# This session's agenda

- Physical topology – Clos network (Leaf-Spine)
- Forwarding in Clos network
- Protocols inside the ACI fabric – IS-IS
- Traditional MPLS service provider networks vs ACI
- Endpoint connections
- ACI endpoint learning
- Protocols inside the ACI fabric – COOP
- Bridge domains, Dataplane in ACI - VXLAN
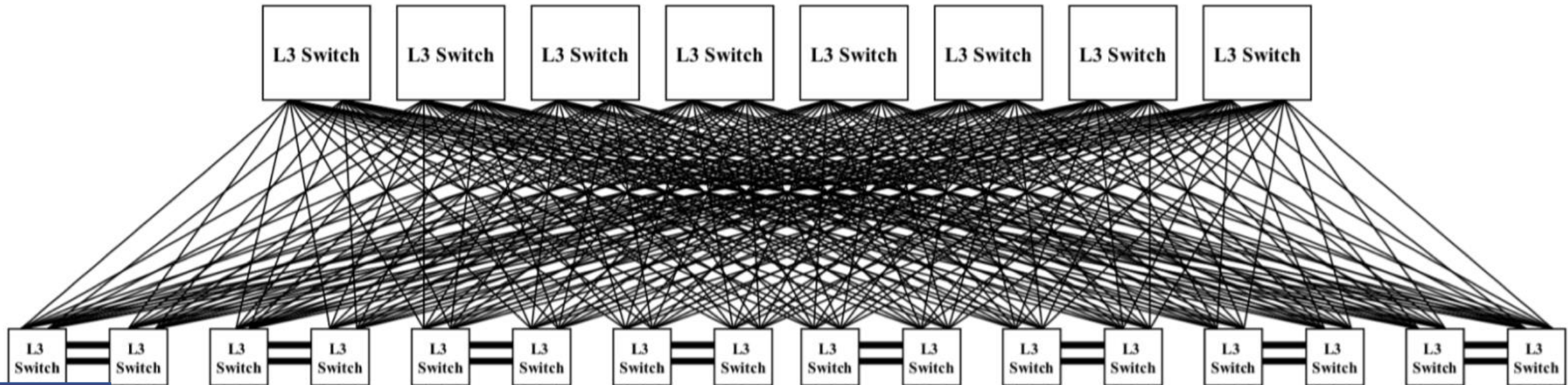- Some design options

# Next session's agenda

- ACI External connections – L3Out
- Protocols inside the ACI fabric – MP-BGP
- Connecting multiple datacentres

- Physical topology – Clos network
- Forwarding in Clos network
- Protocols in ACI - IS-IS
- Traditional service provider networks vs ACI
- Undelay and Overlay
- Endpoints, learning in the ACI fabric
- Protocols in ACI - COOP
- Bridge domains, VXLAN
- Pervasive (Anycast) gateway, VRFs

**Oversimplified**

# Physical layer - Clos topology (leaf-spine)

# Let's start with basics

- Take two routers

- Connect with two cables of the same speed

- Configure any random subnets on these links on both routers, IPv4 or IPv6

- Configure loopbacks

- Run any routing protocols with statement 'redistribute connected'

Routing
protocol

Int Loopback 0
10.10.10.1/32

L3 Link1   x.x.x.x/30

L3 Link2   y.y.y.y/30

Int Loopback 0
20.20.20.1/32

conf t
router bgp / ofsp / isis   x
redistribute connected

# What we'll get

Routing
protocol

Int Loopback 0
10.10.10.1/32

Link1
Link2

Int Loopback 0
20.20.20.1/32

show ip route:
    20.20.20.1/32 :
        via link1 cost 1
        via link2 cost 1

show ip route:
    10.10.10.1/32 :
        via link1 cost 1
        via link2 cost 1

# Simple add more links to increase bandwidth

Routing protocol

Int Loopback 0
10.10.10.1/32

Link1
Link2
Link3

Int Loopback 0
20.20.20.1/32

show ip route:
   20.20.20.1/32 :
      via link1 cost 1
      via link2 cost 1
      via link3 cost 1

show ip route:
   10.10.10.1/32 :
      via link1 cost 1
      via link2 cost 1
      via link2 cost 1

# Add intermediate router

- Take one more router

- Insert between two existing routers

- Run the same routing protocol with default settings

- Call routers on the left and right, border or leaf routers

- Call router in centre – spine router
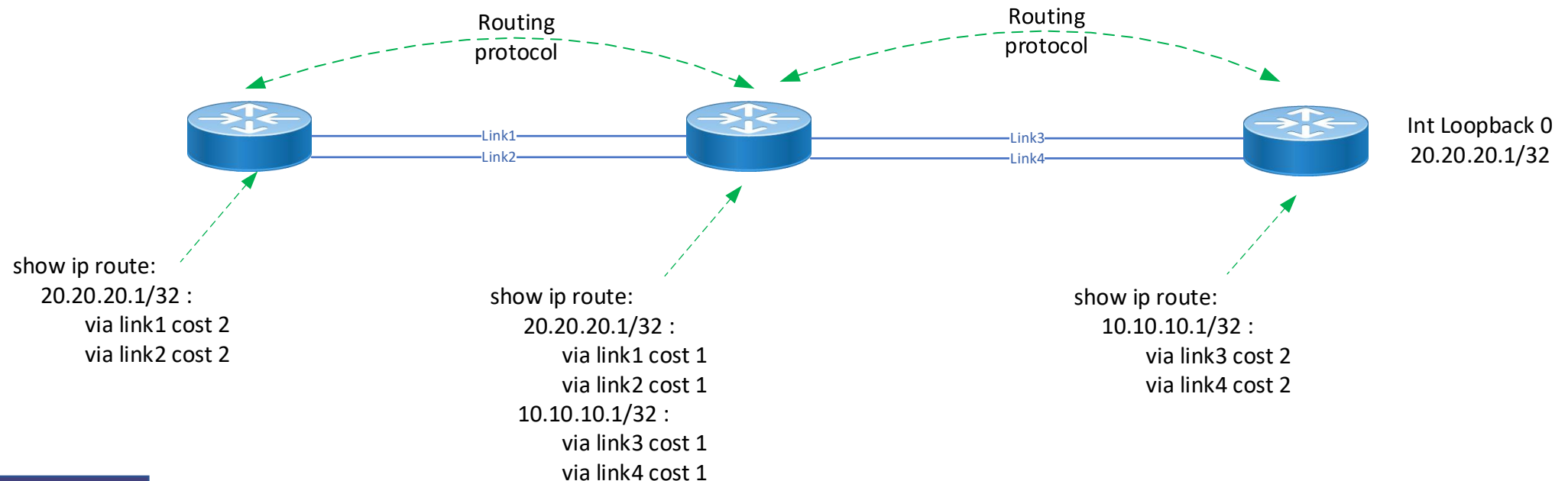
What we've got:

- Routing tables on leaf routers as still the same, the only thing has changed is the cost from 1 to 2



Routing protocol

Routing protocol

Link1
Link2
Link3
Link4

Int Loopback 0
20.20.20.1/32

show ip route:
    20.20.20.1/32 :
        via link1 cost 2
        via link2 cost 2

show ip route:
    20.20.20.1/32 :
        via link1 cost 1
        via link2 cost 1
    10.10.10.1/32 :
        via link3 cost 1
        via link4 cost 1

