



# ACI Design #8 – Endpoint Learning

Cisco Systems G.K.  
2018/1

# Endpoint Learning 基本編

# ACIにおけるEndpoint学習動作

従来のネットワークにおけるMACアドレスとIPアドレスの学習動作とは異なり、ACIではMACアドレステーブルとARPテーブルではなくEndpointテーブルを利用する。(外部接続を除く)。

Table	含まれる情報	概要
Endpoint	MACアドレス・IPアドレス	1つのエントリーにMACアドレスは1つ、IPアドレスは1つもしくは複数 IPアドレスは /32 (IPv4) もしくは /128 (IPv6) でのみ学習する
ARP	L3out 接続については、ARPに基づいて MACアドレスに紐づくIPアドレスを学習する	

デフォルトでは、ACIはARPに依存せずに、Leafスイッチに届いたパケットの送信元MACアドレス・IPアドレスを元にEndpointを学習する。  
(DownlinkからはLocal Endpoint、UplinkからはRemote Endpointとして扱う)

# ACIにおける外部アドレス学習動作

ACI範囲外の外部ネットワークとの接続においては、従来通りRIBテーブルとARPテーブルを利用したアドレス解決を行う (Endpointテーブルは利用しない)。  
※EndpointテーブルはMACアドレスに紐付け可能なIPアドレス数に制限がある (3.1では1024まで)

また、内部Endpointとは異なり、外部ネットワークについては /32 (IPv4)や /128 (IPv6)でのアドレス管理は行わず、SubnetベースのPrefixに基づく宛先管理を行う。

ACIは宛先Prefixへ到達するためのNext-Hopのみ学習するために、従来スイッチと同じくデータプレーンからMACアドレスの学習を行い、ARPによりMACアドレスとIPアドレスの解決を行う(従来スイッチと同じ動作)。

# Local Endpoint と Remote Endpoint

各Leafスイッチは、自身に直接接続したEndpoint (=Local Endpoint) を学習し、Local Endpointの通信相手となった別Leaf (TEP)に接続したEndpoint (=Remote Endpoint)をキャッシュ用として学習する。

Endpoint	学習対象	補足
Local Endpoint	MACアドレスとIPアドレス (IPアドレスの学習は構成次第) ※1 MAC + 必要数IP を紐付けて管理	COOPによるSpineスイッチへの通知対象 (Spineスイッチは全てのEndpointのLeaf配置を管理)  •Default Retention Timer : 900秒 ※IPアドレスは個別にAging処理される
Remote Endpoint	単一のMACアドレスか単一のIPアドレス (どちらを学習するかはiVXLANに含まれる 情報に基づく[BD=MAC, VRF=IP])	キャッシュ用途(Spine Proxyの負荷軽減と最適化) Ingress側でのPolicy適用(Fabric通信の最適化)  •Default Retention Timer : 300秒

# なぜ Spine Proxy があるのに Remote Endpoint を学習する?

各Leafスイッチは、Local Endpointが通信相手とするRemote Endpointの情報をキャッシュすることにより、都度Spine Proxyに問い合わせる時間と動作を削減する。

キャッシュしたRemote Endpoint宛の通信は、Spine Proxyへの問い合わせを行わずに直接宛先Leafスイッチに送信する。学習済の相手に対する通信は送信元側のLeafスイッチでPolicyを適用することが可能(Ingress Policy Enforcement)。

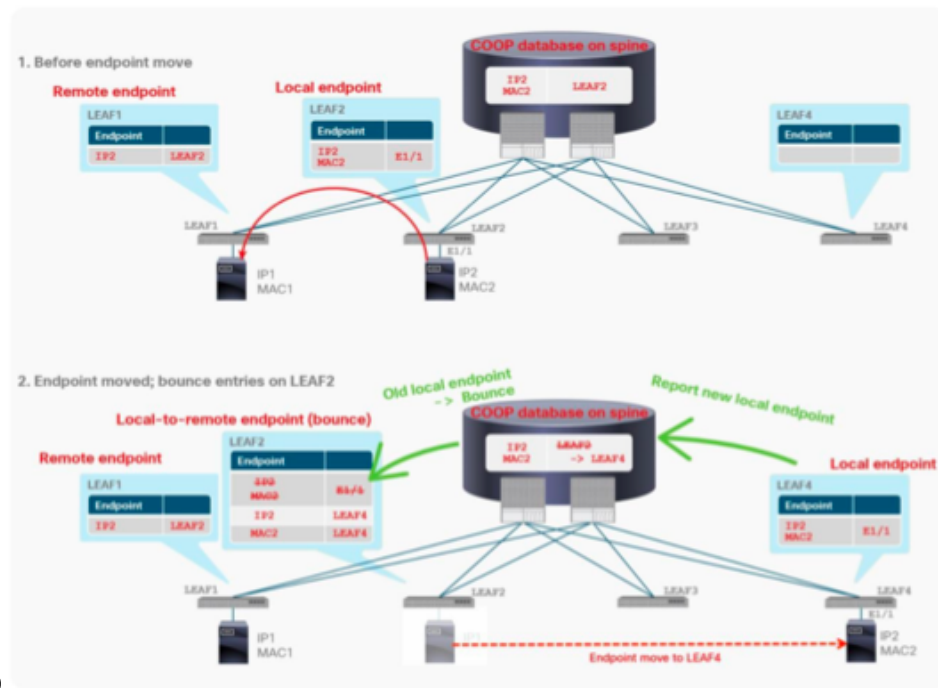
Remote Endpointについては、Local Endpointの様にMACアドレスとIPアドレスの紐付けは管理せず、宛先IPアドレスもしくはMACアドレスを個別に宛先Leafと紐付ける(L2通信相手ならMACアドレスのみ、L3通信相手ならIPアドレスのみ)。

