



Network Fundamentals for Cloud

BITS Pilani
Pilani Campus

Nishit Narang
WILPD-CSIS



CC ZG503: Network Fundamentals for Cloud

Lecture No. 11: Data Center Networks (Contd.)

Story so far...

- Older DCN used L2 networks (Acc-Agg-Core)
 - Issues with L2 networks w.r.t scalability and east-west traffic handling
 - Broadcast storm, bandwidth under-utilization....
 - xSTP, Virtual Chassis technologies introduced – had own limitations
 - L2MP was introduced.....
-
- In the meanwhile, DCN are evolving towards Clos Topology (leaf-spine) for various advantages in scalability and cost economics.....

RECAP: DCN Networking Technology Evolution (contd.)



- 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: DCN Networking Technology Evolution (contd.)



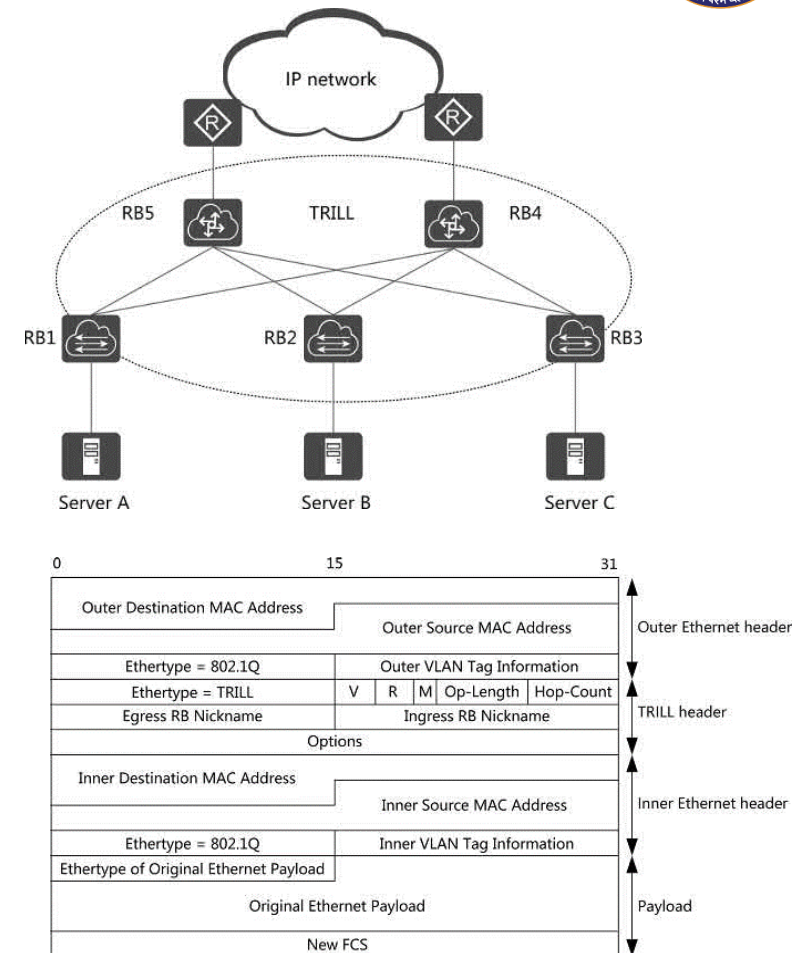
- Layer 2 Multi-Pathing (L2MP) Technologies (contd.)
 - A link-state routing protocol usually requires that each node on a network be addressable
 - Each node uses the link-state routing protocol to calculate the network topology and then calculates the forwarding database based on the network topology
 - Therefore, L2MP technologies need to add an addressable identifier, which is similar to an IP address on an IP network, to each device on a network
 - TRILL is a standard L2MP protocol defined by the Internet Engineering Task Force (IETF) that became popular.



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

Disadvantages of L2MP

- 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



Meanwhile, what is changing in the DCN?