show ip route:
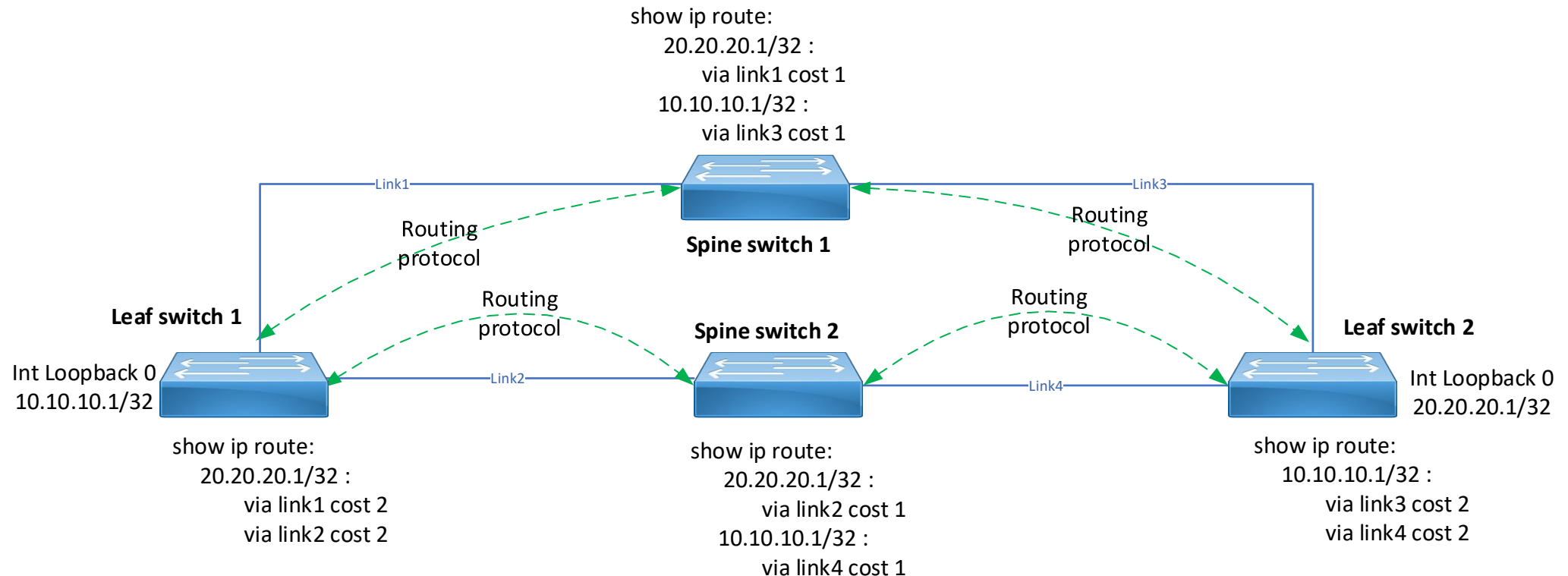    10.10.10.1/32 :
        via link3 cost 2
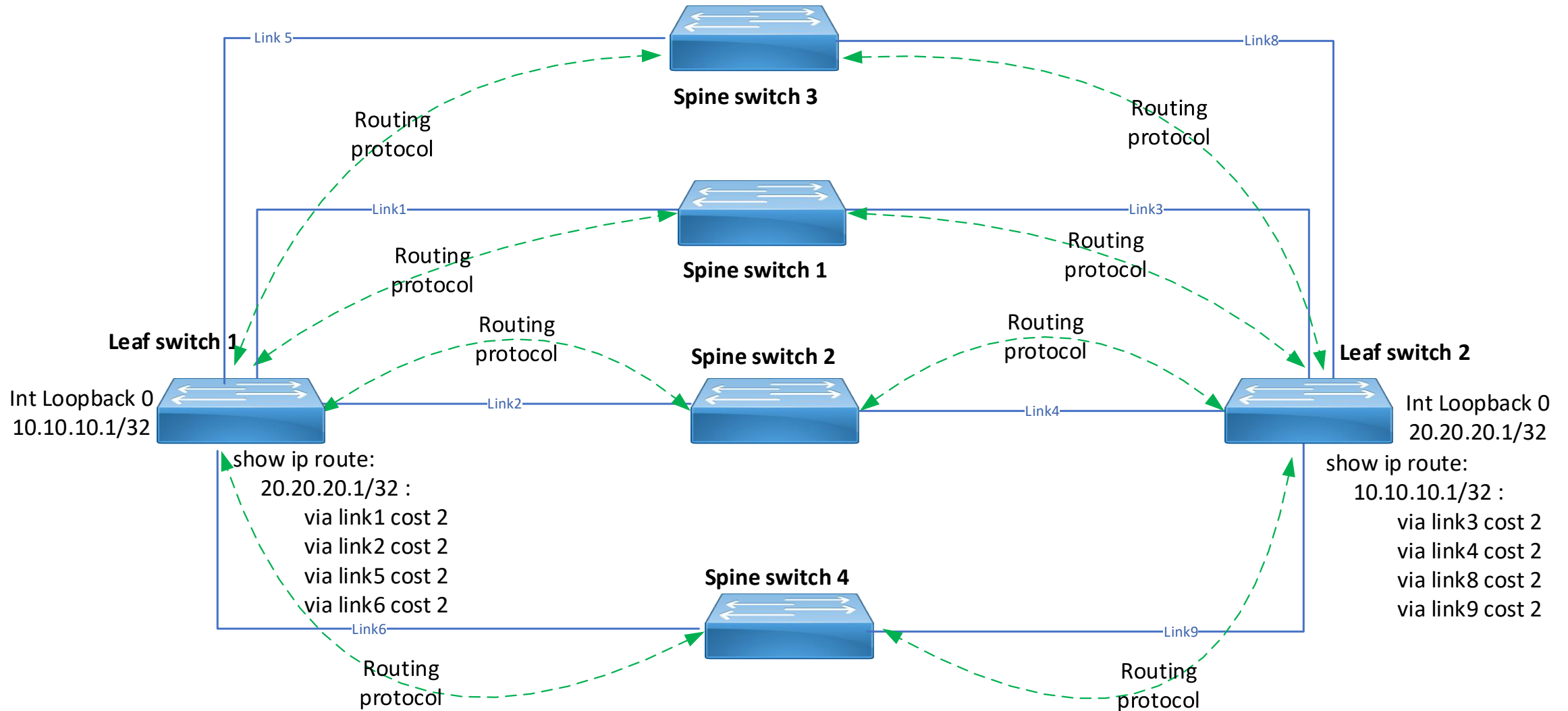        via link4 cost 2

# Let's add more ingredients

- Replace routers with L3 switches – leaf switches

- Add more intermediate routers or L3 switches – spine switches

- Connect all of them with the same cables, run the same routing protocol

What we've got:

- Spine switches send routing updates - Routing tables on leaf switches are still the same, showing **equal paths via different L3 links**

show ip route:
    20.20.20.1/32 :
        via link1 cost 1
    10.10.10.1/32 :
        via link3 cost 1

**Spine switch 1**

Link1

Routing
protocol

Routing
protocol

Link3

Routing
protocol

**Leaf switch 1**

**Spine switch 2**

Routing
protocol

**Leaf switch 2**

Int Loopback 0
10.10.10.1/32

Link2

Link4

Int Loopback 0
20.20.20.1/32

show ip route:
    20.20.20.1/32 :
        via link1 cost 2
        via link2 cost 2

show ip route:
    20.20.20.1/32 :
        via link2 cost 1
    10.10.10.1/32 :
        via link4 cost 1

show ip route:
    10.10.10.1/32 :
        via link3 cost 2
        via link4 cost 2

# And even more



**Spine switch 3**

Routing protocol

Routing protocol

Link 5

Link8

**Spine switch 1**

Link1

Link3

Routing protocol

Routing protocol

**Leaf switch 1**

**Leaf switch 2**

Int Loopback 0
10.10.10.1/32

Int Loopback 0
20.20.20.1/32

Routing protocol

Routing protocol

**Spine switch 2**

Link2

Link4

show ip route:
    20.20.20.1/32 :
        via link1 cost 2
        via link2 cost 2
        via link5 cost 2
        via link6 cost 2

show ip route:
    10.10.10.1/32 :
        via link3 cost 2
        via link4 cost 2
        via link8 cost 2
        via link9 cost 2

