





ACI Troubleshooting

Layer 3 Out (L3Out)

Takuya Kishida - Technical Marketing, DCBU ACI

BRKACI-2642





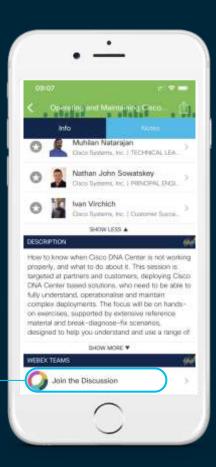
Cisco Webex Teams

Questions?

Use Cisco Webex Teams to chat with the speaker after the session

How

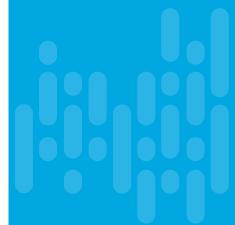
- 1 Find this session in the Cisco Events Mobile App
- 2 Click "Join the Discussion"
- 3 Install Webex Teams or go directly to the team space
- 4 Enter messages/questions in the team space





Agenda

- L3Out Key Components
 - Learning routes (Routing Protocol)
 - Distributing routes within ACI (MP-BGP)
 - Advertising ACI subnet
 - Contract on L3Out (prefix based EPG)
- L3Out Subnet scope options
 - Summary of all options
 - Export Route Control Subnet example

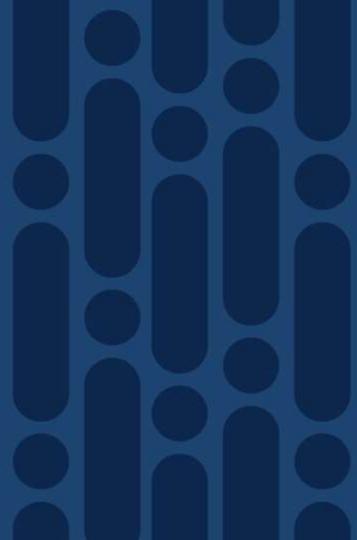


Acronyms/Definitions



Acronyms	Definitions
ACI	Application Centric Infrastructure
APIC	Application Policy Infrastructure Controller
EP	Endpoint
EPG	Endpoint Group
BD	Bridge Domain
VRF	Virtual Routing and Forwarding
L3Out	Layer 3 Out (External Routed Network)
L3Out EPG	Layer 3 Out EPG, Prefix Based EPG (External Network Instance)
MP-BGP	Multi Protocol BGP
VPNv4	Virtual Private Network Version 4
RT	Route Target
RD	Route Distinguisher







What is EPG for ?

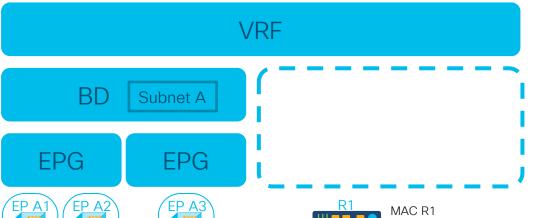
BRKACI-2642

> Endpoint (EP) = MAC & /32 IP (or /128)

What is BD Subnet for ?

- > To be a default gateway
- For ACI Fabric to know a subnet for EPs in a BD

This is for Spine-Proxy Please check BRKACI-3545 for details

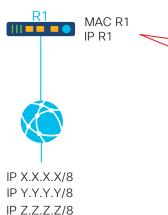






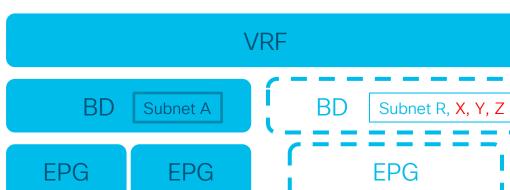


IP A4



A network device (ex. router, loadbalancer) as an endpoint?





MAC A3

IP A3 IP A4

- What is EPG for ?
 - Endpoint (EP) = MAC & /32 IP (or /128)
- What is BD Subnet for ?
 - > To be a default gateway
 - For ACI Fabric to know a subnet for EPs in a BD

MAC R1

IP Z1 - Z999

IP_{R1} IP X1 - X999 IP Y1 - Y999

IP X.X.X.X/8

IP Y.Y.Y.Y/8 IP Z.Z.Z.Z/8

A network device as an endpoint?

➤ All IPs as /32 in a single endpoint

endpoint	vlan	84		
vlan-5		0000.0000.R1R1	L	po3
vlan-5		R.R.R.1	L	po3
vlan-5		X.X.X.1	L	po3
vlan-5		X.X.X.2	L	po3
vlan-5		X.X.X.3	L	po3
vlan-5		Y.Y.Y.1	L	po3
vlan-5		Z.Z.Z.1	L	po3
	vlan-5 vlan-5 vlan-5 vlan-5 vlan-5	vlan-5 vlan-5 vlan-5 vlan-5 vlan-5	vlan-5 R.R.R.1 vlan-5 X.X.X.1 vlan-5 X.X.X.2 vlan-5 X.X.X.3 vlan-5 Y.Y.Y.1	vlan-5 0000.0000.R1R1 L vlan-5 R.R.R.1 L vlan-5 X.X.X.1 L vlan-5 X.X.X.2 L vlan-5 X.X.X.3 L vlan-5 Y.Y.Y.1 L vlan-5 Y.Y.Y.1 L

X One endpoint can have up to 1024 IPs in ACI

This does not scale and efficient. No need to learn each IP as /32.



MAC A2

IP A2

MAC A1

IP A1



po3

What is L3Out for ?

To connect ACI with other network domain = devices with multiple subnet behind it

How is L3Out different from FPG?

- Speak Routing Protocol
- ➤ No IP learning as endpoint
- ➤ Next-hop IP is stored in ARP table
 - = Same as normal routers



FPG

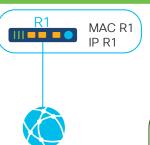
FPG

Subnet A





L3Out EPG



IP X.X.X.X/8

IP Y.Y.Y.Y/8

IP Z.Z.Z.Z/8

Next-hop MAC in endpoint table

leaf1# show endpoint vlan 84

84/TK:VRF1 vxlan-14876665

Next-hop IP in ARP table (only for L3Out)

leaf1# show ip arp vlan 84

Address MAC Address Age Interface R.R.R.1 00:07:51 0000.0000.R1R1 vlan84

0000,0000,R1R1

Other routes via Routing Protocol

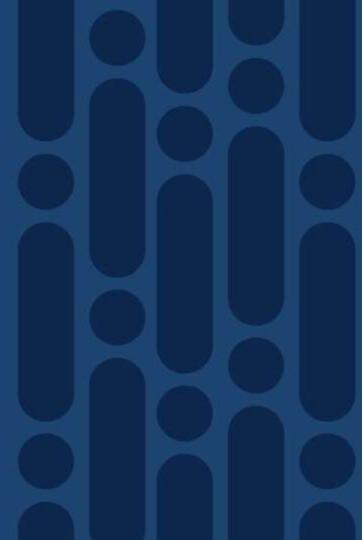
leaf1# show ip route vrf TK:VRF1

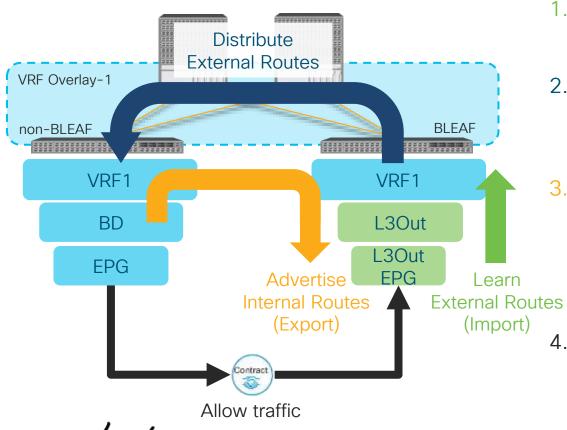
X.0.0.0/8, ubest/mbest: 1/0 *via R.R.R.1, vlan84, [110/5], 2d00h, ospf-default, intra Y.0.0.0/8, ubest/mbest: 1/0

*via R.R.R.1, vlan84, [110/5], 2d00h, ospf-default, intra Z.0.0.0/8, ubest/mbest: 1/0

*via R.R.R.1, vlan84, [110/5], 2d00h, ospf-default, intra







- 1 Learn external routes
 - Routing Protocol in L3Out
- 2. Distribute external routes to other leaves
 - ➤ MP-BGP
- 3. Advertise internal routes (BD subnet) to outside
 - Redistribution and
 - Contract
- 4 Allow traffic with contracts
 - ➤ L3Out EPG (Prefix Based EPG) and
 - Contract

1. Learn External Routes = Routing Protocol

Configurations

External Routed Networks (L3Out)

- VRF to deploy Routing Protocol
- Routing Protocol parameters ex. OSPF area 0.0.0.1 nssa

Node Profile

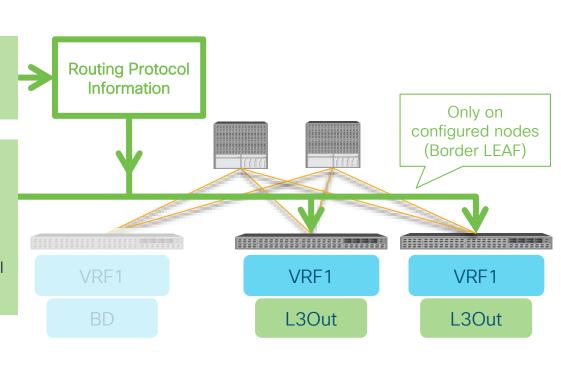
- Node(s) to deploy Routing Protocol
- Static Route (if any)

Interface Profile

- I/F(s) to deploy Routing Protocol
- Routing Protocol I/F parameters ex. OSPF hello interval

Networks (L3Out EPG)

- Contract
- Advanced Route Control ex. route-map



Details for L3Out EPG are in later sections

Verification Examples (OSPF)

1. Is OSPF enabled on a correct I/F?

```
border-leaf# show ip ospf int bri vrf TK:VRF1
 Interface
                                            Cost State Neighbors Status
                       ΙD
 Vlan58
                       134
                                                    BDR
                              backbone
                                                                      up
border-leaf# show vlan id 58 extended
 VLAN Name
                                     Encap
                                              Ports
     TK:VRF1:13out-
                                    vxlan-15695748, Eth1/3, Po2
     L3OUT OSPF:vlan-1425
                                    vlan-1425
```

Same CLI verifications are as useful in ACI too

If anything is not as expected, check config or any faults in APIC GUI.

2. Are OSPF parameters matching with neighbors?

```
border-leaf# show int vlan 58 | grep MTU

MTU 1500 bytes, BW 10000000 Kbit, DLY 1 usec

border-leaf# show ip ospf int vlan 58 | egrep 'IP|State|Timer|auth'

IP address 15.0.0.3/24, Process ID default VRF TK:VRF1, area backbone State BDR, Network type BROADCAST, cost 4

Timer intervals: Hello 10, Dead 40, Wait 40, Retransmit 5

No authentication
```

Is MTU matching?

Is Network Mask matching?

Is Area matching?

Is Timer matching?

Is Network Type expected? etc.

3. Are OSPF neighbors established correctly?

border-leaf#	show ip ospf neighbors	vrf TK:VR	F1	
Neighbor ID	Pri State	Up Time	Address	Interface
4.4.4.4	1 FULL/DR	2d06h	15.0.0.4	Vlan58
9.9.9.9	1 FULL/DROTHER	2d06h	15.0.0.1	Vlan58

Can they ping to each other?

leaf# iping -V <VRF> <target IP>

**OSPF DBD requires unicast reachability etc.

Verification Examples (EIGRP)



1. Is EIGRP enabled on a correct I/F?

border-leaf# show ip eigrp int bri vrf TK:VRF1									
		Xmit Queue	Mean	Pacing Time	Multicast	Pending			
Interface	Peers	Un/Reliable	SRTT	Un/Reliable	Flow Timer	Routes			
vlan92	2	0/0	1	0/0	50	0			
border-leaf# s	show vlan	id 92 extende	ed						
VLAN Name			Enca	ip	Ports	<u>/</u>			

Same CLI verifications are as useful in ACI too

If anything is not as expected, check config or any faults in APIC GUI.

2. Are EIGRP parameters matching with neighbors?

```
border-leaf# show int vlan 92 | grep MTU
    MTU 1500 bytes, BW 10000000 Kbit, DLY 1 usec

border-leaf# show ip int vlan 92 | grep 'IP addr'
    IP address: 16.0.0.3, IP subnet: 16.0.0.0/24
```

border-leaf# show ip eigrp vrf TK:VRF1 | egrep 'AS|K'
IP-EIGRP AS 1 ID 3.3.3.3 VRF TK:VRF1
Metric weights: K1=1 K2=0 K3=1 K4=0 K5=0

Is MTU matching?
Is Network Mask matching?
Is AS matching?
Is K value matching?

3. Are EIGRP neighbors established correctly?

bor	der-leaf#	show ip	eigrp	neighbors vrf	TK:VRF1					
Н	Address			Interface	Hold	Uptime	SRTT	RTO	Q	Seq
					(sec)		(ms)		Cnt	Num
0	16.0.0.4			vlan92	12	2d06h	1	50	0	10
1	16.0.0.1			vlan92	13	2d06h	1	50	0	346

etc.

Verification Examples (BGP)



1. Is BGP neighbor session configured as expected?

```
border-leaf# show ip bgp neighbors vrf TK:VRF1 | egrep 'BGP nei|Using|Opens|hops'
BGP neighbor is 17.0.0.1, remote AS 65001, ebgp link, Peer index 1
  Using Loopback6 as update source for this peer
  External BGP peer might be upto to 2 hops away
  Opens:
border-leaf# show ip int lo6 | grep 'IP addr'
  IP address: 3.3.3.3, IP subnet: 3.3.3.3/32
```

Is it correct remote AS? Is it using correct source I/F with correct IP? Is enough multi-hop configured for eBGP? Is Open message exchanged?

2. Is there IP reachability?

```
border-leaf# iping -V TK:VRF1 17.0.0.1 -S 3.3.3.3
PING 17.0.0.1 (17.0.0.1) from 3.3.3.3: 56 data bytes
64 bytes from 17.0.0.1: icmp seq=0 ttl=255 time=0.76 ms
64 bytes from 17.0.0.1: icmp seq=1 ttl=255 time=0.639 ms
   === snip ===
--- 17.0.0.1 ping statistics ---
5 packets transmitted, 5 packets received, 0.00% packet loss
```

Is there an IP reachability to the BGP neighbor from the correct source IP?

Is it receiving BGP routes?