※ただしvPC Peer同士はOrphan Portを含む全てのEndpointの情報を同期するため動作が異なる

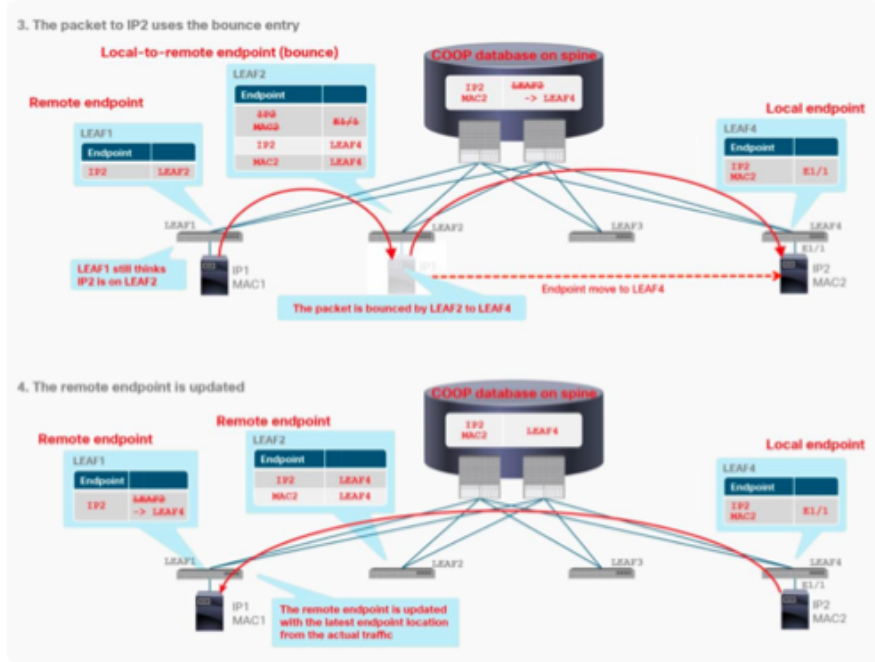
# Endpointの移動とBounce Entry (1)

EndpointがLeafスイッチをまたいで移動した場合、移動先Leafスイッチは当該Endpointの情報をSpineスイッチに通知し、Spineスイッチは更新情報を移動元Leafスイッチにのみ通知する。

この時点では、移動元・移動先のLeafスイッチ以外のLeafスイッチは当該Endpointの移動を知らない。  
→キャッシュに基づいて当該Endpoint宛通信を移動元に送信する可能性がある



# Endpointの移動とBounce Entry (2)



移動元Leafスイッチは移動先Leafスイッチを知っているため、当該パケットを移動先Leafに転送する(Bounce)。

戻りの通信については、移動先Leafは送信元Leafに直接返信するため、送信元Leafはこの時点で当該Endpointの移動を認識し、自身のキャッシュを更新する。

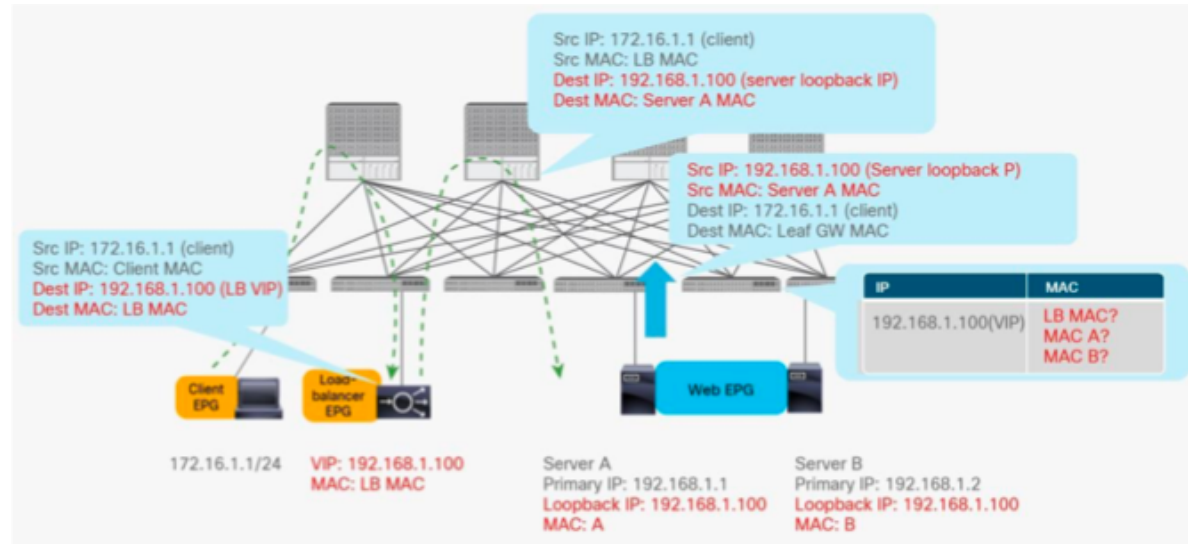
Endpoint移動時にFabricサイズに依存しないSpine処理を実現し、必要Leafのみ必要な時点で情報を更新する為の仕組み



