

The Cisco Live! logo features the word "CISCO" in a dark blue, sans-serif font, followed by "Live!" in a dark blue, cursive script font. The background of the entire image is a vibrant, multi-colored abstract pattern of overlapping, wavy lines and geometric shapes, transitioning from dark blue on the left to bright yellow and white in the center, and then to various shades of blue and green on the right.

CISCO *Live!*

Let's go



The bridge to possible

Cisco ACI Multi-Pod

Design and Deployment

John Weston, Technical Marketing Engineer,
Data Center Networking

CISCO *Live!*

BRKDCN-2949

Session Objectives

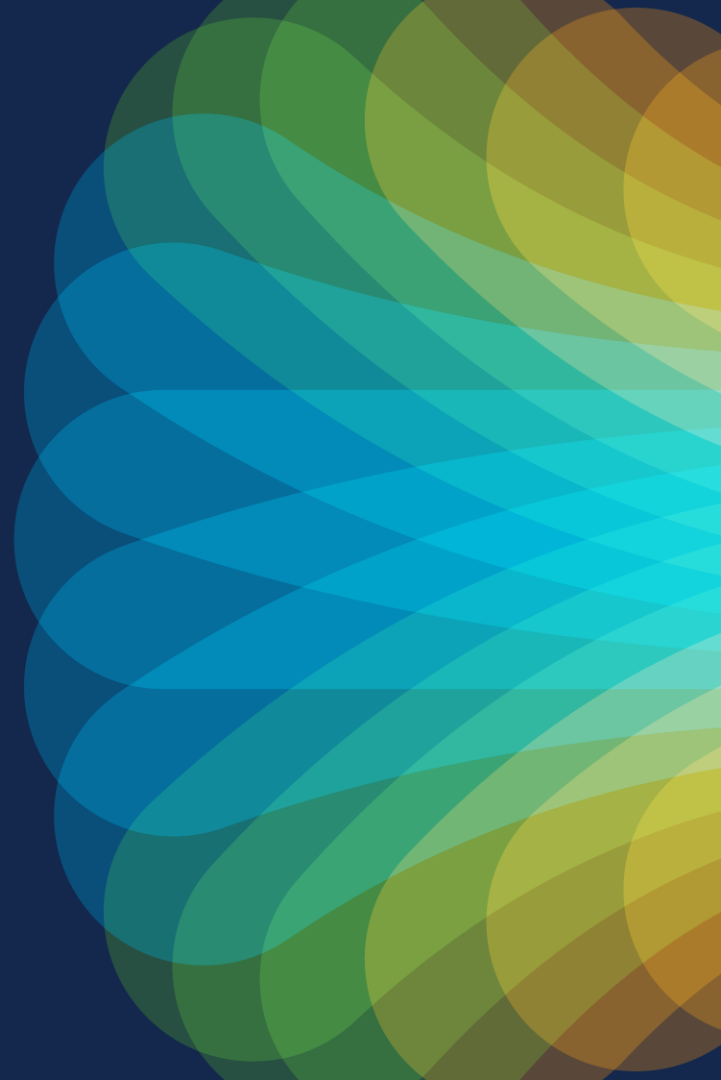


- At the end of the session, the participants should be able to:
 - Articulate the different deployment options to interconnect Cisco ACI networks (Multi-Pod and Multi-Site) and when to choose one vs. the other
 - Understand the functionalities and specific design considerations associated to the ACI Multi-Pod architecture
- Initial assumption:
 - The audience already has a good knowledge of ACI main concepts (Tenant, BD, EPG, L3Out, etc.)

Agenda

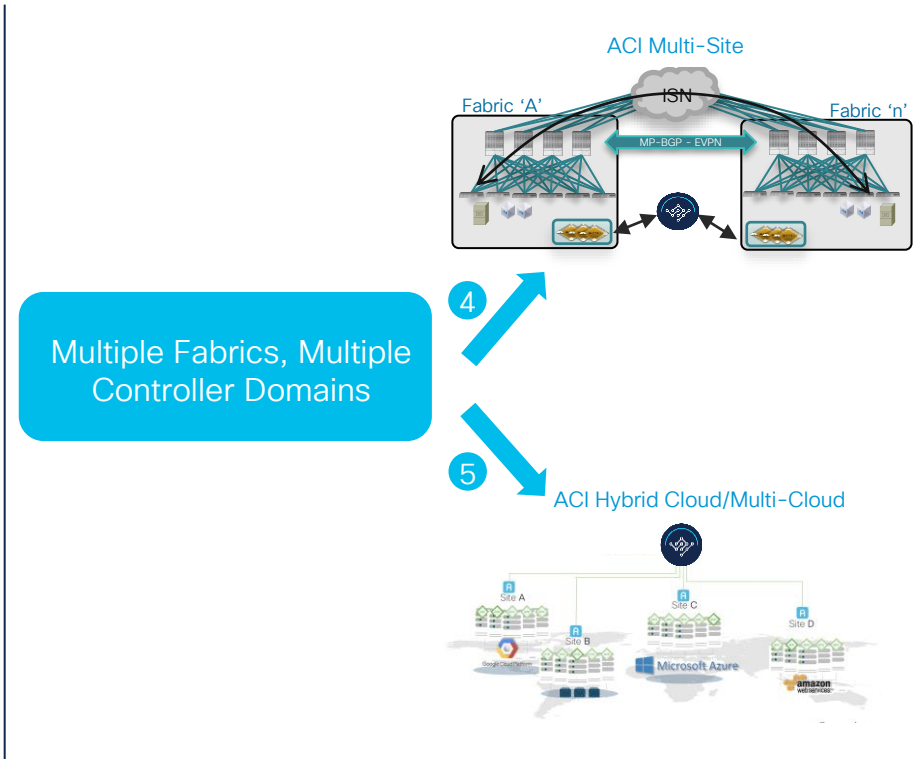
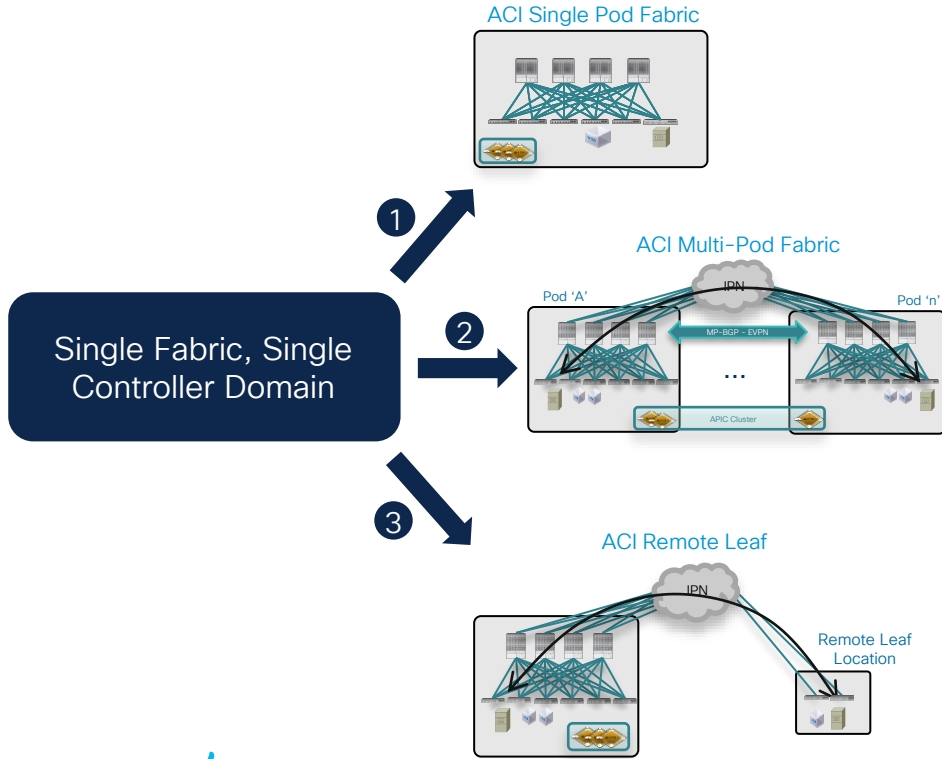
- Introduction
- Supported Topologies
- APIC Cluster Deployment
- Inter-Pod Connectivity
- Control and Data Planes
- Connecting to External Networks
- Network Services Integration
- Multi-Pod and Remote Leaf

Introduction



ACI Architectural Options

Fabric and Policy Domain Evolution



Multi-Pod or Multi-Site?

Where to Get More Information

- ACI Multi-Pod White Paper

<https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/white-paper-c11-737855.html>

- ACI Multi-Site White Paper

<https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/white-paper-c11-739609.html>

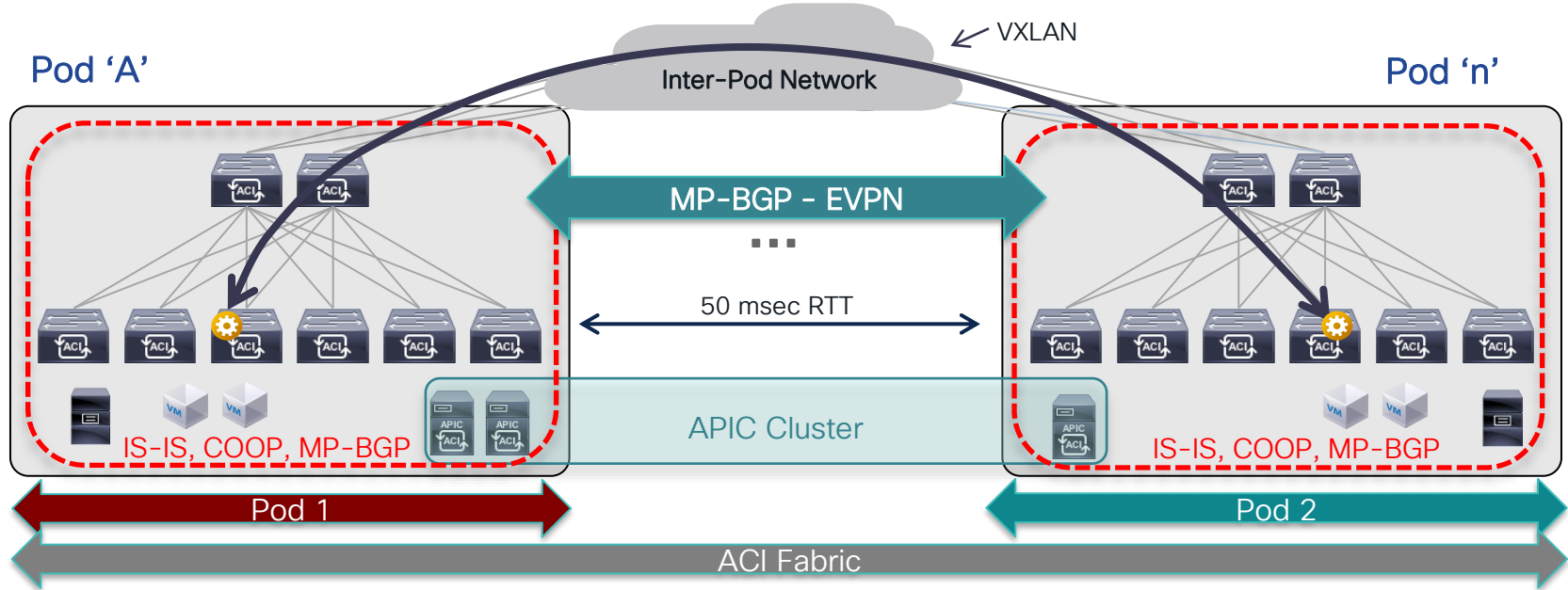
- ACI Multi-Site Cisco Live 2020 Digital Breakout Session

<https://www.ciscolive.com/on-demand/on-demand-library.html?search=ardica&search=ardica#/video/1636411349156002rlx8>

Want to know how to provision Multi-Pod and Multi-Site from scratch? Come to BRKDCN-2919 (Wed @ 10.30 am)

ACI Multi-Pod

The Ideal Architecture for Active/Active DC Deployments



- Multiple ACI Pods connected by an IP Inter-Pod L3 network, each Pod consists of leaf and spine nodes
- Managed by a single APIC Cluster
- Single Management and Policy Domain

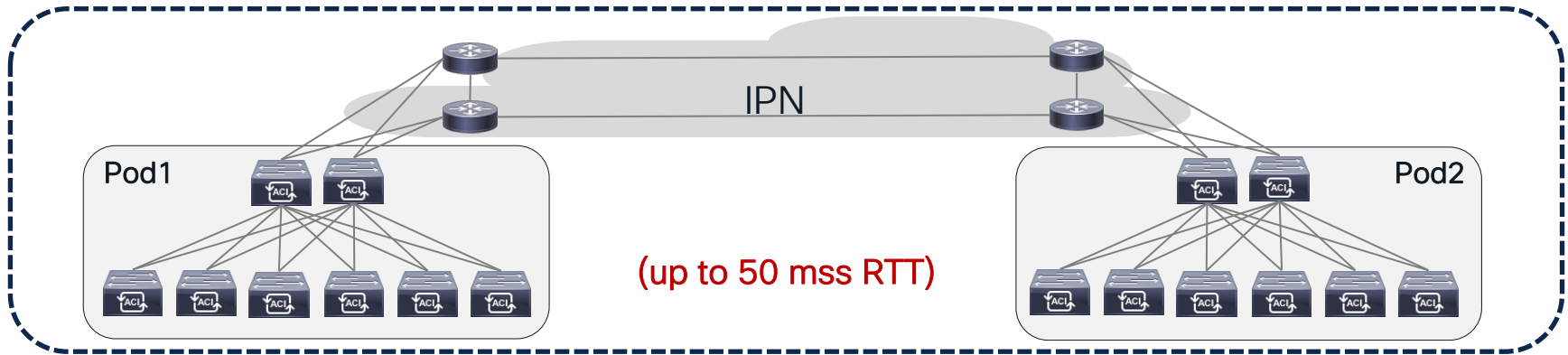
- Forwarding control plane (IS-IS, COOP) fault isolation
- Data Plane VXLAN encapsulation between Pods
- End-to-end policy enforcement

Supported Topologies

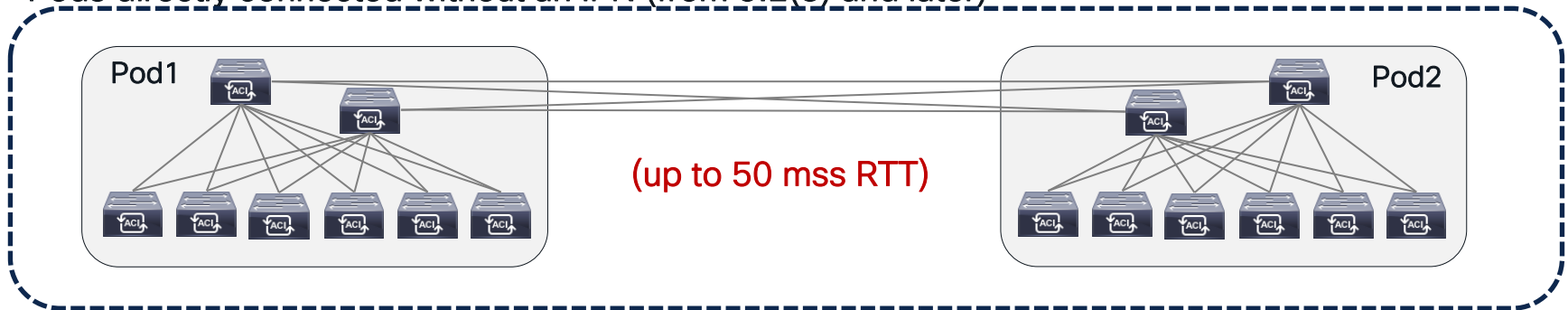


Multi-Pod Supported Topologies

Pods connected via an Inter-Pod Network (IPN)

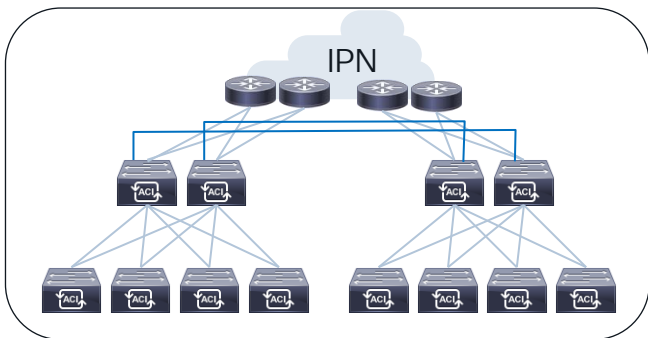
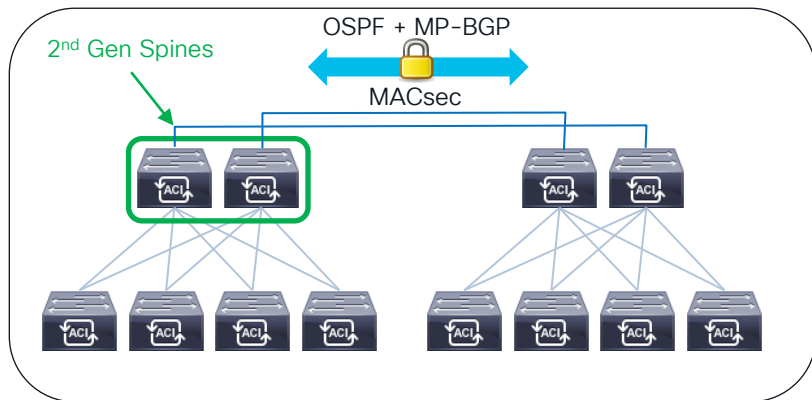


Pods directly connected without an IPN (from 5.2(3) and later)



Multi-Pod Spines Back-to-Back

Guidelines and Restrictions



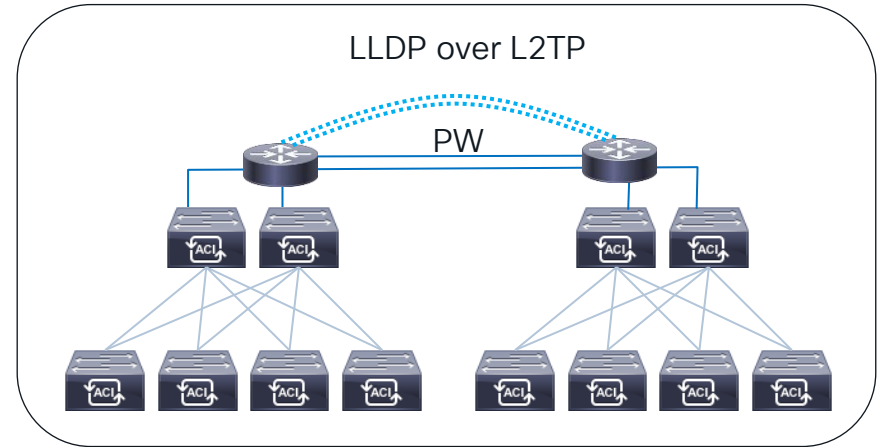
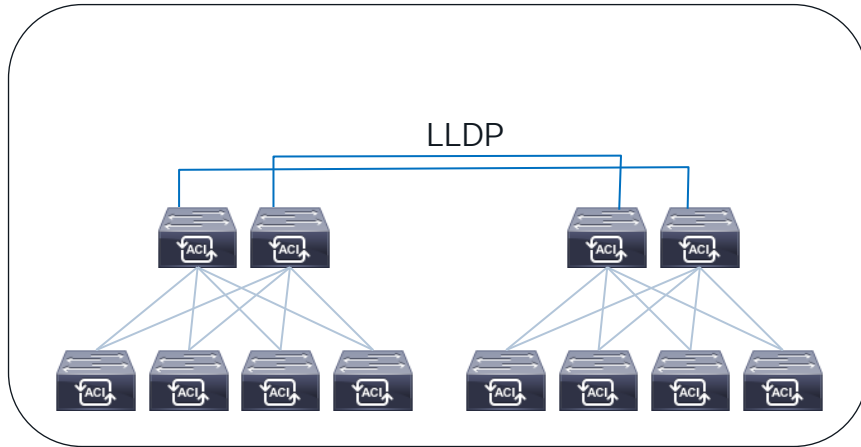
Back-to-Back + IPN (Migration Only)

- Support is limited to a topology with 2 Pods leveraging 2nd generation spines only
- OSPF underlay peering, MP-BGP overlay peering between the spines in separate Pods
 - No need for PIM-Bidir (spines do not run PIM)
- MACsec encryption supported across Pods
- Not compatible functions (any feature requiring an IPN connection)
 - ACI Multi-Site
 - Remote Leaf
 - GOLF
 - Public Cloud Connections
 - APIC connectivity via L3 network
- Back-to-Back + IPN only supported for migration purposes (migration is disruptive)

