



それ本当にEVPNでやるの？ EVPNの使いどころがわかる150分

2016年11月29日
古河ネットワークソリューション株式会社
神谷 尚秀

kamitani@fnsc.co.jp

- 神谷 尚秀 (古河ネットワークソリューション株式会社)
- 三宅 正浩 (ソフトバンク株式会社)
- 大久保 修一 (さくらインターネット株式会社)
- 高澤 信宏 (ヤフー株式会社)

自己紹介

■ 名前

- 神谷 尚秀

■ 所属

- 古河ネットワークソリューション株式会社
ビジネス推進部

■ 業務

- 2012 年入社
- 2015 年までルータのソフト開発に従事
- 現在は最新の通信技術調査、マーケティング、製品企画等に従事

■ コミュニティ活動

- JANOG 運営委員、wakamonog 運営委員

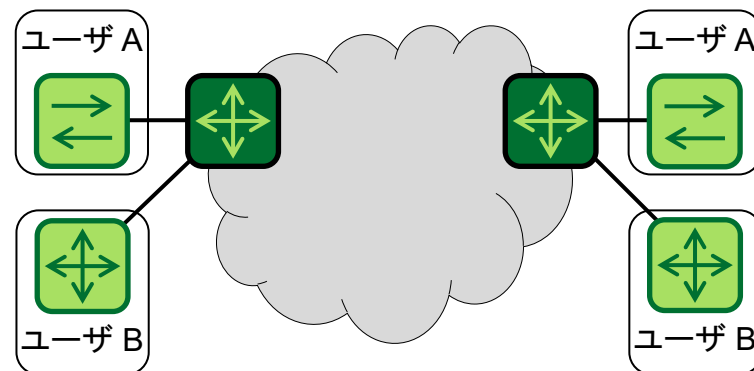
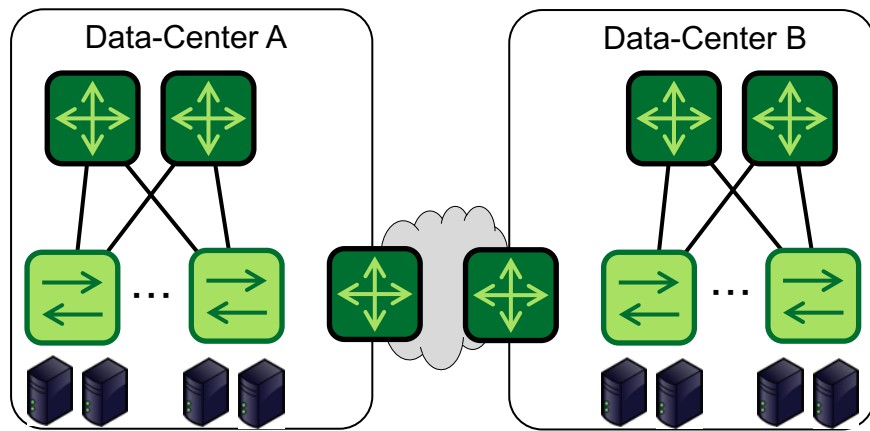
- Ethernet VPN (EVPN) が次世代の Layer-2/Layer-3 VPN 技術として注目を集めている。Interop Tokyo 2015, 2016 の ShowNet でも相互接続検証が実施されている
- EVPN を使えば従来の L2VPN と比べて下記が改善される
 - All-active による Multihoming
 - BUM トラフィックの抑制
 - IP-VPN like なオペレーションによるスケーラビリティ向上
 - L2/L3VPN 網の統合
- これまでカンファレンス等では EVPN の検討話が多かった
- 本プログラムでは EVPN の検討、検証、運用の実際を知ること、EVPN の適用箇所を見極め、ネットワーク改善のきっかけを作りたい

- 1. プログラム導入・EVPN 概要紹介**
 - 神谷 尚秀 (古河ネットワークソリューション株式会社)
- 2. EVPN って本当に良いの？ソフトバンクがキャリア EVPN を考えてみた**
 - 三宅 正浩 (ソフトバンク株式会社)
- 3. マルチベンダ環境における EVPN 構築のノウハウ ～Interop Tokyo 2016 ShowNet での相互接続検証を元に～**
 - 大久保 修一 (さくらインターネット株式会社)
- 4. コンテンツ事業者での EVPN 話**
 - 高澤 信宏 (ヤフー株式会社)
- 5. パネルディスカッション**