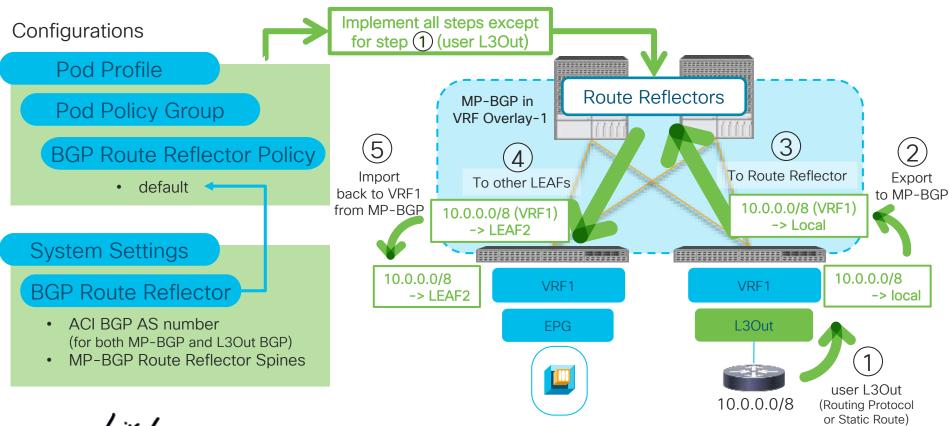
Is ACI BGP using expected local AS?

3. Are BGP neighbors established correctly?

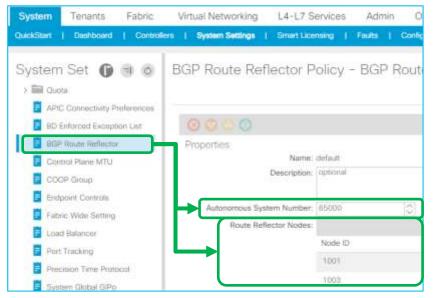
```
border-leaf# show ip bgp summary vrf TK:VRF1
BGP router identifier 3.3.3.3, local AS number 65003
```

Neighbor V AS MsqRcvd MsqSent TblVer InQ OutQ Up/Down State/PfxRcd 17.0.0.1 3300 3302 2d06h

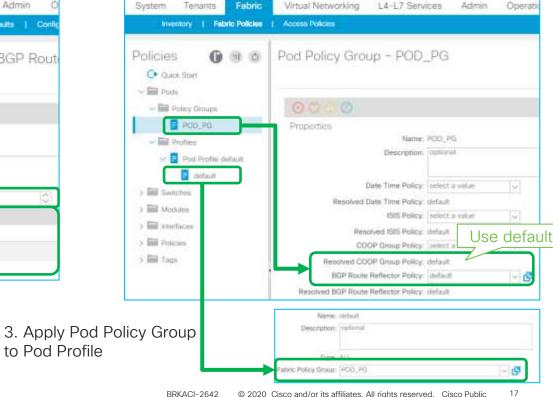
2. Distribute External Routes = MP-BGP in infra



- 2 Distribute External Routes = MP-BGP in infra
- 1. Select ACI BGP AS and Route Reflector SPINEs



2. Apply Route Reflector policy to Pod Policy Group



X L3Out BGP share this same AS with the internal MP-BGP

to Pod Profile

CLI Verification

1. Do both border leaf and non-border leaf have BGP sessions with RR spines?

```
leaf# show bgp sessions vrf overlay-1
Neighbor
                       Flaps LastUpDn|LastRead|LastWrit St Port(L/R) Notif(S/R)
                ASN
10.0.184.65
                             2d07h
                                      Inever
                                               Inever
                                                         E 37850/179
10.0.184.66
                             2d07h
                                                         E 45089/179 0/0
                                               Inever
                                      Inever
leaf# acidiag fnvread | grep spine
    1001
                                    FGE10000000
                                                     10.0.184.65/32
                                                                      spine
                         spine1
                                                                                    active
    1002
                                    SAL10000000
                                                                      spine
                        spine2
                                                                                    active
```

2. Is the external route learned on a border leaf?

```
border-leaf# show ip route vrf TK:VRF1

10.0.0.0/8, ubest/mbest: 1/0

*via 15.0.0.1, Vlan58, [110/5], 2d08h, ospf-default, intra
```

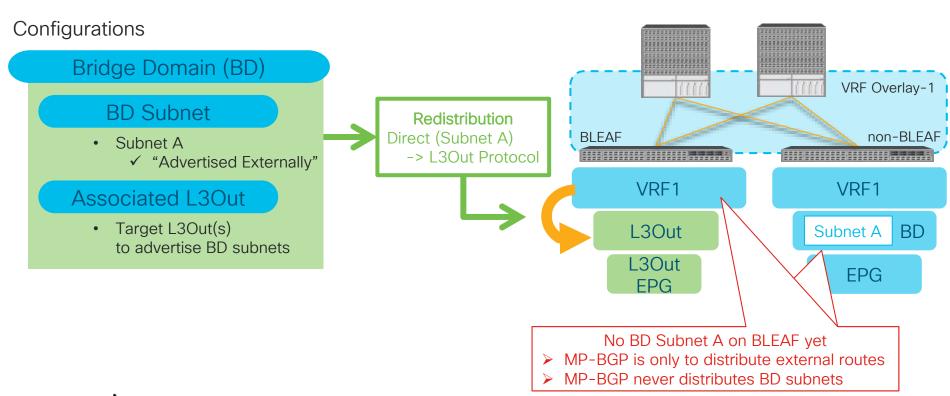
✓ BGP neighbors are RR spines TEP IPs

3. Does non-border leaf show the expected border leaf as next-hop?

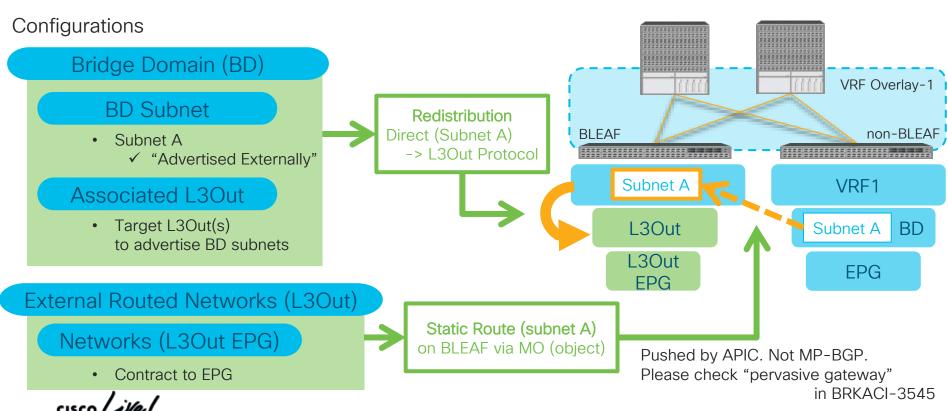
✓ Next-hops are border Leaf TEP IPs✓ Learned via iBGP in ACI AS# (65003)

```
non-border-leaf# show ip route vrf TK:VRF1
10.0.0.0/8, ubest/mbest: 2/0
    *via 10.0.184.67%overlay-1, [200/5], 2d08h, bqp-65003,
                                                                       tag 65003
                                                             internal,
    *via 10.0.184.64%overlay-1, [200/5], 2d08h, bgp-65003
                                                                       tag 65003
non-border-leaf# acidiag fnvread
           Pod ID
                                           Serial Number
                                   Name
                                                                 IP Address
                                                                                Role
                                                                                             State
LastUpdMsgId
     103
                                             SAL10000003
                                                              10.0.184.64/32
                                                                                leaf
                                                                                              active
     104
                                             SAL10000004
                                                                                leaf
                                                                                              active
```

3. Advertise BD subnet

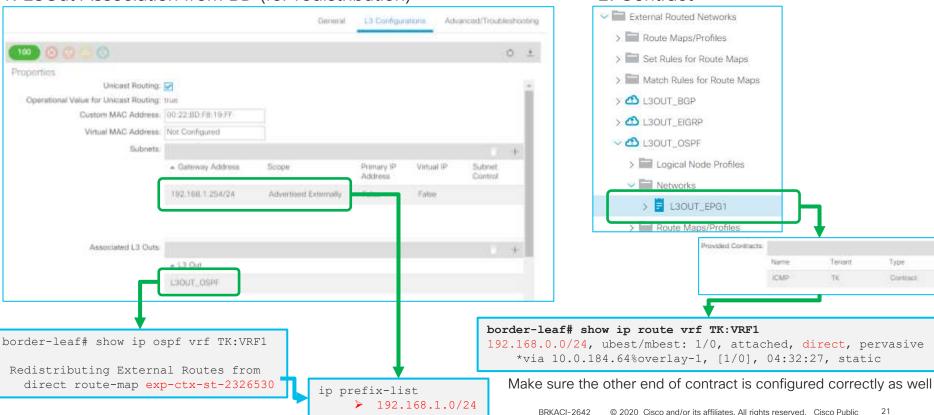


3. Advertise BD subnet



3 Advertise BD subnet

1. L3Out Association from BD (for redistribution)



2. Contract

CLI Verification (OSPF, EIGRP)

1. Does the border leaf have BD subnet to advertise?

```
border-leaf# show ip route vrf TK:VRF1

192.168.1.0/24, ubest/mbest: 1/0, attached, direct,
pervasive

*via 10.0.184.64%overlay-1, [1/0], 04:32:27, static
```

If not, check the contract between the L3Out EPG and the EPG for the BD.

This should be pushed by APIC. Not via MP-BGP.

2. Check a route-map name used by the routing protocol on the border leaf for redistribution

```
border-leaf# show ip ospf vrf TK:VRF1

Redistributing External Routes from direct route-map exp-ctx-st-2097152
```

```
border-leaf# show ip eigrp vrf TK:VRF1

Redistributing:
direct route-map exp-ctx-st-2097152
```

Check next page for BGP

3. Does the route-map have expected BD subnet?

```
border-leaf# show route-map exp-ctx-st-2097152
route-map exp-ctx-st-2097152, deny, sequence 1
Match clauses:
   tag: 4294967295
Set clauses:
route-map exp-ctx-st-2097152, permit, sequence 15804
Match clauses:
   ip address prefix-lists: IPv4-st49158-2097152-exc-int-inferred-export-dst
   ipv6 address prefix-lists: IPv6-deny-all
Set clauses:
```

IP prefix-list should have the BD subnet.

If not, check APIC config and any faults.

- ✓ Is "Advertise Externally" on the BD subnet checked?
- ✓ Is L3Out associated to the BD?

cisco Live!

border-leaf# show ip prefix-list IPv4-st49158-2097152-exc-int-inferred-export-dst ip prefix-list IPv4-st49158-2097152-exc-int-inferred-export-dst: 1 entries

seq 1 permit 192.168.1.254/24

CLI Verification (BGP)

1. Does the border leaf have BD subnet to advertise?

```
--- snip ---
```

2. Check a route-map name used by BGP outbound rule for each neighbor

```
border-leaf# show bgp process vrf TK:VRF1
Information for address family IPv4 Unicast in VRF TK:VRF1
Redistribution
direct, route-map permit-all
```

BGP redistributes all direct routes first, then limit the routes with an outbound route-map.

```
border-leaf# show ip bgp neighbors vrf TK:VRF1 | egrep '^BGP|Out'

BGP neighbor is 17.0.0.1, remote AS 65001, ebgp link, Peer index 1
   Outbound route-map configured is exp-l3out-L3OUT_BGP-peer-2097152, handle obtained
```

3. Does the BGP outbound route-map have the expected BD subnet?

```
border-leaf# show route-map exp-l3out-L3OUT_BGP-peer-2097152
route-map exp-l3out-L3OUT_BGP-peer-2097152, permit, sequence 15801
Match clauses:
   ip address prefix-lists: IPv4-peer49157-2097152-exc-int-inferred-export-dst
   ipv6 address prefix-lists: IPv6-deny-all
   Set clauses:
route-map exp-l3out-L3OUT_BGP-peer-2097152, deny, sequence 16000
Match clauses:
   route-type: direct
Set clauses:
```

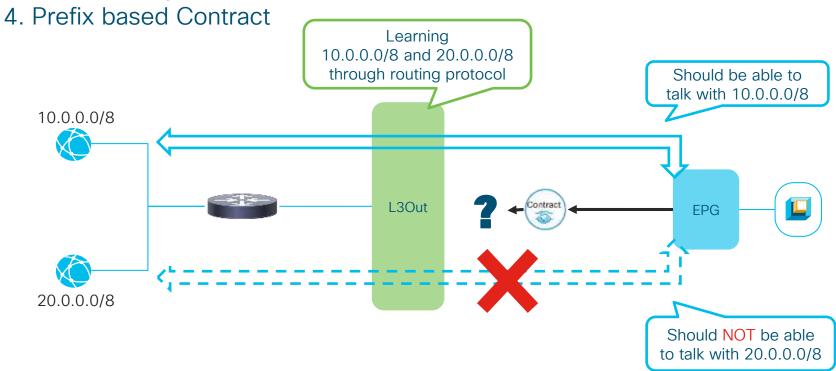
IP prefix-list should have the BD subnet.

If not, check APIC config and any faults.

- ✓ Is "Advertise Externally" on the BD subnet checked?
- ✓ Is L3Out associated to the BD?

cisco Live!

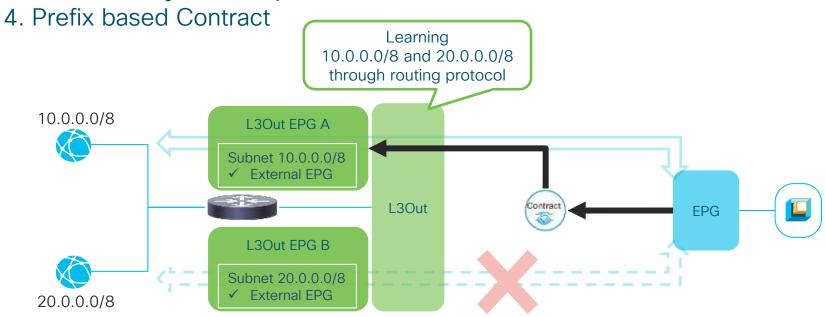
border-leaf# show ip prefix-list IPv4-peer49157-2097152-exc-int-inferred-export-dst
ip prefix-list IPv4-peer49157-2097152-exc-int-inferred-export-dst: 1 entries
 seq 1 permit 192.168.1.254/24



How do we accomplish this ??