Multi-Pod Spines Back-to-Back

Supported Topologies

- Back-to-back spine connectivity must be point-to-point (physical or logical)
- Spines discover back-to-back connections via LLDP
- Links can be directly connected or must support tunneling of LLDP packets



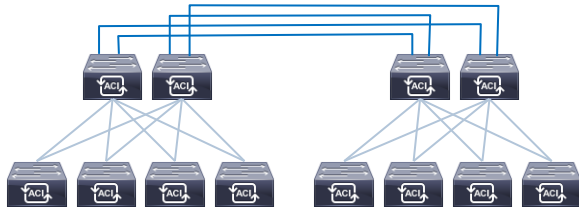
Multi-Pod Spines Back-to-Back

Supported Topologies

It is not mandatory for all spines in a Pod to connect to all the spines in the other Pod, the design decision must be made based on resiliency/bandwidth considerations

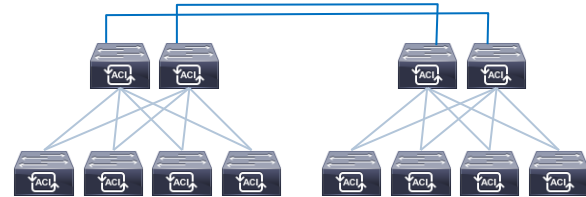
Recommended

Full mesh between spines



Supported

Partial mesh between spines



ACI Multi-Pod

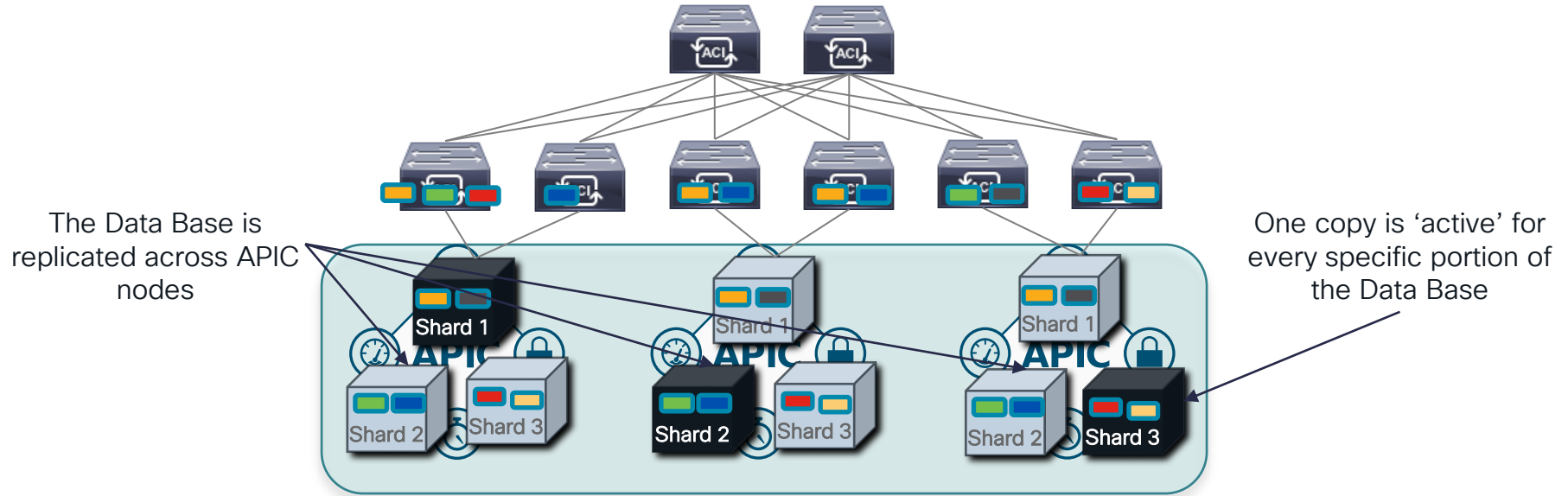
SW/HW Support and Scalability Values



- All existing Nexus 9000 HW supported as leaf and spine nodes*
- Maximum number of supported ACI leaf nodes (across all Pods)
 - Up to 80 leaf nodes supported with a **3 node** APIC cluster
 - 200 leaf nodes (across Pods) with a **4 node** APIC cluster (from ACI release 4.1)
 - 300 leaf nodes (across Pods) with a **5 node** APIC Cluster
 - 400 leaf nodes (across Pods) with a **7 node** APIC Cluster (from ACI release 2.2(2e))
 - 500 leaf nodes (across Pods) with a **7 node** APIC Cluster (from ACI release 4.2(4))
 - Maximum 400 leaf nodes per Pod (from ACI release 4.2(4))
 - Up to 6 spines per Pod, 50 spines per Fabric (from ACI release 6.0(1))
- Maximum number of supported Pods
 - 4 in 2.0(1)/2.0(2) releases
 - 6 in 2.1(1) release
 - 10 in 2.2(2e) release
 - 12 in 3.0(1) release
 - 25 in 6.0(1) release

APIC Cluster Deployment Considerations

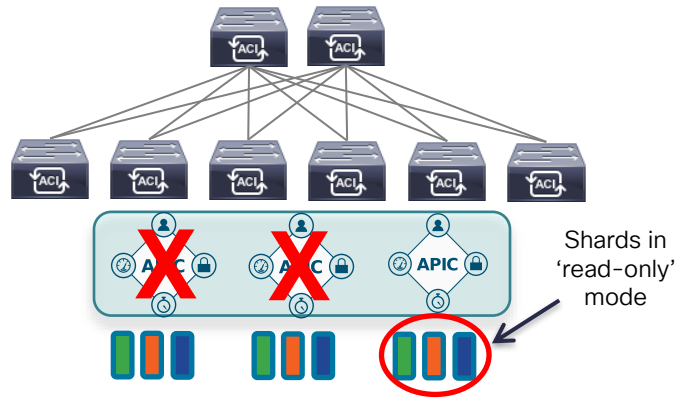
APIC - Distributed Multi-Active Data Base



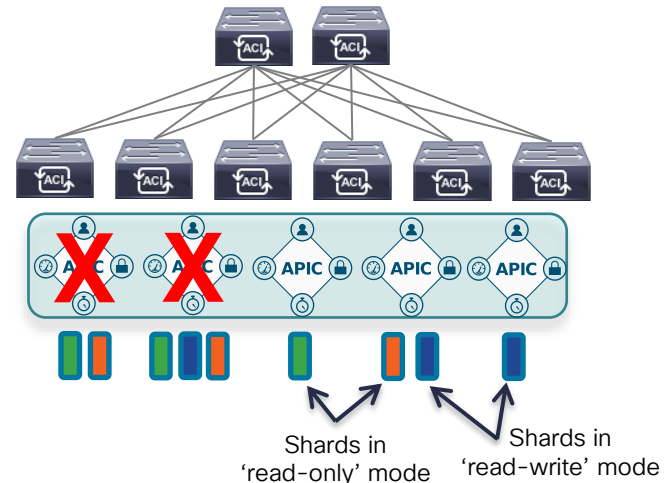
- Processes are active on all nodes (not active/standby)
- The Data Base is distributed as active + 2 backup instances (shards) for every attribute

APIC Cluster Deployment Considerations

Single Pod Scenario



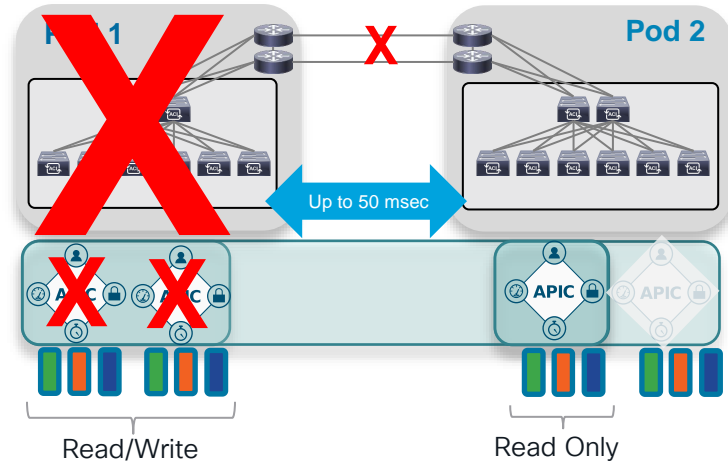
- APIC will allow read-only access to the DB when only one node remains active (standard DB quorum)
- Hard failure of two nodes cause all shards to be in 'read-only' mode (of course reboot etc. heals the cluster after APIC nodes are up)



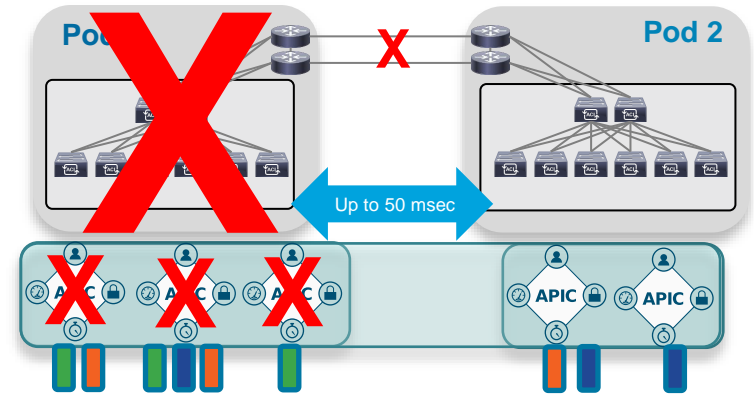
- Additional APIC will increase the system scale (up to 7* nodes supported) but does not add more redundancy
- Hard failure of two nodes would cause inconsistent behaviour across shards (some will be in 'read-only' mode, some in 'read-write' mode)

APIC Cluster Deployment Considerations

Multi-Pod – 2 Pods Scenario



- **Pod isolation scenario:** changes still possible on APIC nodes in Pod1 but not in Pod2
- **Pod hard failure scenario:** recommendation is to activate a standby node to make the cluster fully functional again

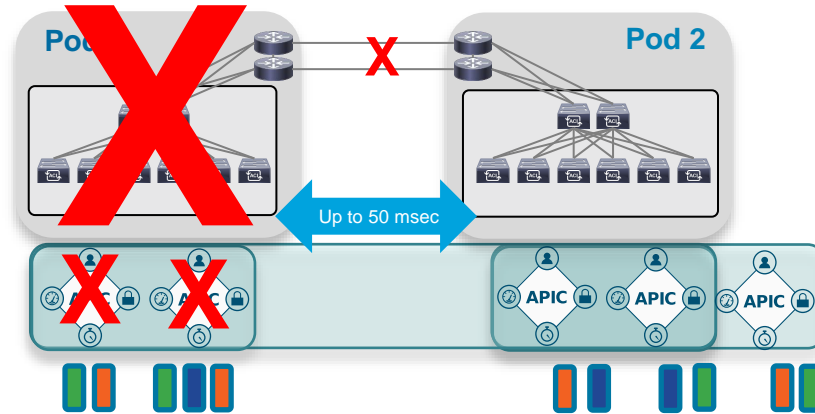


- **Pod isolation scenario:** same considerations as with single Pod (different behaviour across shards)
- **Pod hard failure scenario:** may cause the loss of information for the shards replicated across APIC nodes in the failed Pod

Possible to restore the whole fabric state to the latest taken configuration snapshot ('ID Recovery' procedure – **needs BU and TAC involvement**)

APIC Cluster Deployment Considerations

What about a 4 Nodes APIC Cluster?



- Intermediate scalability values compared to a 3 or 5 nodes cluster scenario (up to 200 leaf nodes supported)
- **Pod isolation scenario:** same considerations as with 5 nodes (different behaviour across shards)
- **Pod hard failure scenario**
 - No chance of total loss of information for any shard
 - Can bring up a standby node in the second site to regain full majority for all the shards

APIC Cluster Deployment Considerations

Deployment Recommendations

- **Main recommendation:** Deploy a 3-node APIC cluster when fewer than 85 leaf nodes are deployed across Pods
- From 4.1(1) can deploy 4 nodes if the scalability requirements are met
- When 5 (or 7) nodes are really needed for scalability reasons, follow the rule of thumb of never placing more than two APIC nodes in the same Pod (when possible):

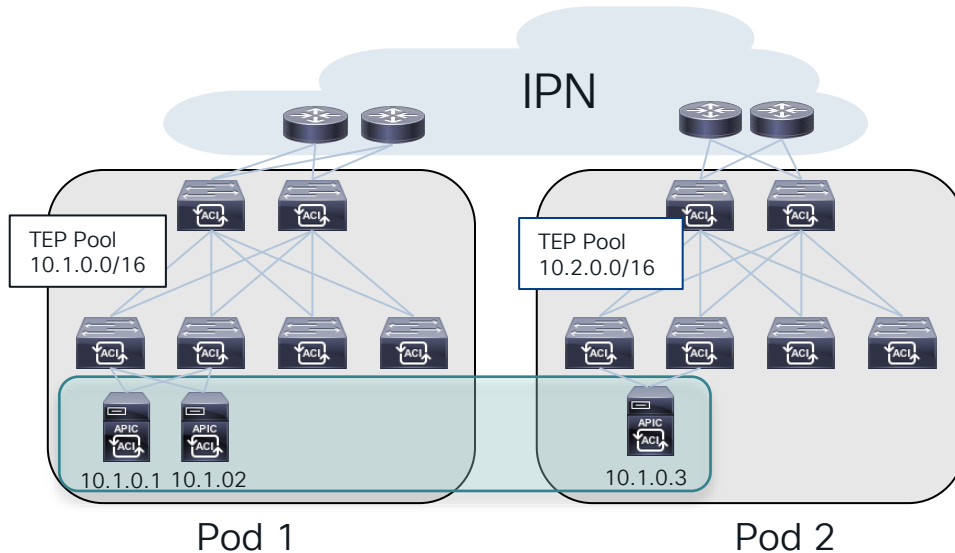
	Pod1	Pod2	Pod3	Pod4	Pod5	Pod6
2 Pods*						
3 Pods						
4 Pods						
5 Pods						
6+ Pods						

APIC Connectivity over L3 Network



APIC Connectivity Options

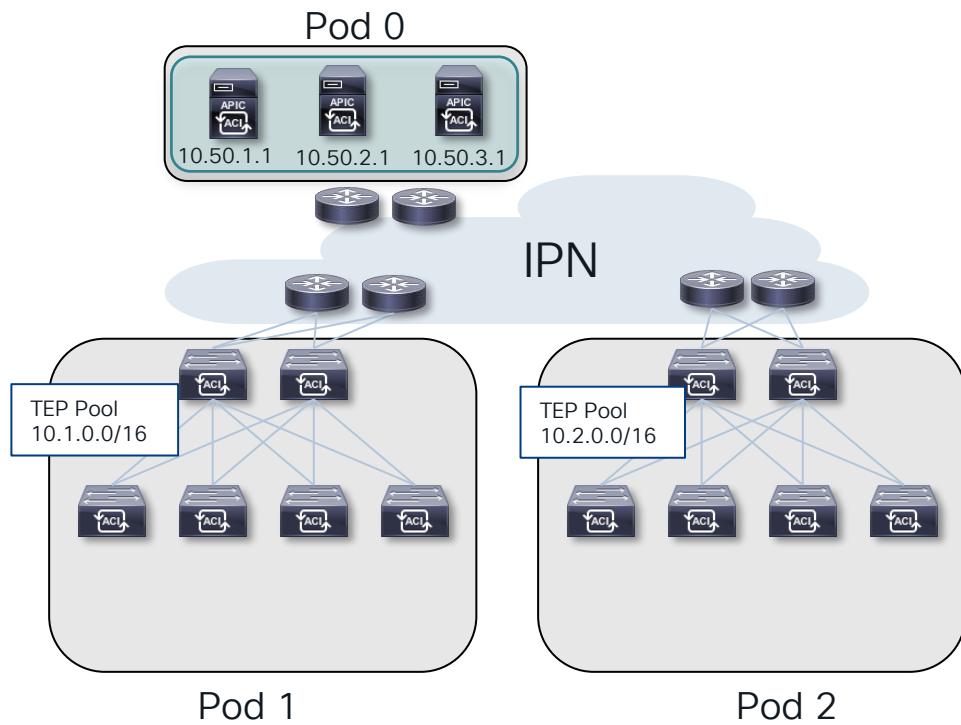
APIC Cluster directly connected to fabric



- APICs can be placed in any pod
- APIC fabric IP addresses are always assigned from pod 1 TEP pool
- Recommended to distribute APICs across pods so loss of a pod does not bring down the entire cluster

APIC Connectivity Options

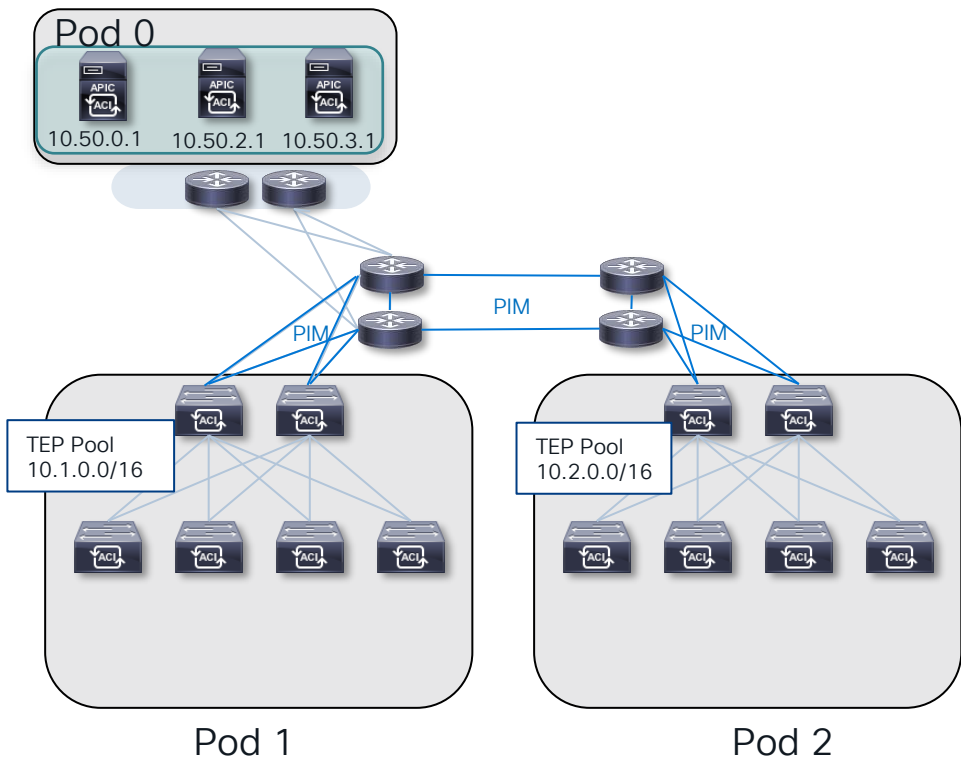
APIC cluster connected over L3 Network



- APICs do not need to connect to a leaf switch. Can be connected via the IPN
- APICs connected over the L3 network are considered part of Pod 0
- APIC fabric IP addresses are user configurable. Not assigned from any pod TEP range
- APIC fabric IPs can be in the same or different subnet per APIC
- APICs can be geographically distributed within the Multi-Pod 50 msec distance limitation

APIC Connectivity Options

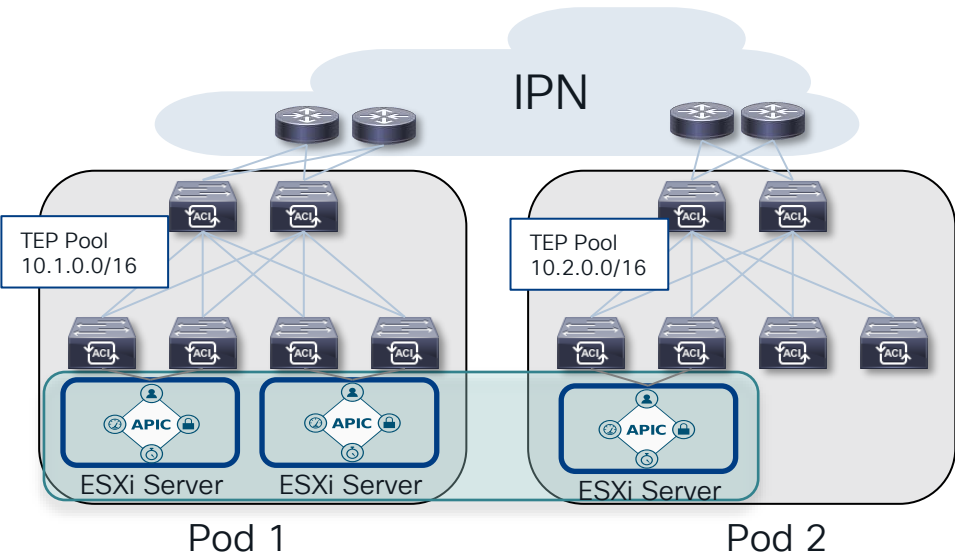
APIC cluster connected over L3 Network, IPN Multicast Requirement



