix carried in the route.
IP Prefix	IP prefix carried in the route.
GW IP Address	Default gateway address.
MPLS Label	Layer 3 VNI carried in the route.

## 6.9 NLRI format in a type 5 route.

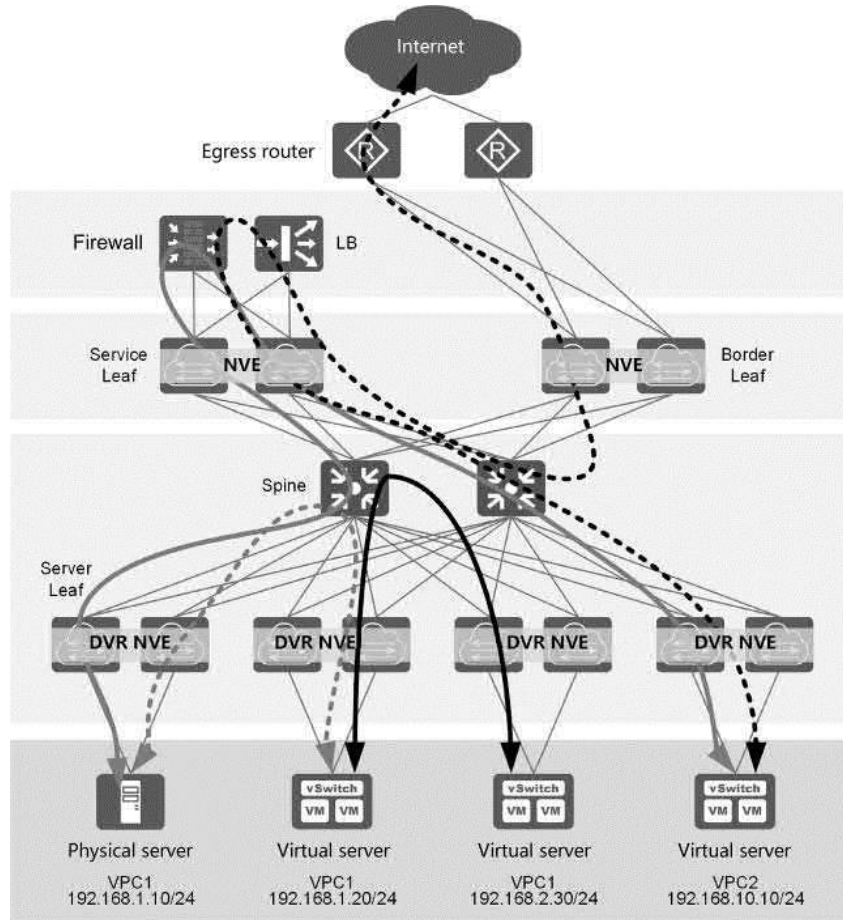
# For completeness – the VXLAN Data Plane!



- Depending on the traffic flow direction and scope, DCN traffic can be classified into east-west traffic (transmitted within a DC) and north-south traffic (sent across the DC)
  - Traffic transmitted within the same subnet of a VPC is forwarded by a TOR switch after Layer 2 VXLAN encapsulation.
  - Traffic transmitted between subnets of the same VPC is forwarded by a TOR switch based on Layer 3 routes. This is done after Layer 3 VXLAN encapsulation.
  - Traffic transmitted between VPCs is forwarded across subnets, and isolation for security purposes is required. Therefore, to meet this, the Traffic needs to pass through a firewall and reach the Layer 3 VXLAN gateway.
  - Traffic sent from a user outside the DC to a server in a VPC passes through the Intrusion Prevention System (IPS) or firewall, LB, VXLAN gateway, and TOR switch before reaching the server.



# VXLAN Data Plane



Source: Lei Zhang, Le Chen. Cloud Data Center Network Architectures and Technologies, CRC Press 2021

**Thank You!**