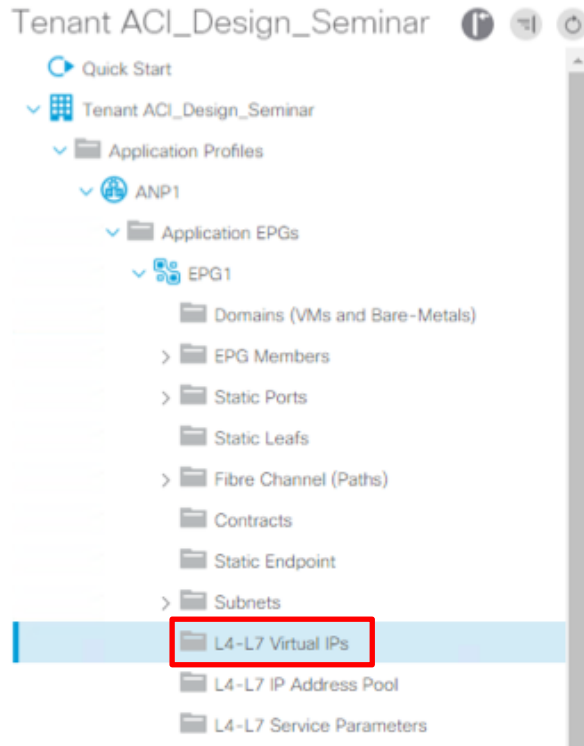
# Endpoint Learning パラメータチューニング編

# EPG - L4-L7 Virtual IPs

LBのDirect Server Return (L2 DSR)動作のために、DSR用仮想IPアドレスについてEndpoint Data-plane Learningを無効化する(フラップ抑制)。



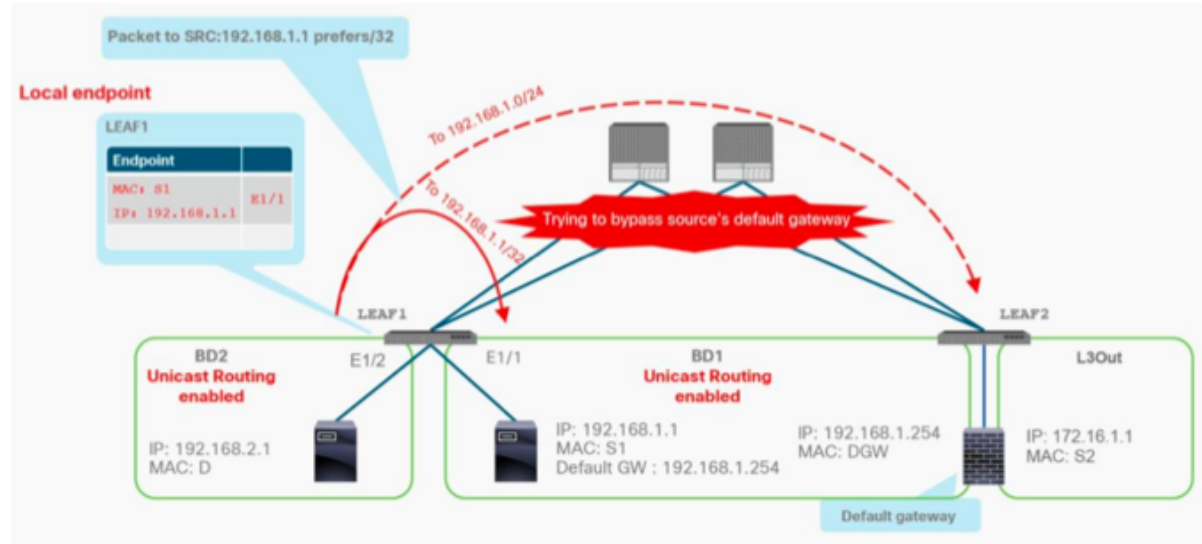
# EPG - L4-L7 Virtual IPs



# BD – Unicast Routing

BDに対して構成したSubnet設定をLeafスイッチに適用する。  
(Routing動作およびEndpointのIPアドレス学習の両方に対して機能する)


BDをL2利用する場合は、  
想定外のGWバイパス  
動作の発生に注意する。



# BD – Unicast Routing

Tenant ACI\_Design\_Seminar   

 Quick Start

▼  Tenant ACI\_Design\_Seminar

>  Application Profiles

▼  Networking

▼  Bridge Domains

>  BD1

>  BD2

>  BD3

>  BD4

>  L3-BD1

>  L3-BD2

>  L3-BD3

## Bridge Domain - BD1



### Properties

Unicast Routing:

Operational Value for Unicast Routing: true

Custom MAC Address:

Virtual MAC Address:

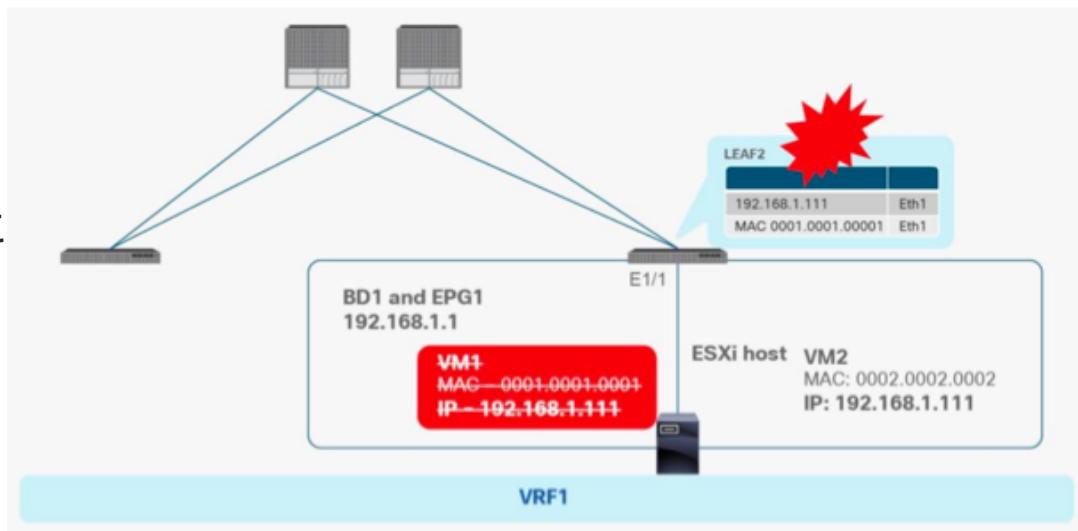
Subnets:

▲ Gateway Address

# BD – GARP-based EP Move Detection

Unicast Routing有効・ARP Flooding有効に加えて有効化することで、同一ポートの配下でIPアドレスが新しいMACアドレスに紐付いたGARP通知に基づいたEndpointの学習を行う。

LeafスイッチがVM1からの通信で一旦学習したIPアドレスを別のVMが同一ホスト上で再利用した場合などに必要となる。



# BD – GARP-based EP Move Detection

The screenshot displays the Cisco ACI configuration interface for a Bridge Domain (BD1) within the Tenant ACI\_Design\_Seminar. The left-hand navigation pane shows a tree structure with 'Bridge Domains' expanded to 'BD1'. The right-hand pane shows the configuration for 'Bridge Domain - BD1'. Under the 'Properties' section, the 'Unicast Routing' checkbox is checked, and the 'Operational Value for Unicast Routing' is 'true'. The 'Custom MAC Address' is set to '00:22:BD:F8:19:FF' and the 'Virtual MAC Address' is 'Not Configured'. The 'Subnets' section is currently empty. At the bottom of the configuration pane, the 'EP Move Detection Mode' is set to 'GARP based detection', which is highlighted with a red rectangular box. Below this, the 'Associated L3 Outs' section is also empty.

Tenant ACI\_Design\_Seminar

- Quick Start
- Tenant ACI\_Design\_Seminar
  - Application Profiles
  - Networking
    - Bridge Domains
      - BD1**
      - BD2
      - BD3
      - BD4
      - L3-BD1
      - L3-BD2
      - L3-BD3
    - VRFs
    - External Bridged Networks
    - External Routed Networks
    - Dot1Q Tunnels
    - Contracts
    - Policies

Bridge Domain - BD1

Properties

Unicast Routing:

Operational Value for Unicast Routing: true

Custom MAC Address: 00:22:BD:F8:19:FF

Virtual MAC Address: Not Configured

Subnets:

▲ Gateway Address

EP Move Detection Mode:  GARP based detection

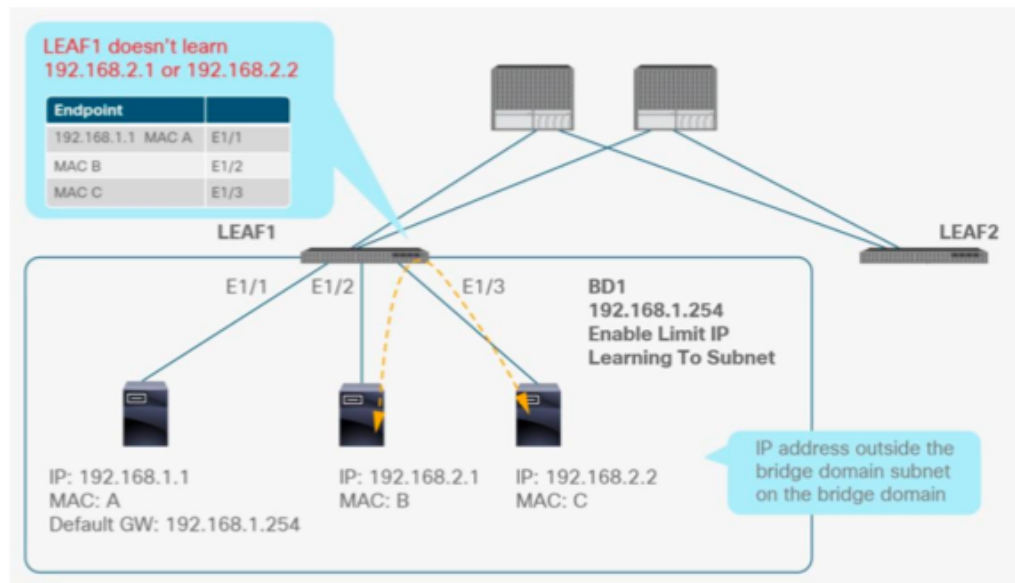
Associated L3 Outs:

▲ L3 Out

# BD – Limit IP Learning To Subnet (旧名称 : Enforce Subnet Check for IP Learning)

- 2.3(1e)および3.0(1k)以降で新規に作成したBDではデフォルト有効(CSCvb16668)
- 3.0(1k)以前のバージョンで有効・無効を切り替える場合はEndpoint学習が120秒間停止するので注意(CSCve29663)

誤ったIPアドレスを設定した場合に  
当該IPアドレスの学習を抑制する。  
(Local Endpointのみ)





# BD – Limit IP Learning To Subnet (旧名称 : Enforce Subnet Check for IP Learning)

Tenant ACI\_Design\_Seminar



Bridge Domain - BD1

Quick Start

Tenant ACI\_Design\_Seminar

Application Profiles

Networking

Bridge Domains

BD1

BD2

BD3

BD4

L3-BD1

L3-BD2

L3-BD3

VRFs

External Bridged Networks

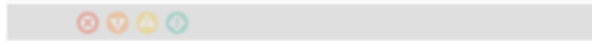
External Routed Networks

Dot1Q Tunnels

Contracts

Policies

Services



Properties

Name: BD1

Alias:

Description: optional

Type: fc **regular**

Global Alias:

Legacy Mode: No

VRF: VRF1

Resolved VRF: ACI\_Design\_Seminar/VRF1

L2 Unknown Unicast: **Flood** Hardware Proxy

L3 Unknown Multicast Flooding: **Flood** Optimized Flood

Multi Destination Flooding: **Flood in BD** Drop Flood in Encapsulation

PIM:

IGMP Policy: select an option

ARP Flooding:

Endpoint Dataplane Learning:

Clear Remote MAC Entries:

**Limit IP Learning To Subnet:**

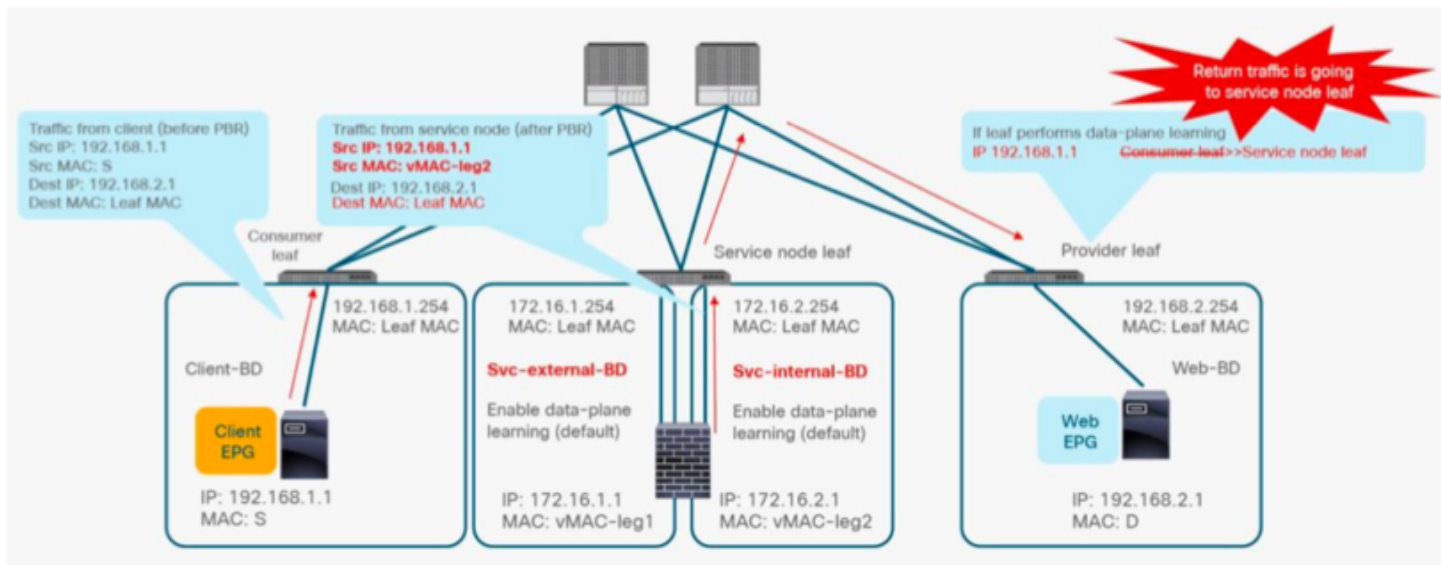
Endpoint Retention Policy: select a value

This policy only applies to local L2, L3, and remote L3 entries

IGMP Snoop Policy: select a value

# BD – Endpoint Dataplane Learning

PBRで利用するService Nodeを接続するためのBDでは、Dataplane Learningは無効化する必要がある。



# BD – Endpoint Dataplane Learning

Tenant ACI\_Design\_Seminar

Quick Start

Tenant ACI\_Design\_Seminar

Application Profiles

Networking

Bridge Domains

BD1

BD2

BD3

BD4

L3-BD1

L3-BD2

L3-BD3

VRFs

External Bridged Networks

External Routed Networks

Dot1Q Tunnels

Contracts

Policies

Services

Bridge Domain - BD1

Properties

Name: BD1

Alias:

Description: optional

Type: fc regular

Global Alias:

Legacy Mode: No

VRF: VRF1

Resolved VRF: ACI\_Design\_Seminar/VRF1

L2 Unknown Unicast: Flood Hardware Proxy

L3 Unknown Multicast Flooding: Flood Optimized Flood

Multi Destination Flooding: Flood in BD Drop Flood in Encapsulation

PIM:

IGMP Policy: select an option

ADP Flooding:

Endpoint Dataplane Learning:

Clear Remote MAC Entries:

Limit IP Learning To Subnet:

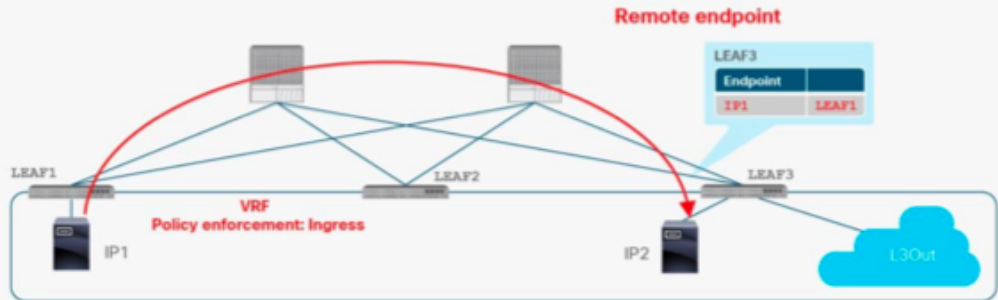
Endpoint Retention Policy: select a value

This policy only applies to local L2, L3, and remote L3 entries

IGMP Snoop Policy: select a value

# Fabric - Disable Remote EP Learn (on border Leaf) ①

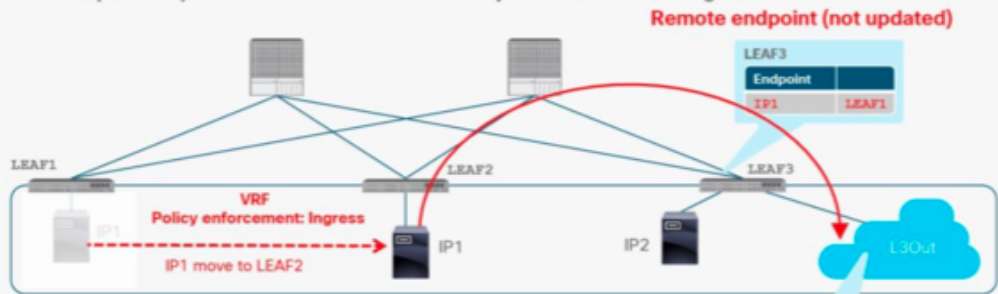
1. IP1 is learned on border leaf LEAF3 because of IP1-to-IP2 traffic



2. IP1-to-IP2 traffic ceases

3. IP1 moves to LEAF2

4. IP1-to-L3Out (external) traffic starts before remote endpoint on border leaf ages out



Remote endpoint IP1 won't be updated with new source LEAF2, nor will it age out because of packet-to-L3Out behavior in the VRF instance with ingress policy enforcement mode

Endpoint接続LeafをBorder Leafとしても利用する注意点

Leaf1でBounce Entryが期限切れとなった後にLeaf3からIP1宛の通信が失敗する(Retention Timer経過時間まで)

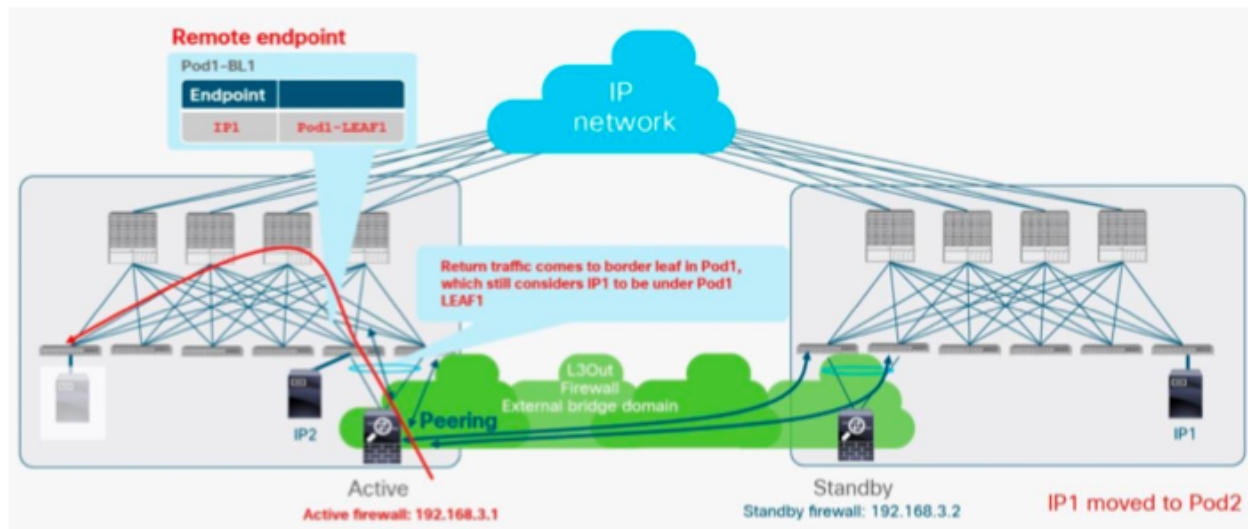
Leaf3でRemote Endpoint情報を削除することにより事象は解決する(要CLI操作)