**Spine switch 4**

Link6

Link9

Routing protocol

Routing protocol

# Re-arrange the diagram

# Add more leaf switches



Spine switch 1

Spine switch 2

L3 link1

L3 link2

L3 link4

L3 link3

Leaf switch 1

Leaf switch 2

Leaf switch 3

Leaf switch 4

# Or add more spine switches

# All switches are connected with single physical link



Spine switch 1

Spine switch 2

Spine switch 3

Link 1 - 100G

Link 2 – 200G

Link 3 – 200G

Leaf switch 1

Leaf switch 2

Leaf switch 3

Leaf switch 4

show ip route:
    all leafs :
        via link2 cost 2
        via link3 cost 2

link1 cost 4 – not in routing table

Spine switch 1    Spine switch 2    Spine switch 3
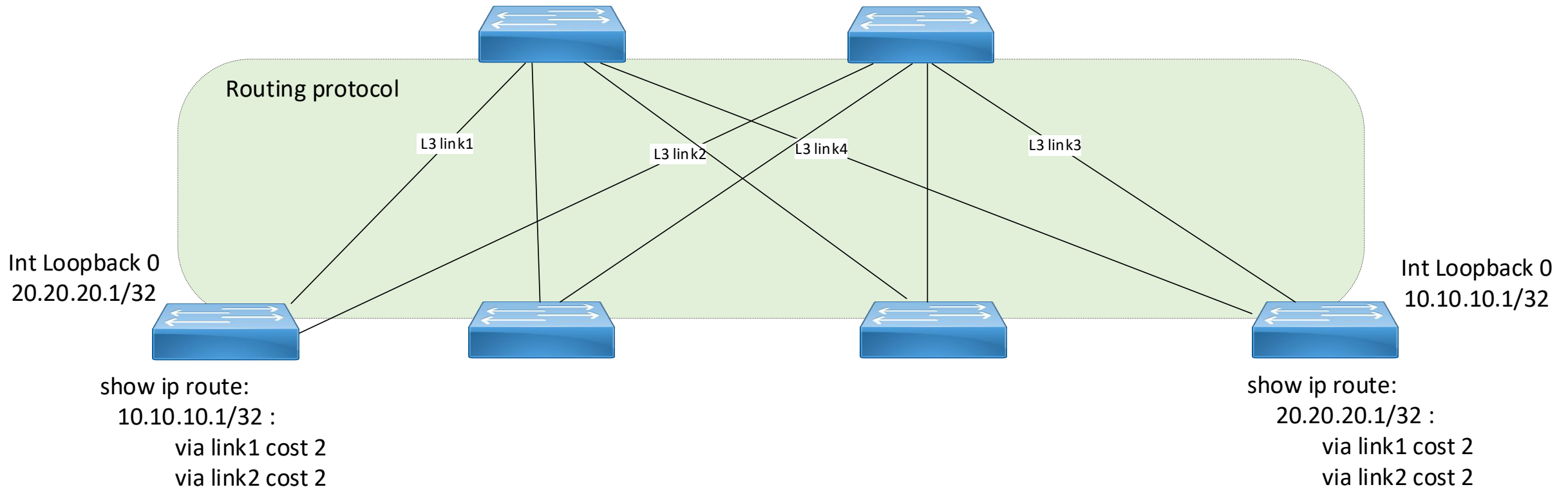
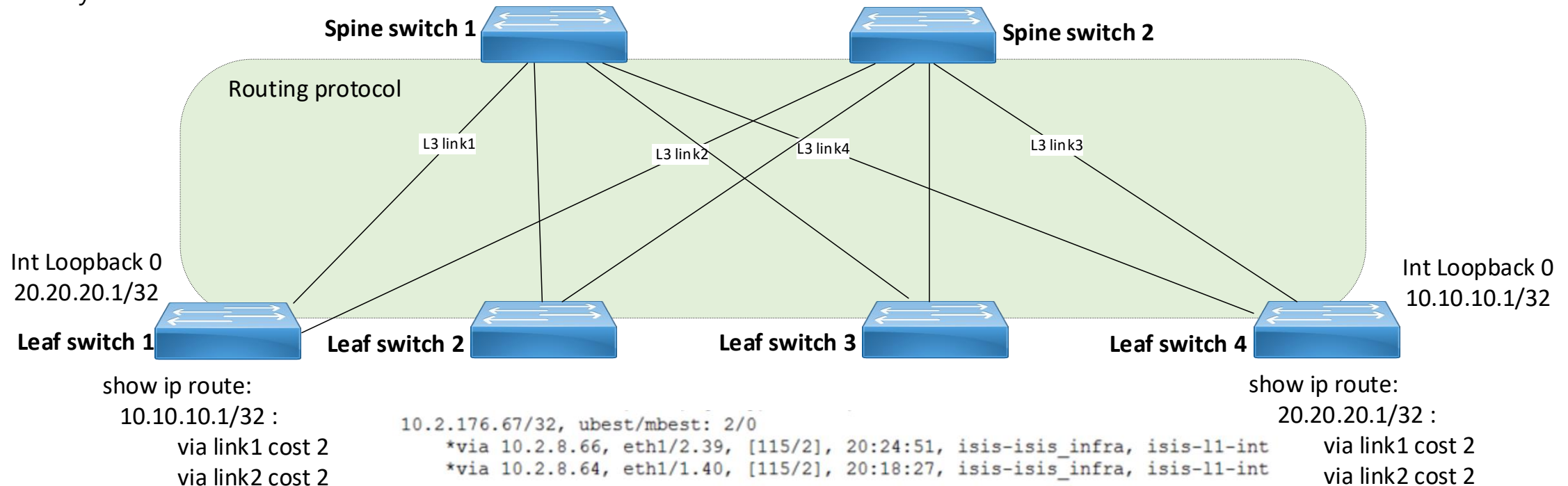Leaf switch 1    Leaf switch 2    Leaf switch 3    Leaf switch 4

# IS-IS in ACI

- ACI uses IS-IS between leaf and spines inside Pod and Site (later about multi-Pod)
- IS-IS has very basic settings – single Level 1 Area
- Not configurable by administrators

Routing protocol

L3 link1

L3 link2

L3 link4

L3 link3

Int Loopback 0
20.20.20.1/32

Int Loopback 0
10.10.10.1/32

show ip route:
    10.10.10.1/32 :
        via link1 cost 2
        via link2 cost 2

show ip route:
    20.20.20.1/32 :
        via link1 cost 2
        via link2 cost 2

# Clos topology – ACI physical layer

- All links are the same speed, Layer 3, **the same metric** from routing protocol's perspective

- All links are active, load-balanced – forwarding based on L3 **Equal-cost** multi-path routing (**ECMP**)

- IP addressing on the links doesn't matter – can be IPv4, IPv6 link-local, anything – as long as both sides can reach each other

- What's important is reachability of Loopbacks – they are called VTEPs (remember this definition, will be explained later):

- Each leaf switch knows how to reach other leaf switch (more specifically, VTEP) via **multiple equal** paths

- Easy to scale

**Spine switch 1**　　　　　　　　　　　　**Spine switch 2**

Routing protocol

L3 link1　　　　　　　L3 link2　　　L3 link4　　　　　L3 link3

Int Loopback 0
20.20.20.1/32

Int Loopback 0
10.10.10.1/32

**Leaf switch 1**　　**Leaf switch 2**　　　　**Leaf switch 3**　　　**Leaf switch 4**

show ip route:
　10.10.10.1/32 :
　　via link1 cost 2
　　via link2 cost 2

```
10.2.176.67/32, ubest/mbest: 2/0
   *via 10.2.8.66, eth1/2.39, [115/2], 20:24:51, isis-isis_infra, isis-l1-int
   *via 10.2.8.64, eth1/1.40, [115/2], 20:18:27, isis-isis_infra, isis-l1-int
```