- Multicast (PIM Bidir) is only required for inter-pod BUM traffic
- If APIC cluster over L3 network is managing only one pod, multicast is not required in the IPN
- If it is a Multi-Pod fabric, multicast is only required on the links interconnecting the pods

APIC Connectivity Options

Virtual APIC cluster

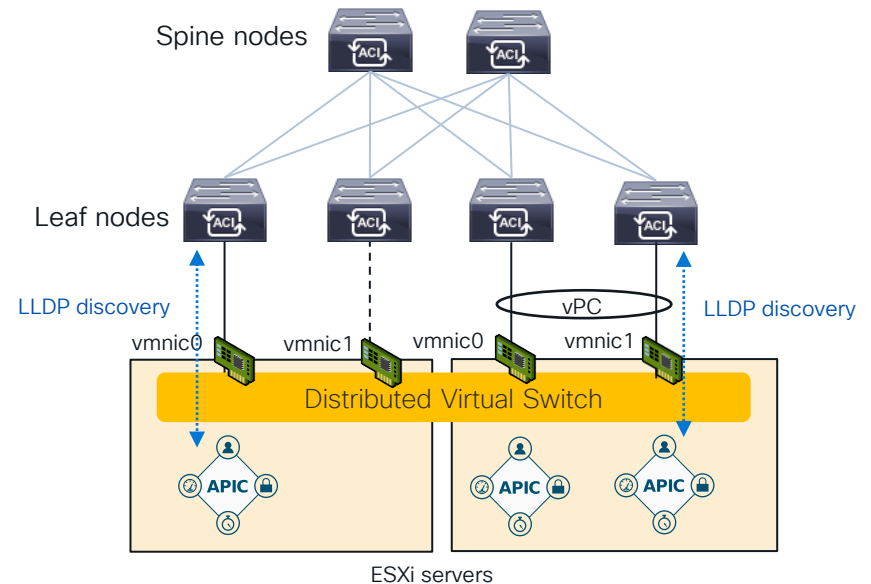
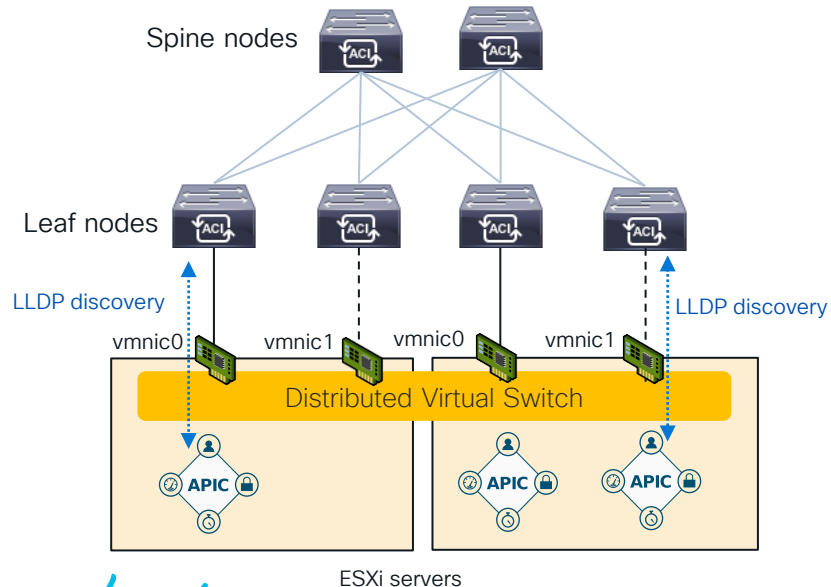


- Virtual APIC cluster (all virtual APICs)
- Runs as a VM on an ESXi hypervisor
- ESXi server directly connected to fabric
- No mixed cluster support. Must be all virtual or all physical
- Supports all types of deployments, Remote Leaf, Multi-Pod, Multi-Site.

Topology considerations for virtual APIC on ESXi

Directly Attached

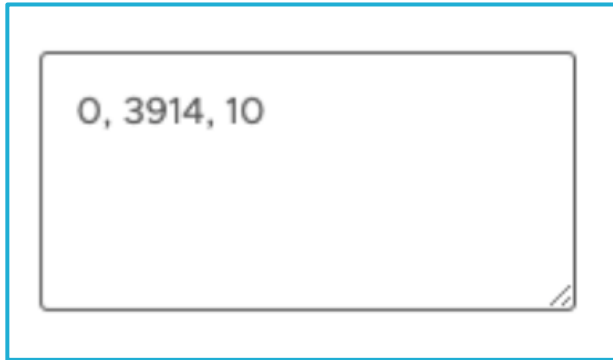
- ESXi servers need to be directly connected to ACI leaf nodes via individual links or vPC. (APIC1 must use Active-Standby instead of Active-Active with vPC)
- LLDP must be disabled on the virtual switch for LLDP discovery between leaf nodes and vAPICs.



Distributed Switch configuration

Port Group VLAN configuration

- VLAN type: VLAN Trunking
- VLAN trunk range:
 - VLAN 0 (VLAN 0 is required for APIC LLDP discovery)
 - ACI Infra VLAN (for example, 3914 is used as the default value during APIC initial setup)
 - Inband VLAN(s) (VLAN 10 in the example)

A screenshot of the "New Distributed Port Group" configuration interface. The interface is divided into two main sections: "Name and location" and "Configure settings". The "Configure settings" section is active and shows various configuration options. The "VLAN" section is expanded, showing "VLAN type" set to "VLAN trunking" and "VLAN trunk range" set to "0, 3914, 10". The "VLAN trunk range" field is highlighted with a blue box, and a blue arrow points from this box to the "0, 3914, 10" text in the previous image. At the bottom right, there are "CANCEL", "BACK", and "NEXT" buttons.

New Distributed Port Group

Configure settings

Set general properties of the new port group.

1 Name and location

2 Configure settings

3 Ready to complete

Port binding: Static binding

Port allocation: Elastic

Number of ports: 8

Network resource pool: (default)

VLAN

VLAN type: VLAN trunking

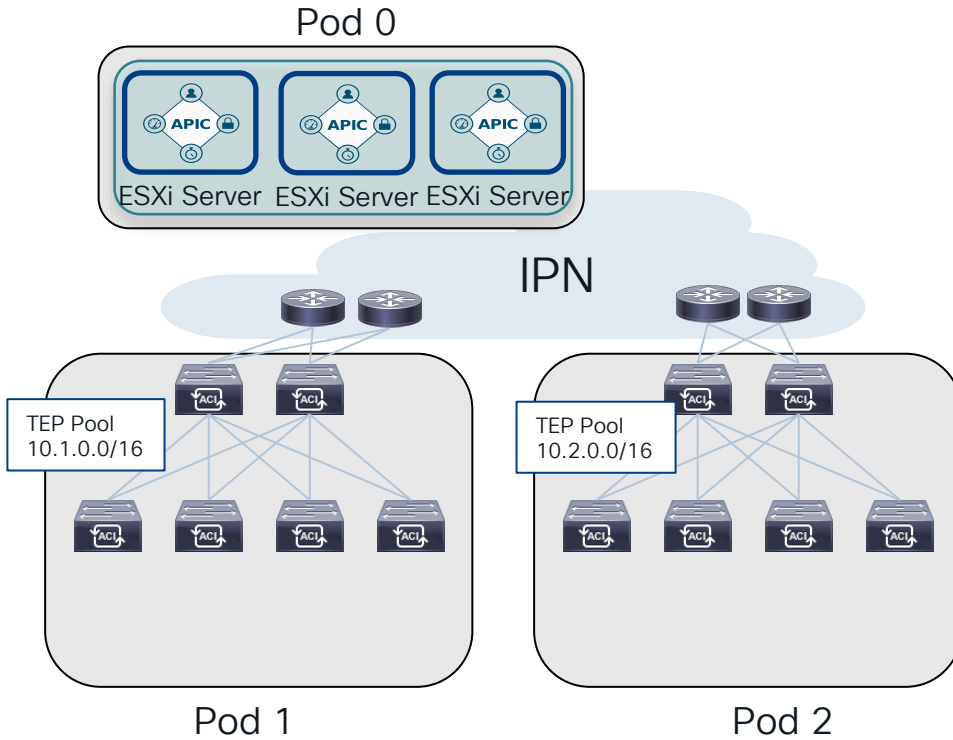
VLAN trunk range: 0, 3914, 10

Advanced

CANCEL BACK NEXT

APIC Connectivity Options

Virtual APIC cluster over L3 Network



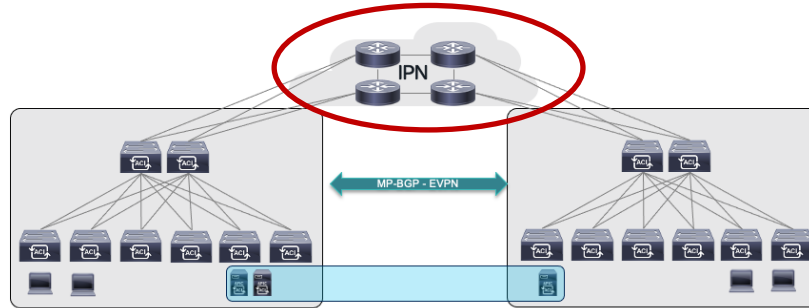
- Virtual APIC over L3 Network
- Same or different IP addresses per APIC same as physical APIC over L3 network
- Cannot mix virtual APIC over L3 Network with directly connected virtual APIC

Inter-Pod Connectivity Deployment Considerations



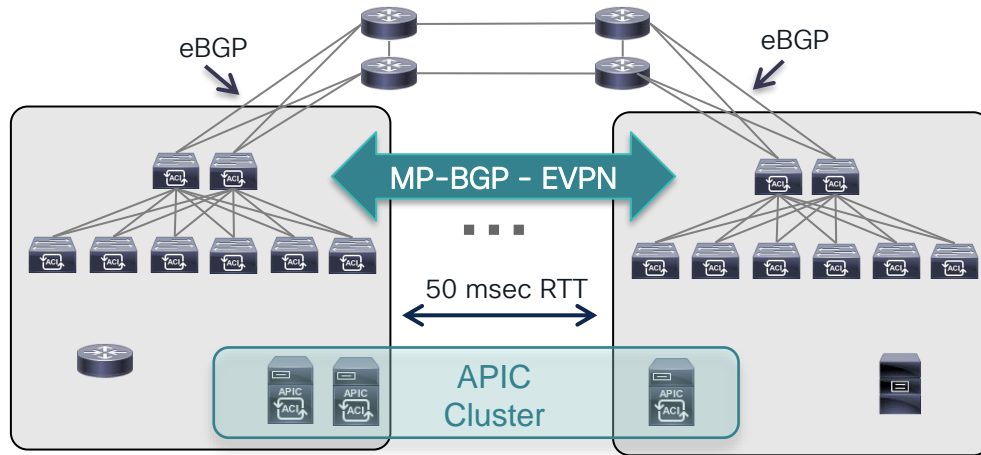
ACI Multi-Pod

Inter-Pod Network (IPN) Requirements



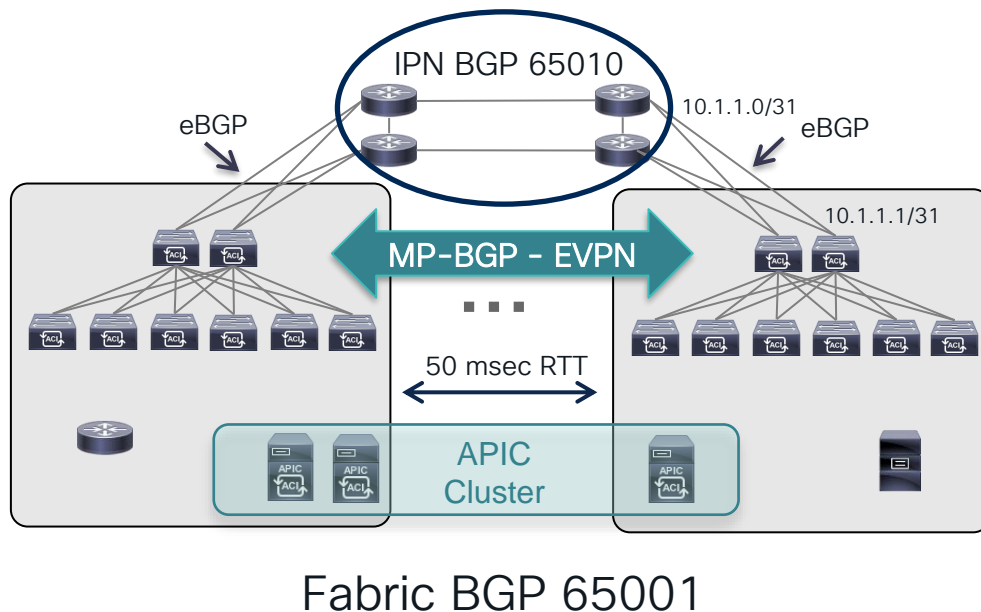
- Not managed by APIC, must be separately configured (day-0 configuration)
- IPN topology can be arbitrary, not mandatory to connect to all spine nodes
- Main requirements:
 - Multicast BiDir PIM → needed to handle Layer 2 BUM* traffic
 - OSPF or BGP to peer with the spine nodes and learn VTEP reachability
 - Increase MTU support to handle VXLAN encapsulated traffic
 - DHCP-Relay

BGP Underlay Support for IPN links



- From ACI 5.2(3) you can use either OSPF and/or BGP for IPN connectivity
- Only eBGP is supported
- Infra L3Out interfaces can be configured with OSPF, BGP, or both protocols at the same time (typically used for migration)
- When both protocols are configured, BGP routes will be preferred due to lower admin distance
- Supported with Multi-Pod, Remote Leaf, Multi-Site, and APIC over L3 Network (not supported with GOLF or cloud sites)
- Enabling BGP underlay may cause spine BGP router-id change

BGP Underlay Support for IPN links



- Configure BGP 'disable-peer-as-check' if Nexus switches are used for IPN
- Nexus switches will not advertise prefixes to peer if peer AS is already in the AS PATH. 'disable-peer-as-check' turns off this behavior

Sample IPN configuration (Nexus 9000)

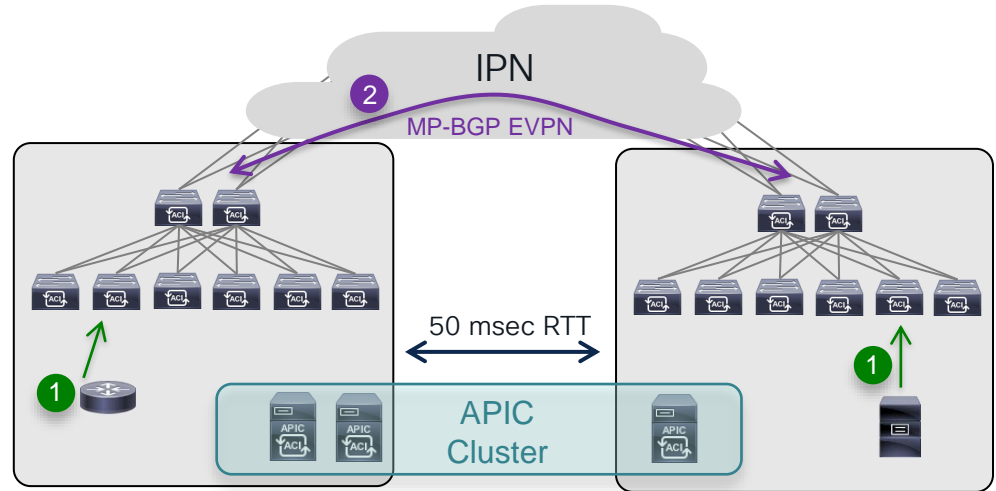
```
feature bgp

router bgp 65010
  router-id 10.10.10.1
  vrf IPN
    address-family ipv4 unicast
    neighbor 10.1.1.1
    remote-as 65001
    address-family ipv4 unicast
    disable-peer-as-check
```


ACI Multi-Pod and MTU

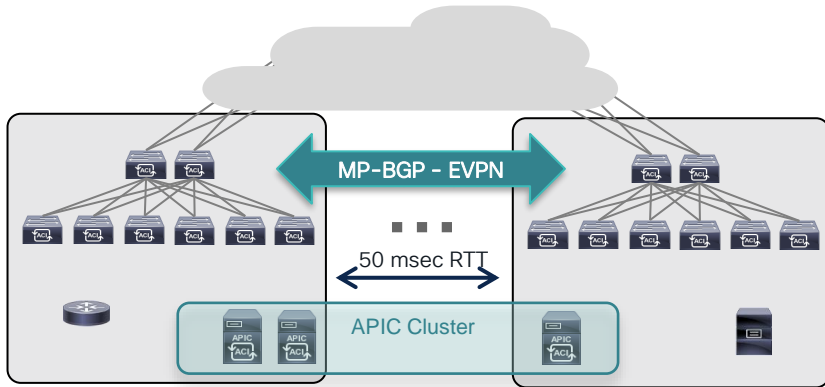
Different MTU Meanings

1. **Data Plane MTU:** MTU of the traffic generate by endpoints (servers, routers, service nodes, etc.) connected to ACI leaf nodes
 - Need to account for 50B of overhead (VXLAN encapsulation) for inter-Pod communication
2. **Control Plane MTU:** for CPU generated traffic like EVPN across sites
 - The default value is **9000B**, can be tuned to the maximum MTU value supported in the IPN



ACI Multi-Pod and MTU

Tuning CP MTU for EVPN Traffic across Pods

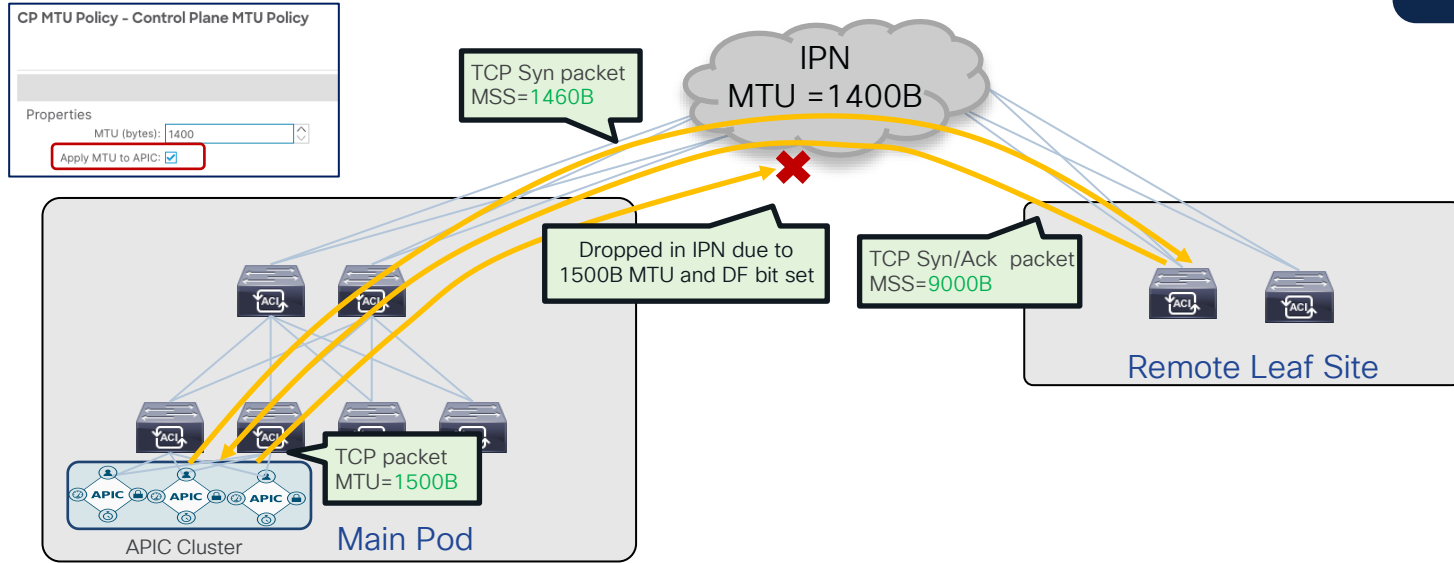


- Control Plane MTU can be set leveraging the “CP MTU Policy” on APIC
- The required MTU in the IPN would then depend on this setting and on the Data Plane MTU configuration
 - Always need to consider the VXLAN encapsulation overhead for data plane traffic (50/54 bytes)

The screenshot shows the APIC System Settings page. The 'System Settings' menu is open, and the 'CP MTU Policy - Control Plane MTU Policy' page is selected. The 'Properties' section shows the 'MTU (bytes)' field set to 9000. A red dashed box highlights this field, and an arrow points to it with the text 'Modify the default 9000B MTU value'.