Leaf3がBorder Leaf専用である場合や、Disable Remote EP Learnが構成されている場合も防止できる

## Fabric – Disable Remote EP Learn (on border leaf) ②

複数のLeafスイッチに渡って同一VLAN SVIでL3outを構成する場合に構成する。

例：Multi-PodでL3outをActive/Standby構成の環境でEndpointが移動した場合



# Fabric - Disable Remote EP Learn (on border Leaf)

## System Settings

- > Quota
  - APIC Connectivity Preferences
  - BD Enforced Exception List
  - Control Plane MTU
  - Endpoint Controls
  - Fabric Wide Setting**
  - System Global GIPo
  - BGP Route Reflector
  - COOP Group
  - Load Balancer
  - Precision Time Protocol



## Fabric Wide Setting Policy

### Properties

- Disable Remote EP Learning:**  To disable remote endpoint learning in VRFs containing external bridged/routed domains
- Enforce Subnet Check:**  To disable IP address learning on the outside of subnets configured in a VRF, for all VRFs
- Reallocate Gipo:**  Reallocate some non-stretched BD gipos to make room for stretched BDs.
- Enforce Domain Validation:**  Validation check if a static path is added but no domain is associated to an EPG
- Opflex Client Authentication:**  To enforce Opflex client certificate authentication for GOLF and Linux

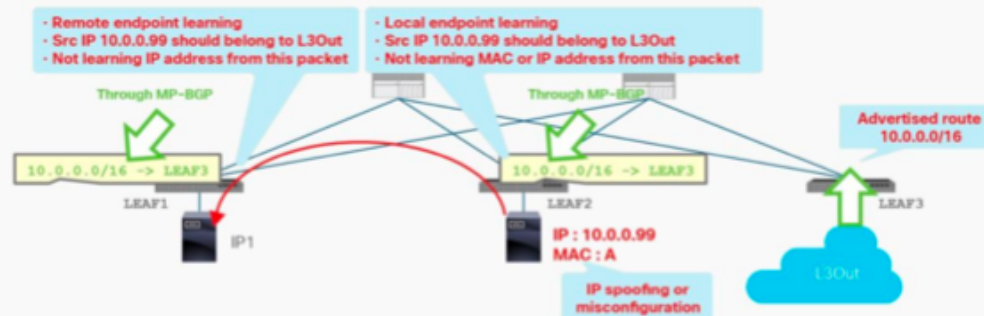
# Fabric - Enforce Subnet Check ①

外部Route範囲のIPアドレスを設定したEndpointがパケットを誤送信した場合、L3outから受信している経路情報がPrefixの場合は登録は防止される(第2世代Leafスイッチの場合のみ)が、0.0.0.0/0範囲だった場合には防止できない。

第2世代Leafスイッチの場合は [Enforce Subnet Check] を有効化して防止する。

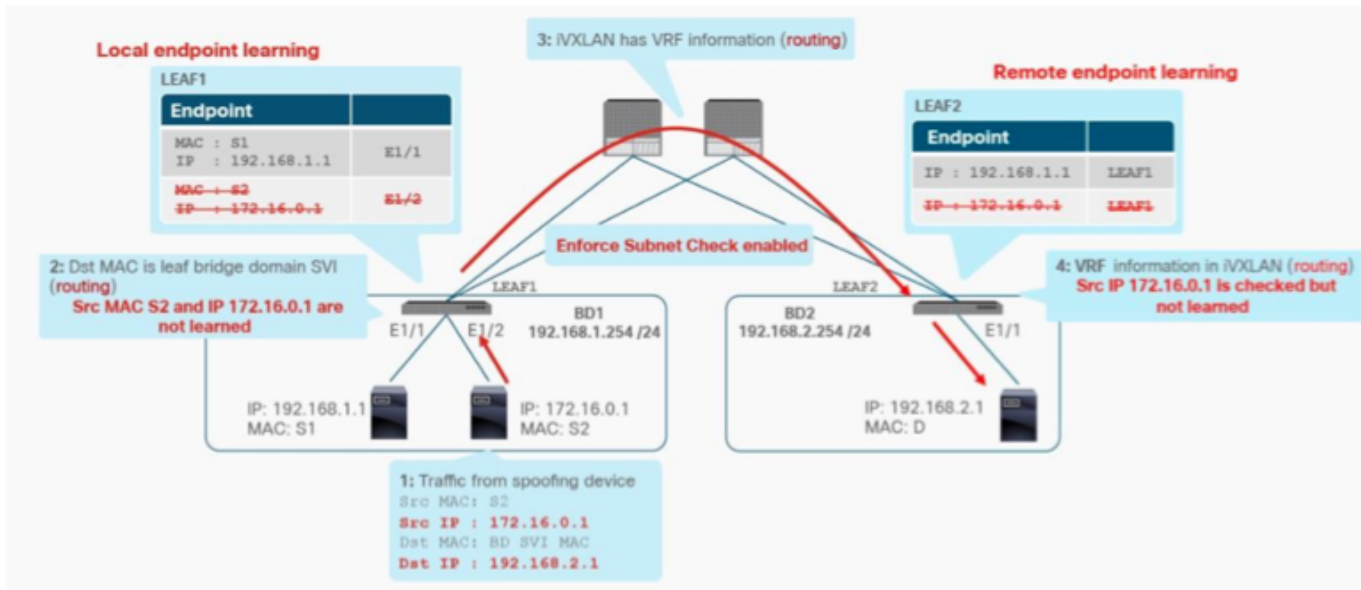
第1世代Leafスイッチの場合は [Limit IP Learn To Subnet] と [Disable Remote EP Learn] を併用設定する。