BRKACI-2642

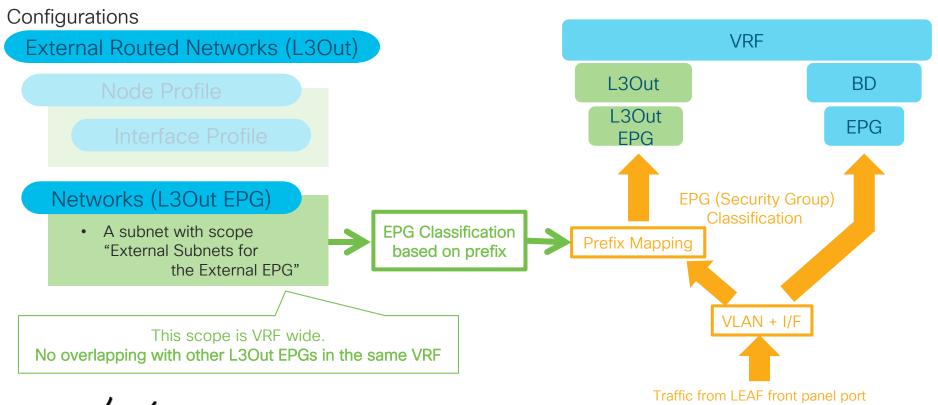




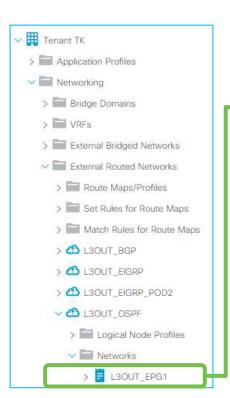
Prefix Based EPG (= L3Out EPG)

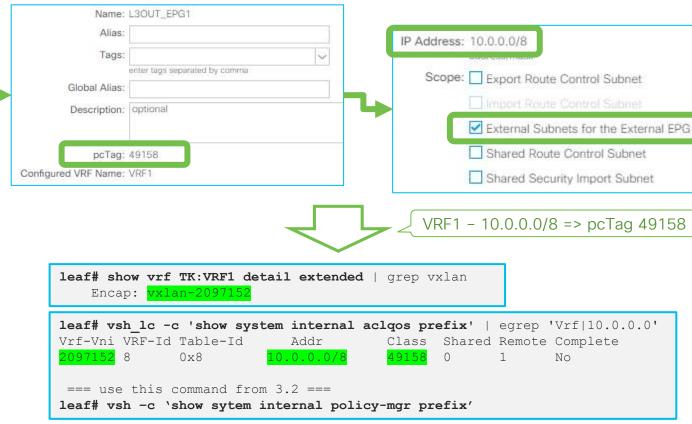


4. Prefix based Contract



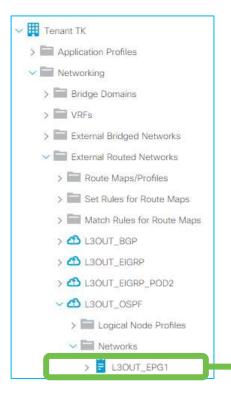
4. Prefix based Contract

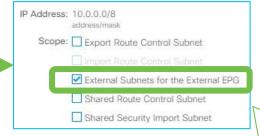






4. Prefix based Contract







> To create prefix to pcTag mapping

NOTF:

It has nothing to do with routing table or routing protocol behavior unlike other Route Control Subnet scopes

A common mistake is selecting both "External Subnets for the External EPG" and "Export Route Control Subnet" for the same subnet, which implies a conflicting situation where the subnet behind the L3Out but the same L3Out is also expected to advertise/redistribute the subnet back to where it came from. It may not cause an immediate issue but unnecessary redistribution should always be avoided.

Check L3Out Subnet scope section for details.



CLI Verifications

Contract Drop on this leaf shows up in this command.

1. Check if there is any contract drops

Check both ingress/egress leaf just in case, or see appendix for Policy Control Enforcement Direction

192.168.1.1

leaf# show logging ip access-list internal packet-log deny [Wed May 8 18:34:31 2019 155907 usecs]: CName: TK:VRF1 (VXLAN: 2719744), VlanType: FD VLAN, Vlan-Id: 26, SMac: 0x0050569185d1, DMac:0x0022bdf819ff, SIP: 192.168.1.1, DIP: 10.0.0.1, SPort: 58968, DPort: 80, Src Intf: port-channel1, Proto: 6, PktLen: 74

2. Check VRF VNID

leaf# show vrf TK:VRF1 detail extended | grep vxlan Encap: vxlan-2097152

3. Check source (or destination) EPG pcTag

```
Vlan id : 30 ::: Vlan vnid : 9025 ::: VRF name : TK:VRF1
BD vnid : 16318374 ::: VRF vnid : 2097152
Flags: 0x80005c04 ::: sclass: 49100 ::: Ref count: 5
EP Flags: local|IP|MAC|host-tracked|sclass|timer|
```

pcTag/contract is per VRF except for shared service (VRF route leaking)

If your source/destination is an endpoint, it should be in here

sclass = pcTag = EPG ID for contract

Make sure the external IP is not here. This pcTag takes precedence over "prefix-pcTag mapping table". If it is, check the traffic path that caused ACI to learn the external IP as an endpoint.

4. Check destination (or source) L3Out prefix based EPG pcTag

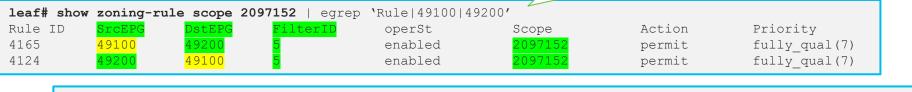
leaf# show system internal epm endpoint ip 192.168.1.1 | egrep

```
leaf# vsh lc -c 'show system internal aclgos prefix' | egrep 'Vrf|10.0.0.0'
Vrf-Vni VRF-Id Table-Id
                              Addr
                                          Class Shared Remote Complete
                                                                                "External Subnet for the External EPG"
                          10.0.0.0/8
               0x8
                                                                No
                                                                                       config is reflected here.
 === use this command from 3.2 ===
                                                                                    This is Longest Prefix Match.
leaf# vsh -c 'show sytem internal policy-mgr prefix'
```

CLI Verifications

L3Out **EPG** 192.168.1.1 10.0.0.0/8 49100 49200

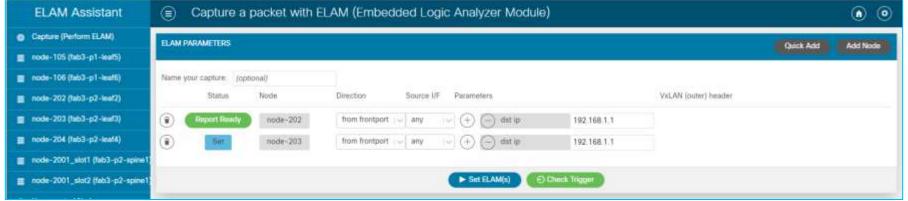
5. Check contracts between two pcTags



```
leaf# show zoning-filter filter 5
FilterId
         Name
                EtherT
                                            ~snip~ SFromPort
                                                                                                   ~snip~
                       Arp0pc
                                    Prot
                                                               SToPort
                                                                           DFromPort.
                                                                                       DToPort.
                                                                                                   ~snip~
                                            ~snip~ =====
                                            ~snip~ unspecified unspecified unspecified ~snip~
          5 0
                ip
                        unspecified icmp
```

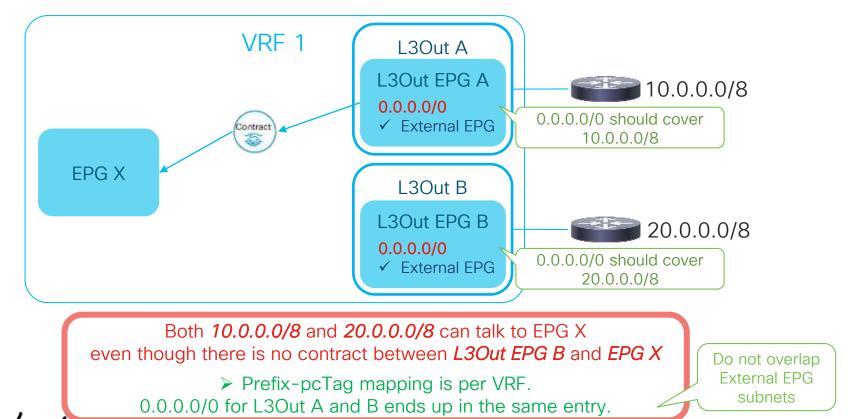
scope = VRF VNID

6. Check ELAM to see if the traffic is using correct src pcTag and dst pcTag



L3Out Contract

Common Issue (L3Out EPGs with 0.0.0.0/0)



L3Out Contract

Common Issue (L3Out EPGs with 0.0.0.0/0)



1. Check VRF VNID

leaf# show vrf TK:VRF1 detail extended | grep vxlan
Encap: vxlan-2097152

2. Check source (or destination) EPG pcTag

leaf# show system internal epm endpoint ip 192.168.1.1 | egrep 'VR
Vlan id : 30 ::: Vlan vnid : 9025 ::: VRF name : TK:VRF1
BD vnid : 16318374 ::: VRF vnid : 2097152
Flags : 0x80005c04 ::: sclass : 49100 ::: Ref count : 5

3. Check destination L3Out 0.0.0.0/0 EPG pcTag

"0.0.0.0/0 -> 15" is the only pcTag entry in this VRF.

➤ Both L3Out A & B will share it since there is no other granular LPM entries

NOTE:

Scope

2097152

- 0.0.0.0/0 always use pcTag 15
- This is not a routing table. It doesn't matter even if the routing table has more granular routes

leaf# vsh_lc -c 'show system internalaclqos prefix' | egrep 'Vrf|2097152'Vrf-Vni VRF-Id Table-IdAddrClass Shared Remote Complete209715280x80.0.0.0/01500No

4. Check contracts between pcTags

leaf# showzoning-rulescope2097152egrep'Rule|49162'Rule IDSrcEPGDstEPGFilterIDoperSt416549100155enabled

This contract is due to "EPG X <-> L3Out A"

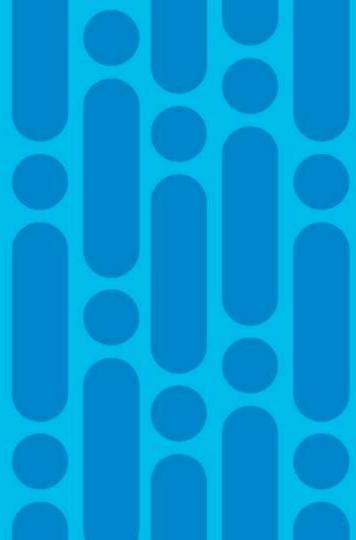
But any traffic that hits 0.0.0.0/0 in the prefix table

can use this rule

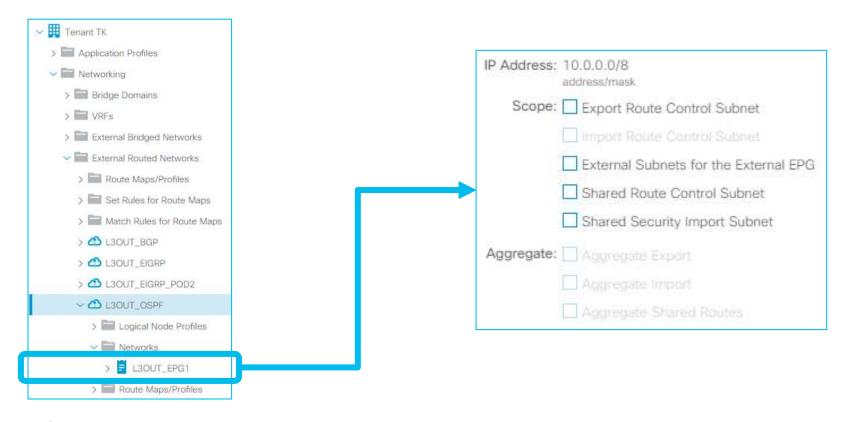
Action Priority permit fully_qual(7)



L3Out Subnet Scope



L3Out Subnet Scope





L3Out Subnet Scope

Route Control for Routing Protocol

- **Export Route Control Subnet**
- Import Route Control Subnet
- Shared Route Control Subnet

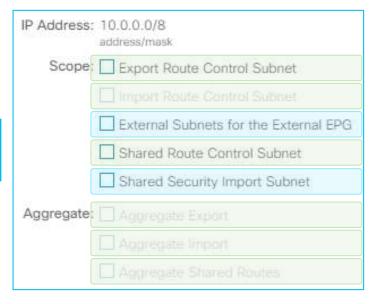
Traffic Classification for Contract

- External Subnets for the External EPG
- Shared Security Import Subnet

<u>Aggregate</u>

- Aggregate Export
- Aggregate Import
- Aggregate Shared Routes

Grouping by functionality



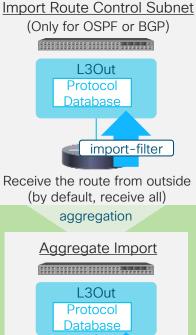
L3Out Subnet Scope Summary

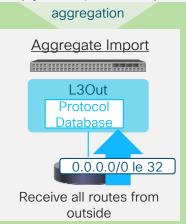
Route Control for Routing Protocol

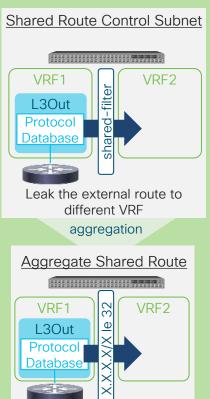
Export Route Control Subnet (Mainly for Transit Routing) L3Out Protocol Database export-filter Advertise the route from ACI

to outside aggregation





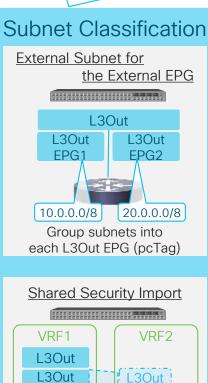




Leak multiple external

routes to different VRF





EPG1

Leak prefix-pcTag mapping

to different VRF

EPG₁

10.0.0.0/8

Route Control Enforcement

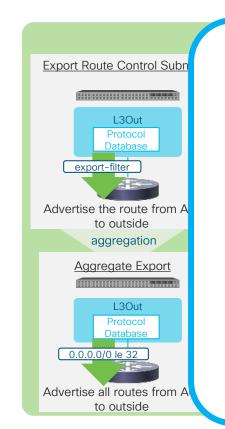


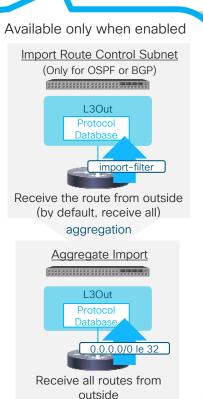
Import is disabled by default.

Receive all routes by default.

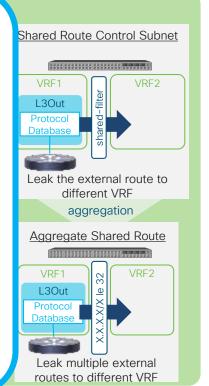
No import route control. Export is always enabled.