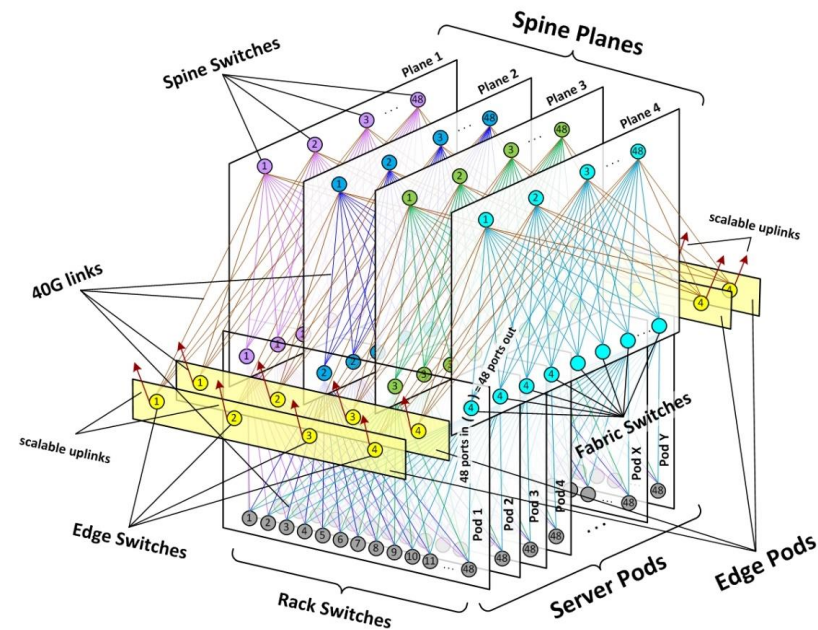
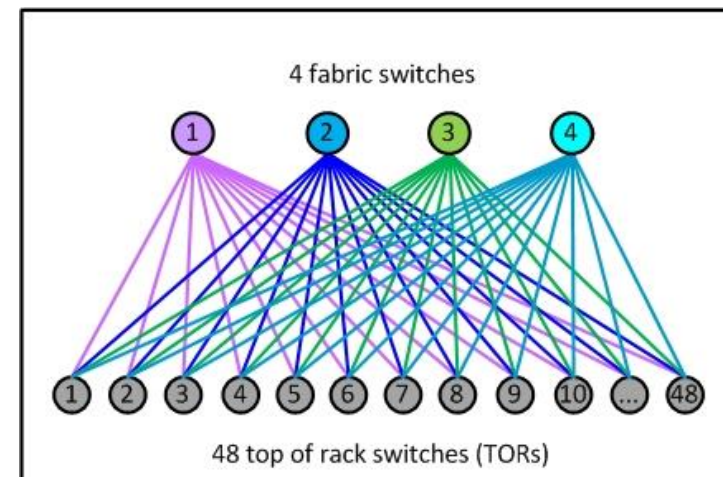
- As DCNs scale and traffic patterns change to more east-west traffic....
- ...DCN topologies evolve towards Clos Topology and its variants....
-DCN protocols evolve towards use of L3 technologies (IP & associated networking/control plane protocols)
 - Refer following case studies

Case Studies

- [A Scalable, Commodity Data Center Network Architecture](http://ccr.sigcomm.org/online/files/p63-alfares.pdf) (Research Paper):
(<http://ccr.sigcomm.org/online/files/p63-alfares.pdf>)
 - Introduction (Sec 1): DC Applications and their traffic patterns
 - Section 2.1: DC Network Topologies and Bandwidth Oversubscription
 - Section 2.2: Clos Networks and Fat-Tree Topology

Case Studies

- [Introducing data center fabric, the next-generation Facebook data center network - Engineering at Meta](#)
 - Facebook is a good example application to explain cloud application traffic patterns and networking needs
 - The example shows performance limitations encountered by FB in M2M traffic when using cluster-design. And the migration to the next-generation fabric design to overcome this challenge. This introduces modularity, leading to gradual scalability.
 - Role of BGP4 (distributed routing) alongside a centralized BGP controller (for centralized override). FB calls this hybrid approach as “**distributed control, centralized override**”.
 - High-capacity 40G links connecting fabric switches, TOR switches and Spine switches. Rack servers connected to TOR via 10G links.