Apply Control Plane MTU to APIC

ACI Release 6.0(4)F

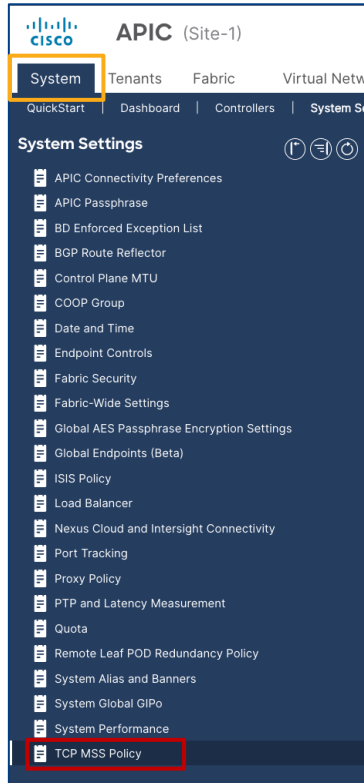


- Leaf discovery may fail for Remote Leaf and Multi-Pod switches if the IPN network supports less than 1500-byte MTU
- From ACI release 6.0(4) you can set the APIC fabric interface to the configured CP MTU value (if less than 1500 bytes)

ACI Multi-Site and MTU Size

Introducing the TCP-MSS Adjust Functionality

ACI Release 6.0(3)F



The image shows the 'TCP MSS Policy' configuration page. The 'Type' is set to 'Global'. The 'IPv4' field is set to '8888' and the 'IPv6' field is set to '8868'. A red box highlights these two fields, with a blue arrow pointing to the explanatory text below.

Supported values are
688-9104 bytes

- TCP MSS adjust policy is enabled at System Settings level
- Supports different TCP MSS adjust setting for IPv4 and IPv6
- Supports three different options:
 1. **Global**: applies to all flows (Multi-Pod, Multi-Site, RLS to LLS/RLs to RLS)
 2. **RL and Msite**: applies to Multi-Site and RLS to LLS/RLs to RLS flows
 3. **RL Only**: applies only to RLS to LLS/RLs to RLS flows

TCP-MSS Adjust Functionality

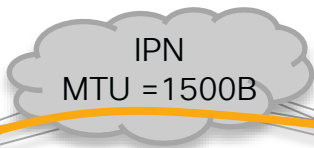
ACI Release 6.0(3)F

SYN Packet

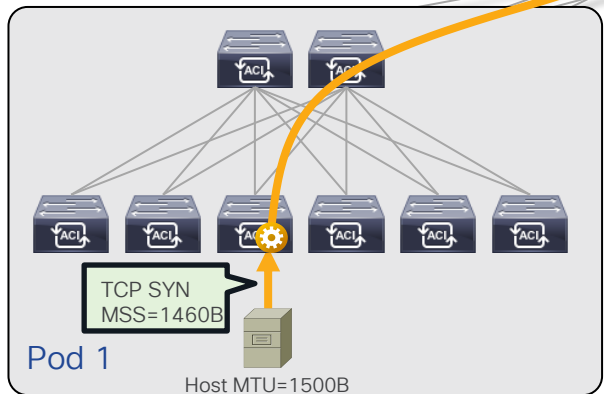
Type: Global RL and Multi-Site RL Only Disable

IPv4: 1400

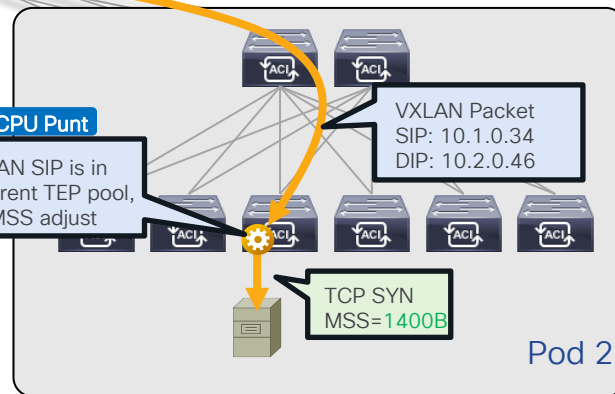
IPv6: 8868



$$\text{MSS} = \text{MTU} - \text{IP Header Size} - \text{TCP Header Size}$$



TEP Pool 10.1.0.0/16



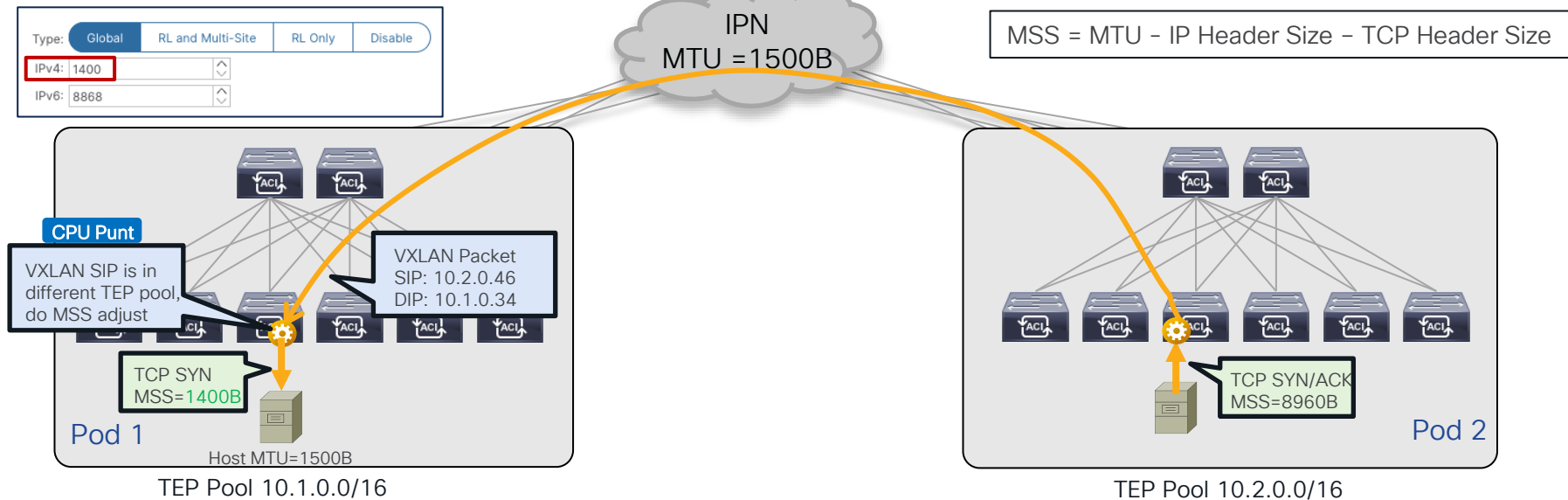
TEP Pool 10.2.0.0/16

- TCP MSS adjust is always performed on the egress leaf node
- Adjusts TCP MSS value on SYN and SYN/ACK packets
- Checks for Source IP in the VXLAN header → TCP-MSS adjusts performed if the source IP is not part of the pods local internal TEP pool

TCP-MSS Adjust Functionality

ACI Release 6.0(3)F

SYN/ACK Packet



- TCP MSS adjust is always performed on the egress leaf node
- Adjusts TCP MSS value on SYN and SYN/ACK packets
- Checks for Source IP in the VXLAN header → TCP-MSS adjusts performed if the source IP is not part of the pods local internal TEP pool

TCP-MSS Adjust Functionality

ACI Release 6.0(3)F

Inter-Pod Data Packets

Type: Global RL and Multi-Site RL Only Disable

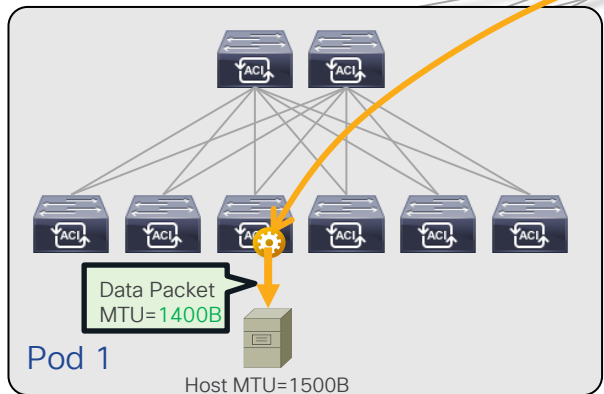
IPv4:

IPv6:

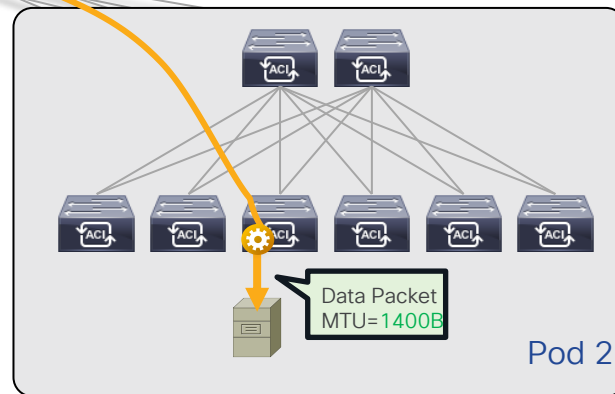
IPN
MTU = 1500B

VXLAN Packet
MTU=1450B

$MSS = MTU - IP\ Header\ Size - TCP\ Header\ Size$



TEP Pool 10.1.0.0/16



TEP Pool 10.2.0.0/16

- As a result of the MSS negotiation, the endpoints generate packets for that TCP communication with MTU 1400B (irrespective of the local Host MTU)
- The VXLAN encapsulated traffic can be successfully forwarded across the IPN

ACI Multi-Pod and QoS

Inter-Pod QoS Behavior

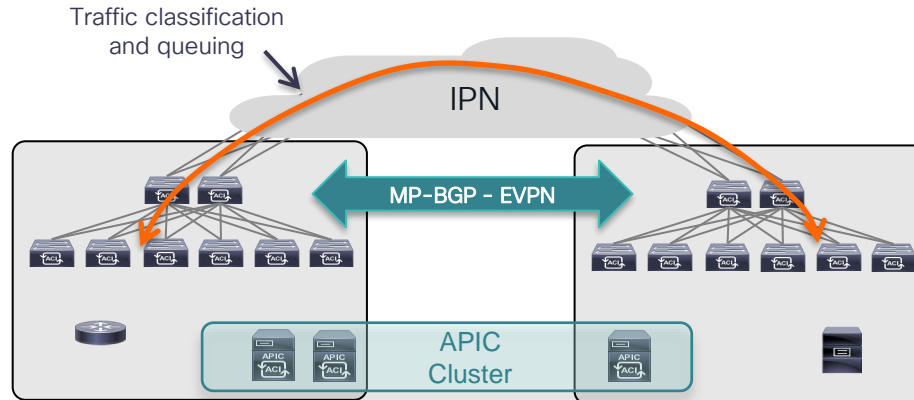
- Traffic across sites should be consistently prioritized (as it happens intra-site)
- To achieve this end-to-end consistent behavior, it is required to configure DSCP-to-CoS mapping in the 'infra' Tenant
 - Allows to classify traffic received on the spines from the IPN based on outer DSCP value
 - Without the DSCP-to-CoS mapping configuration, classification for the same traffic will be CoS based (preserving CoS value in the IPN is harder)
- The traffic can also then be properly treated inside the IPN (classification/queuing)
 - Recommended to always prioritize at least Policy and Control Plane traffic

DSCP class-CoS translation policy for L3 traffic

Properties

Translation Policy State: Disabled Enabled

User Level 1:	CS1	▼
User Level 2:	CS2	▼
User Level 3:	CS3	▼
User Level 4:	AF11 low drop	▼
User Level 5:	AF21 low drop	▼
User Level 6:	AF31 low drop	▼
Control Plane Traffic:	CS0	▼
Policy Plane Traffic:	CS4	▼
Span Traffic:	CS5	▼
Traceroute Traffic:	CS6	▼



DSCP class-CoS translation policy for L3 traffic

Properties

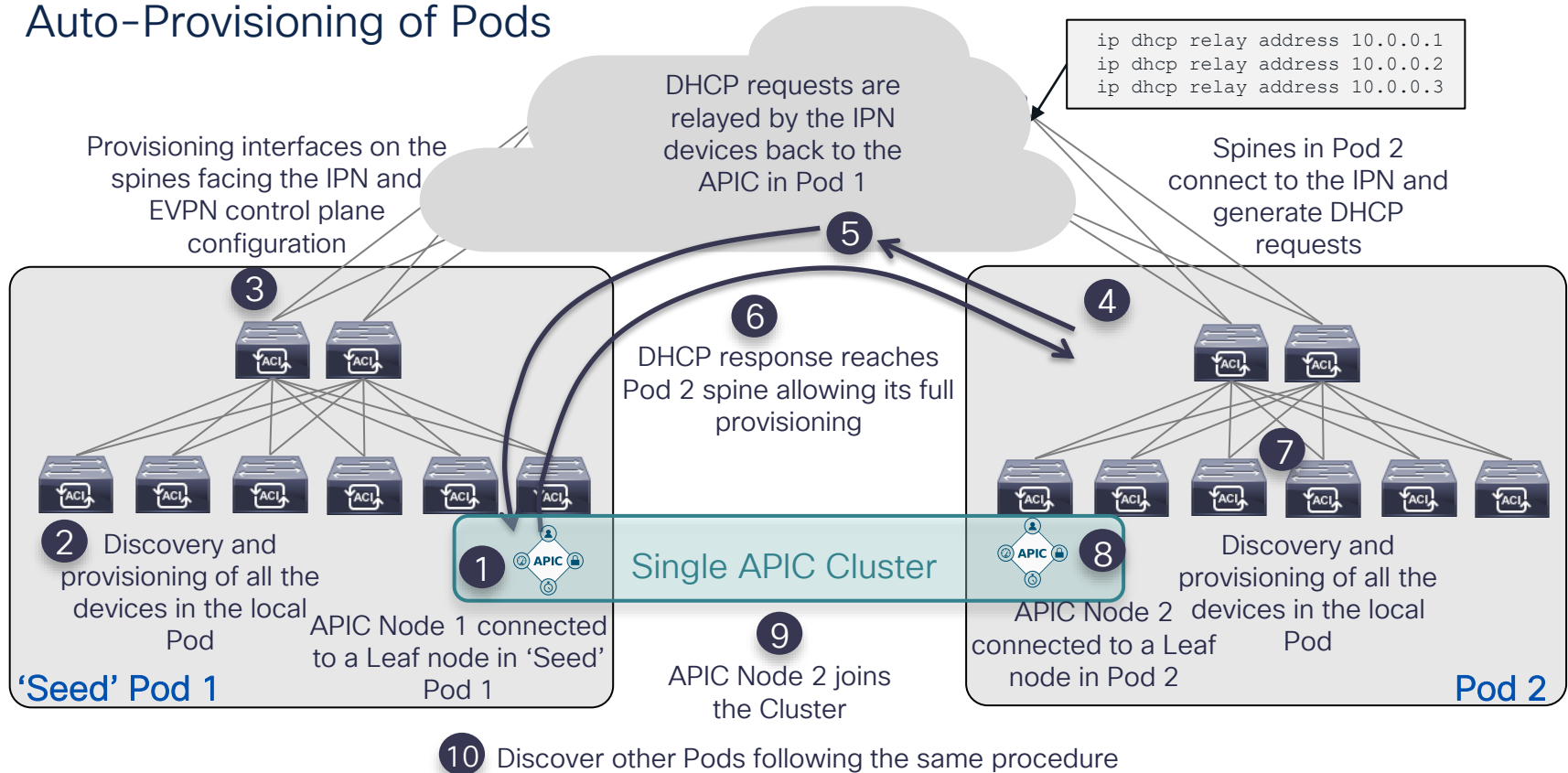
Translation Policy State: Disabled Enabled

User Level 1:	CS1	▼
User Level 2:	CS2	▼
User Level 3:	CS3	▼
User Level 4:	AF11 low drop	▼
User Level 5:	AF21 low drop	▼
User Level 6:	AF31 low drop	▼
Control Plane Traffic:	CS0	▼
Policy Plane Traffic:	CS4	▼
Span Traffic:	CS5	▼
Traceroute Traffic:	CS6	▼

Control and Data Planes

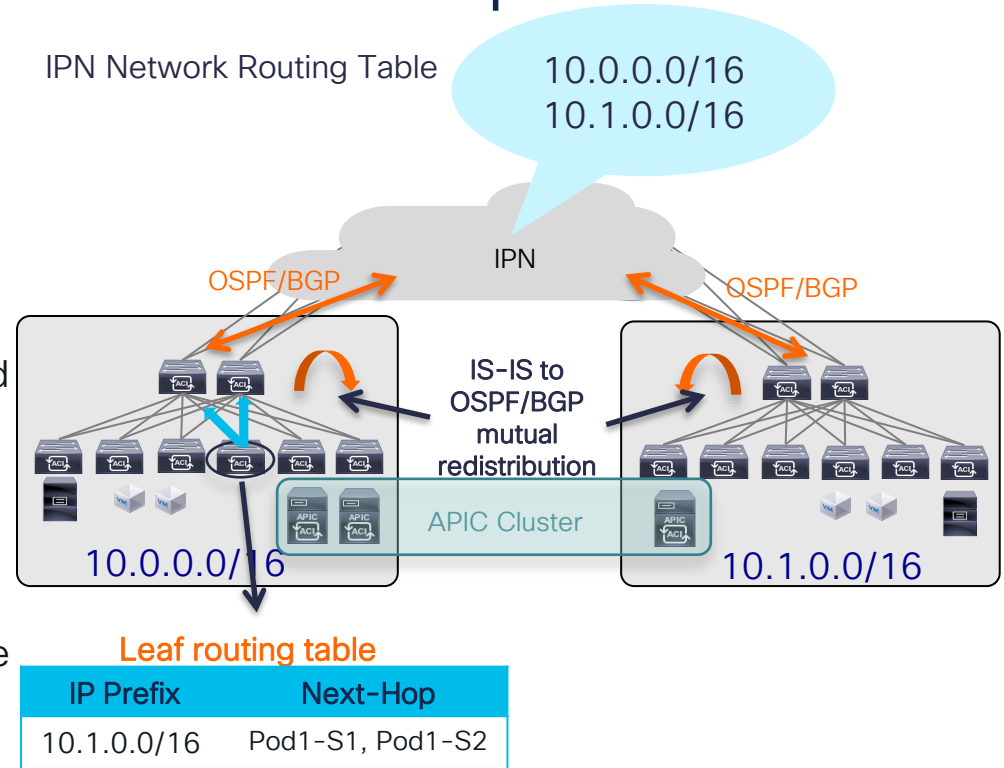
ACI Multi-Pod

Auto-Provisioning of Pods



Exchanging TEP information across pods

- Separate IP address pools for VTEPs assigned by APIC to each Pod
 - Summary routes advertised toward the IPN via OSPF or BGP routing
 - IS-IS convergence events local to a Pod not propagated to remote Pods
- Spine nodes redistribute other Pods summary routes into the local IS-IS process
 - Needed for local VTEPs to communicate with remote VTEPs