show ip route:
　20.20.20.1/32 :
　　via link1 cost 2
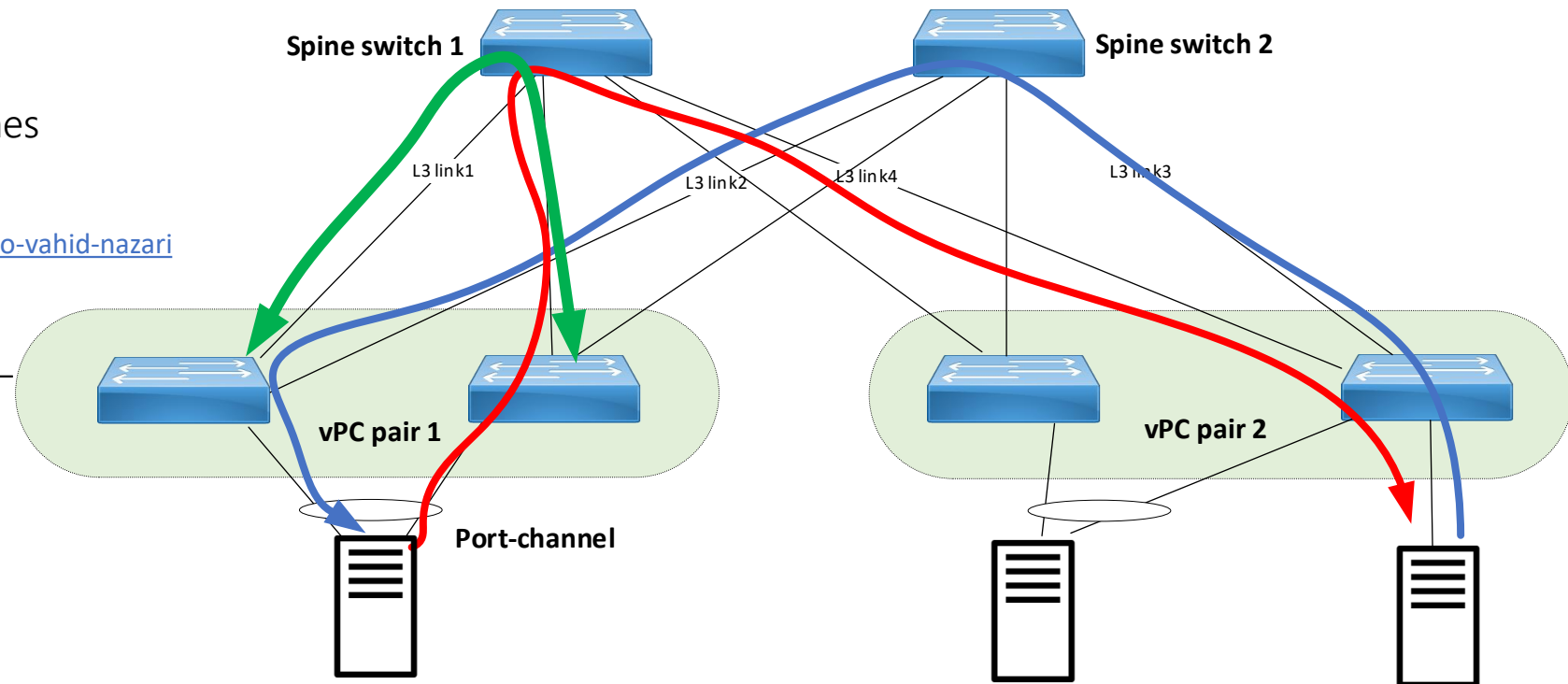　　via link2 cost 2

# Clos topology – consequences

- VPC pair are no longer require VPC peer-link and PKL  https://www.cisco.com/c/dam/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/aci-guide-vpc.pdf

- VPC peers communicate over spine

- No more issues with orphan (single-homed) connections in VPC

- Traffic flows can be asymmetric, but it's OK, as **paths are equal**    **(see picture below)**

- Leaf switches can be ToR or shared between racks (Middle-of-row topology)

- No FEXes (normally),

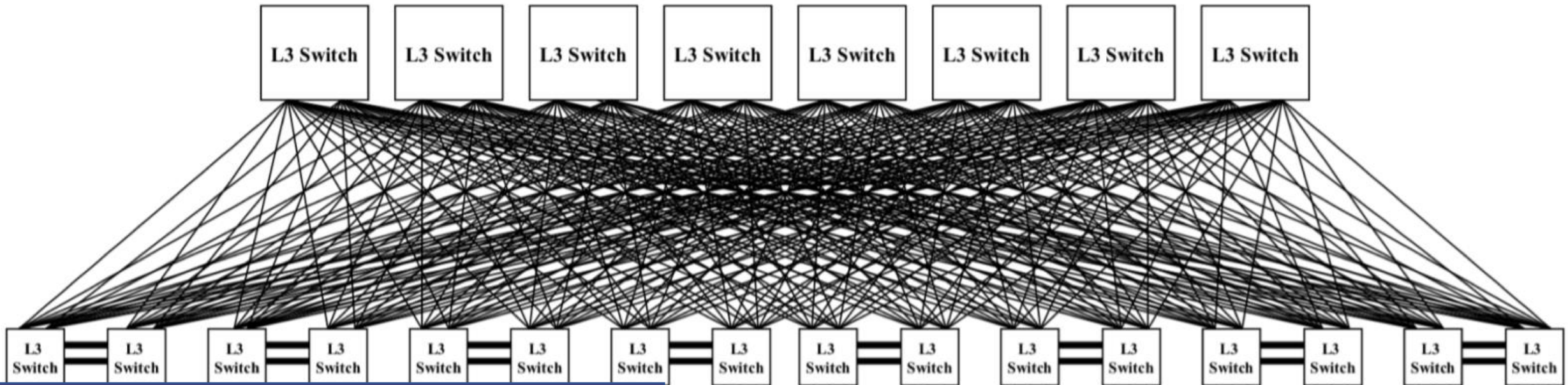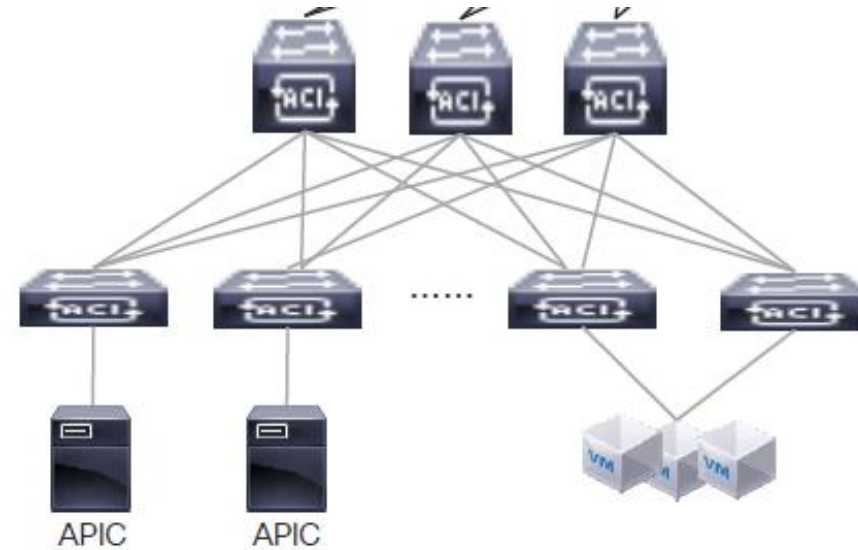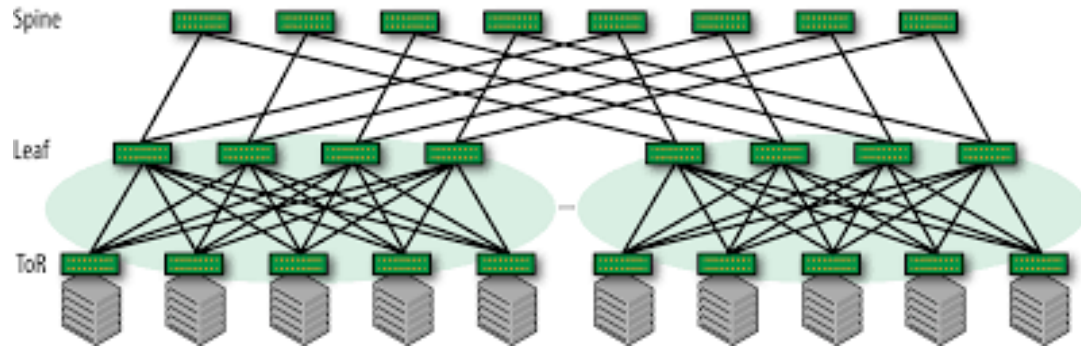but still can be connected to leaf switches