EVPN 概要紹介

1. L2VPN の背景
2. EVPN 機能紹介
3. EVPN プロトコル紹介
4. EVPN 動作例
5. まとめ



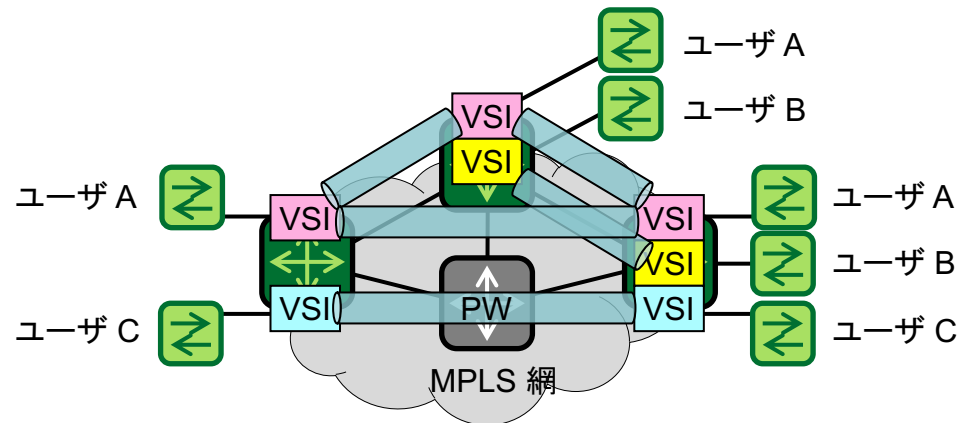
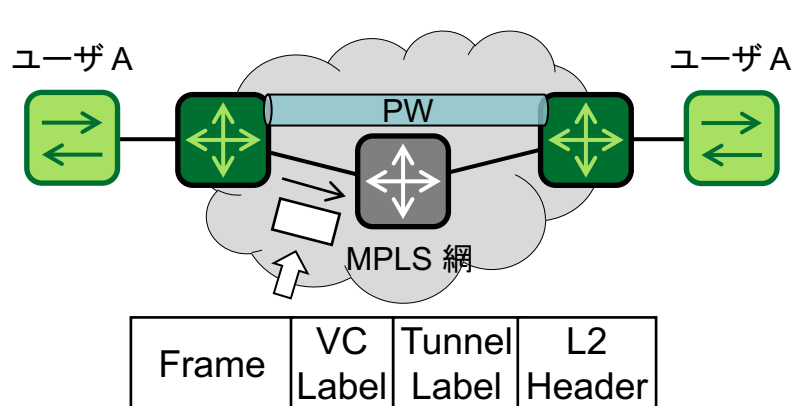
■ Data-Center の L2 網延伸

- DC 間を L2VPN によりフラットに接続することで、DC 内・DC 間での仮想マシンのライブマイグレーション等を実現

■ 広域イーサネットサービス

- ユーザ間をキャリアが提供する L2VPN サービスによりフラットに接続することで、ユーザ間で Ether フレームのやり取りを実現

これまでの L2VPN 技術一例



■ VPWS (Virtual Private Wire Service)