Exchanging TEP information across pods

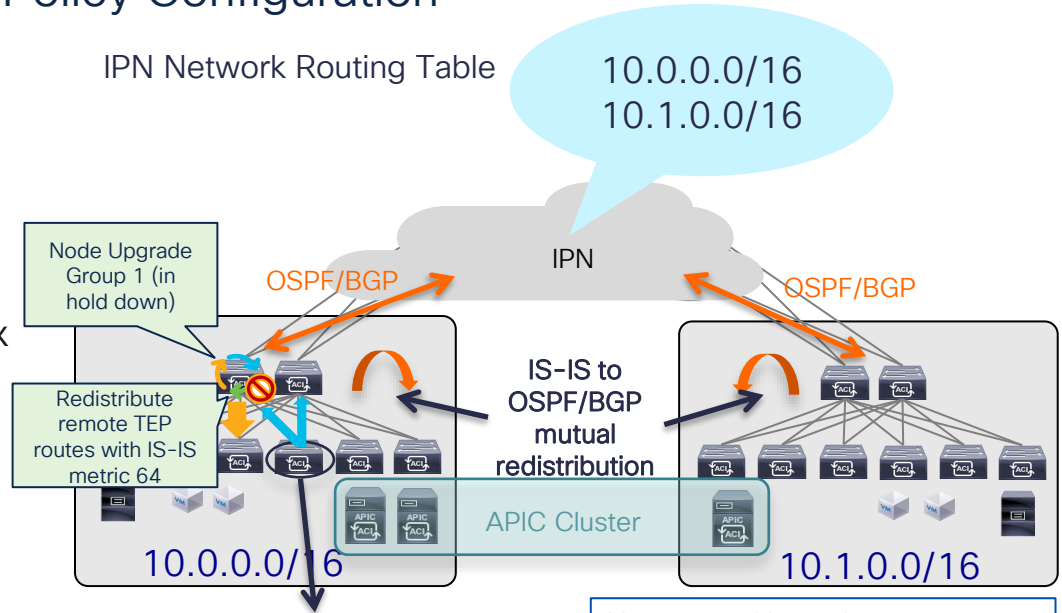
Issue with Default IS-IS Metric Policy Configuration

Default fabric wide IS-IS metric is set at 63 (max value)

During upgrade, spines set the overload mode while policy is being downloaded

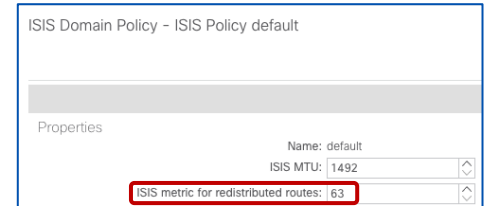
If fabric-wide value is already using the max value for routes redistributed into IS-IS, the overload functionality is ineffective

This can create unexpected traffic interruption if leaf sends traffic to a spine which is not fully upgraded (and ready to forward traffic)



Leaf routing table

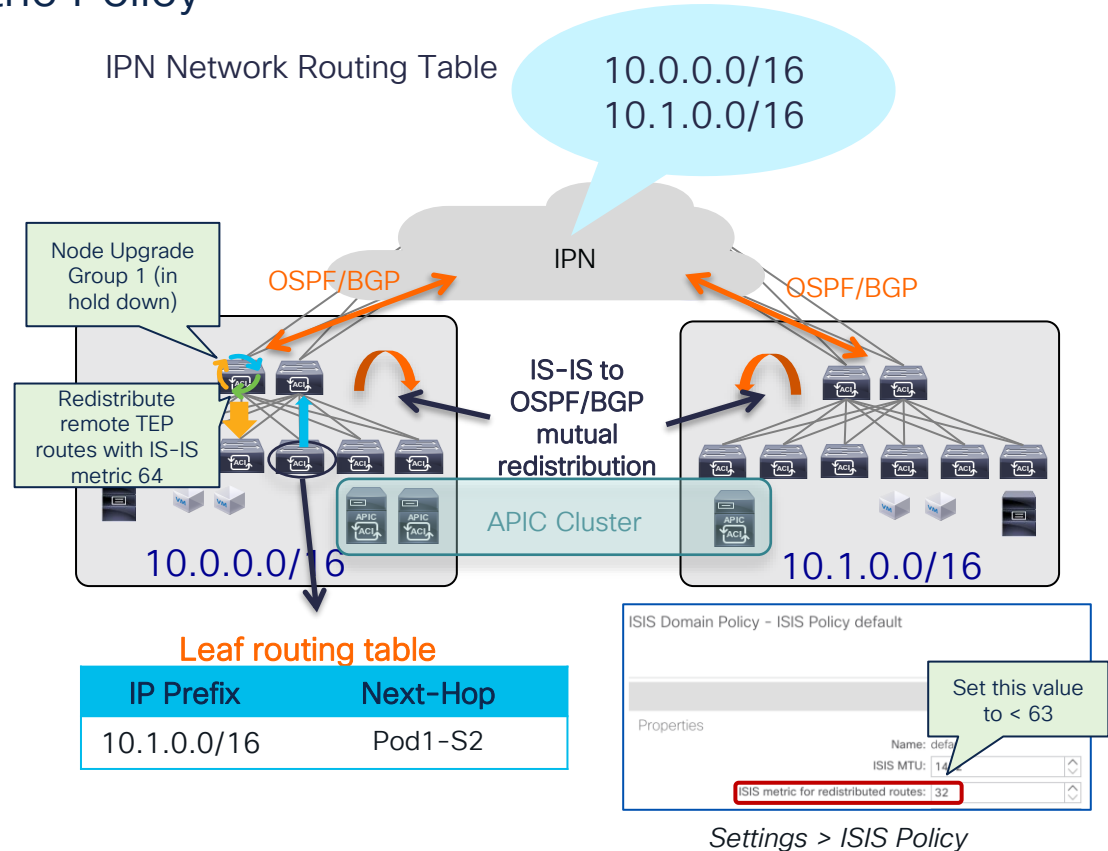
IP Prefix	Next-Hop
10.1.0.0/16	Pod1-S1 Pod1-S2



Exchanging TEP information across pods

Lowering the Default IS-IS Metric Policy

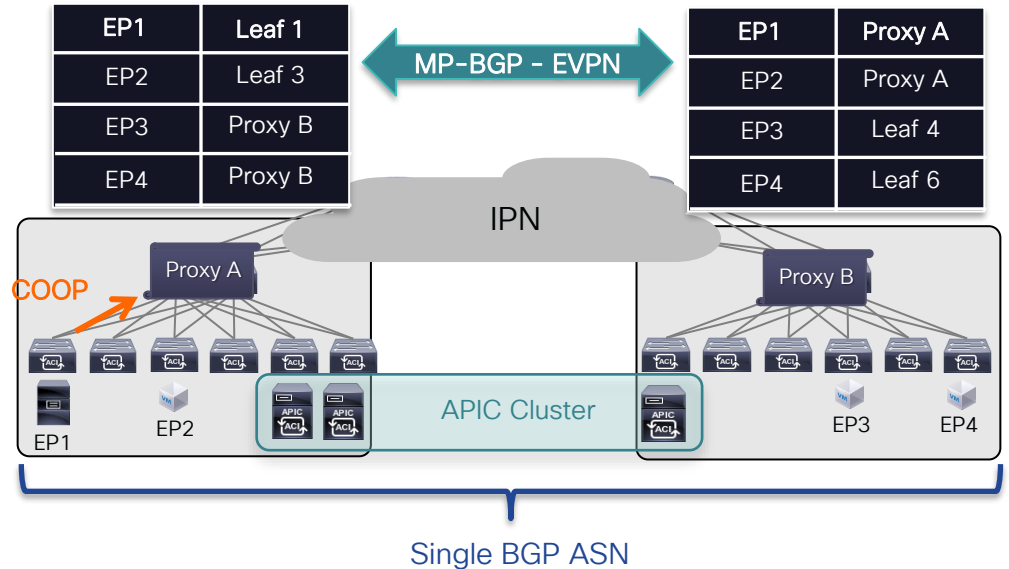
- By lowering the default ISIS metric value, connectivity to TEP prefixes received from the remote site will be preferred through the remaining spines
- This behavior gives time to the spine for completing the upgrade



ACI Multi-Pod

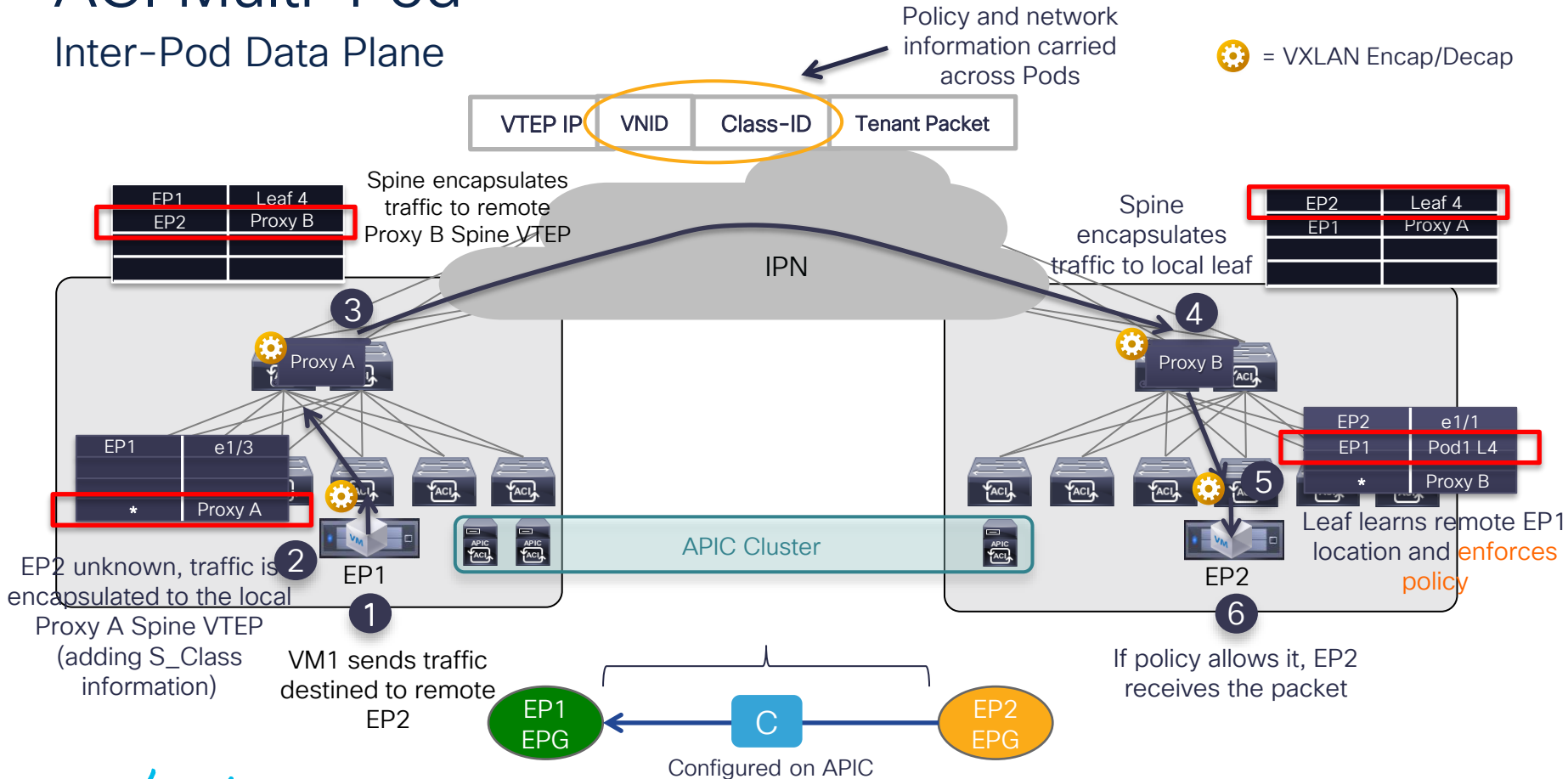
Inter-Pod MP-BGP EVPN Control Plane

- MP-BGP EVPN to sync Endpoint (EP) and Multicast Group information
 - All remote Pod entries associated to a Proxy VTEP next-hop address (not part of local TEP Pool)
 - Same BGP AS across all the Pods
- iBGP EVPN sessions between spines in separate Pods
 - Full mesh MP-iBGP EVPN sessions between local and remote spines (default behavior)
 - Optional RR deployment (recommended one RR in each Pod for resiliency)



ACI Multi-Pod

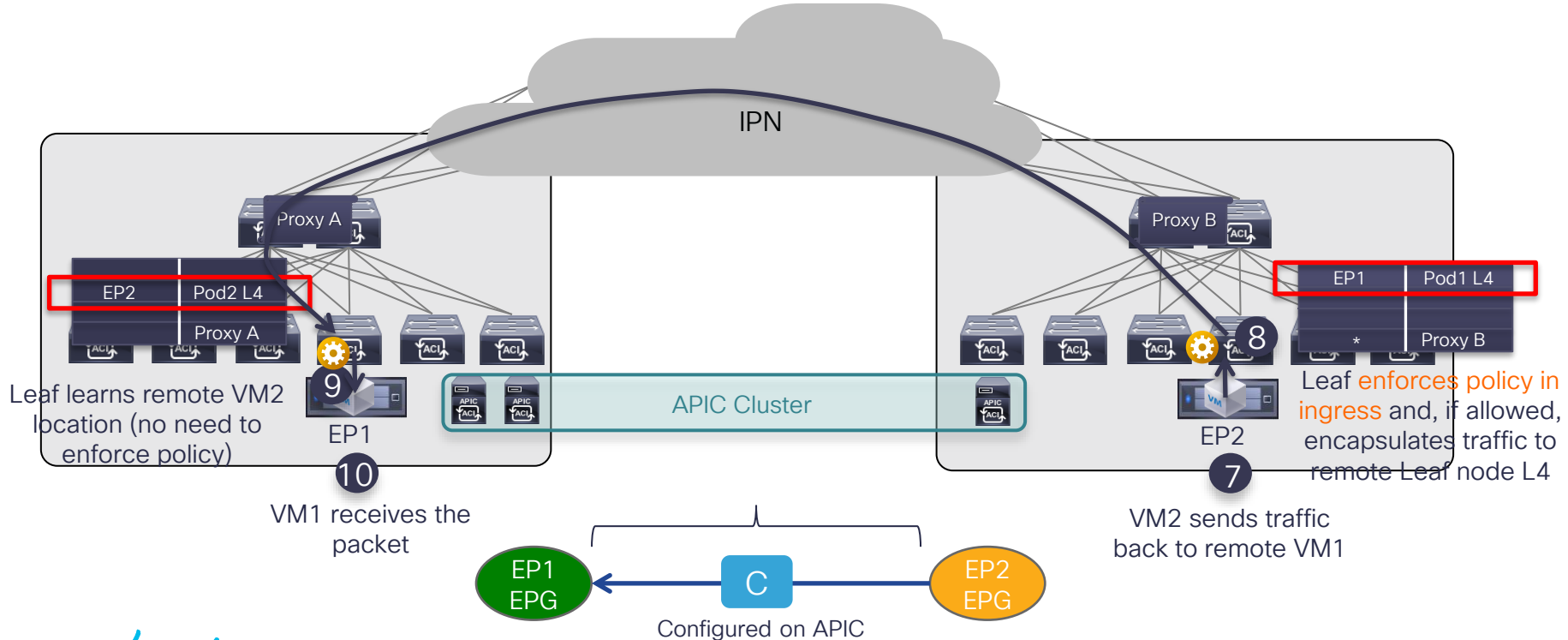
Inter-Pod Data Plane



ACI Multi-Pod

Inter-Pod Data Plane (2)

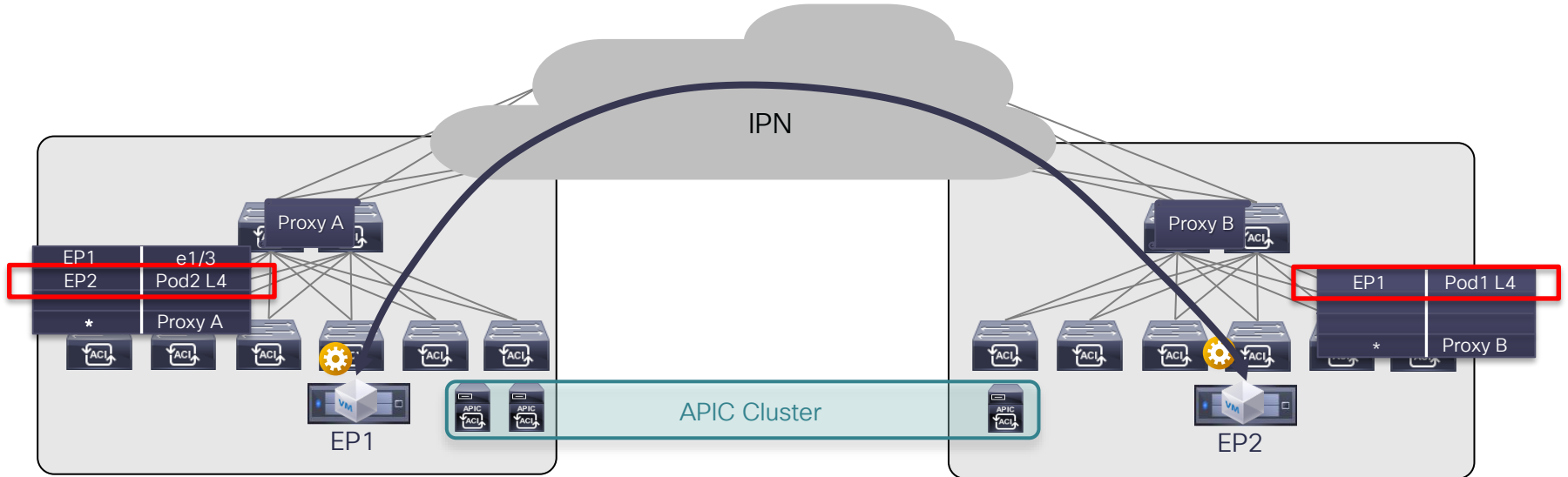
 = VXLAN Encap/Decap



ACI Multi-Pod

Inter-Pod Data Plane (3)

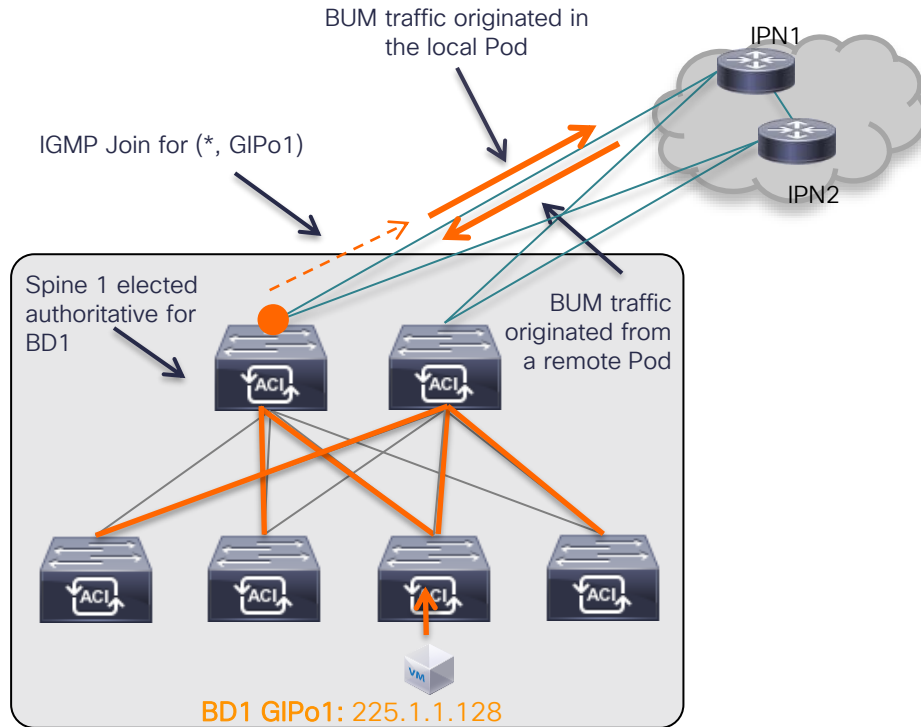
 = VXLAN Encap/Decap



- 11 From this point EP1 to EP2 communication is encapsulated Leaf to Leaf (VTEP to VTEP) and policy always applied at the ingress leaf (applies to both L2 and L3 communication)