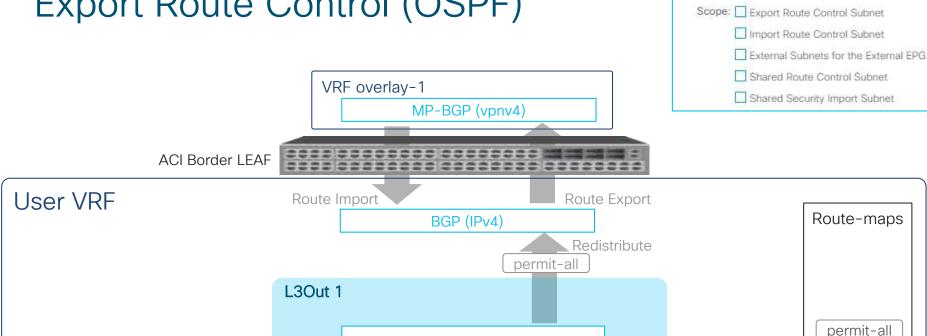


Route Control Enforcement: Import





Export Route Control (OSPF)



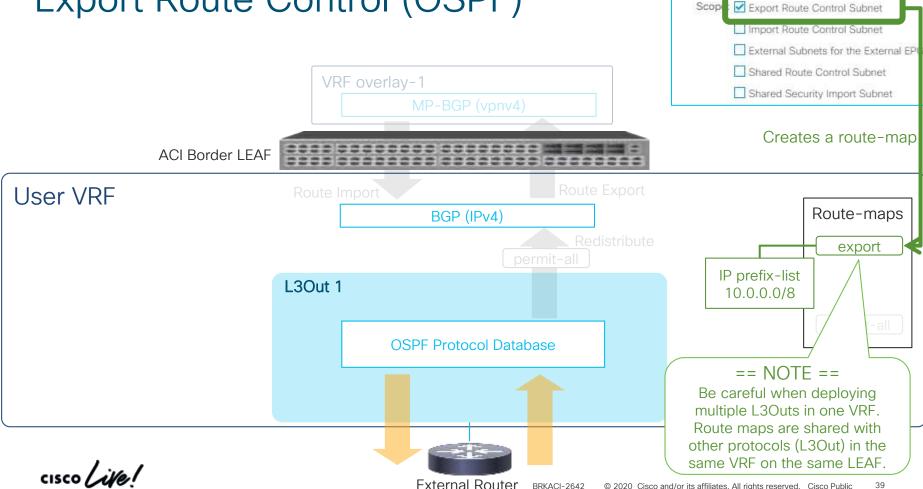
BRKACI-2642

OSPF Protocol Database

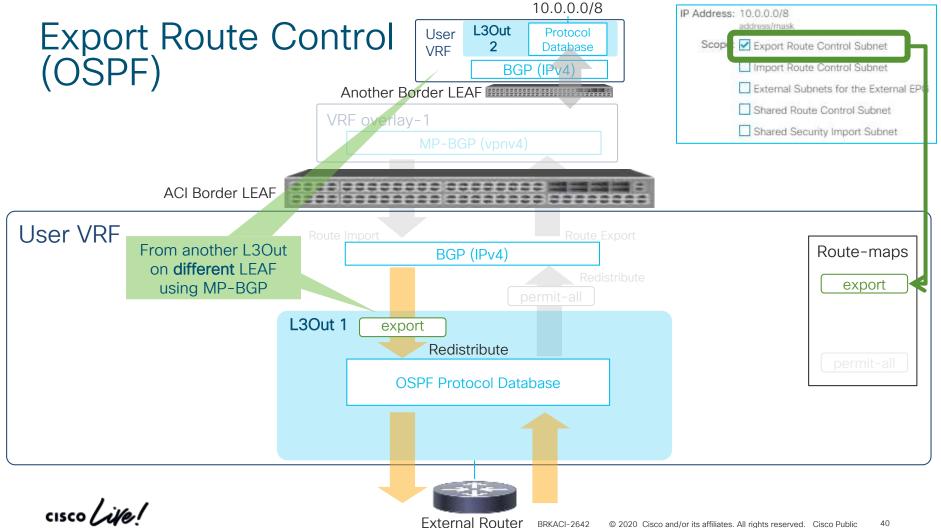
IP Address: 10.0.0.0/8

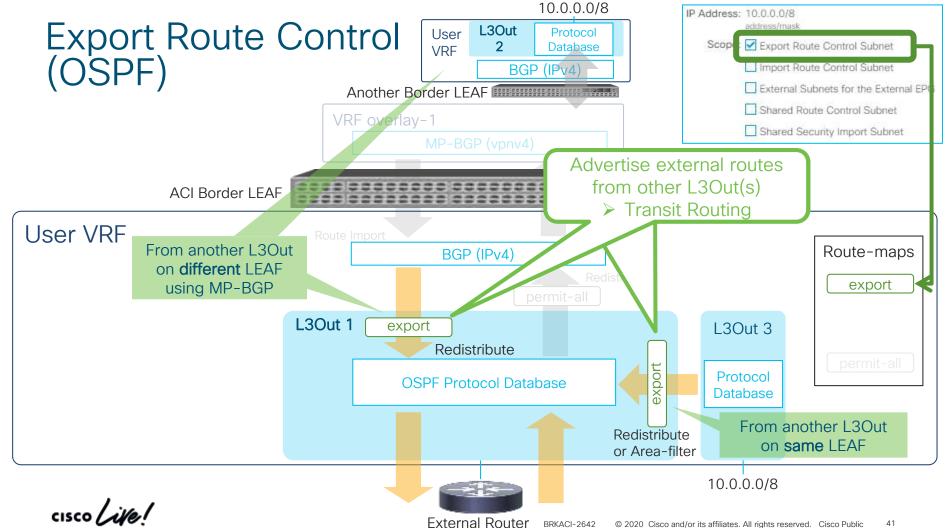
address/mask

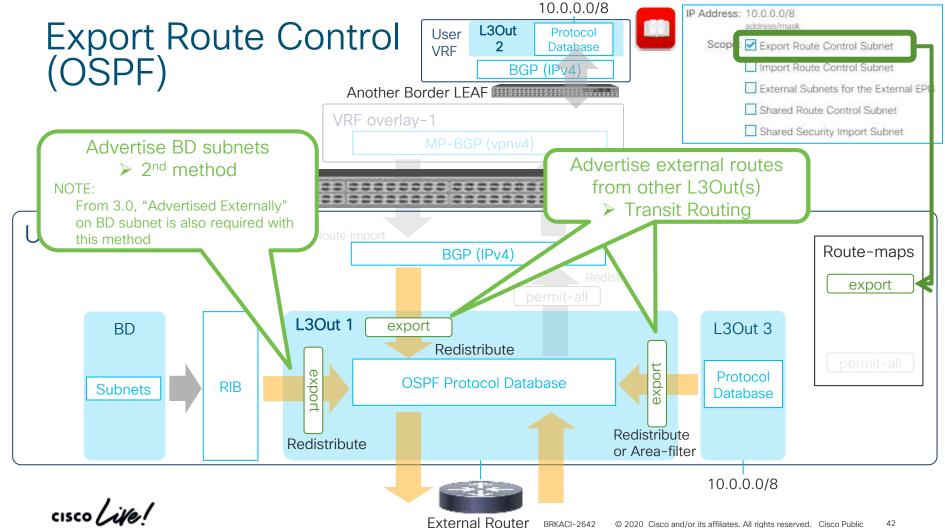
Export Route Control (OSPF)



IP Address: 10.0.0.0/8







CLI Verification (OSPF/EIGRP)

1. OSPF/EIGRP Redistribution route-map

```
border-leaf# show ip ospf vrf TK:VRF1
 Redistributing External Routes from
   static route-map exp-ctx-st-2097152
   direct route-map exp-ctx-st-2097152
   bgp route-map exp-ctx-proto-2097152
   eigrp route-map exp-ctx-proto-2097152
   Area (backbone)
        Area-filter in 'exp-ctx-proto-2097152'
border-leaf# show ip eigrp vrf TK:VRF1
  Redistributing:
    static route-map exp-ctx-st-2097152
    ospf-default route-map exp-ctx-proto-2097152
    direct route-map exp-ctx-st-2097152
    bgp-65003 route-map exp-ctx-proto-2097152
```

It shares the same route-map with other protocols in the same VRF on the same LEAF route-map naming:

exp-ctx-st-<vrf vnid> or

exp-ctx-proto-<vrf vnid>

EIGRP doesn't support Transit Routing on a same LEAF.

No equivalent filter like OSPF area-filter in EIGRP

All Export Route Control subnet on a same LEAF is added here

Same goes to exp-cxt-st-2097152

```
2. route-map and ip prefix-list
```

tag 4294967295

```
border-leaf# show route-map exp-ctx-proto-2097152
route-map exp-ctx-proto-2097152, permit, sequence 15801
Match clauses:
ip address prefix-lists: IPv4-proto49158-2097152-exc-ext-inferred-export-dst
ipv6 address prefix-lists: IPv6-deny-all
Set clauses:
```

border-leaf# show ip prefix-list IPv4-proto49158-2097152-exc-ext-inferred-export-dst ip prefix-list IPv4-proto49158-2097152-exc-ext-inferred-export-dst: 1 entries seq 1 permit 10.0.0.0/8

CLI Verification (BGP)

BGP has a route-map per L3Out

A bit more granular control

1. BGP outbound route-map

route-map naming: exp-l3out-
bgp l3out name>-peer-<vrf vnid>

All Export Route Control subnets from

the same BGP I 3Out is added here

```
border-leaf# show ip bgp neighbors vrf TK:VRF1

BGP neighbor is 17.0.0.1, remote AS 65001, ebgp link, Peer index 1
```

Outbound route-map configured is exp-13out-L30UT_BGP-peer-2097152, handle obtained

2. route-map and ip prefix-list

```
border-leaf# show route-map exp-13out-L3OUT_BGP-peer-2097152
```

route-map exp-13out-L3OUT_BGP-peer-2097152, permit, sequence 15804

Match clauses:

ip address prefix-lists: IPv4-peer49157-2097152-exc-ext-inferred-export-dst

ipv6 address prefix-lists: IPv6-deny-all

Set clauses:

tag 4294967295

route-map exp-13out-L3OUT BGP-peer-2097152, deny, sequence 16000

Match clauses:

route-type: direct

Set clauses:

border-leaf# show ip prefix-list IPv4-peer49157-2097152-exc-ext-inferred-export-dst

ip prefix-list IPv4-peer49157-2097152-exc-ext-inferred-export-dst: 4 entries

seq 1 permit 10.0.0.0/8



BD Subnet and Export Route Control



border-leaf# show ip route vrf TK:VRF1

192.168.1.0/24, ubest/mbest: 1/0, attached, direct, pervasive *via 11.0.248.0%overlay-1, [1/0], 00:00:05, static, tag 4294967295

"Advertised Externally" removes VRF tag from BD subnet



border-leaf# show ip route vrf TK:VRF1

192.168.1.0/24, ubest/mbest: 1/0, attached, direct, pervasive *via 11.0.248.0%overlay-1, [1/0], 00:00:05, static

Prior to 3.0

From 3.0

border-leaf# show route-map exp-ctx-st-2097152

route-map exp-ctx-st-2097152, permit, sequence 15804

Match clauses:

ip address prefix-lists: IPv4-st49158-2097152-exc-int-inferred-export-dst ipv6 address prefix-lists: IPv6-deny-all

Set clauses:

IP prefix-list from "Export Route Control" for 192.168.1.0/24

border-leaf# show route-map exp-ctx-st-2097152

route-map exp-ctx-st-2097152, deny, sequence 1
Match clauses:

tag: 4294967295

Set clauses:

route-map exp-ctx-st-2097152, permit, sequence 15804

Match clauses:

ip address prefix-lists: IPv4-st49158-2097152-exc-int-inferred-export-dst

ipv6 address prefix-lists: IPv6-deny-all

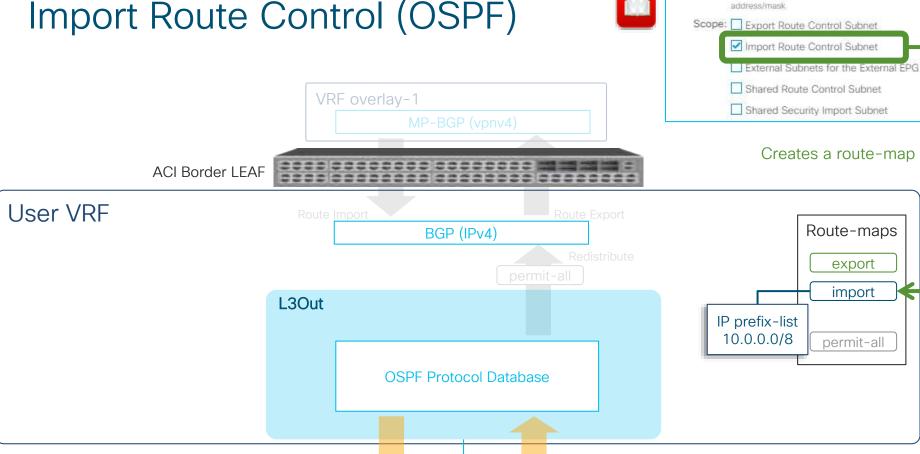
Set clauses:

New rule to prevent BD subnets without "Advertised Externally" from being advertised

cisco Live!

IP prefix-list from "Export Route Control" for 192.168.1.0/24

Import Route Control (OSPF)



External Router

IP Address: 10.0.0.0/8

Import Route Control (OSPF)

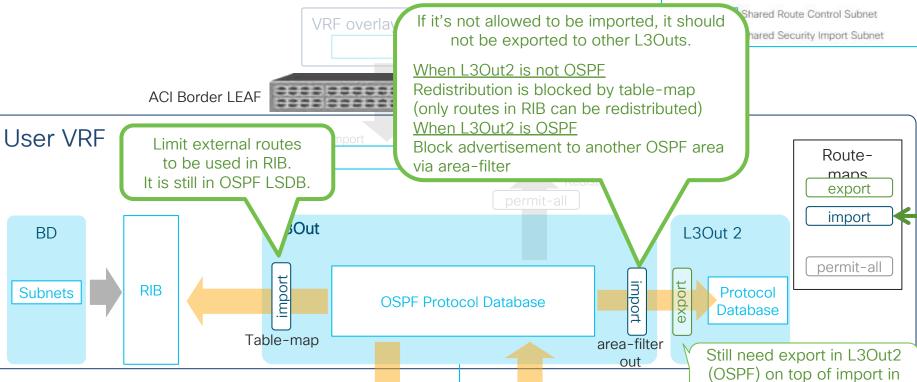


Scope: Export Route Control Subnet

IP Address: 10.0.0.0/8

✓ Import Route Control Subnet

External Subnets for the External EPG



BD

External Router BRKACI-2642

© 2020 Cisco and/or its affiliates. All rights reserved. Cisco Public

L3Out1 (OSPF)

CLI Verification (OSPF)

```
border-leaf# show ip ospf vrf TK:VRF1
Table-map using route-map exp-ctx-2097152-deny-external-tag
Area (backbone)
         Area-filter out 'imp-ctx-ospf-area20971520'
```

- Table-map to prevent the routes from being used in RIB
- "Area-filter out" to prevent the routes from being advertised to another OSPF area on a same LEAF (Transit Routing)

```
border-leaf# show route-map exp-ctx-2097152-deny-external-tag
route-map exp-ctx-2097152-deny-external-tag, deny, sequence 1
  Match clauses:
    tag: 4294967295
  Set clauses:
route-map exp-ctx-2097152-deny-external-tag, permit, sequence 15801
  Match clauses:
    ip address prefix-lists: IPv4-ospf-49158-2097152-exc-ext-inferred-import-dst-rtpfx
    ipv6 address prefix-lists: IPv6-deny-all
    ospf-area: backbone
  Set clauses:
route-map exp-ctx-2097152-deny-external-tag, deny, sequence 19999
  Match clauses:
    ospf-area: backbone
  Set clauses:
route-map exp-ctx-2097152-deny-external-tag, permit, sequence 20000
  Match clauses:
  Set clauses:
```

route-map for table-map

- 1. blocks any routes with VRF tag
- 2. allow routes with Import Route Control subnet in OSPF area X
- 3. block any routes from OSPF area X