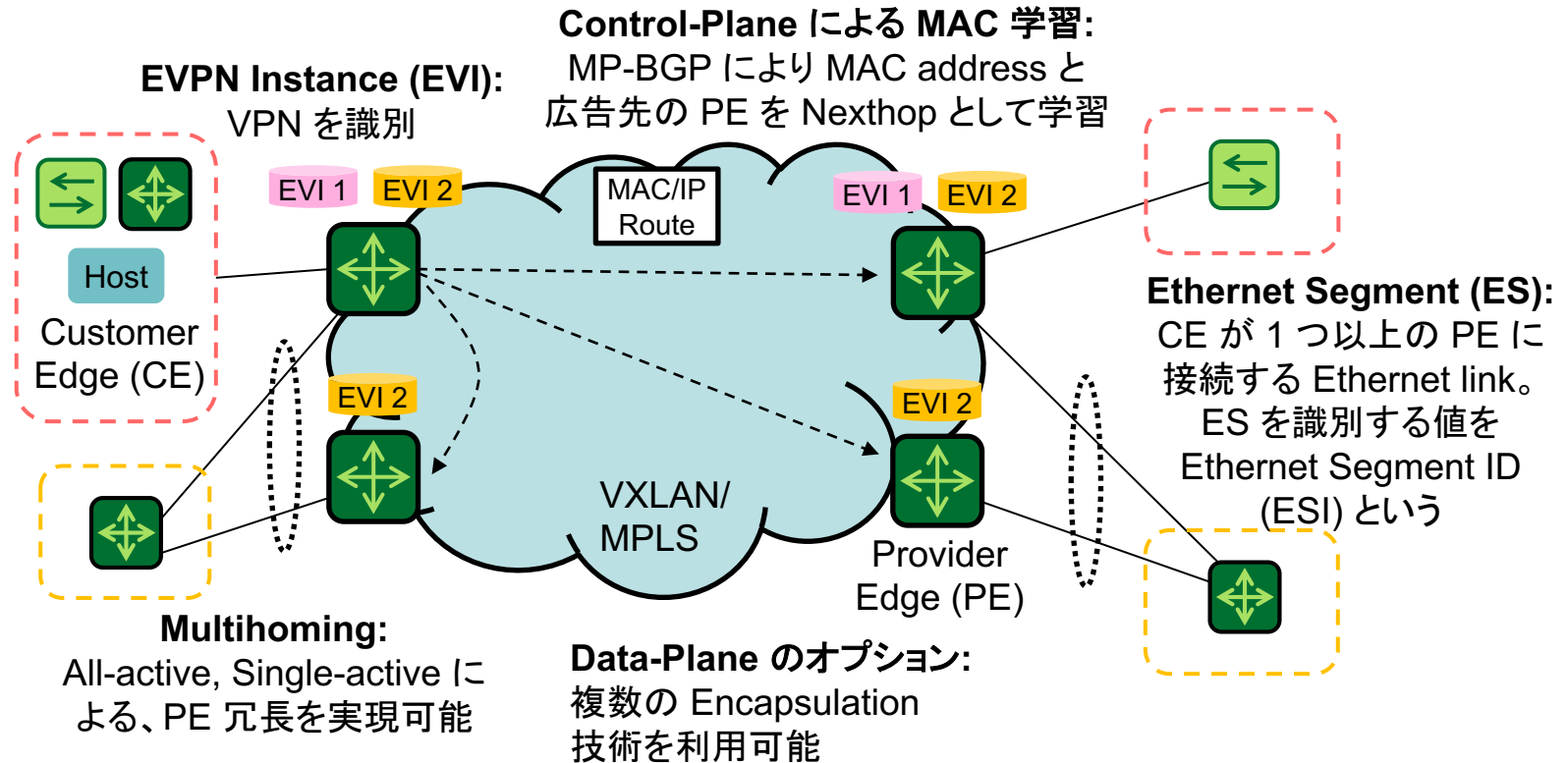
- 二つのユーザネットワークを P2P の論理回線 PW (Pseudo Wire) で接続する技術

■ VPLS (Virtual Private LAN Service)

- 複数のユーザネットワークを MP2MP で接続し、仮想的な単一 LAN を実現する技術

- これまでの L2VPN 技術には以下のような制限があった
 - Multihoming と冗長機能
 - Multicast 最適化
 - シンプルな Provisioning
 - フローベースの load balancing
 - マルチパス化
- EVPN は上記のような制限を解決できると期待されている

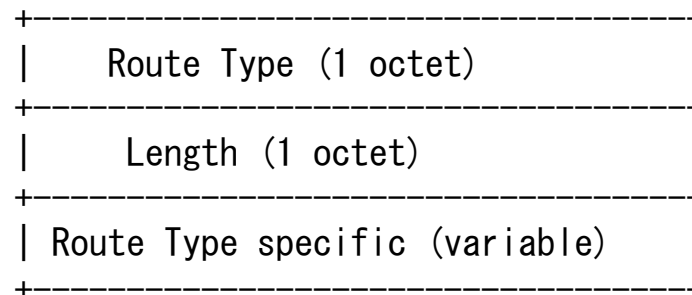
EVPN 概要



- EVPN は L2VPN を実現するための技術である
- BGP EVPN Route として MP_REACH_NLRI/MP_UNREACH_NLRI に EVPN NLRI (AFI 25, SAFI 70) を追加し、EVPN NLRI と付随する Ext-comm 等により、L2VPN を制御する
- 各 VPN は EVPN Instance (EVI) によって区別され、IP-VPN like な L2VPN の構築が可能