ACI Multi-Pod

Use of Multicast for Inter-Pod Layer 2 BUM Traffic

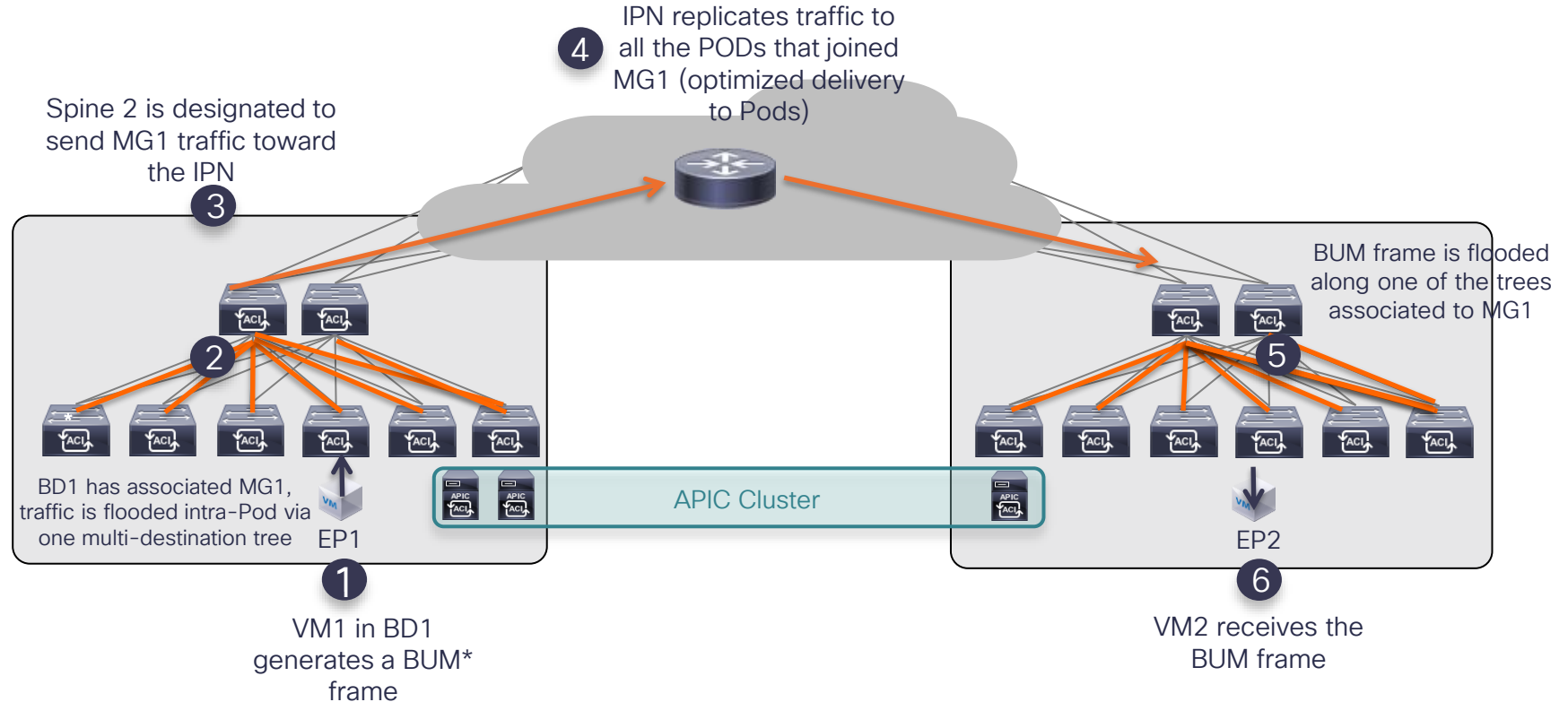


BUM: Broadcast, Unknown Unicast, Multicast

- Ingress replication for BUM* traffic not supported with Multi-Pod
- PIM Bidir is the only validated and supported option
 - Scalable: only a single (*,G) entry is created in the IPN for each BD
 - Fast-convergent: no requirement for data-driven multicast state creation
- A spine is elected authoritative for each Bridge Domain:
 - Generates an IGMP Join on a specific link toward the IPN
 - Always sends/receives BUM traffic on that link

ACI Multi-Pod

Use of Multicast for Inter-Pod BUM Traffic

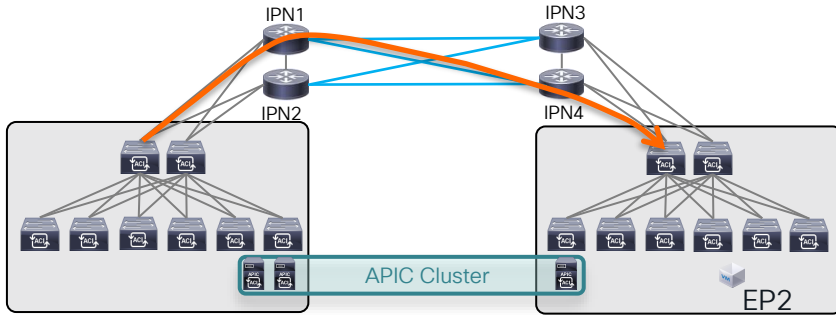


BUM: Layer 2 Broadcast, Unknown Unicast, Multicast

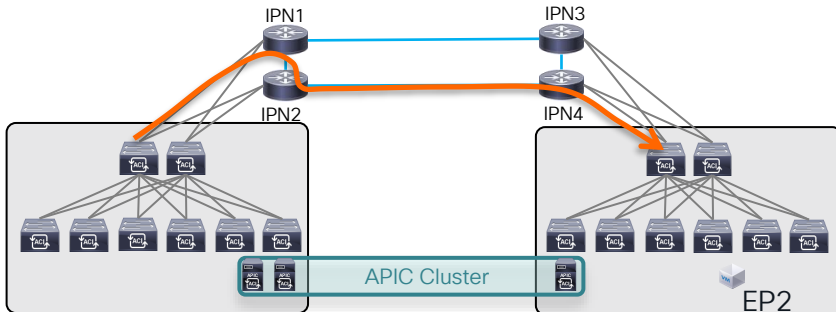
ACI Multi-Pod

PIM Bidir for BUM – Supported Topologies

Full Mesh between remote IPN devices

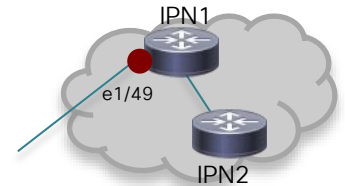


Directly connect local IPN devices

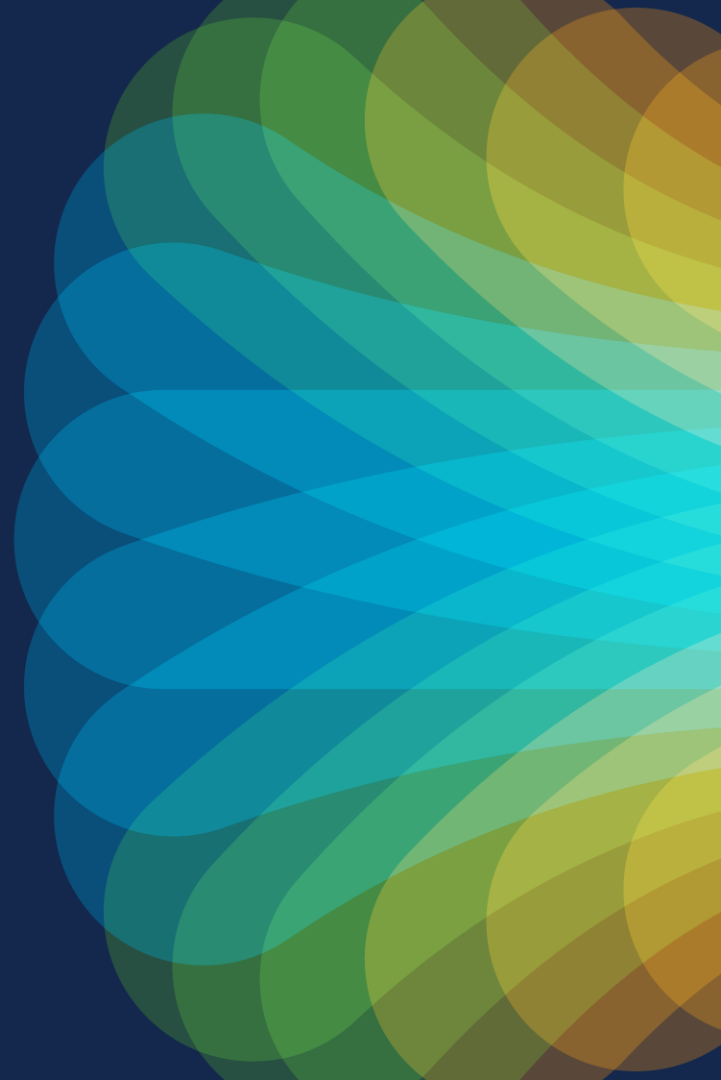


- Create full-mesh connections between IPN devices
- More costly for geo-dispersed Pods, as it requires more links between sites
- Alternatively, connect local IPN devices with a port-channel interface (for resiliency)
- In both cases, it is **critical** to ensure that the preferred path toward the RP from any IPN devices is not via a spine
- Recommendation is to increase the OSPF cost of the interfaces between IPN and spines

```
interface Ethernet1/49.4
description L3 Link to Pod1-Spine1
mtu 9150
encapsulation dot1q 4
ip address 192.168.1.1/31
ip ospf cost 100
ip ospf network point-to-point
ip router ospf IPN area 0.0.0.0
ip pim sparse-mode
ip dhcp relay address 10.1.0.2
ip dhcp relay address 10.1.0.3
```

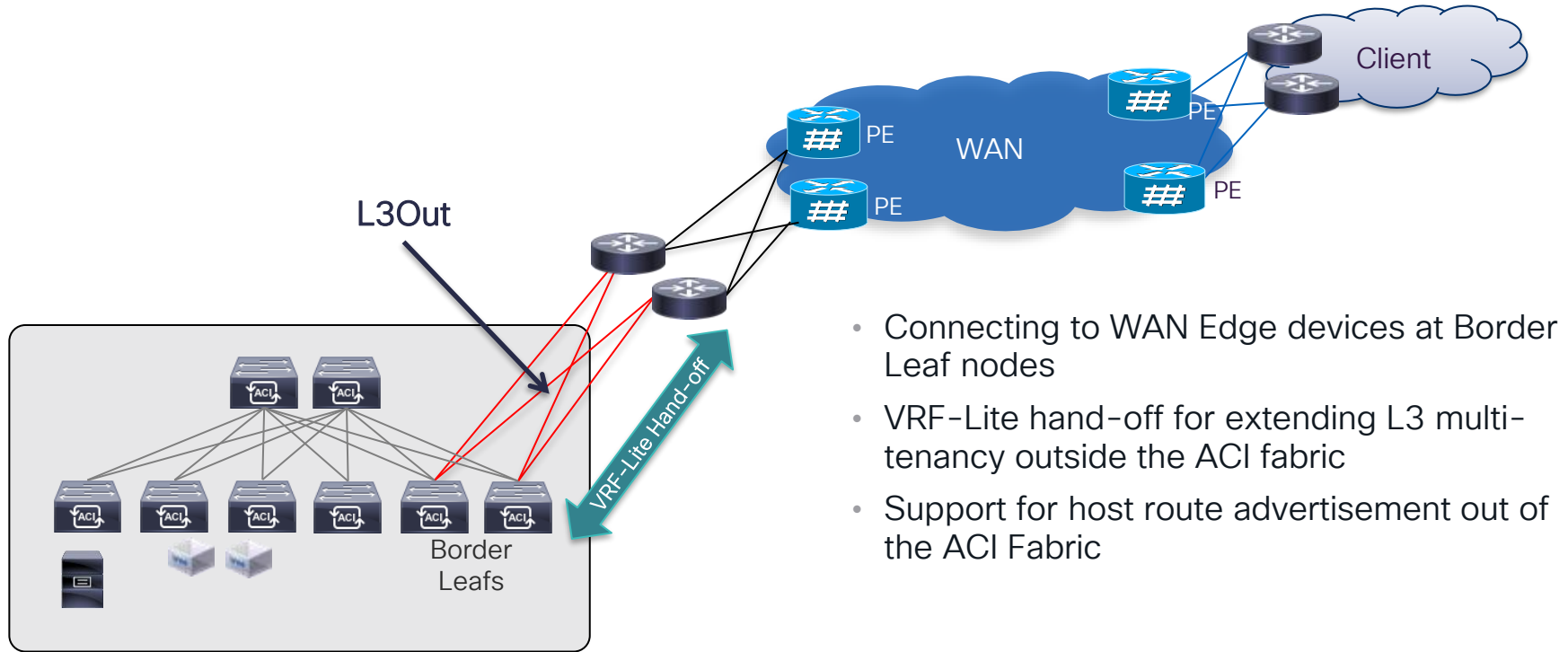


Connecting to the External Layer 3 Domain



Connecting ACI to Layer 3 Domain

'Traditional' L3Out on the BL Nodes

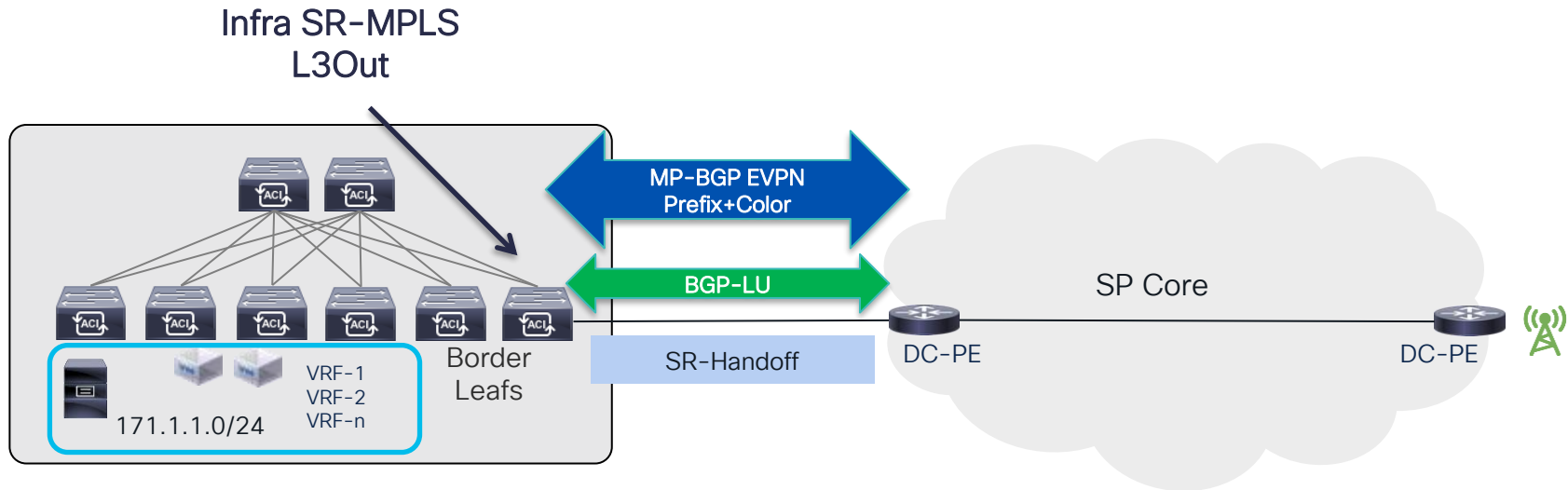


- Connecting to WAN Edge devices at Border Leaf nodes
- VRF-Lite hand-off for extending L3 multi-tenancy outside the ACI fabric
- Support for host route advertisement out of the ACI Fabric

Connecting ACI to Layer 3 Domain

'SR-MPLS Handoff'

- Border Leafs connect to PE router in SP core
- Single BGP EVPN session for all VRFs
- ACI BL is advertising EVPN type-5 routes with BGP color community



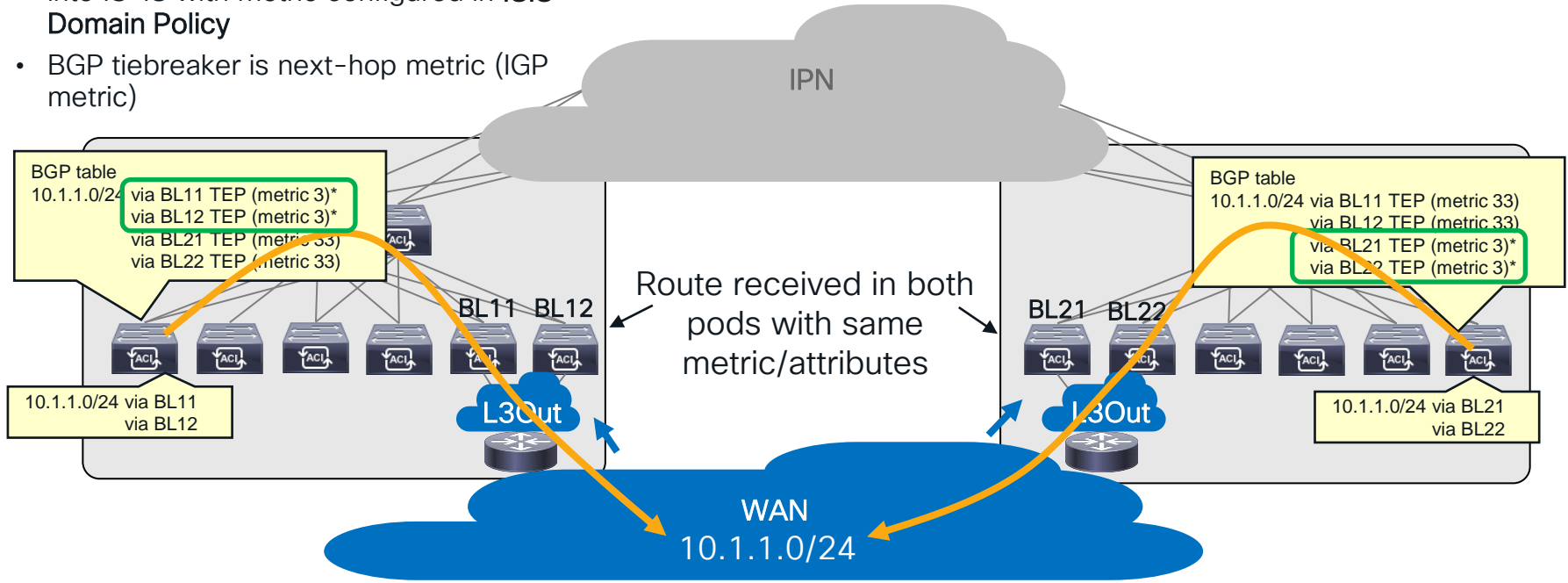
Connecting to the External L3 Domain

Local L3Outs preferred over L3Outs in remote pods

- Remote pod TEP routes are redistributed into IS-IS with metric configured in **ISIS Domain Policy**
- BGP tiebreaker is next-hop metric (IGP metric)

ISIS Domain Policy - ISIS Policy default
ISIS metric for redistributed routes: 32