https://www.linkedin.com/pulse/

why-you-shouldnt-think-fabric-extenders-fex-along-cisco-vahid-nazari

- **No any other connections to spines –**
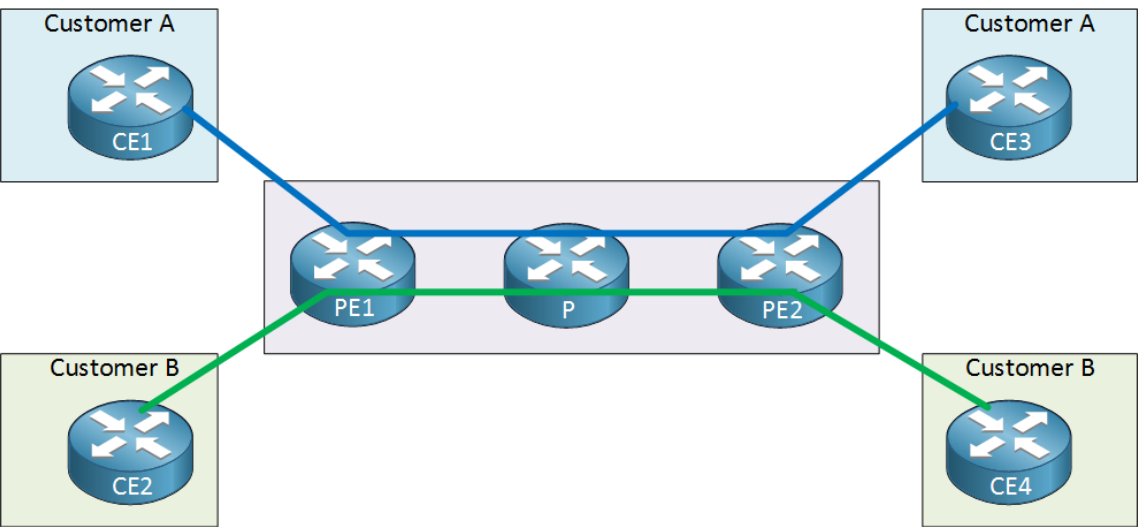
(the only exception is DCI)

# Clos topology

- Physical topology – Clos network
- Forwarding in Clos network
- Protocols in ACI - IS-IS
- Traditional service provider networks vs ACI
- Undelay and Overlay
- Endpoints, learning in the ACI fabric
- COOP
- Bridge domains, VXLAN
- Pervasive (Anycast) gateway, VRFs

Oversimplified

# Traditional IP-VPN networks – connecting L3 customer segments

How it works – connecting different customer sites

Site 1 --> Site 3:

- PE1 knows how to reach PE2 via SP backbone (via IGP – OSPF or IS-IS)

- PE1 learns customer routes at Site 1

- PE3 learns customer routes at Site 2

- PE1 establishes iBGP session with PE2



- PE1 and PE2 send each other these learned customer routers with MP-BGP and assign them a 'VPN ID' – Route target and assign Inner Label

- As a result PE1 and PE3 know each others' routes and which customer these prefixes belong to

- When packet arrives at PE1 from Site 1 with dest at Site 2, PE1 takes the original packet adds MPLS header (with Inner Label) and sends via the SP backbone to PE2

- Intermediate P router don't about customers' routes, it only responsible to deliver packets between PE router

- PE2 receives the packet, examines **inner label**, showing which **customer** the packet belongs to, and forward the original packet to Site2

Edge Routers can also exchange customers' L2 MAC in a special MP-BGP address family - EVPN (Ethernet VPN)

# Service Provider networks vs ACI



Two completely different networks, not visible to each other:

– SP Transport network – **underlay**

- Customers' networks – **customer overlays**

Important points about underlay (SP transport) network:

- Provider edge routers (Leaf) know how to reach each other

- Provider edge routers (Leaf) exchange customer L3 prefixes or if it's L2 VPN – MAC addresses

- The edge routers (Leaf) set some kind of **label** to identify what customer VPN the route belongs to and send this labelled packet to remote edge router

- The edge router don't exchange customer VLANs – customers can configure any VLANs, they have local significance

- Transit P routers (**Spine**) **don't handle customer** (tenant) packets – only forwards packets between edge routers

Cisco ACI fabric is a Service Provider network – Transport underlay

In Cisco ACI this transport network is called Infra VRF

PE == Leaf switches, P router == Spine switches

ACI Integration with WAN at Scale
'Project GOLF' Overview

- Physical topology – Clos network
- Forwarding in Clos network
- Protocols in ACI - IS-IS
- Traditional service provider networks vs ACI
- Undelay and Overlay
- Endpoints, learning in the ACI fabric
- Protocols in ACI - COOP
- Bridge domains, VXLAN
- Pervasive (Anycast) gateway, VRFs

Oversimplified

# Connecting 'our customers' – Endpoints

## End Point (EP)

### What is an EP?
- It stands for hosts, in other words MAC address with IP(s)
  - ➤ sometimes MAC only
  - ➤ IP in EP is always /32

### What Forwarding Table is used?
- End Point Table
  - ➤ host information (MAC and /32 IP address)

- LPM(Longest Prefix Match) Table
  - ➤ non /32 IP route information (exception: /32 for SVI or L3OUT route)



These are End Points

| Legacy | ACI |
|---|---|
| RIB ( non-/32 & /32 ) | RIB ( non /32 ) |
| MAC | EndPoint ( mac & /32 ip ) |
| ARP | ARP (only for L3OUT) |

**Forwarding table lookup order**
1. EndPoint Table (show endpoint)
2. RIB (show ip route)

RIB : Routing Information Base

https://www.ciscolive.com/c/dam/r/ciscolive/emea/docs/2019/pdf/BRKACI-3545.pdf

# How Leaf switches know about Endpoints

**Local Endpoint (MAC)**
A leaf learns MAC A as local if a packet with src MAC A comes in from its front panel port.

**Local Endpoint (/32 host IP)**
A leaf learns IP A /32 as local
- if a packet with src IP A comes in from its front panel port AND IP lookup is done on ACI.
  (which means IP addr is learned only when a leaf handles L3 traffic)
    or
- if ARP request with sender IP A comes in from its front panel port. (regardless of ARP Flooding setup)

**Remote Endpoint (MAC)**
A leaf learns MAC A as remote when L2 traffic with src MAC A comes in from SPINE.

**Remote Endpoint (/32 host IP)**
A leaf learns IP A as remote when L3 traffic with src IP A comes in from SPINE.

# Protocol inside fabric – COOP

## COOP (End Point Learning on Spine)

**SPINEs do NOT learn EP from data plane like LEAF**

**SPINEs receive all EP data from Leafs**
1. LEAF learns EP (either MAC or/and IP) as local
2. LEAF reports local EP to Spine via COOP process
3. SPINE stores these in COOP DB and synchronize with other SPINEs

**What is the purpose of COOP?**
When Leaf doesn't know dst EP, LEAF can forward packet to Spine in order to let Spine decide where to send. This behavior is called Spine-Proxy.

**Note :**
- Normally SPINE doesn't push COOP DB entries to each LEAF. It just receives and stores. The exception is for bounce entries.

- Remote Endpoints are stored on each Leaf nodes as cache. This is not reported to Spine COOP.

COOP Table
MAC/IP A -> LEAF 1
MAC/IP B -> LEAF 2
MAC/IP C -> LEAF 3

I know EP B
I know EP C
I know EP A

LEAF 1   LEAF 2   LEAF 3

MAC A IP A   MAC B IP B   MAC C IP C

https://www.ciscolive.com/c/dam/r/ciscolive/emea/docs/2019/pdf/BRKACI-3545.pdf

COOP table
| MAC | IP | VTEP | |
|---|---|---|---|
| 0000.0000.0001 | 10.10.10.1 | 1.1.1.1 | (Leaf 1) |
| 0000.0000.0020 | 10.20.10.2 | 2.2.2.2 | (Leaf 2) |
| 0000.0000.0030 | 10.30.30.1 | 3.3.3.3 | (Leaf 3) |
| 0000.0000.0021 | 10.10.10.1 | 2.2.2.2 | (Leaf 2) |
| 0000.0000.0040 | 10.40.10.1 | 4.4.4.4 | (Leaf 4) |
| 0000.0000.0041 | 10.40.10.2 | 4.4.4.4 | (Leaf 4) |