Source: [Introducing data center fabric, the next-generation Facebook data center network - Engineering at Meta](#)



Case Studies

- Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google's Datacenter Network
 - CLOS Topologies (Leaf'n'spine) in data centers
 - Google datacenter networks run at dozens of sites across the planet, scaling in capacity by 100x over ten years to more than 1Pbps of bisection bandwidth.
 - Much of the general, but complex, decentralized network routing and management protocols supporting arbitrary deployment scenarios are overkill for single-operator, pre-planned datacenter networks. Here, a centralized control and management mechanism is discussed, based on a global configuration that is pushed to all datacenter switches.
 - Granular control over ECMP tables with proprietary, scalable in-house IGP
 - Use standard BGP between Cluster Border Routers and external vendor gear
 - Proprietary Neighbor Discovery (ND) protocol for online liveness and peer correctness checking - used for correcting cabling errors, one of the key challenges in a large data center
 - Data Center Challenges (viz. Fabric Congestion and Outages) discussed in section 6.

Back to DCN Networking Technology Evolution: NVO3: & Overlay Networks



- Virtual chassis, L2MP, and multi-chassis link aggregation technologies can solve problems of xSTP technologies. However, these technologies are fundamentally traditional network technologies and are still hardware-centric.
- NVO3 technologies are overlay network technologies driven by IT vendors and aim to get rid of the dependency on the traditional physical network architecture

Overlay Network:

- A software-defined logical network built over an existing underlay network
- Is completely decoupled from the underlay network
 - Allows the underlay network to be flexibly expanded
 - Facilitates SDN architecture deployment
 - SDN controller is not required to consider the underlay network architecture and can flexibly deploy services on the overlay network
- Created using NVO3 technology
 - Overlay / NVO3 technology is a tunnel encapsulation technology that encapsulates Layer 2 packets over tunnels and transparently transmits the encapsulated packets
 - In a DCN environment, it enables layer 2 communication between large-scale VMs on the DCN

NVO3 technologies include VXLAN and NVGRE. VXLAN is used by the majority of enterprises



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

Where else have you seen the use of tunnels??



- A tunnel is often used to encapsulate a packet of one protocol into a packet of another protocol to carry it over an intermediate network that supports the latter protocol



DCN Networking Technology Evolution (contd.)



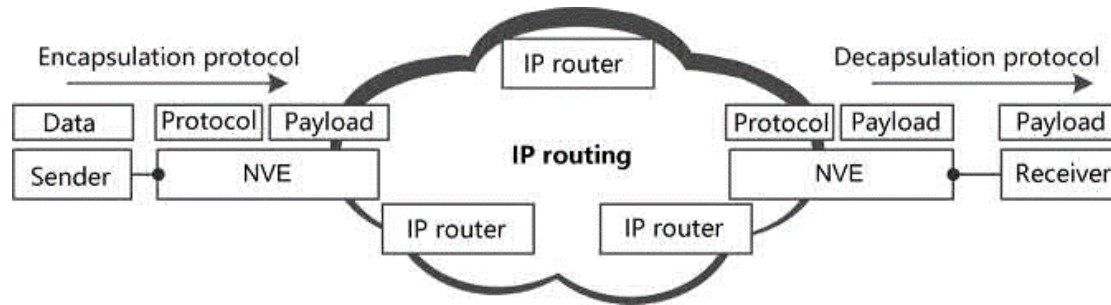
- Network Virtualization Overlays (NVO3) Technologies
 - Since an overlay network is a virtual network topology constructed on top of a physical network, thus...
 - Each virtual network instance that is implemented as an overlay, requires that an original frame is encapsulated on a Network Virtualization Edge (NVE).
 - The encapsulation identifies the device that will perform decapsulation.
 - Before sending the frame to the destination endpoint, the device decapsulates the frame to obtain the original frame.
 - Intermediate network devices forward the encapsulated frame based on the outer encapsulation header and are oblivious to the original frame carried in the encapsulated frame.
 - The NVE can be a traditional switch or router, or a virtual switch in a hypervisor.
 - The endpoint can be a VM or physical server.
 - A VXLAN network identifier (VNI) can be encapsulated into an overlay header to identify a virtual network to which a data frame belongs.
 - Because a virtual DC supports both routing and bridging, the original frame in the overlay header can be a complete Ethernet frame containing a MAC address or an IP packet.



DCN Networking Technology Evolution (contd.)



- 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

DCN Networking Technology Evolution (contd.)