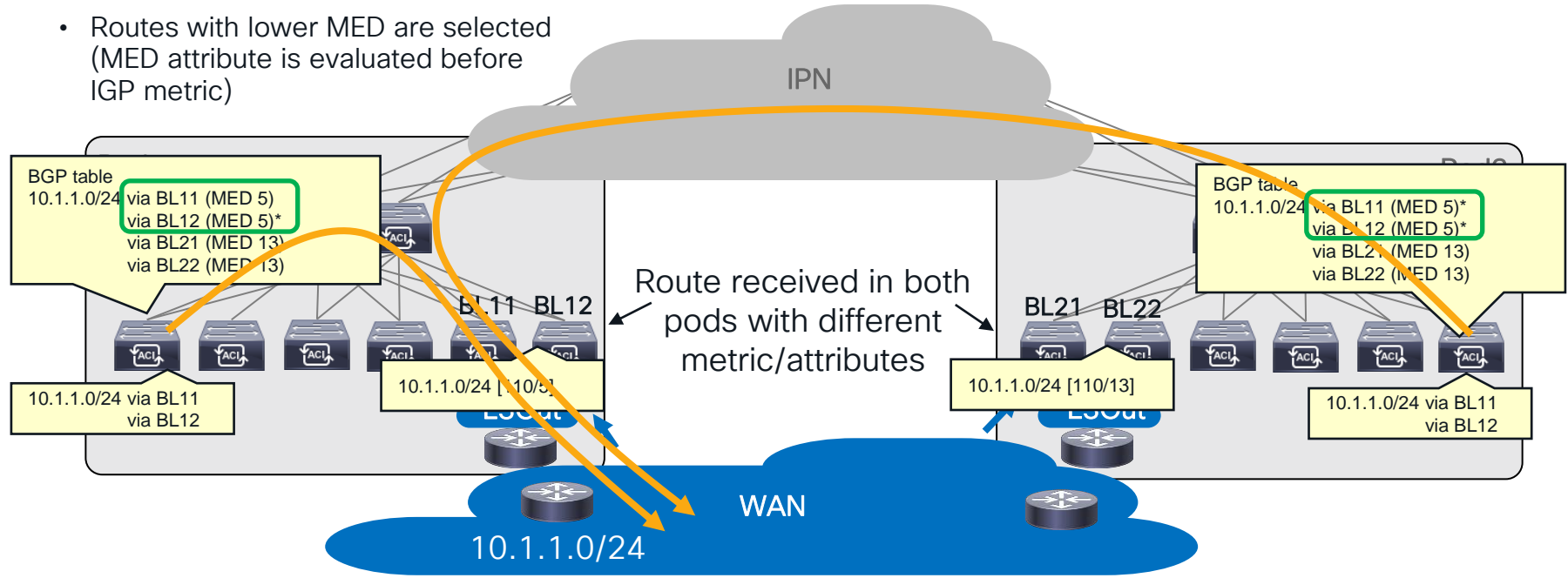
Best Practice



Connecting to the External L3 Domain

Remote pod L3Out may be used if it has a better external metric

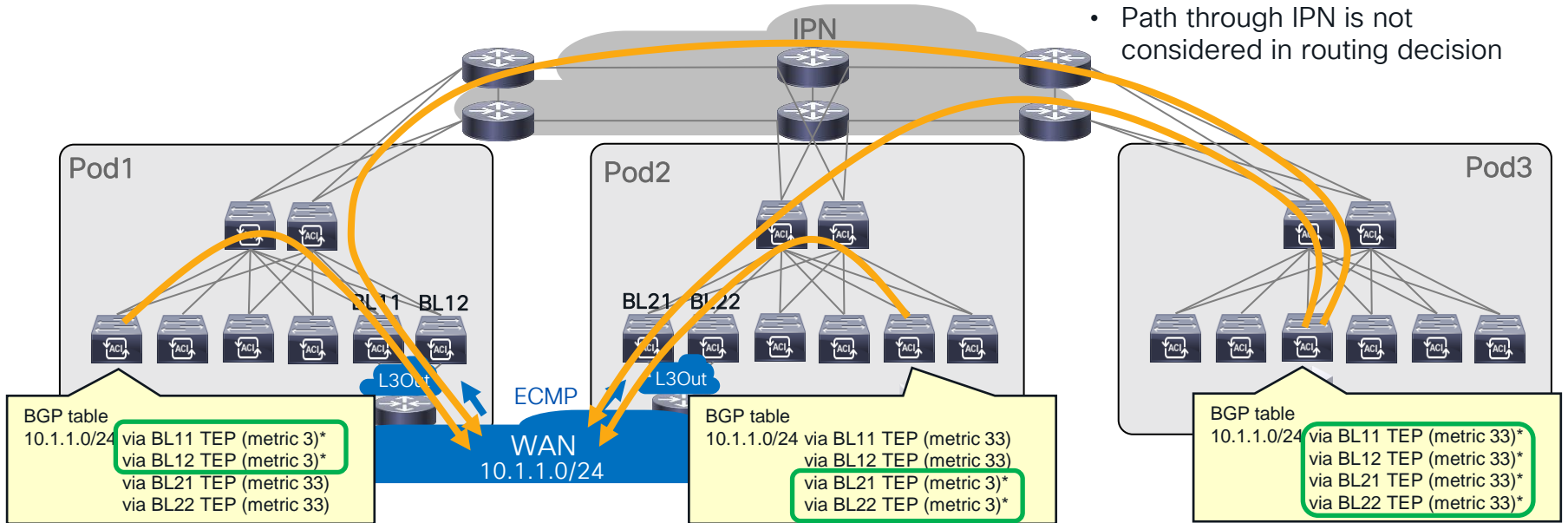
- BGP MED is set to OSPF metric when redistributed into fabric
- Routes with lower MED are selected (MED attribute is evaluated before IGP metric)



Connecting Multi-Pod to the Layer 3 Domain

What happens when there are more than two pods?

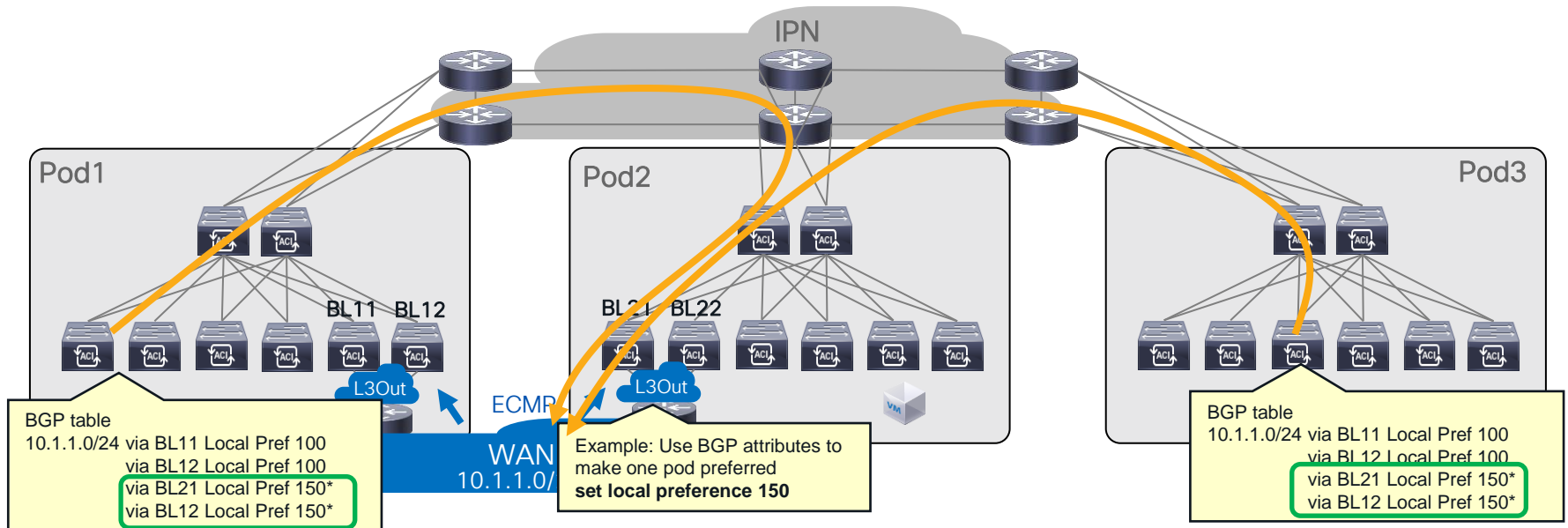
- Traffic flows are load balanced across all remote pods
- Path through IPN is not considered in routing decision



- A pod does not need a dedicated L3Out. Flows to external destinations can use an L3Out in another pod

Connecting Multi-Pod to the Layer 3 Domain

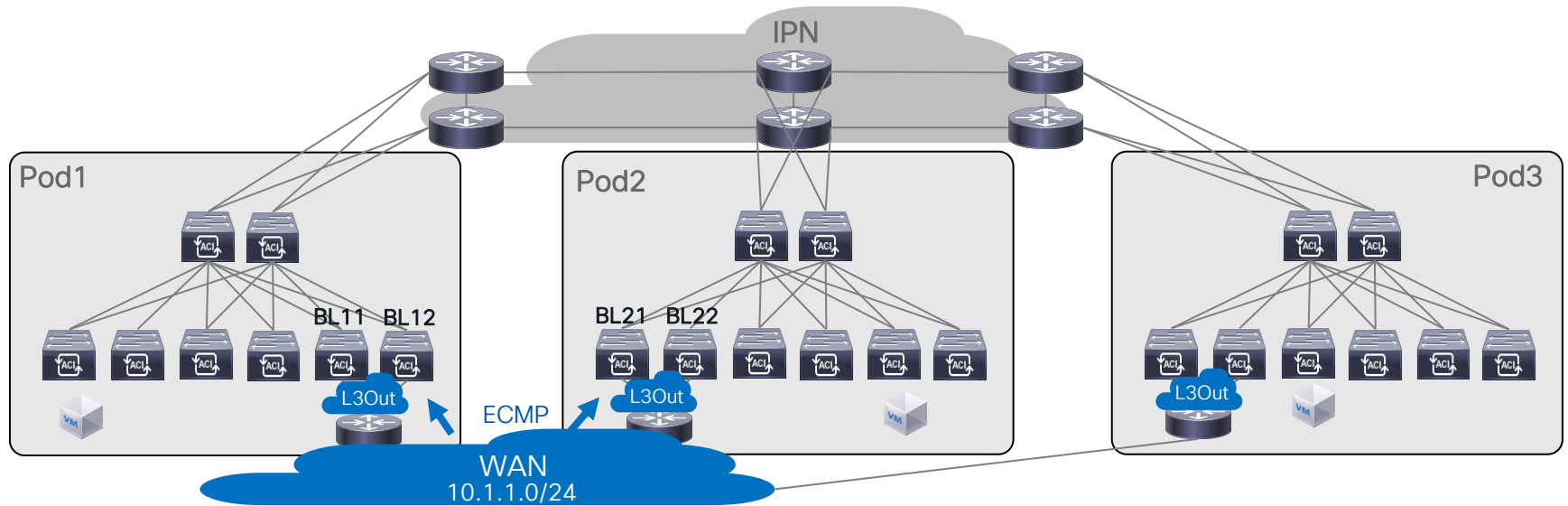
How to prefer one remote pod over another?



But change will affect all pods!

Connecting Multi-Pod to the Layer 3 Domain

How to prefer one remote pod over another?

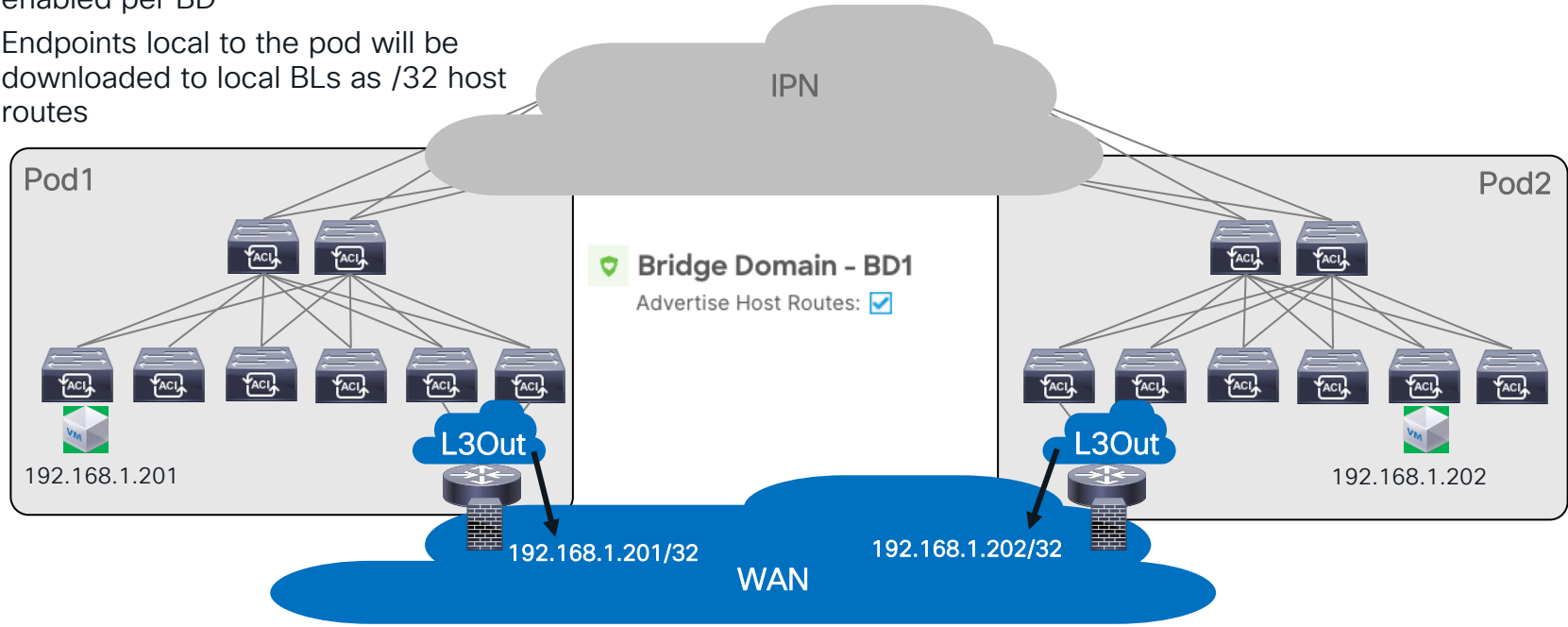


Adding a local L3out may be a better option

Connecting to the External L3 Domain

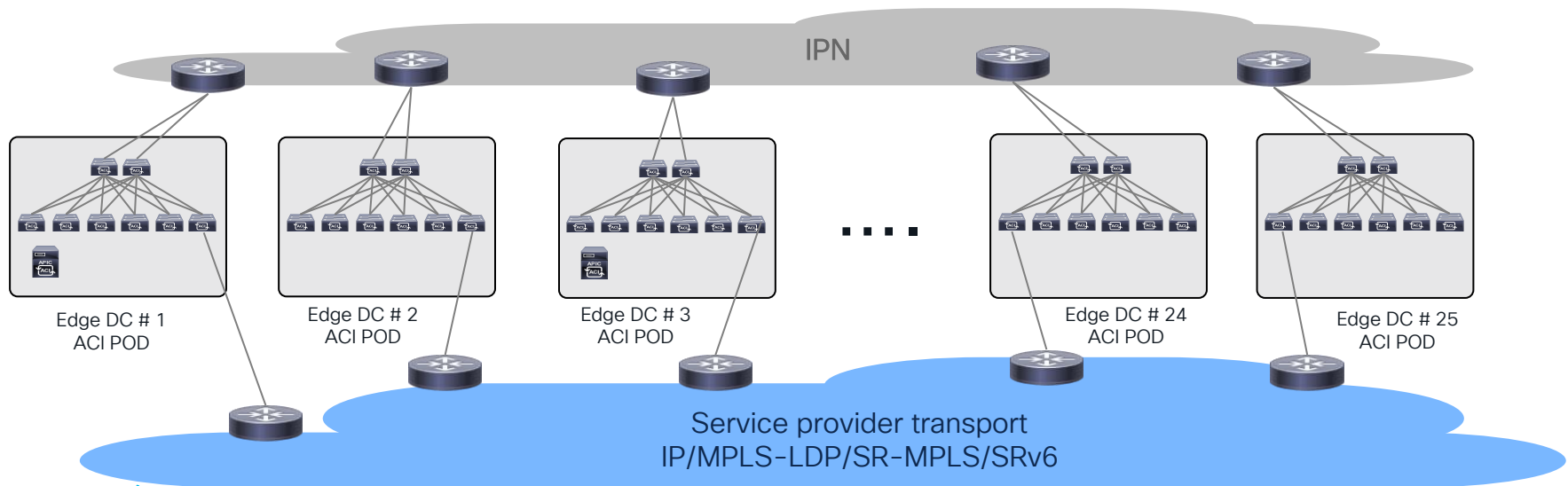
Influencing inbound path: Host route advertisement

- Host route advertisement can be enabled per BD
- Endpoints local to the pod will be downloaded to local BLs as /32 host routes



Edge DCs with Multi-Pod architecture

- APIC controllers are needed only in some Pods
- Communication across Pods is typically through SR-MPLS L3out
- 25 Pods per fabric is supported starting 6.0(1) release
- Leaf scale per fabric remains same. 2 Spines per Pod is supported
- Latency requirement remains same - 50 msec RTT requirement across APIC clusters and between switches and APIC
- No need to enable PIM-Bidir in IPN if L2 extension across Pod is not required

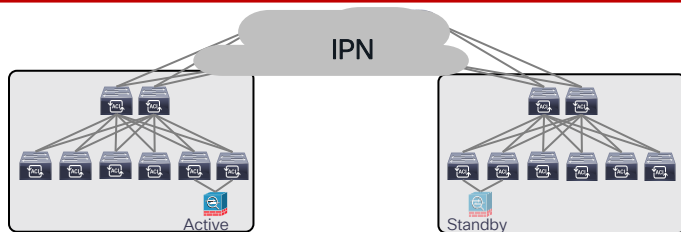


Network Services Integration

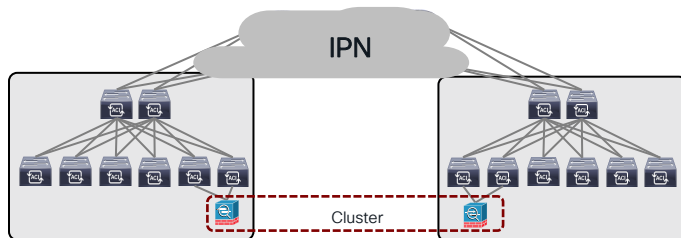
ACI Multi-Pod

Design options

Typical options for an Active/Active DC use case



- Active and Standby pair deployed across Pods
- No issues with asymmetric flows



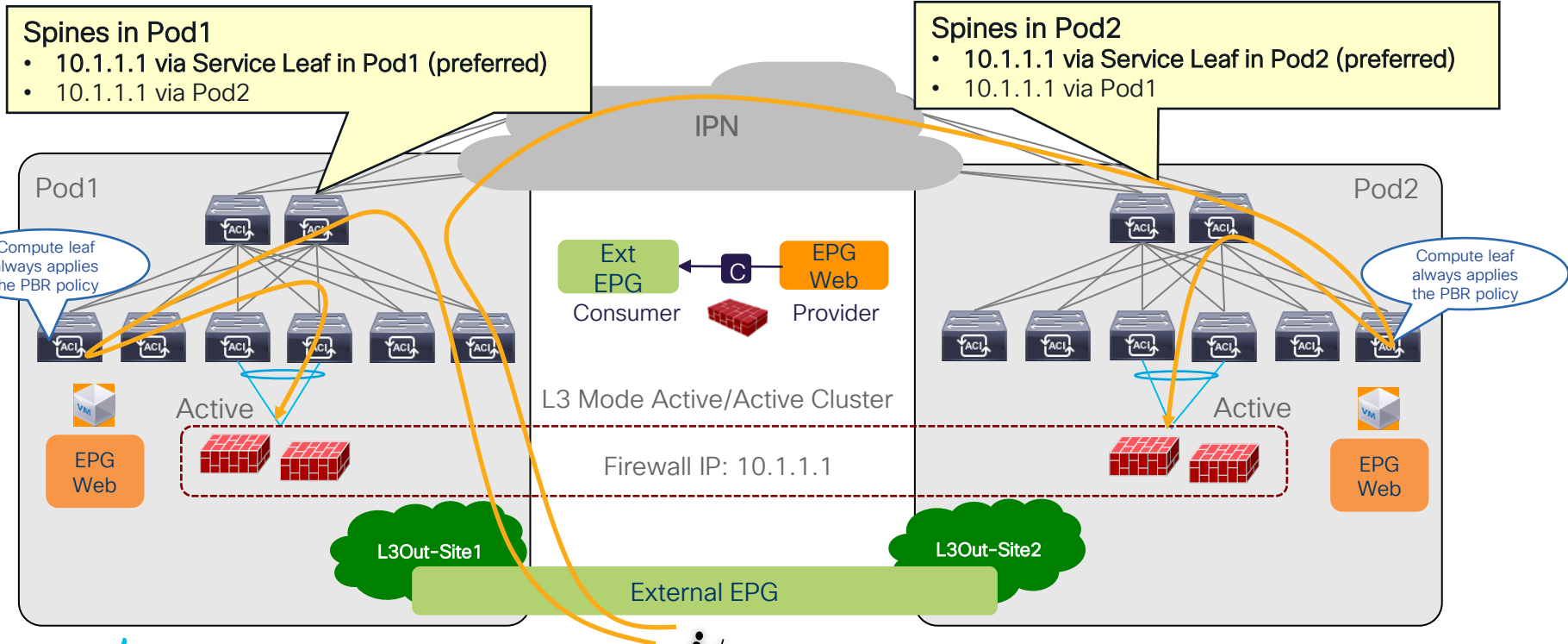
- Active/Active FW cluster nodes stretched across Sites (single logical FW)
- Requires the ability of discovering the same MAC/IP info in separate sites at the same time
- Supported from ACI release 3.2(4d) with the use of Service-Graph with PBR



- Independent Active/Standby pairs deployed in separate Pods
- Use of Symmetric PBR to avoid the creation of asymmetric paths crossing different active FW nodes