A prefix configured by "Import Route Control Subnet"

border-leaf# show ip prefix-list IPv4-ospf-49158-2097152-exc-ext-inferred-import-dst-rtpfx ip prefix-list IPv4-ospf-49158-2097152-exc-ext-inferred-import-dst-rtpfx: 1 entries seq 1 permit 10.0.0.0/8

BRKACI-2642



CLI Verification (OSPF) cont.



border-leaf# show ip ospf vrf TK:VRF1

Table-map using route-map exp-ctx-2097152-deny-external-tag 4

Area (backbone)

Area-filter out 'imp-ctx-ospf-area20971520'

- Table-map to prevent the routes from being used in RIB
- "Area-filter out" to prevent the routes from being advertised to another OSPF area on a same LEAF (Transit Routing)

route-map for area-filter

border-leaf# show route-map imp-ctx-ospf-area20971520

route-map imp-ctx-ospf-area20971520, permit, sequence 15801

Match clauses:

ip address prefix-lists: IPv4-ospf-rt-ospf-import49158-2097152-exc-ext-inferred-import-dst-

ipv6 address prefix-lists: IPv6-deny-all

Set clauses:

border-leaf# show ip prefix-list IPv4-ospf-rt-ospf-import49158-2097152-exc-ext-inferred-import-dst-

ip prefix-list IPv4-ospf-rt-ospf-import49158-2097152-exc-ext-inferred-import-dst-: 1 entries seq 1 permit 10.0.0.0/8

A prefix configured by "Import Route Control Subnet"



CLI Verification (BGP)



BGP uses an inbound route-map (per L3Out) instead of table-map

border-leaf# show ip bgp neighbors vrf TK:VRF1

BGP neighbor is 17.0.0.1, remote AS 65001, ebgp link, Peer index 1

Inbound route-map configured is imp-l3out-L3OUT_BGP-peer-2097152, handle obtained

```
border-leaf1# show route-map imp-l3out-L3OUT_BGP-peer-2097152
```

route-map imp-l3out-L3OUT_BGP-peer-2097152, permit, sequence 15801

Match clauses:

ip address prefix-lists: IPv4-peer49157-2097152-exc-ext-inferred-import-dst

ipv6 address prefix-lists: IPv6-deny-all

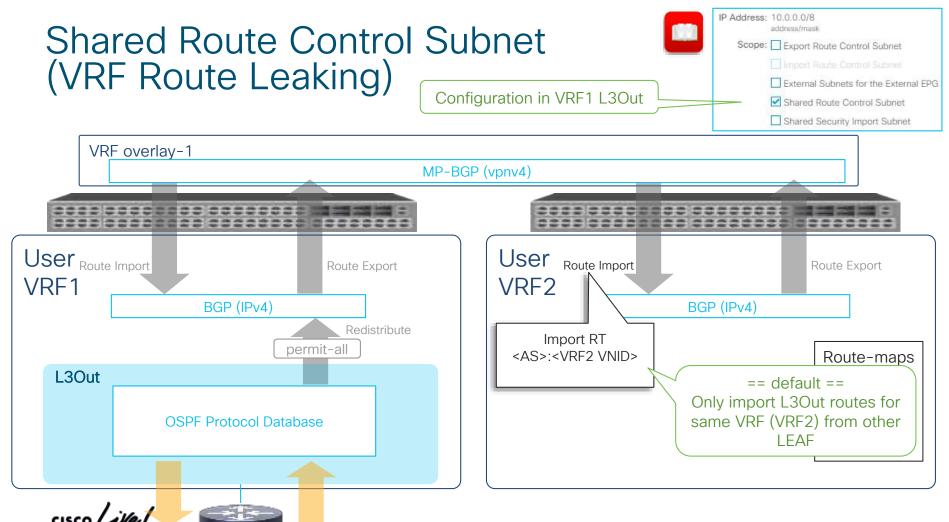
Set clauses:

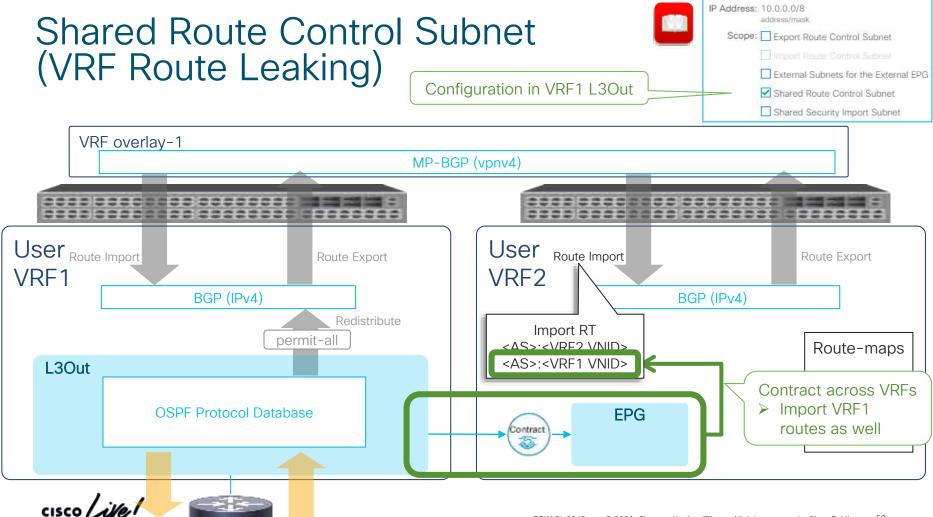
border-leaf# show ip prefix-list IPv4-peer49157-2097152-exc-ext-inferred-import-dst

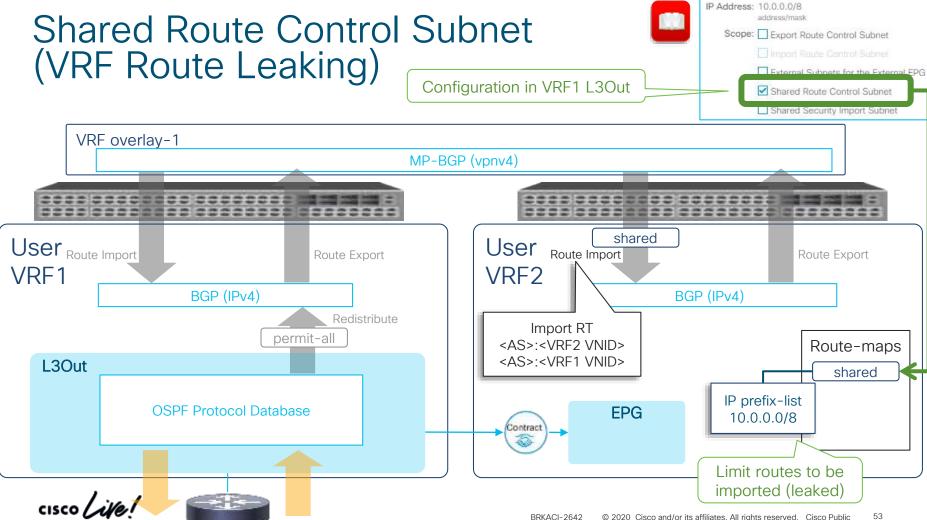
seq 1 permit 10.0.0.0/8

A prefix configured by "Import Route Control Subnet"









CLI Verification



1. MP-BGP Import rule with another VRF VNID route-target and a route-map

```
leaf# show bgp process vrf TK:VRF2

Information for address family IPv4 Unicast in VRF TK:VRF2

Import route-map 2588672-shared-svc-leak
Export RT list:
    65003:2588672
Import RT list:
    65003:2097152
    65003:2588672
Label mode: per-prefix
```

- It always has Import and Export RT for its own VRF2 VNID (65003:2588672)
- VRF Route Leaking is handled by Import RT and Import route-map (highlighted ones)

```
leaf# show vrf TK:VRF1 detail extended | egrep 'RD|vxl'
    RD: 10.0.184.64:2
    Encap: vxlan-2097152

leaf# show vrf TK:VRF2 detail extended | egrep 'RD|vxl'
    RD: 10.0.184.64:13
    Encap: vxlan-2588672
```

VRF VNID can be checked with this command to confirm Import RT is correct



CLI Verification



2. A route-map for shared service (VRF Route Leaking)

```
leaf# show route-map 2588672-shared-svc-leak
route-map 2588672-shared-svc-leak, deny, sequence 1
  Match clauses:
    pervasive: 2
  Set clauses:
route-map 2588672-shared-svc-leak, permit, sequence 2
  Match clauses:
    extcommunity (extcommunity-list filter): 2588672-shared-svc-leak
  Set clauses:
route-map 2588672-shared-svc-leak, permit, sequence 1000
  Match clauses:
  ip address prefix-lists: TPv4-2097152-32771-18-2588672-shared-svc-leak
  ipv6 address prefix-lists: IPv6-deny-all
  Set clauses:
```

Prevent BD subnet (pervasive route) from being imported via MP-BGP.
 BD subnet distribution should be done by APIC instead of MP-BGP.

 Allow importing any routes from the same VRF.
 Extended community list has RT for the same VRF VNID.

3. Allow importing certain routes from another VRF

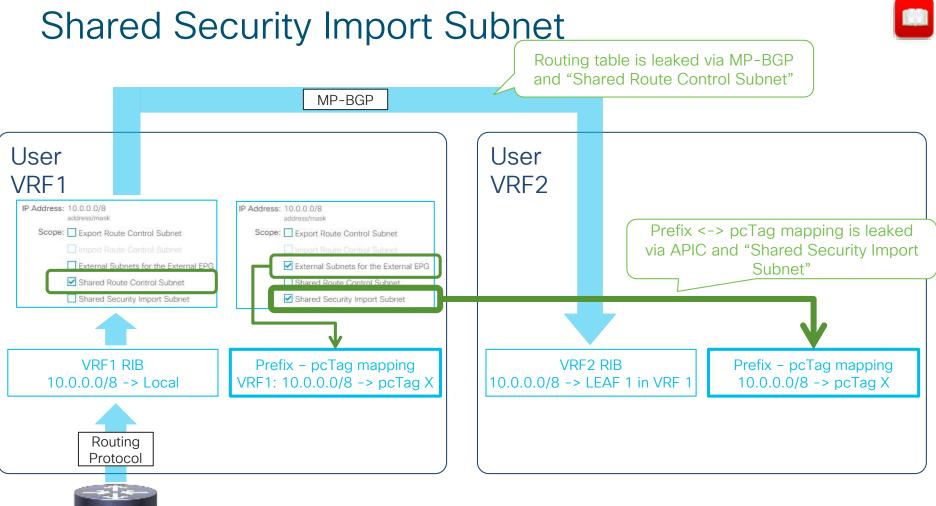
leaf# show ip extcommunity-list 2588672-shared-svc-leak
Standard Extended Community List 2588672-shared-svc-leak
 permit RT:65003:2588672

RT for the same VRF VNID Not for VRF Route Leaking

leaf# show ip prefix-list IPv4-2097152-32771-18-2588672-shared-svc-leak
ip prefix-list IPv4-2097152-32771-18-2588672-shared-svc-leak: 1 entries
 seq 1 permit 10.10.10.0/8

IP Prefix-List from Shared Route Control Subnet





CLI Verification



1. Check VRF VNID

leaf# show vrf TK:VRF1 detail extended | grep vxlan Encap: vxlan-2097152 leaf# show vrf TK:VRF2 detail extended | grep vxlan Encap: vxlan-2588672

pcTag (class) for shared route

prefix-pcTag mapping is sahred to VRF2 (VNID 2588672)

2. Prefix - pcTag mapping table

```
1st-gen-leaf# vsh lc -c 'show system internal aclgos prefix' | egrep '^Shared|52.52.52'
Shared Addr
               Mask
                                     RefCnt
                        Scope Class
10.0.0.0
               ffffff
```

2nd-gen-leaf# vsh_l	lc -c 'show system i	internal aclqos	prefix' egr	rep 'Shared	52.52.52'
Vrf-Vni VRF-Id Tabl	le-Id Addr		<mark>Class</mark> Sh	nared Remote	Complete
2097152 8 0x8	10.0.	.0.0/8	18	0 1	No
<mark>2588672</mark> 11 0xb	10.0.	.0.0/8	18	1 1	No

```
leaf# vsh -c 'show system internal policy-mgr prefix'
Vrf-Vni VRF-Id Table-Id Table-State VRF-Name
                                                         Addr
                                                                      Shared Remote Complete
2097152 8
               0x8
                                     TK:VRF1
                                                   10.0.0.0/8
                                                                                      False
                              Uρ
                                                                               True
                                                                        True
2588672 11
               0xb
                                     TK:VRF2
                                                   10.0.0.0/8
                                                                                      False
                              Uρ
                                                                               True
                                                                        True
```



From 3.2 release, use this command regardless of leaf generations

Aggregate Route Control



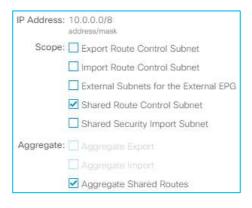


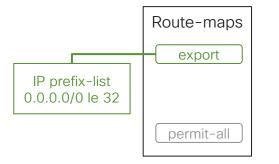


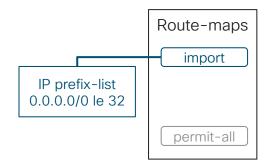
== Import ==

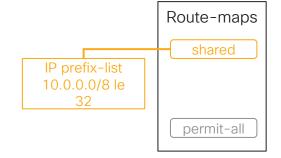


== Shared ==









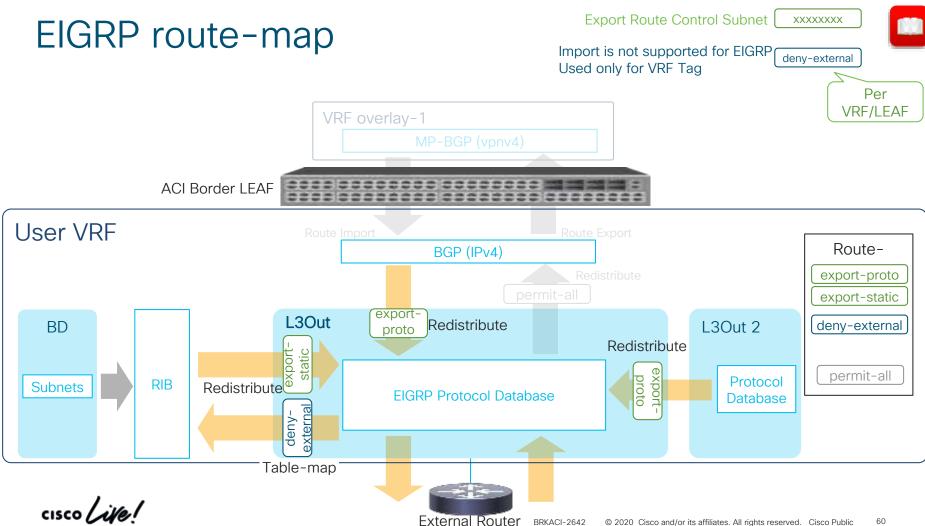
Only "Aggregate Shared Routes" support non-0.0.0.0/0 aggregation