- L2MP Vs NVO3 Technologies

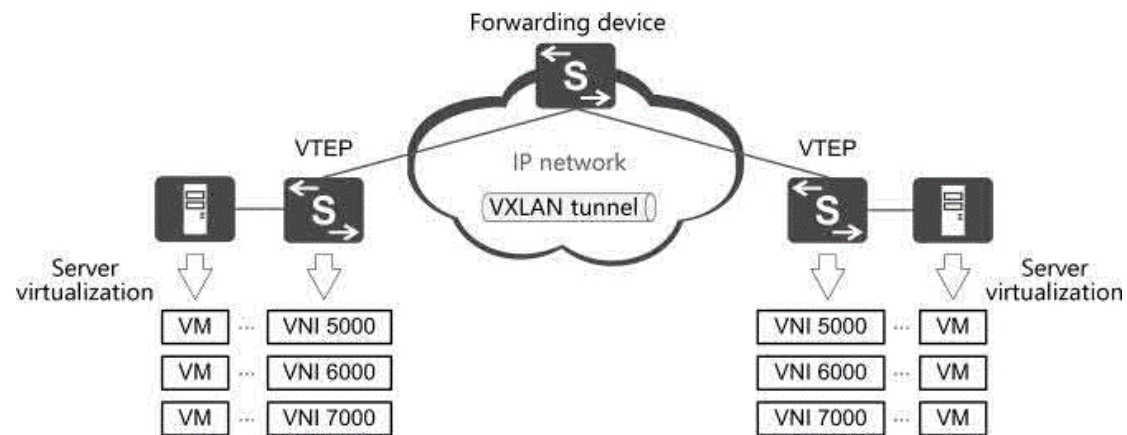
- To some extent, NVO3 and L2MP technologies are similar.
- They both build an overlay network on the physical network
- The difference is that L2MP technologies add a new forwarding identifier to the original Layer 2 network, thereby requiring that chips on hardware devices support L2MP technologies.
- In contrast, NVO3 technologies reuse the current IP forwarding mechanism and only add a new logical network that does not depend on the physical network environment on the traditional IP network.
- The logical network is not perceived by physical devices, and its forwarding mechanism is the same as the IP forwarding mechanism.
- In this way, the threshold of NVO3 technologies is greatly lowered, and this is why NVO3 technologies have become popular on DCNs in a few years

Typical NVO3 technologies include VXLAN, Network Virtualization Using Generic Routing Encapsulation (NVGRE), and Stateless Transport Tunneling (STT), among which VXLAN is the most popular one.



VXLAN Basics

- 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



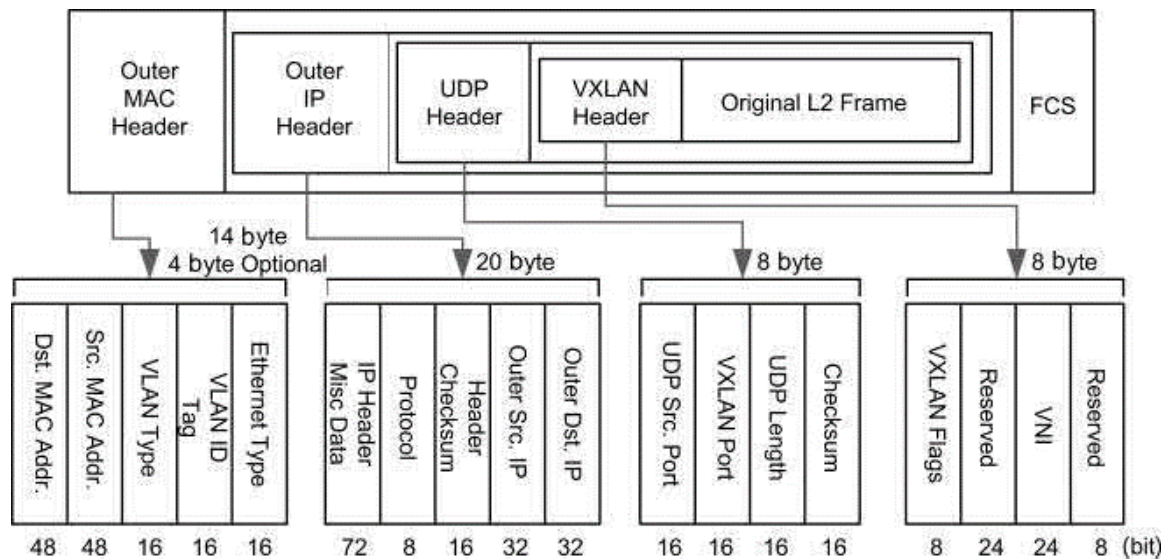
Source: Lei Zhang, Le Chen. Cloud Data Center Network Architectures and Technologies, CRC Press 2021