Spine switch 1

Spine switch 2

Infra VRF

L3 link1

L3 link2

L3 link4

L3 link3

Leaf switch 1

Leaf switch 2

Leaf switch 3

Leaf switch 4

Lo0
VTEP IP 1.1.1.1

Lo0
VTEP IP 2.2.2.2

Lo0
VTEP IP 3.3.3.3

Lo0
VTEP IP 4.4.4.4

EP table
| MAC | IP | VTEP | |
|---|---|---|---|
| 0000.0000.0001 | 10.10.10.1 | local (Eth 1/20) | |
| 0000.0000.0020 | 10.20.10.2 | 2.2.2.2 | (Leaf 2) |
| 0000.0000.0030 | 10.30.30.1 | 3.3.3.3 | (Leaf 3) |

EP table
| MAC | IP | VTEP |
|---|---|---|
| 0000.0000.0030 | 10.10.10.1 | local (Eth 1/2) |

MAC 0000.0000.0001

MAC 0000.0000.0020
MAC 0000.0000.0021

MAC 0000.0000.0030

MAC 0000.0000.0040
MAC 0000.0000.0041

```
apic1# fabric 2101 show coop internal info repo ep | egrep -i "mac|real|-"
--------------------------------------------------
EP mac :   00:50:56:8A:F8:32
MAC Tunnel  : 10.2.176.67
Ep vpc-id : 685
Ep vpc virtual switch-id : 10.2.176.67
Real IPv4 EP : 10.208.12.200
MAC Tunnel  : 10.2.176.67
--------------------------------------------------
EP mac :   B4:96:91:89:16:5F
MAC Tunnel  : 10.2.8.66
Real IPv4 EP : 10.210.12.10
MAC Tunnel  : 10.2.8.66
```

```
apic1# fabric 2101 show coop internal info ip-db
---------------------------------------------------
 Node 2101 (Spine2101)
---------------------------------------------------

IP address : 10.208.12.1
Vrf : 2686976
Flags : 0
EP vrf vnid : 2686976
EP IP :  10.208.12.1
Publisher Id : 10.2.8.66
Record timestamp : 06 10 2021 10:34:18 93121693
Publish timestamp : 06 10 2021 10:34:18 95126317
Seq No: 0
Remote publish timestamp: 01 01 1970 10:00:00 0
URIB Tunnel Info
Num tunnels : 1
        Tunnel address : 10.2.8.66
        Tunnel ref count : 1
```

```
apic1# fabric 2101 show coop internal info repo ep | egrep -i "mac|real|-"
--------------------------------------------
EP mac :  00:50:56:8A:F8:32
MAC Tunnel  : 10.2.176.67
Ep vpc-id : 685
Ep vpc virtual switch-id : 10.2.176.67           apic1# fabric 2101 show coop internal info ip-db
Real IPv4 EP : 10.208.12.200                     --------------------------------------------------
MAC Tunnel  : 10.2.176.67                          Node 2101 (Spine2101)
--------------------------------------------       --------------------------------------------------
EP mac :  B4:96:91:89:16:5F
MAC Tunnel  : 10.2.8.66
Real IPv4 EP : 10.210.12.10                       IP address : 10.208.12.1
MAC Tunnel  : 10.2.8.66                            Vrf : 2686976
                                                   Flags : 0
                                                   EP vrf vnid : 2686976
                                                   EP IP :  10.208.12.1
                                                   Publisher Id : 10.2.8.66
                                                   Record timestamp : 06 10 2021 10:34:18 93121693
                                                   Publish timestamp : 06 10 2021 10:34:18 95126317
                                                   Seq No: 0
                                                   Remote publish timestamp: 01 01 1970 10:00:00 0
                                                   URIB Tunnel Info
                                                   Num tunnels : 1
                                                           Tunnel address : 10.2.8.66
                                                           Tunnel ref count : 1
```

# ACI Integration with WAN at Scale
## 'Project GOLF' Overview

- Physical topology – Clos network
- Forwarding in Clos network
- Protocols in ACI - IS-IS
- Traditional service provider networks vs ACI
- Undelay and Overlay
- Endpoints, learning in the ACI fabric
- Protocols in ACI - COOP
- Bridge domains, VXLAN
- Pervasive (Anycast) gateway, VRFs

**Oversimplified**

# Bridge Domains



**Mapping:**

Eth 1/20 VLAN 20 –
Bridge Domain 1

Eth 1/7 VLAN 50 –
Bridge Domain 2

Spine switch 1   Spine switch 2

Infra VRF
**L3 links**
Routing: IS-IS

L3 link1   L3 link2   L3 link4   L3 link3

**Mapping:**

Eth 1/1  VLAN 20 –
Bridge Domain 1

Eth1/22 VLAN 40 –
Bridge domain 2

**Lo0**
**VTEP IP 1.1.1.1**
**Leaf switch1**

**Leaf switch2**
Lo0
**VTEP IP 2.2.2.2**

**Leaf switch3**
Lo0
**VTEP IP 3.3.3.3**

Lo0
**VTEP IP 4.4.4.4**
**Leaf switch4**

MAC
0000.0000.0100
IP 20.20.20.1

MAC
0000.0000.0001
IP 10.10.10.1

MAC
0000.0000.0020
IP 20.20.20.2

MAC
0000.0000.0021
IP 10.10.10.2

MAC 0000.0000.0030
IP 10.10.10.3

MAC
0000.0000.0040
IP 20.20.20.4

MAC
0000.0000.0041
IP 10.10.10.4

Bridge Domain 1  - Subnet 10.10.10.0/24

Bridge Domain 2  - Subnet 20.20.20.0/24

# Virtual Extensible LAN  - VXLAN

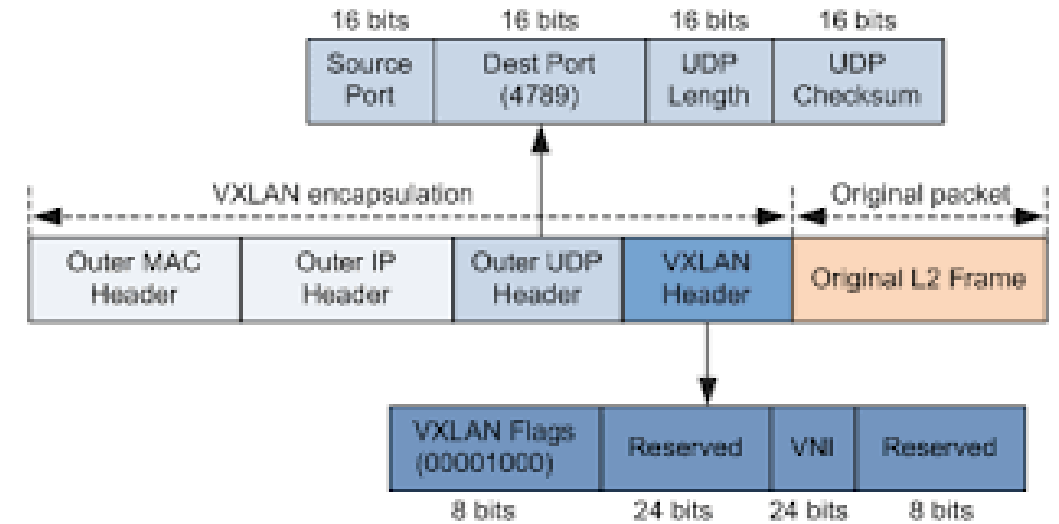- Main purpose to deliver L2 frames over L3 networks
- Standard-based (but Cisco uses proprietary iVXLAN)
- Uses the MAC-in-UDP.  UDP port 4789
- Requires MTU to be at least 1574 bytes, standard setting is 9000 bytes
- It uses a VLAN-like encapsulation, but instead of 12-bit VLAN ID
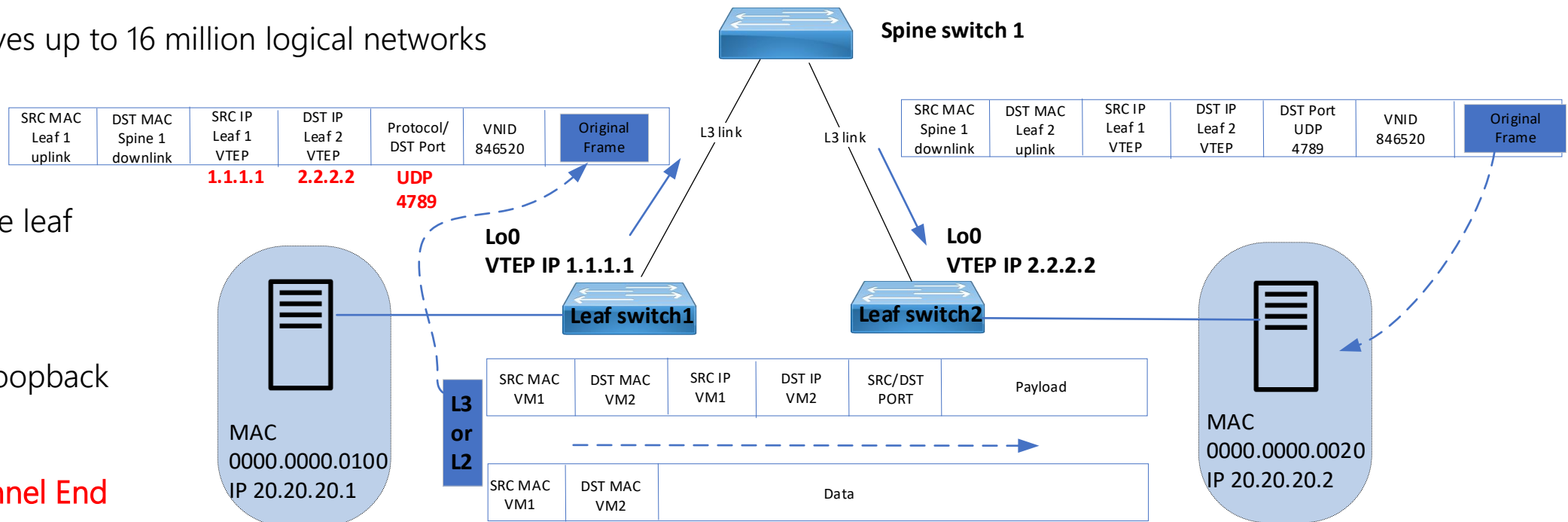
Uses 24-bit VNID – Virtual Network ID