Export Route Control Subnet XXXXXXXX OSPF route-map Import Route Control Subnet XXXXXXXX Per VRF/LEAF VRF overlay-1 **ACI Border LEAF** User VRF Route-BGP (IPv4) export-proto export-static Redistribute or export-L3Out area-filter in Redistribute L3Out 2 BD deny-external proto Redistribute export-proto import-ospf permit-all Protocol **RIB** Subnets **OSPF** Protocol Database Database importexport deny-external Still needs export in area-filter Table-map L3Out2 on top of out import in L3Out1

External Router

BRKACI-2642

© 2020 Cisco and/or its affiliates. All rights reserved. Cisco Public

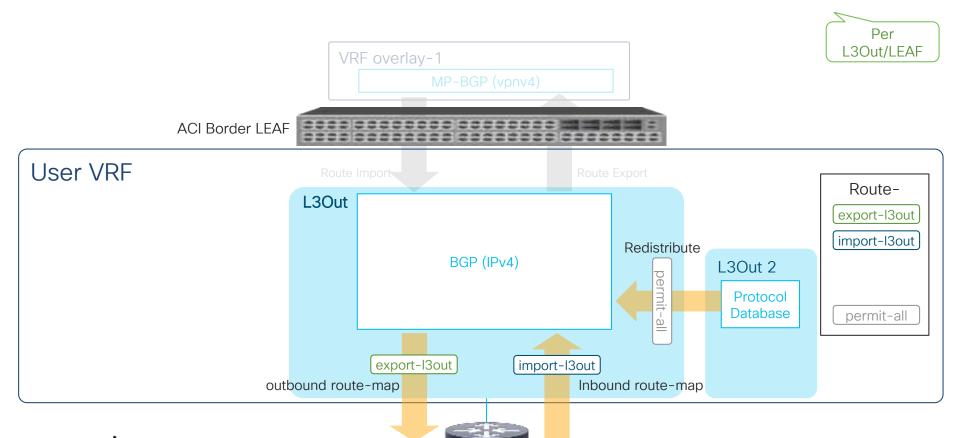


BGP route-map

Export Route Control Subnet export-I3out

Import Route Control Subnet __ir





External Router

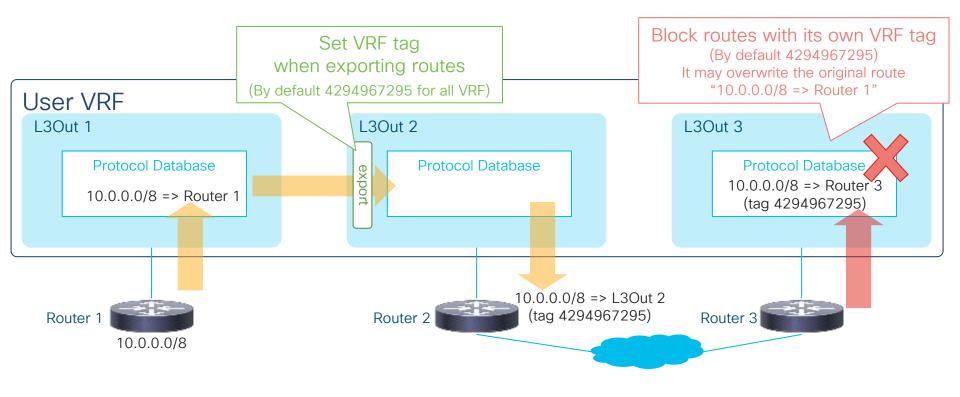
Routing Loop

Loop Avoidance for OSPF/EIGRP - VRF Tag

EIGRP and MP-BGP Redistribution Issue



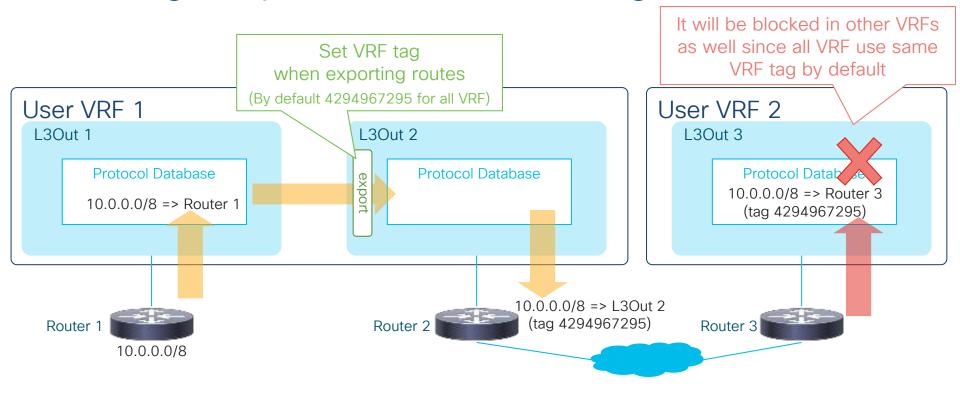






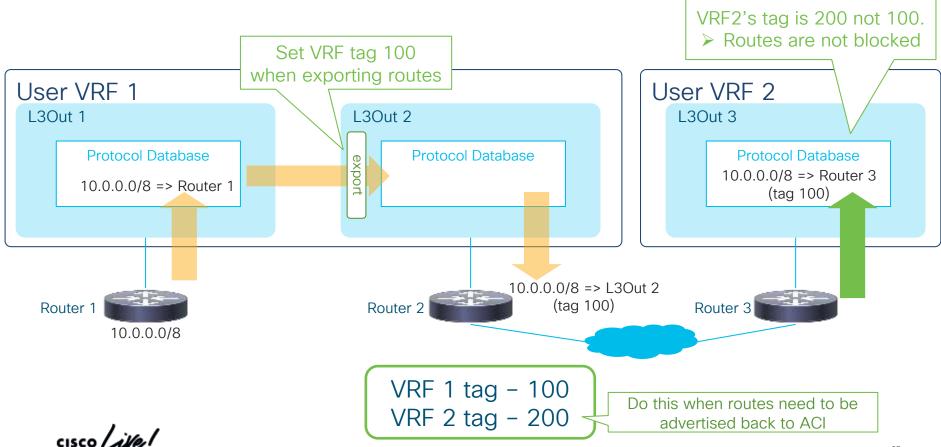
X VRF tagging for exported routes and blocking routes with VRF tag are always enabled



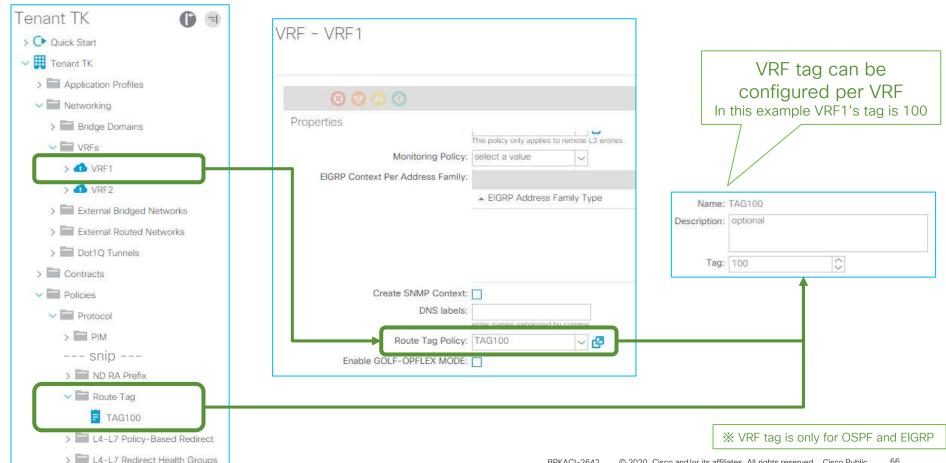




X VRF tagging for exported routes and blocking routes with VRF tag are always enabled











leaf# show ip prefix-list IPv4-proto49158-2097152-exc-ext-inferred-export-dst ip prefix-list IPv4-proto49158-2097152-exc-ext-inferred-export-dst: 1 entries seq 1 permit 10.0.0.0/8