VXLAN Benefits

- **VLAN flexibility in multitenant segments:** It provides a solution to extend Layer 2 segments over the underlying network infrastructure so that tenant workload can be placed across physical pods in the data center.
- **Higher scalability:** VXLAN uses a 24-bit segment ID known as the VXLAN network identifier (VNID), which enables up to 16 million VXLAN segments to coexist in the same administrative domain.
- **Improved network utilization:** VXLAN solves the Layer 2 STP limitations. VXLAN packets are transferred through the underlying network based on its Layer 3 header and can take complete advantage of Layer 3 routing, equal-cost multipath (ECMP) routing, and link aggregation protocols to use all available paths.

VXLAN Packet Format

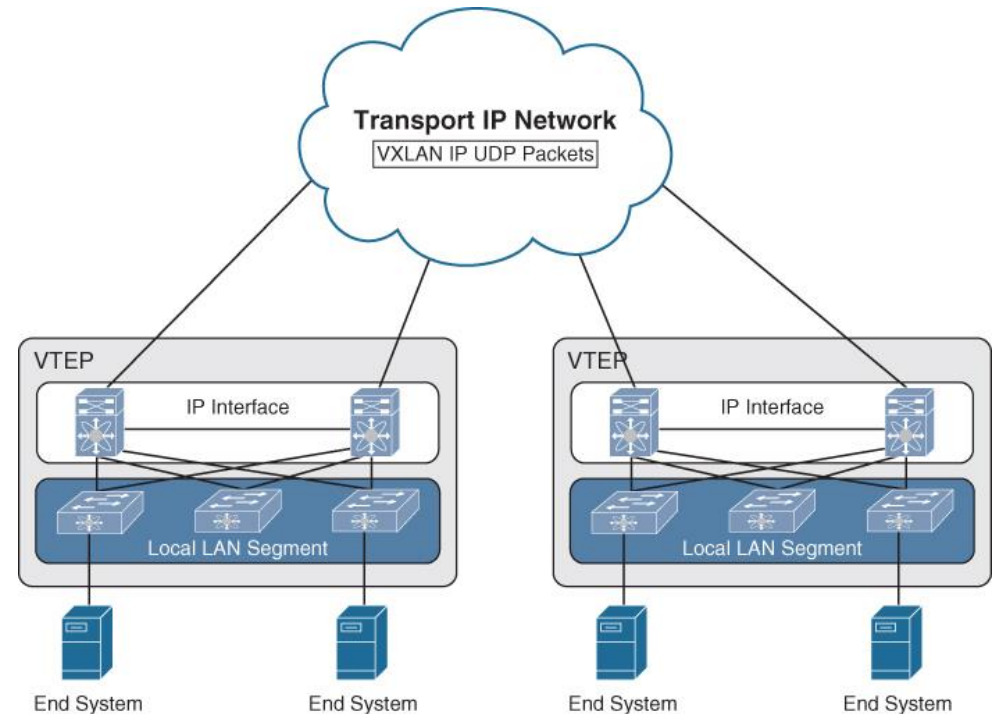
- The transport protocol over the physical data center network is UDP/IP
- With MAC-in-UDP encapsulation, VXLAN tunnels the Layer 2 network over the Layer 3 network.



- In the UDP header,
 - the destination port number has a fixed value of 4789
 - The source port number is the hash value of the original Ethernet frame
- Outer IP (/MAC) header
 - source IP (/MAC) address specifies the IP (/MAC) address of the VTEP, where the source VM belongs.
 - destination IP address indicates the IP address of the VTEP where the destination VM belongs
 - destination MAC address is the MAC address of the next-hop device on the path to the destination VTEP.