ACI Multi-Pod: Active/Active cluster across pods

North-South Traffic Flow



ACI Multi-Pod: Active/Active cluster across pods

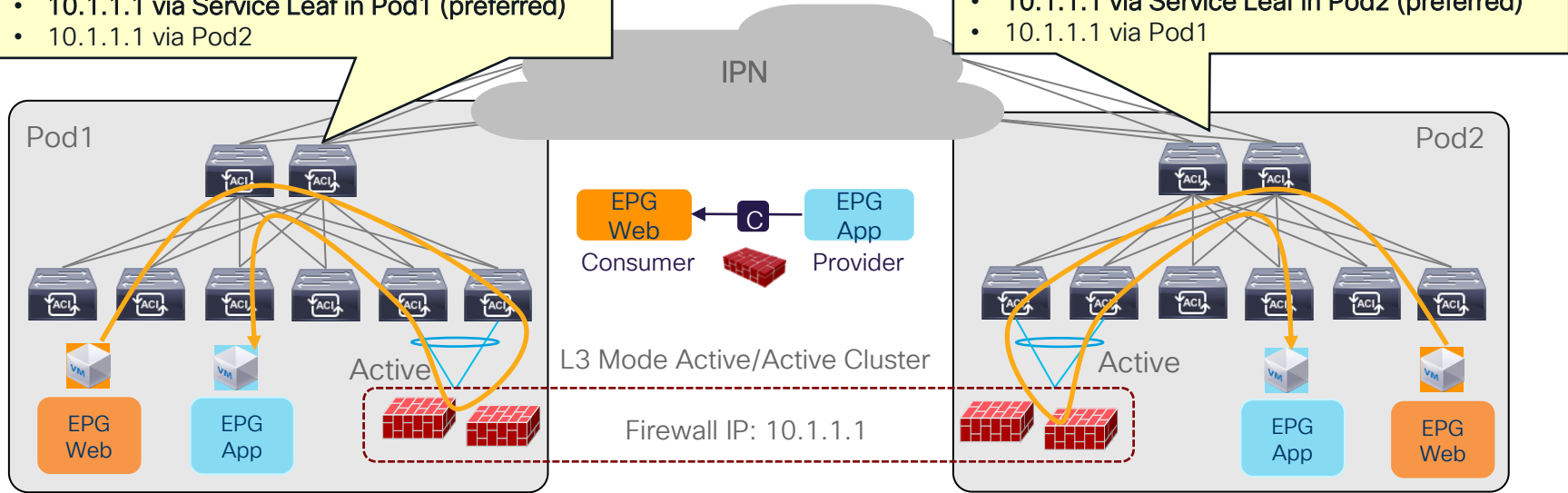
East-West Traffic Flow (Intra-Pod)

Spines in Pod1

- 10.1.1.1 via Service Leaf in Pod1 (preferred)
- 10.1.1.1 via Pod2

Spines in Pod2

- 10.1.1.1 via Service Leaf in Pod2 (preferred)
- 10.1.1.1 via Pod1

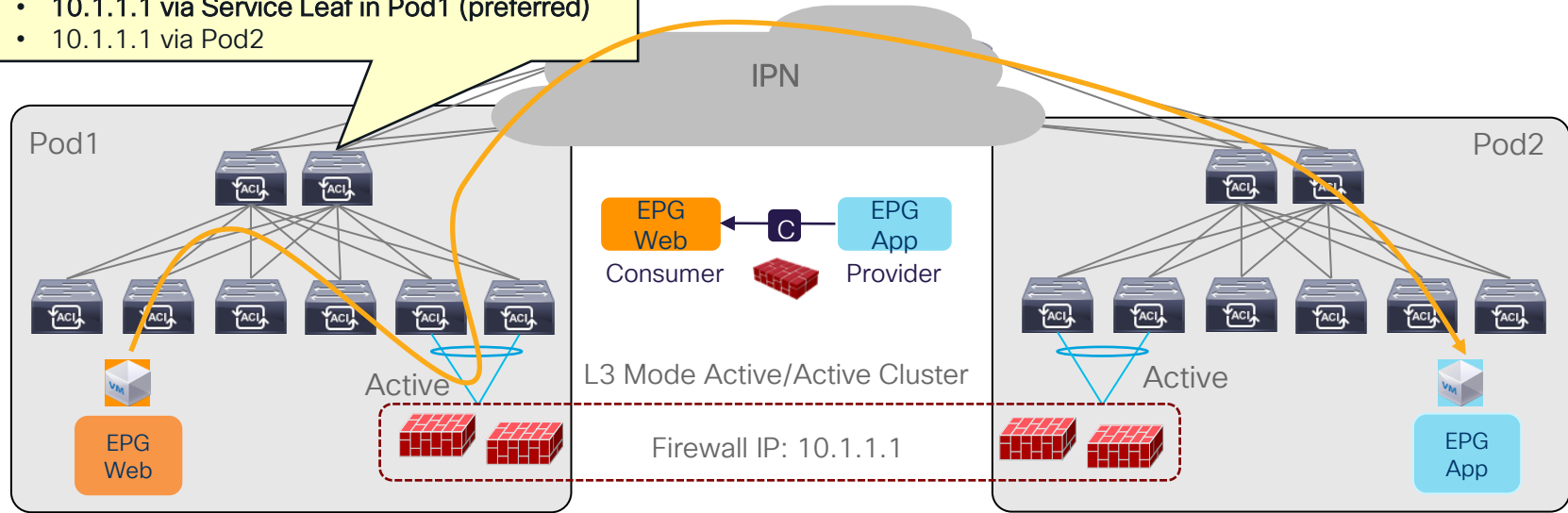


ACI Multi-Pod: Active/Active cluster across pods

East-West Traffic Flow (Inter-Pod) incoming traffic

Spines in Pod1

- 10.1.1.1 via Service Leaf in Pod1 (preferred)
- 10.1.1.1 via Pod2



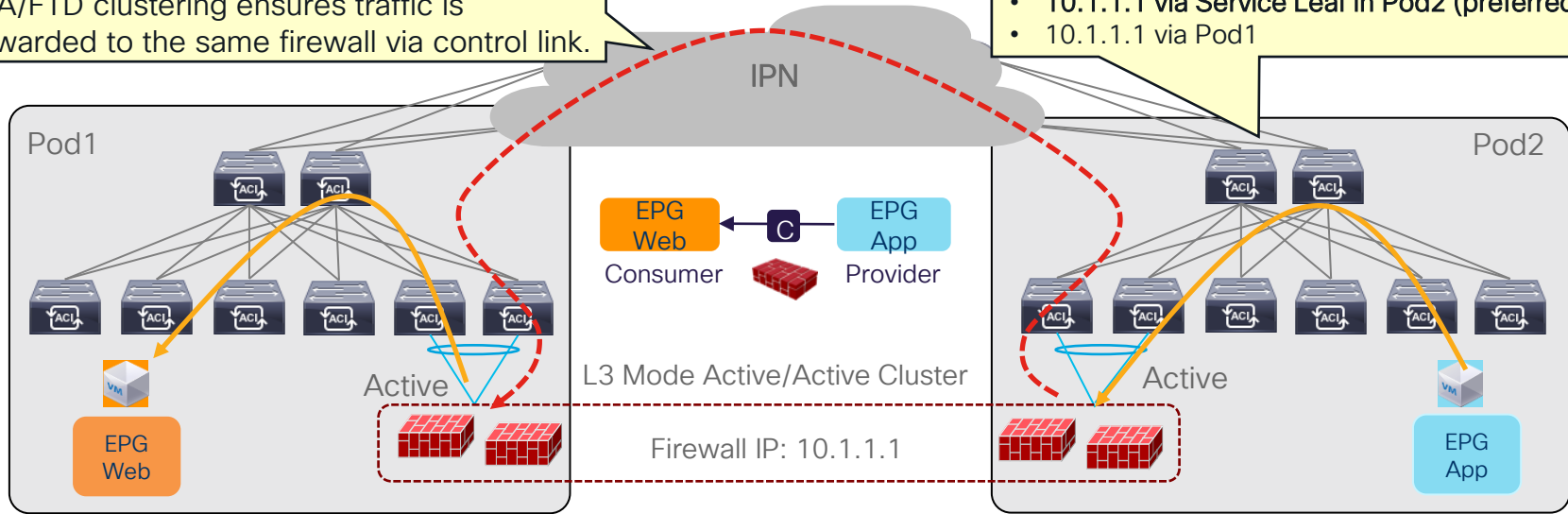
ACI Multi-Pod: Active/Active cluster across pods

East-West Traffic Flow (Inter-Pod) return traffic

Even if asymmetric redirection happens, ASA/FTD clustering ensures traffic is forwarded to the same firewall via control link.

Spines in Pod2

- 10.1.1.1 via Service Leaf in Pod2 (preferred)
- 10.1.1.1 via Pod1

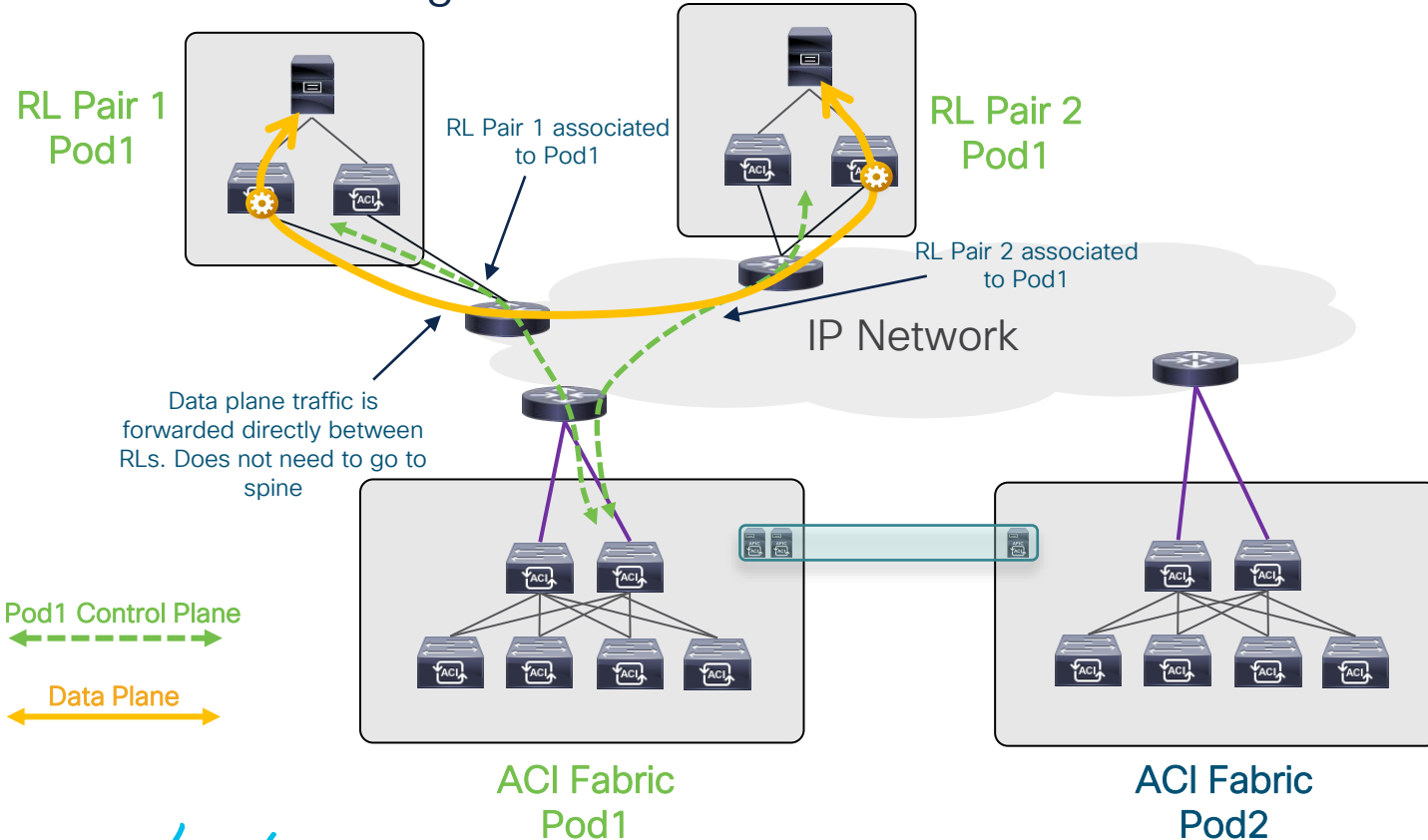


Multi-Pod with Remote Leaf



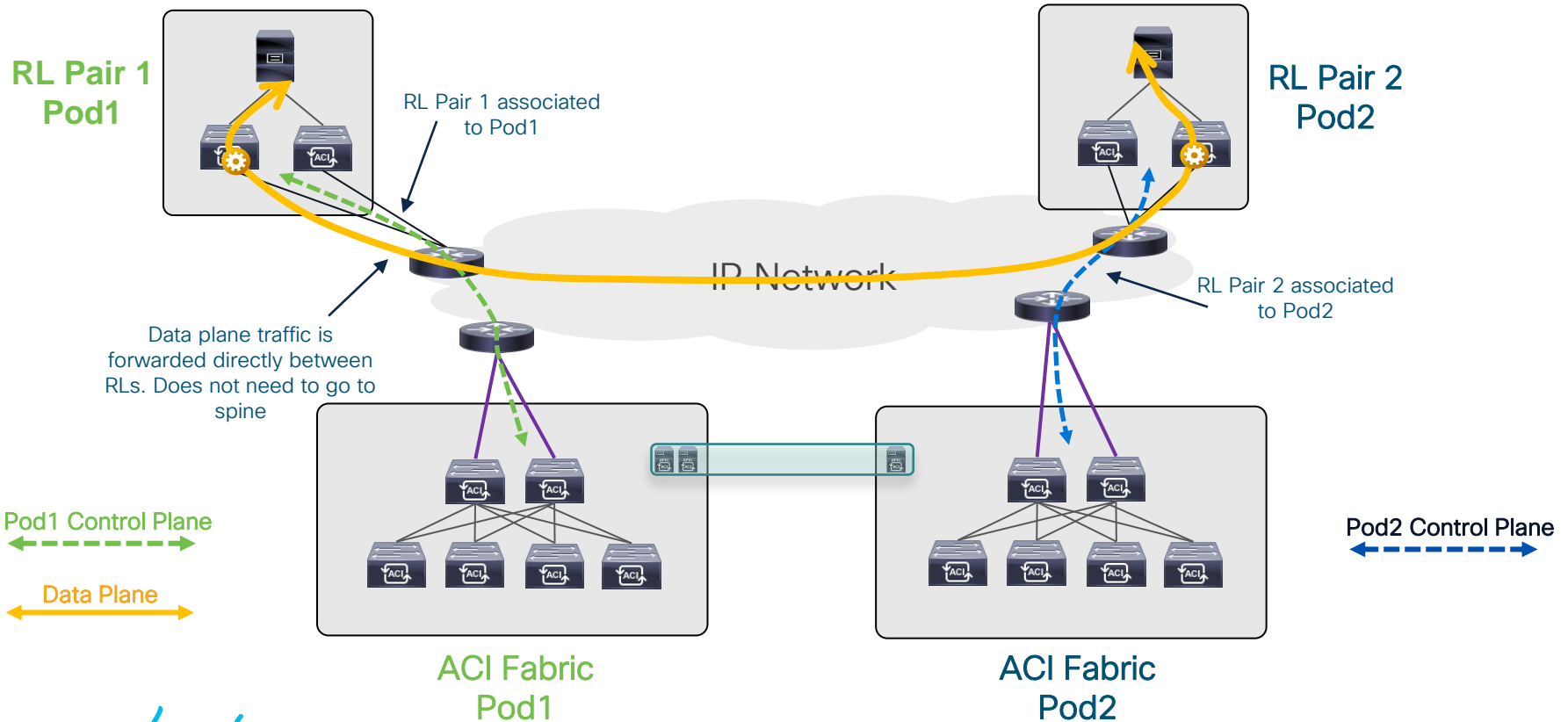
ACI Remote Leaf with Multi-Pod

Direct Forwarding between RL Pairs Part of the Same Pod



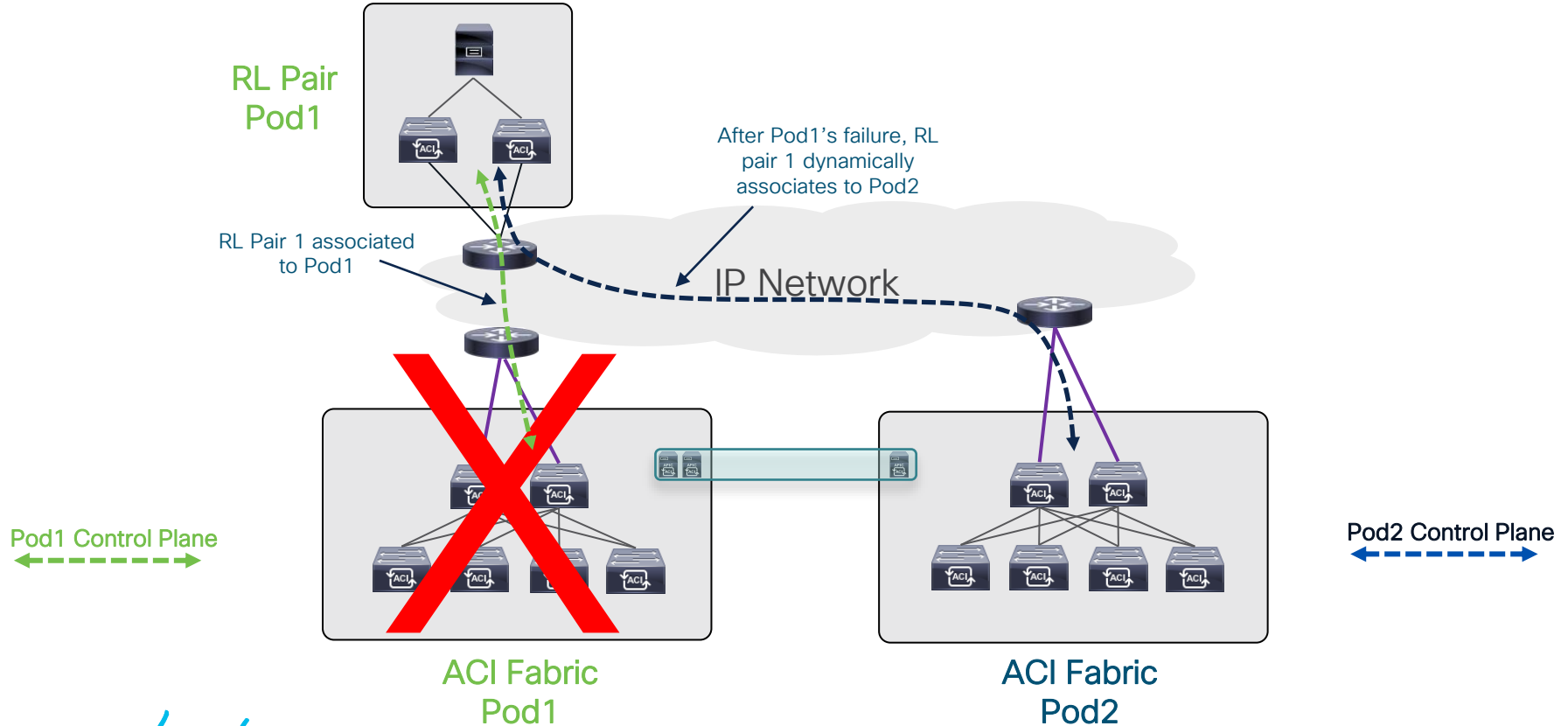
ACI Remote Leaf with Multi-Pod

Direct Forwarding between RL Pairs Part of Different Pods



ACI Remote Physical Leaf

RL Pair Resiliency in a Pod Failure Scenario



Useful Links

- ✓ ACI Multi-Pod White Paper

<https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/white-paper-c11-737855.html>

ACI Multi-Pod Configuration Paper

<https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/white-paper-c11-739714.html>

- ✓ ACI Multi-Pod and Service Node Integration White Paper

<https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/white-paper-c11-739571.html>

- ✓ ACI Remote Leaf Architecture White Paper

<https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/white-paper-c11-740861.html>



The bridge to possible

Thank you

CISCO *Live!*

The Cisco Live! logo features the word "CISCO" in a bold, black, sans-serif font, followed by "Live!" in a black, cursive script font. The background of the entire image is a vibrant, multi-colored abstract pattern of overlapping, wavy bands in shades of red, orange, yellow, green, and blue, creating a sense of motion and energy.

CISCO *Live!*

Let's go