- 24-bit VNID gives up to 16 million logical networks

| | 16 bits | 16 bits | 16 bits | 16 bits |
|---|---|---|---|---|
| | Source Port | Dest Port (4789) | UDP Length | UDP Checksum |

VXLAN encapsulation / Original packet

| Outer MAC Header | Outer IP Header | Outer UDP Header | VXLAN Header | Original L2 Frame |
|---|---|---|---|---|

| VXLAN Flags (00001000) | Reserved | VNI | Reserved |
|---|---|---|---|
| 8 bits | 24 bits | 24 bits | 8 bits |

**Spine switch 1**

| SRC MAC Leaf 1 uplink | DST MAC Spine 1 downlink | SRC IP Leaf 1 VTEP 1.1.1.1 | DST IP Leaf 2 VTEP 2.2.2.2 | Protocol/ DST Port UDP 4789 | VNID 846520 | Original Frame |
|---|---|---|---|---|---|---|

| SRC MAC Spine 1 downlink | DST MAC Leaf 2 uplink | SRC IP Leaf 1 VTEP | DST IP Leaf 2 VTEP | DST Port UDP 4789 | VNID 846520 | Original Frame |
|---|---|---|---|---|---|---|

L3 link        L3 link

**Lo0
VTEP IP 1.1.1.1**

**Lo0
VTEP IP 2.2.2.2**

Source IP is source leaf Loopback (VTEP)

**Leaf switch1**        **Leaf switch2**

Destination IP is destination leaf Loopback (VTEP)

VTEP - Virtual Tunnel End Point

MAC
0000.0000.0100
IP 20.20.20.1

MAC
0000.0000.0020
IP 20.20.20.2

**L3 or L2**

| SRC MAC VM1 | DST MAC VM2 | SRC IP VM1 | DST IP VM2 | SRC/DST PORT | Payload |
|---|---|---|---|---|---|

| SRC MAC VM1 | DST MAC VM2 | Data |
|---|---|---|

# Forwarding in ACI



Source LEAF knows the destination ( on the remote LEAF )

VxLAN has VRF or BD VNID

Anycast TEP is used for proxy not used in this scenario

VRF overlay-1

Anycast TEP

3 Forward based on outer IP (dIPo)

| dMACo | sMACo | sIPo (TEP1) | dIPo (TEP2) | VxLAN | dMACi | sMACi | sIPi | dIPi |

2 Add VxLAN header

TEP1    TEP2    TEP3

1 Send to LEAF2 (TEP2)

| dMAC | sMAC | sIP | dIP |

4 Decapsulate VxLAN

VRF1    VRF1    VRF1

| dMAC | sMAC | sIP | dIP |

BD1    BD1    BD2    5 Send to EP    L3OUT

Remote EndPoint

Packet goes through infra network with additional encaps (iVxLAN)

EPG1    EPG2    EPG 2    EPG 3    EPG 4

EP1-1    EP2-1    EP2-2    EP3-1  EP3-2    EP4-1    WAN

# Design option 1 – ACI as a big L2 switch



Spine switch 1

Spine switch 2

Infra VRF
**L3 links**
Routing: IS-IS

L3 link1

L3 link2

L3 link4

L3 link3

Lo0
VTEP IP 1.1.1.1

Leaf switch1

Leaf switch2

Leaf switch3

Lo0
VTEP IP 4.4.4.4

Leaf switch4

Lo0
VTEP IP 2.2.2.2

Lo0
VTEP IP 3.3.3.3

Default GW for
BD1

Default GW for
BD2

MAC
0000.0000.0100
IP 20.20.20.1

MAC
0000.0000.0001
IP 10.10.10.1

MAC
0000.0000.0020
IP 20.20.20.2

MAC
0000.0000.0021
IP 10.10.10.2

MAC 0000.0000.0030
IP 10.10.10.3

L3 Network

Bridge Domain 1  - Subnet 10.10.10.0/24

Bridge Domain 2  - Subnet 20.20.20.0/24

- Physical topology – Clos network
- Forwarding in Clod network
- Traditional service provider networks vs ACI
- Undelay and Overlay
- Endpoints, learning in the ACI fabric
- COOP
- Bridge domains, VXLAN
- Pervasive (Anycast) gateway, VRFs

Oversimplified

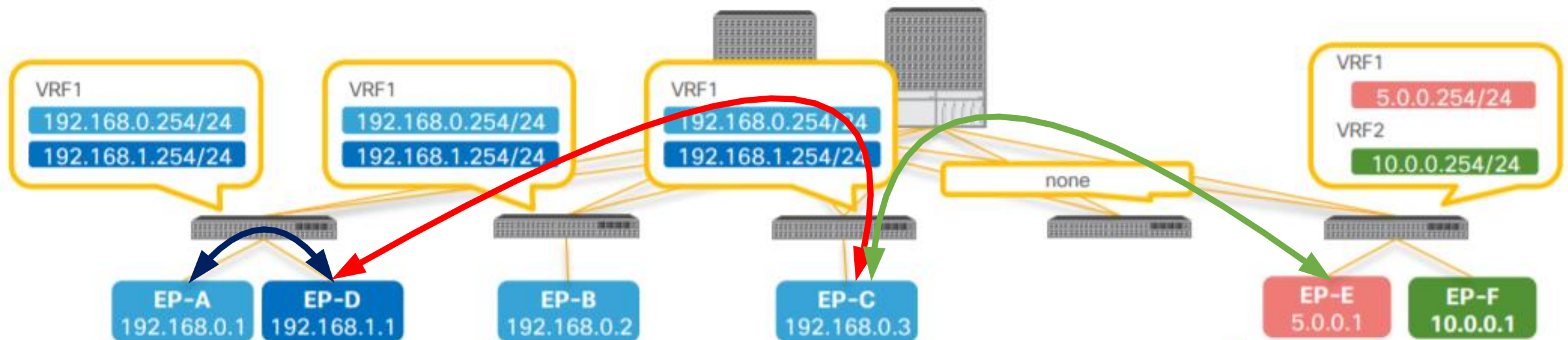# Traditional network – default gateway with HSRP/VRRP