1. Each leaf has 10.0.0.0/16 from L3Out
2. EP2 uses IP 10.0.0.99 because of IP spoofing or misconfiguration
3. Traffic from IP 10.0.0.99 to IP1 is injected by EP2
4. Endpoint learning of spoofed IP 10.0.0.99 is prevented



# Fabric – Enforce Subnet Check ②

Limit IP Learning To Subnetよりも強力的に、Local Endpointだけでなく、Remote Endpointについても誤学習を防止できるが、第2世代Leafスイッチに対してのみ有効。





# Fabric – Enforce Subnet Check

## System Settings

> Quota

APIC Connectivity Preferences

BD Enforced Exception List

Control Plane MTU

Endpoint Controls

Fabric Wide Setting

System Global GIPo

BGP Route Reflector

COOP Group

Load Balancer

Precision Time Protocol



## Fabric Wide Setting Policy

### Properties

**Disable Remote EP Learning:**  To disable remote endpoint learning in VRFs containing external bridged/routed domains

**Enforce Subnet Check:**  To disable IP address learning on the outside of subnets configured in a VRF, for all VRFs

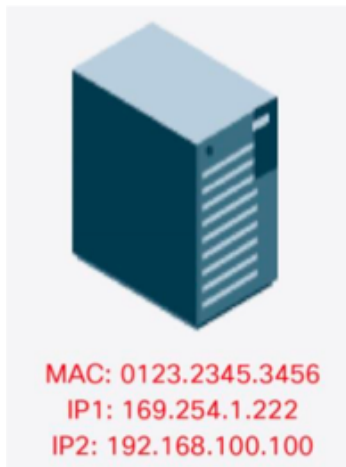
**Reallocate Gipo:**  Reallocate some non-stretched BD gipos to make room for stretched BDs.

**Enforce Domain Validation:**  Validation check if a static path is added but no domain is associated to an EPG

**Opflex Client Authentication:**  To enforce Opflex client certificate authentication for GOLF and Linux

# Fabric – IP Aging

不要になったIPアドレス情報をいつまでも保持することを防止する。  
(DHCPでIPアドレスを取得できなかった際の169.254.x.xアドレスなど)



Aging指定時間の75%が経過した時点で、LeafスイッチはUnicastでARPパケットを送出し、IPアドレスの存在確認を行う  
(応答がない場合は当該IPアドレスエントリが削除される)

# Fabric – IP Aging

## System Settings

> Quota

- APIC Connectivity Preferences
- BD Enforced Exception List
- Control Plane MTU
- Endpoint Controls**
- Fabric Wide Setting
- System Global GIPo
- BGP Route Reflector
- COOP Group
- Load Balancer
- Precision Time Protocol



## Endpoint Controls

### IP Aging Policy

#### Properties

Administrative State:

Disabled

**Enabled**

# 推奨チューニング設定

ACIにおけるEndpointの学習動作を正しく理解した上で、適切な構成を行う。

		第1世代 Leafスイッチのみの ACI Fabric	第2世代 Leafスイッチのみの ACI Fabric	第1世代・第2世代 Leafスイッチ混在の ACI Fabric
BD	Limit IP Learning To Subnet	有効	有効	有効
Fabric	IP Aging	有効 (Default)	有効 (Default)	有効 (Default)
	Disable Remote EP Learn (on border leaf)	有効	有効	有効
	Enforce Subnet Check		有効	有効

# 主なEndpoint学習動作に関連する機能の実装バージョン

機能	実装バージョン	概要
L4-L7 Virtual IP	1.2(1m)	特定IPアドレスのData-plane Learningを無効化する (DSR動作に対応するため)
Unicast Routing	1.0(1e)	BDにおけるL3ルーティングとIP Learning動作
GARP-based EP Move Detection	1.1(1j)	ARP FloodingがEnable時に同一Interface配下におけるEndpointの移動の検出にGARPを利用する (Endpoint DBにおける IP-MAC紐付け管理)
Limit IP Learning To Subnet	1.1(1j)	BD構成Subnet範囲外のIP学習の防止 (Local Endpointのみ)
Endpoint Dataplane Learning	2.0(1m)	BDにおけるEndpoint Data-plane Learningを無効化する (PBR動作に対応するため)
Disable Remote EP Learn (on border leaf)	2.2(2e)	Ingress Policy Enforcement (Default) 動作のBorder LeafにおけるRemote Endpointの学習を無効化する (対象Border LeafスイッチはSpine Proxyのみ利用する)
Enforce Subnet Check	2.2(2q) ※2.3, 3.0(1x)では 利用不可	VRFに紐付いたBD Subnet範囲外のEndpoint学習を防止する (Local / Remote 両方) ※第2世代Leafスイッチに対してのみ有効
IP Aging	2.1(1h)	学習済のEndpointに紐付いた未使用IPアドレスの開放 ※2.1(1h)以降ではデフォルト有効、それ以前からの場合デフォルト無効

詳しくは・・・

ACI Fabric Endpoint Learning ホワイトペーパーを参照ください!!

[https://www.cisco.com/c/en/us/solutions/collateral/  
data-center-virtualization/application-centric-  
infrastructure/white-paper-c11-739989.html](https://www.cisco.com/c/en/us/solutions/collateral/data-center-virtualization/application-centric-infrastructure/white-paper-c11-739989.html)