VXLAN Tunnels and VTEP

- VXLAN tunnel endpoint (VTEP) maps tenants' end devices to VXLAN segments performs VXLAN encapsulation and de-encapsulation
- Each VTEP function has two interfaces:
 - one is a switch interface on the local LAN segment to support local endpoint communication, and
 - the other is an IP interface to the transport IP network
- A VTEP device discovers the remote VTEPs for its VXLAN segments and learns remote MAC Address-to-VTEP mappings through its IP interface



- A virtual network identifier (VNI) is a value that identifies a specific virtual network in the data plane
- It is typically a 24-bit value part of the VXLAN header, which can support up to 16 million individual network segments.
 - Valid VNI values are from 4096 to 16,777,215.
- There are two main VNI scopes:
- **Network-wide scoped VNIs:** The same value is used to identify the specific Layer 3 virtual network across all network edge devices
 - A uniform VNI per VPN is a simple approach → eases network operations
- **Locally assigned VNIs:** In an alternative approach supported as per RFC 4364, the identifier has local significance to the network edge device that advertises the route
 - uses the same existing semantics as an MPLS VPN label

VXLAN Overlay Network Types

- Classified as one of three types:
 - Network overlay: All VTEPs are deployed on physical switches.
 - Host overlay: All VTEPs are deployed on vSwitches.
 - Hybrid overlay: Some VTEPs are deployed on physical switches with others deployed on vSwitches

VXLAN Control Plane

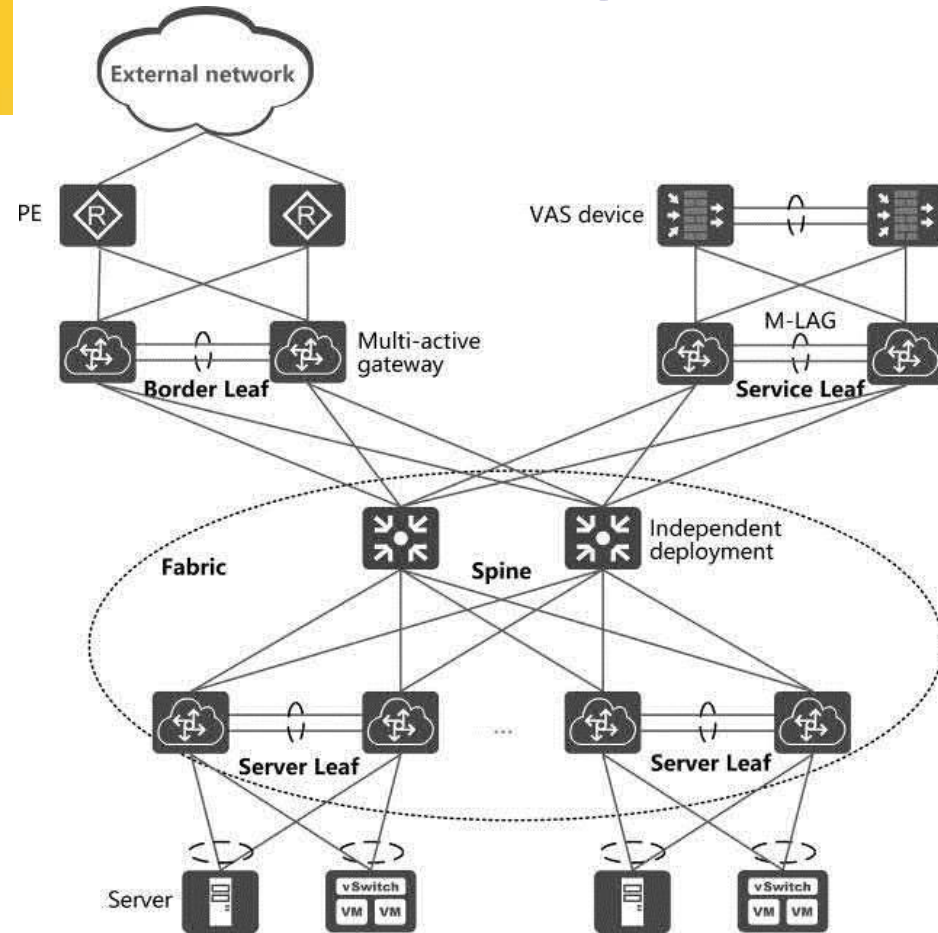
- Two widely adopted control planes are used with VXLAN:
 - the VXLAN Flood and Learn Multicast-Based Control Plane and
 - the VXLAN MPBGP EVPN Control Plane.

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

Thank You!