Set clauses: tag 100

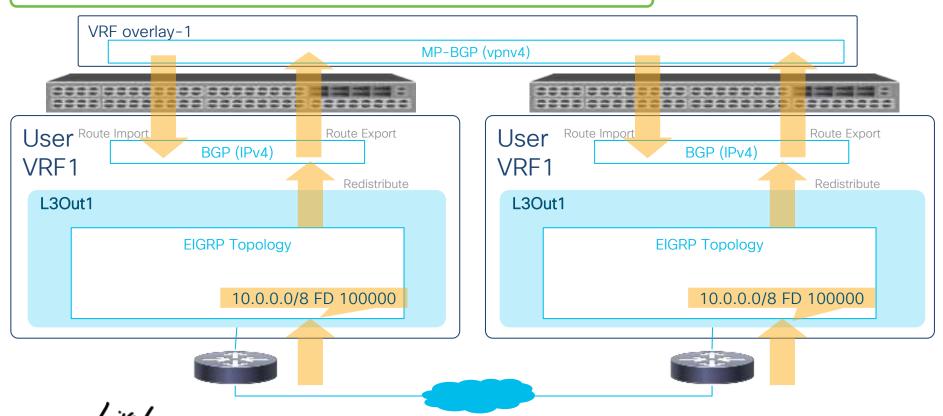




EIGRP & MP-BGP Redistribution Issue



1. 10.0.0.0/8 from L3Out 1 via EIGRP on two Border LEAFs



EIGRP & MP-BGP Redistribution Issue

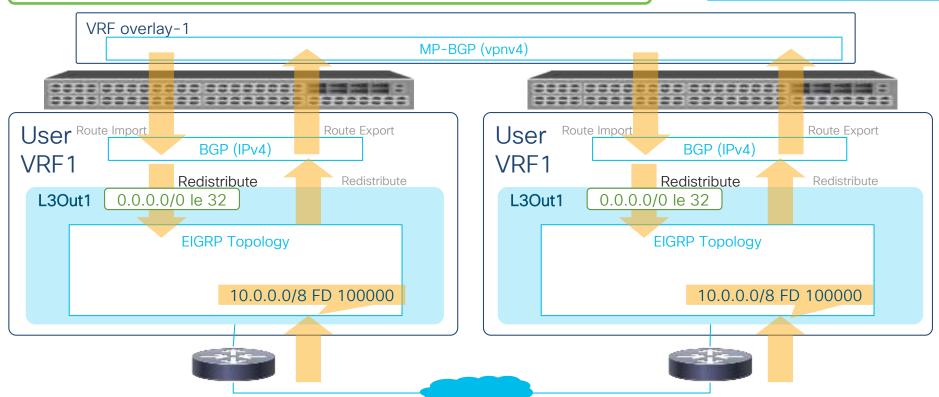
address/mask

Scope: ✓ Export Route Control Subnet

Aggregate: ✓ Aggregate Export

IP Address: 0.0.0.0/0

2. L3Out 1 exports all routes (including 10.0.0.0/8)



BRKACI-2642

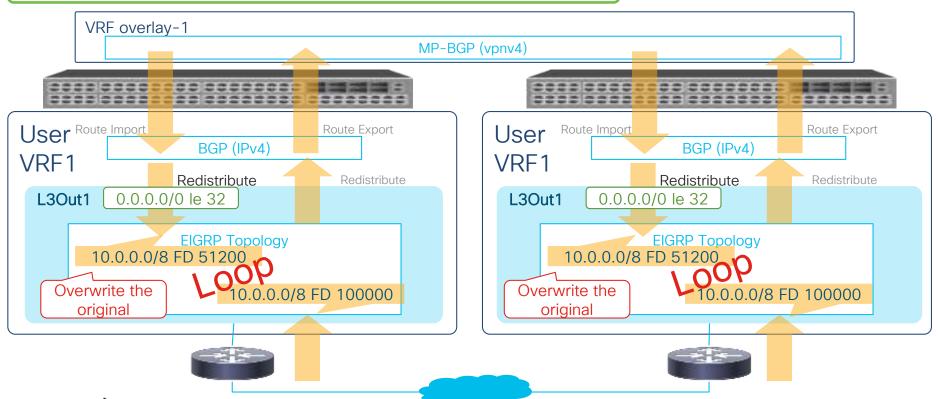
EIGRP & MP-BGP Redistribution Issue

3. Redistributed routes have lower metric than the original

IP Address: 0.0.0.0/0
address/mask

Scope: ✓ Export Route Control Subnet

Aggregate: ✓ Aggregate Export

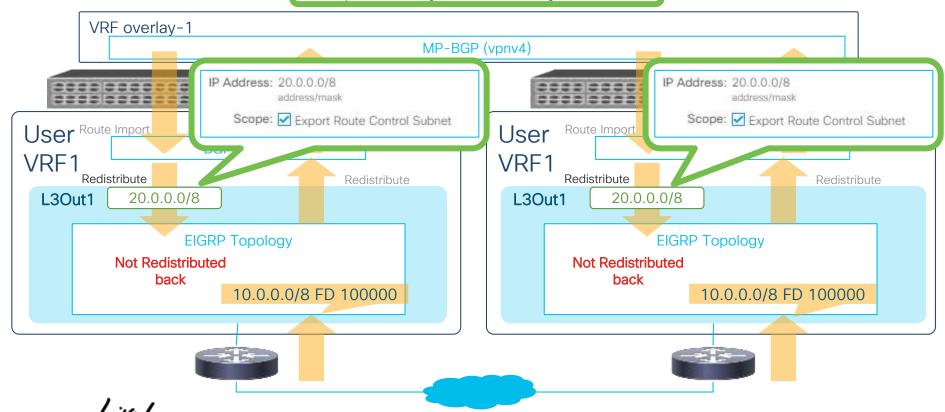


BRKACI-2642

EIGRP & MP-BGP Redistribution Solution 1



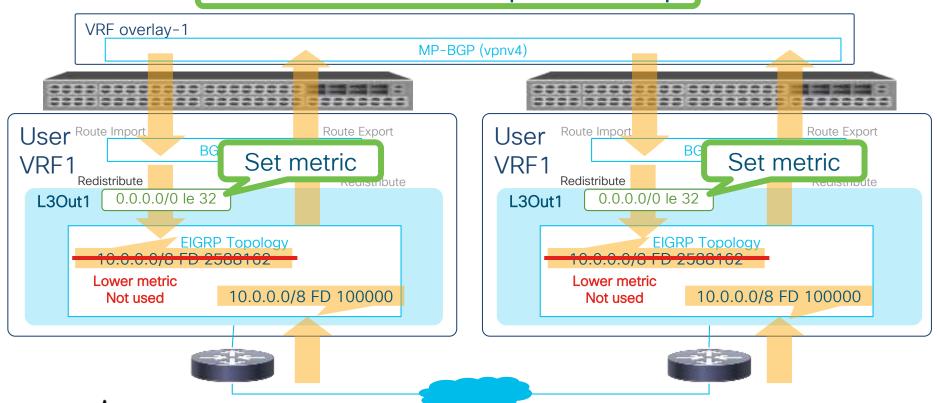
Export only necessary routes



EIGRP & MP-BGP Redistribution Solution 2



Add "set metric" rule to export route-map

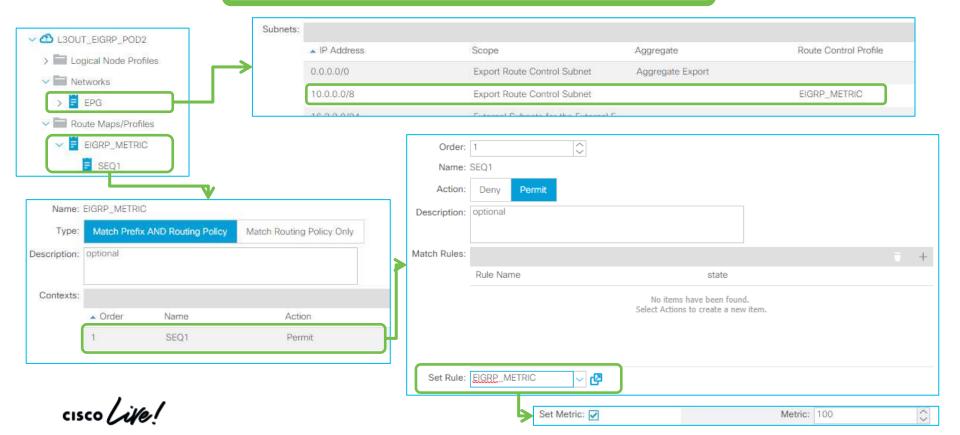


BRKACI-2642

EIGRP & MP-BGP Redistribution Solution 2



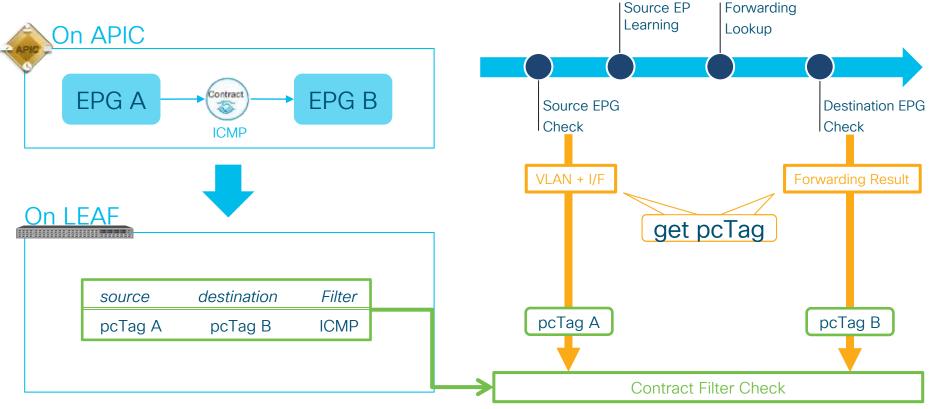
Add "set metric" rule to export route-map



L3Out Contract deep dive



pcTag (policy control Tag) in normal EPG



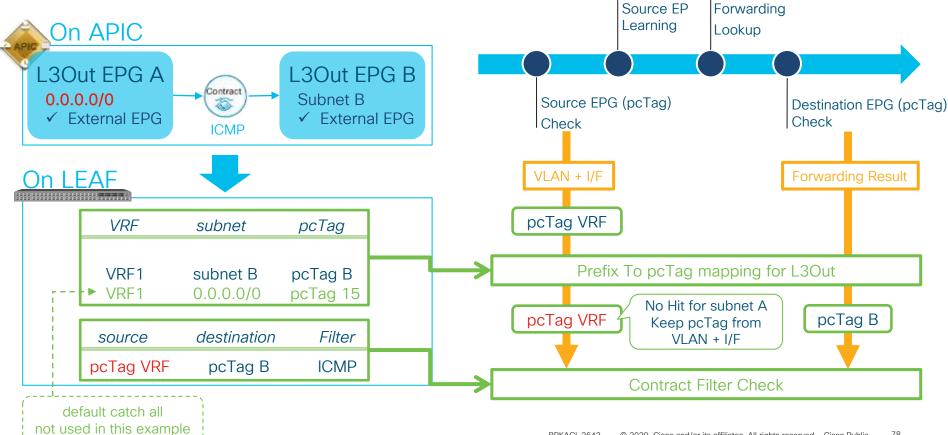
BRKACI-2642

L3Out Contract Src: Subnet A -> Dst: Subnet B pcTag (policy control Tag) in L3Out EPG Source EP | Forwarding Learning Lookup On APIC L3Out EPG A L3Out EPG B Contract Subnet A Subnet B Source EPG (pcTag) Destination EPG (pcTag) ✓ External FPG ✓ External FPG Check Check **ICMP** VLAN + I/F Forwarding Result ----pcTag VRF **VRF** subnet pcTag VRF1 subnet A pcTag A Prefix To pcTag mapping for L3Out VRF1 subnet B pcTag B VRF1 0.0.0.0/0pcTag 15 Hit subnet A pcTag A pcTag B destination Filter source Hit subnet B pcTag A pcTag B **ICMP** Contract Filter Check

default catch all not used in this example

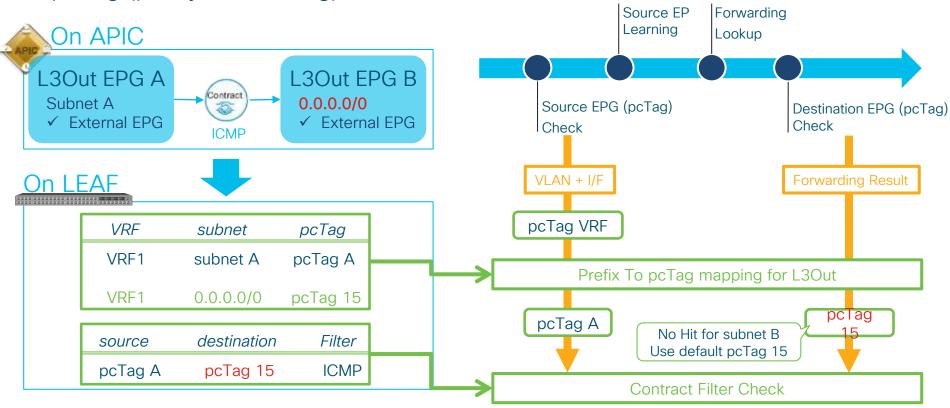


pcTag (policy control Tag) in L3Out EPG with 0.0.0.0/0





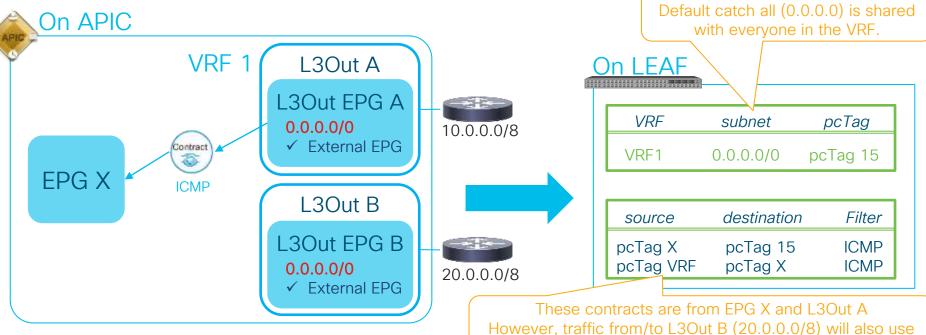
pcTag (policy control Tag) in L3Out EPG with 0.0.0.0/0





Prefix-pcTag entry is per VRF.

Common Issue (L3Out EPGs with 0.0.0.0/0)

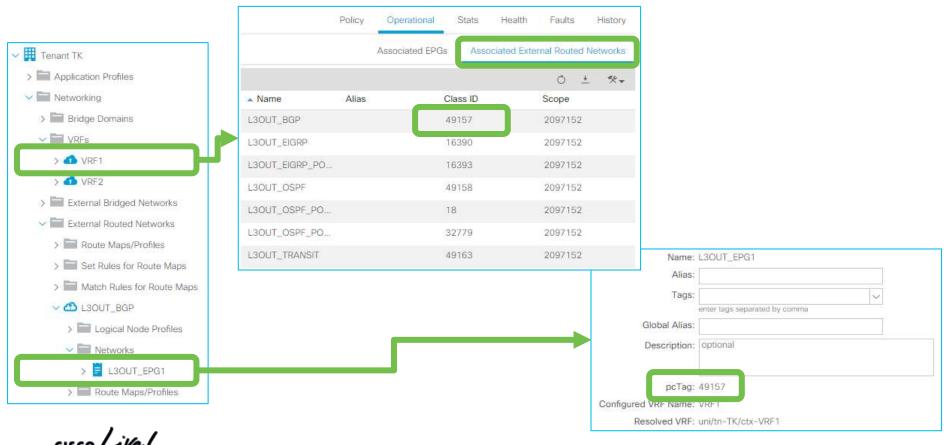


No overlap of External EPG L3Out subnets in same VRF Use 0.0.0.0/0 (External subnet for the external EPG) only for one L3Out EPG per VRF

default pcTag (VRF or 15) due to 0.0.0.0/0 config.

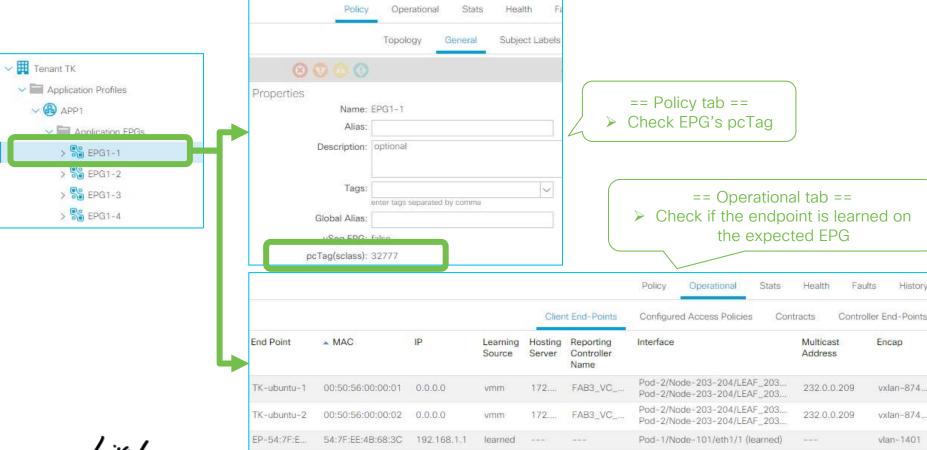
How to get pcTag for L3Out





How to get pcTag for normal EPG





BRKACI-2642

How to get VRF pcTag



From APIC

```
admin@apic1:~> moquery -c fvCtx -f 'fv.Ctx.name=="VRF1"' | egrep '#|dn|pcTag'
# fv.Ctx
dn : uni/tn-TK/ctx-VRF1
pcTag : 49153
```

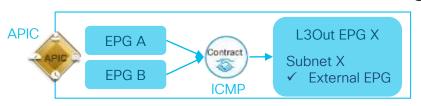
From LEAF



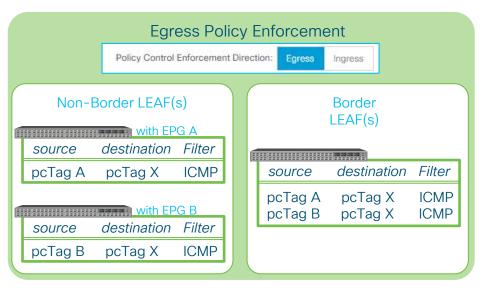
Policy Control Enforcement Direction

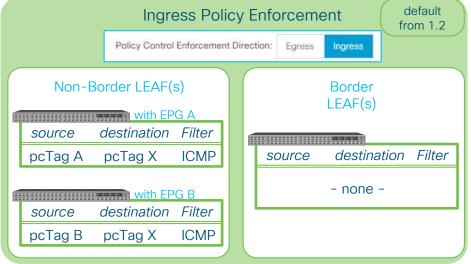


A feature to save contract TCAM usage on border LEAF



No effects on EPG <-> EPG traffic

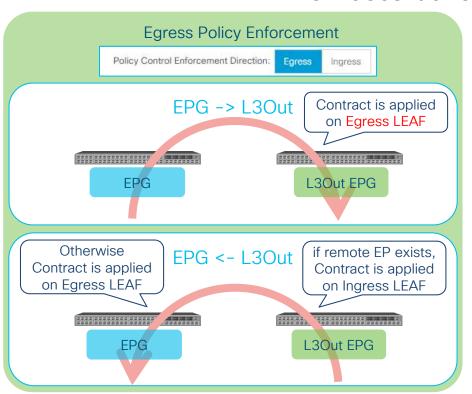


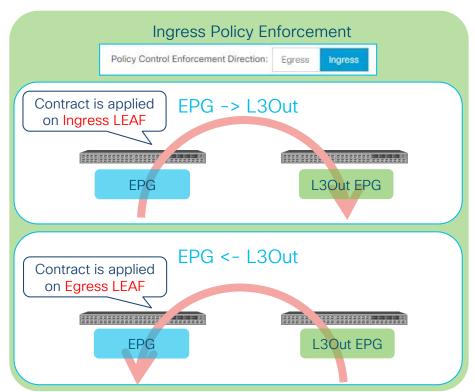


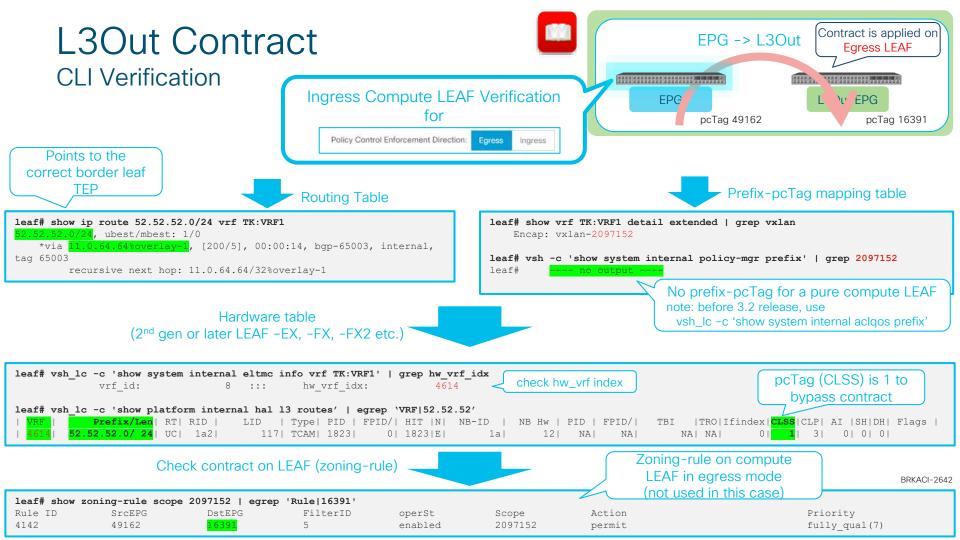
Under VRF Policy Control Enforcement Direction: Egress Ingress