# HSRP/VRRP - Traffic flow between subnets

# Anycast (Pervasive) Gateway

- **Every leaf switch** is configured as a default gateway for all connected L3 endpoint subnets
- SVI with same IP and MAC address on all leaf switches
- No concept as active/standby – all leafs are 'active' default gateway for **their** connected endpoints
- No central default gateway
- Endpoint send traffic to the local leaf - default GW, the leaf then sends traffic to remote leaf directly using VXLAN
- Traffic goes directly between every leaf (via Spines obviously, as they are physically connected via Spines)
- In ACI it is called Pervasive Gateway, in all other vendors implementations it's called Distributed Anycast Gateway

# Configuring Pervasive Gateway

# Pervasive Gateway(BD SVI)

**Tenant TK**
- Quick Start
- Tenant TK
  - Application Profiles
  - Networking
    - Bridge Domains
      - BD1
        - DHCP Relay Labels
        - L4-L7 Service Parameters
        - Subnets
          - 192.168.0.254/24
          - 192.168.1.254/24
        - ND Proxy Subnets
      - BD2
      - BD3
      - BD_SG_PBR1

**Subnet - 192.168.0.254/24**

Properties

IP Address: 192.168.0.254/24
Description: optional

Treat as virtual IP address: ☐
Make this IP address primary: ☐
Scope: ☑ Private to VRF
☐ Advertised Externally
☐ Shared between VRFs
Subnet Control: ☑ ☐
☐ No Default SVI Gateway
☐ Querier IP
L3 Out for Route Profile: select a value
Route Profile: select value

```
leaf1# show ip route vrf TK:VRF1

192.168.0.0/24, ubest/mbest: 1/0, attached, direct, pervasive
    *via 10.0.184.64%overlay-1, [1/0], 04:32:16, static

192.168.0.254/32, ubest/mbest: 1/0, attached
    *via 192.168.0.254, vlan10, [1/0], 04:32:16, local, local
```

Pervasive route

Pervasive SVI

BD SVI with PI-VLAN

## What is pervasive GW for?

- To be a default GW for EPs in the Fabric
  - All EPs can have consistent gateway IP address one hop away

- To represent subnets(IP ranges) for a BD
  - ACI knows which BD may have potential hidden/silent EPs

## How is pervasive GW deployed?

- Installed as an SVI on LEAFs
  - PI-VLAN for BD is used to represent a pervasive GW SVI
  - A pervasive SVI has secondary IP when multiple pervasive GWs are configured on the same BD
    - User can choose a primary address

# Briefly about Tenants and VRFs

Constructs in ACI vs Public Cloud:

VRF = VPC or Vnet
Bridge Domain = Subnet
Tenant = Account or Customer

```
apic1# fabric 2101 show vrf
---------------------------------------------------------------
Node 2101 (Spine2101)
---------------------------------------------------------------
VRF-Name                    VRF-ID State   Reason
black-hole                       3 Up      --
management                       2 Up      --
mgmt:inb                         5 Up      --
overlay-1                        4 Up      --


apic1# fabric 2201 show vrf
---------------------------------------------------------------
Node 2201 (Leaf2201)
---------------------------------------------------------------
VRF-Name                    VRF-ID State   Reason
black-hole                       3 Up      --
common:SharedServices_VRF       14 Up      --
management                       2 Up      --
mgmt:inb                         8 Up      --
overlay-1                        4 Up      --
Tenant01:Production_VRF         23 Up      --
Tenant06:Production_VRF         24 Up      --
Tenant07:Production_VRF         26 Up      --
Tenant08:Production_VRF         27 Up      --
Tenant09:Production_VRF         18 Up      --
Tenant10:Production_VRF         21 Up      --
Tenant11:Production_VRF         17 Up      --
Tenant12:Production_VRF         19 Up      --
```
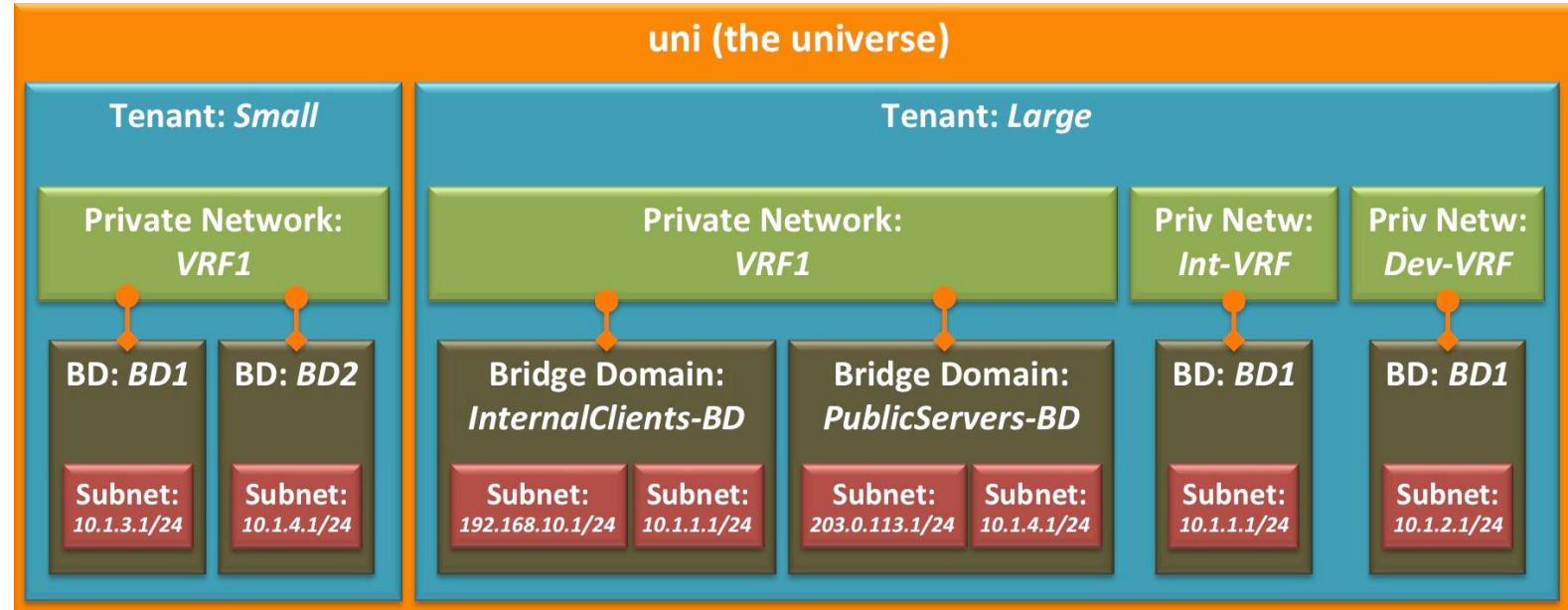


https://rednectar.files.wordpress.com/2015/05/the-universe5.jpg

<----- Note Spine switches don't have information about Tenant VRFs

VRF **overlay-1** is fabric **underlay** (please don't ask why ☺)

# Summary

- **Main definitions:** Clos fabric, VTEP, Endpoint, Bridge Domain, Endpoint tables, COOP Database, Pervasive (Anycast) Gateway, VRF, ACI Tenant, VXLAN

- **Control plane protocols in ACI**: IS-IS, COOP

- **Endpoint learning** – local and remote

- **Traffic forwarding** in ACI/EVPN, VXLAN encapsulation

Next time

- ACI External connections – L3Out

- Protocols inside the ACI fabric – MP-BGP

- Connecting multiple datacenters

Possible topics for further sessions

- Endpoint Groups and Contracts, micro segmentation

- Integration with VMware ESXi

- Policy-based routing

- ACI controllers, main UI sections

# Good reading

- ACI Fabric Endpoint Learning White Paper

- Mastering ACI Forwarding Behaviour - A day in the life of a packet

  https://www.ciscolive.com/c/dam/r/ciscolive/emea/docs/2019/pdf/BRKACI-3545.pdf

- Virtual Port Channel (vPC) in ACI

https://www.cisco.com/c/dam/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/aci-guide-vpc.pdf

- Why You shouldn't Think about Fabric Extenders (FEX) along with Cisco ACI anymore?

https://www.linkedin.com/pulse/why-you-shouldnt-think-fabric-extenders-fex-along-cisco-vahid-nazari

Thanks!