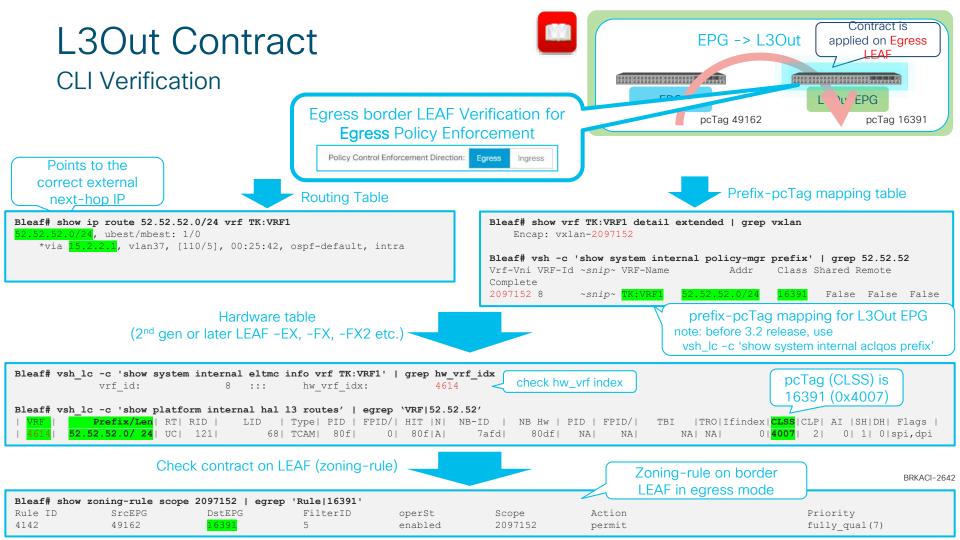
Policy Control Enforcement Direction

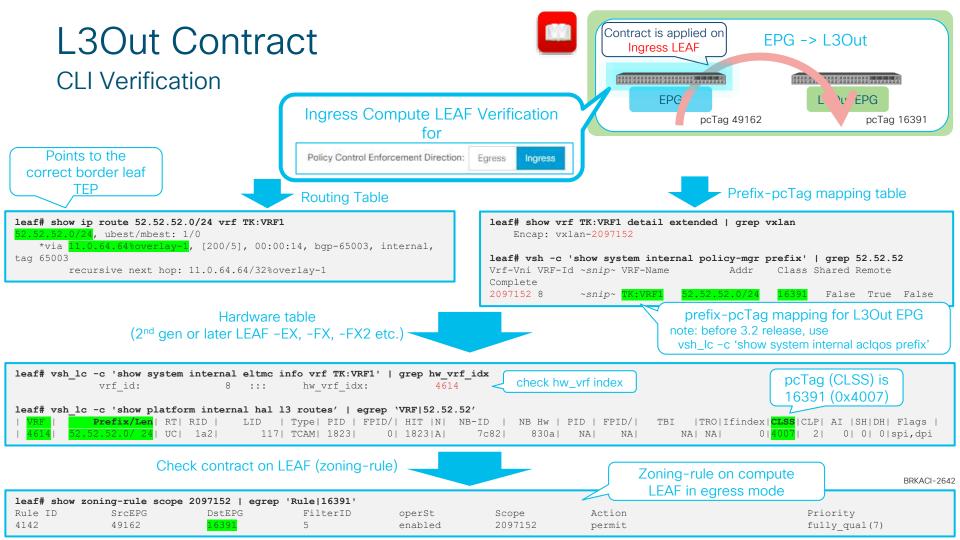
How does it affect traffic flow and contract?

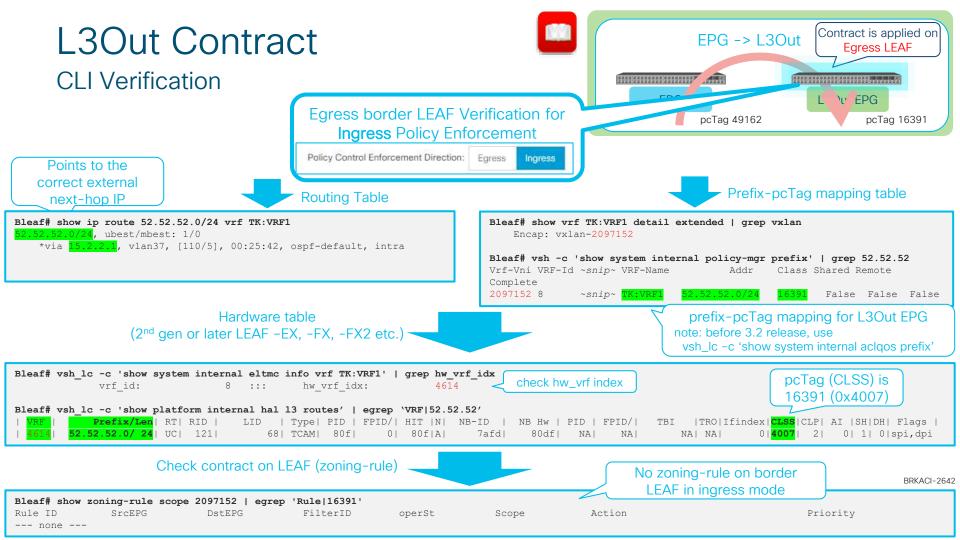












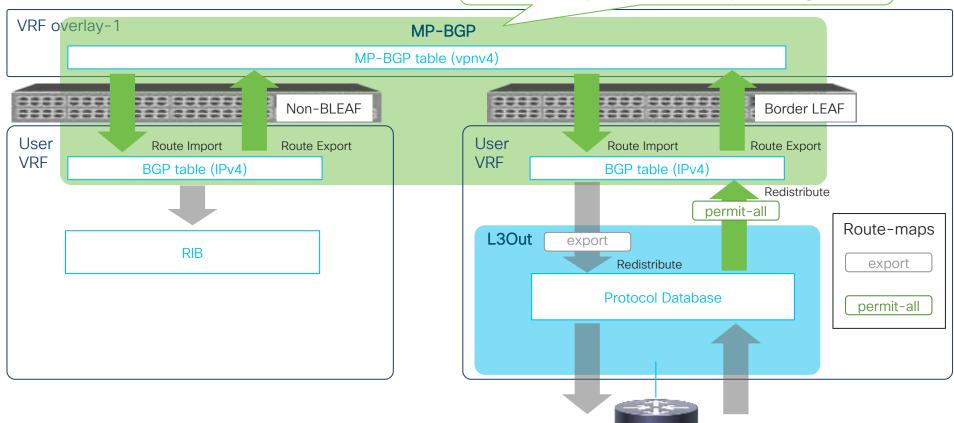
MP-BGP Deep Dive



MP-BGP

MP-BGP is automatically deployed once Route Reflector (and MP-BGP AS) is configured

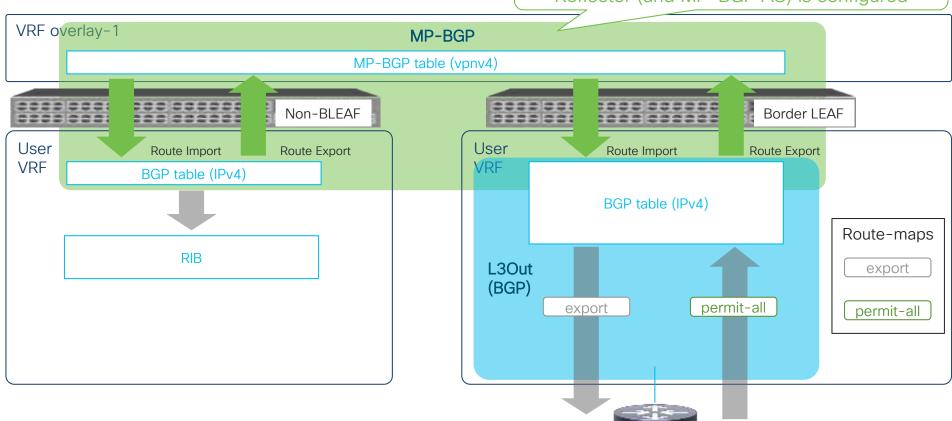




cisco Live!

MP-BGP with L3Out BGP

MP-BGP is automatically deployed once Route Reflector (and MP-BGP AS) is configured



cisco Live!

CLI Verifications



1. BGP process in your VRF with expected Redistribution and Route-Target

```
border-leaf# show bgp process vrf TK:VRF1

VRF RD : 10.0.184.64:2

Information for address family IPv4 Unicast in VRF TK:VRF1

Redistribution
direct, route-map permit-all
static, route-map imp-ctx-bgp-st-interleak-2097152
eigrp, route-map permit-all
ospf, route-map permit-all
export RT list:
65003:2097152

Import RT list:
65003:2097152

Information for address family IPv6 Unicast in VRF TK:VRF1
--- snip ---
```

Automatically created regardless of routing protocol used in L3Out. If not, check Route Reflector policy on APIC

- VRF RD (Route Distinguisher) is based on TEP IP
- BGP redistributes (almost) all external routes to export them into MP-BGP vpnv4 by default. Check a later page for the exception on BD subnets (direct routes).
- RT (Route Target) is based on ACI BGP AS and VRF VNID.

2. External routes are redistributed/exported into VPNv4 in VRF overlay-1

border-leaf# show bgp vpnv4 unicast vrf overlay-1							
Network	Next Hop	Metric	LocPrf	Weight Path			
Route Distingui * i5.5.5.0/24 *>r * i15.0.0.0/24 *>r	sher: 10.0.184.64:2 10.0.184.67 0.0.0.0 10.0.184.67 0.0.0.0	(VRF TK:VRF) 5 5 0 0	1) 100 100 100 100	0 ? 32768 ? 0 ? 32768 ?			

MP-BGP VPNv4 table can be checked via normal CLI in vrf overlay-1

NOTE:

This example shows two routes are learned locally (r with next-hop 0.0.0.0) and also from another leaf with TFP 10.0.184.67.

CLI Verifications



3. MP-BGP on all leaves should have all external routes in VPNv4 format in VRF overlay-1

non-border-leaf# show bgp vpnv4 unicast vrf overlay-1						
Route Distinguisher: *>i5.5.5.0/24 * i	10.0.184.64:2 10.0.184.64 10.0.184.64	5 5	100	0 ?		
Route Distinguisher:		5 5	100	0 3		

```
5.5.5.0/24 is advertised from border-leaf1 (10.0.184.64) and border-leaf2 (10.0.184.67)
```

Two entries with the same next-hop LEAF TEP means there are two Route Reflectors.

```
non-border-leaf# show bgp vpnv4 unicast 5.5.5.0/24 vrf overlay-1
Route Distinguisher: 10.0.184.67:1
BGP routing table entry for 5.5.5.0/24, version 598 dest ptr 0xaa7e840c

AS-Path: NONE, path sourced internal to AS
10.0.184.67 (metric 3) from 10.0.184.65 (1.1.1.101)
Extcommunity:
RT:65003:2097152
```

Each VPNv4 route in VRF overlay-1 has RT with its original VRF VNID

4. BGP process is running also in your VRF on non-border-leaf for MP-BGP

```
non-border-leaf# show bgp process vrf TK:VRF1

Information for address family IPv4 Unicast in VRF TK:VRF1

Export RT list:
65003:2097152

Import RT list:
65003:2097152

Information for address family IPv6 Unicast in VRF TK:VRF1
--- snip ---
```

IPv4 BGP imports all the external routes for its own VRF based on RT from VPNv4 table

If BGP is not running in your VRF on non-border-leaf, check Route Reflector Policy config on APIC

CLI Verifications



5. The external routes are imported in IPv4 BGP table based on RT

```
        non-border-leaf# show bgp ipv4 unicast vrf TK:VRF1

        Network
        Next Hop
        Metric
        LocPrf
        Weight Path

        *>i5.5.5.5.0/24
        10.0.184.64
        5
        100
        0 ?

        *|i
        10.0.184.67
        5
        100
        0 ?
```

6. The routing table shows border leaves as next-hop learned from iBGP

```
non-border-leaf# show ip route 5.5.5.0/24 vrf TK:VRF1
5.5.5.0/24, ubest/mbest: 2/0
   *via 10.0.184.67%overlay-1, [200/5], 2d10h, bgp-65003, internal, tag 65003
        recursive next hop: 10.0.184.67/32%overlay-1
   *via 10.0.184.64%overlay-1, [200/5], 2d10h, bgp-65003, internal, tag 65003
        recursive next hop: 10.0.184.64/32%overlay-1
```



CLI Verifications for BD subnet exception



```
border-leaf# show bgp vpnv4 unicast neighbors vrf overlay-1
BGP neighbor is 10.0.184.65, remote AS 65003, ibgp link, Peer index 1
  For address family: VPNv4 Unicast
  Outbound route-map configured is deny-pervasive, handle obtained
BGP neighbor is 10.0.184.66, remote AS 65003, ibgp link, Peer index 2
  For address family: VPNv4 Unicast
  Outbound route-map configured is deny-pervasive, handle obtained
border-leaf# show route-map deny-pervasive
route-map deny-pervasive, deny, sequence 1
 Match clauses:
    pervasive: 2
  Set clauses:
route-map deny-pervasive, deny, sequence 2
 Match clauses:
    interface: NullO
  Set clauses:
route-map deny-pervasive, permit, sequence 3
  Match clauses:
  Set clauses:
```

BD subnets and Null0 I/F should not be distributed via MP-BGP.

Outbound route-map to BGP Route Reflector (spines) limits BD subnets (pervasive routes) and Nullo.

In this example, 10.0.184.65 and 10.0.184.66 are RR spines.

BD subnets are deployed based on object policies from APIC instead of routing protocol



BRKACI-2642

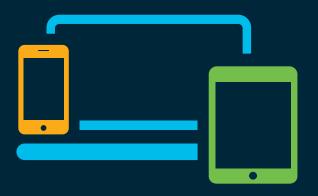
Reference

ACI Fabric L3Out Guide -

https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/guide-c07-743150.html



Complete your online session survey

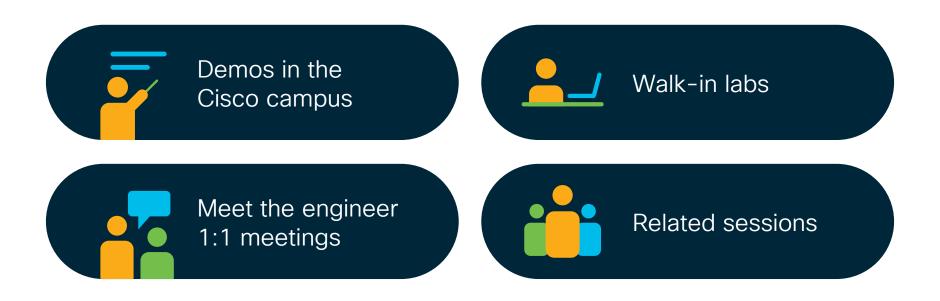


- Please complete your session survey after each session. Your feedback is very important.
- Complete a minimum of 4 session surveys and the Overall Conference survey (starting on Thursday) to receive your Cisco Live t-shirt.
- All surveys can be taken in the Cisco Events Mobile App or by logging in to the Content Catalog on <u>ciscolive.com/emea</u>.

Cisco Live sessions will be available for viewing on demand after the event at ciscolive.com.



Continue your education





illiilli CISCO

Thank you



cisco live!





You make possible