- EVPN NLRI の Route Type 毎に EVPN の各種機能を実現する情報が格納されている

※各 Route Type のフォーマットは Appendix 参照



EVPN NLRI フォーマット

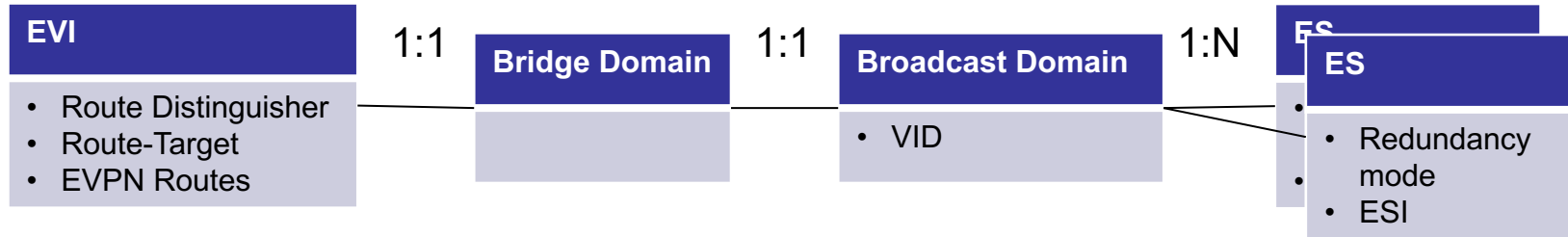
Type	Route Type Name	Usage	Advertised-with
0x1	Ethernet Auto-Discovery per ES (EAD/ES)	- Fast Convergence - Split Horizon	-ESI Label Ext-comm
	Ethernet Auto-Discovery per EVI (EAD/EVI)	- Aliasing and Backup Path	
0x2	MAC/IP Advertisement (MAC/IP)	- Determining Reachability to MAC Address - MAC Mobility - Default Gateway	- MAC Mobility Ext-comm - Default Gateway Ext-comm
0x3	Inclusive Multicast Ethernet Tag (IMET)	- Handling of Multi-Destination Traffic	- PMSI Tunnel
0x4	Ethernet Segment (ES)	- Designated Forwarder Election - Multi-homed ES Auto-discovery	- ES-Import RT Ext-comm

Route Type の名称と役割

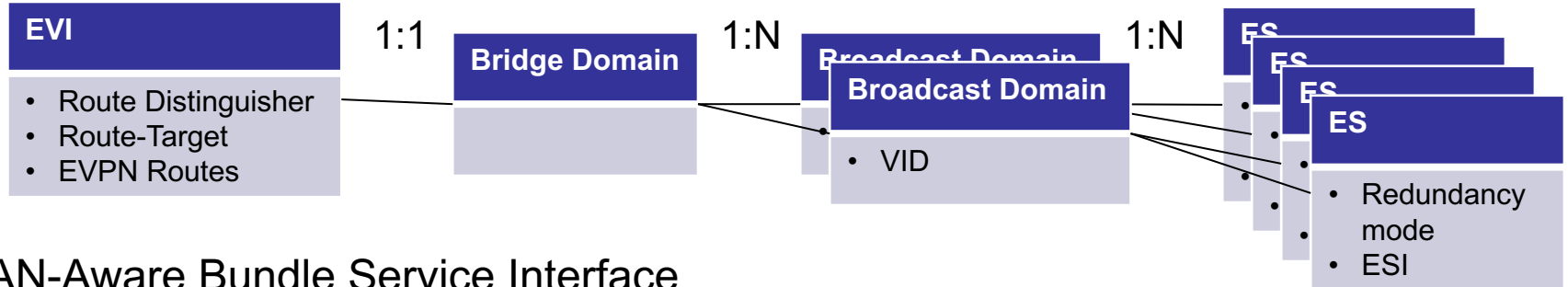
EVPN Service Interfaces

- EVPN service interface は EVI, Bridge Domain, Broadcast Domain (VLAN) の関係により、3つのタイプに区別される

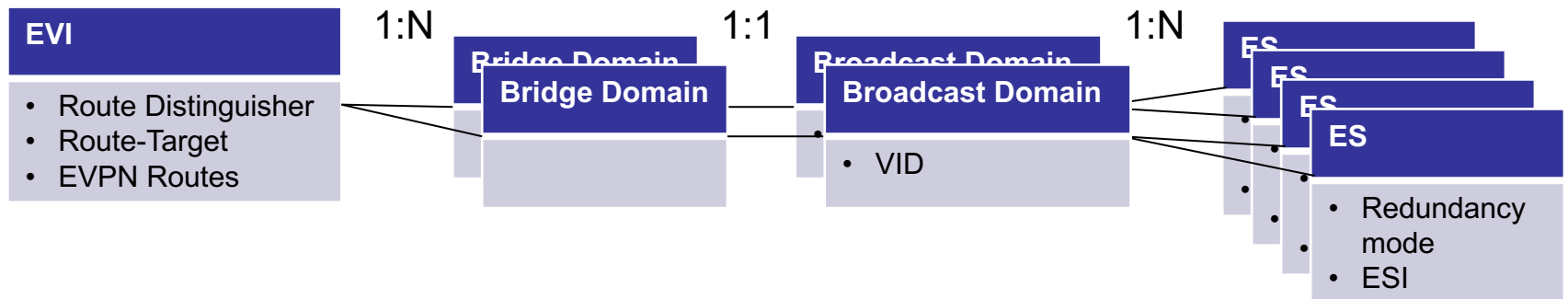
VLAN-Based Service Interface

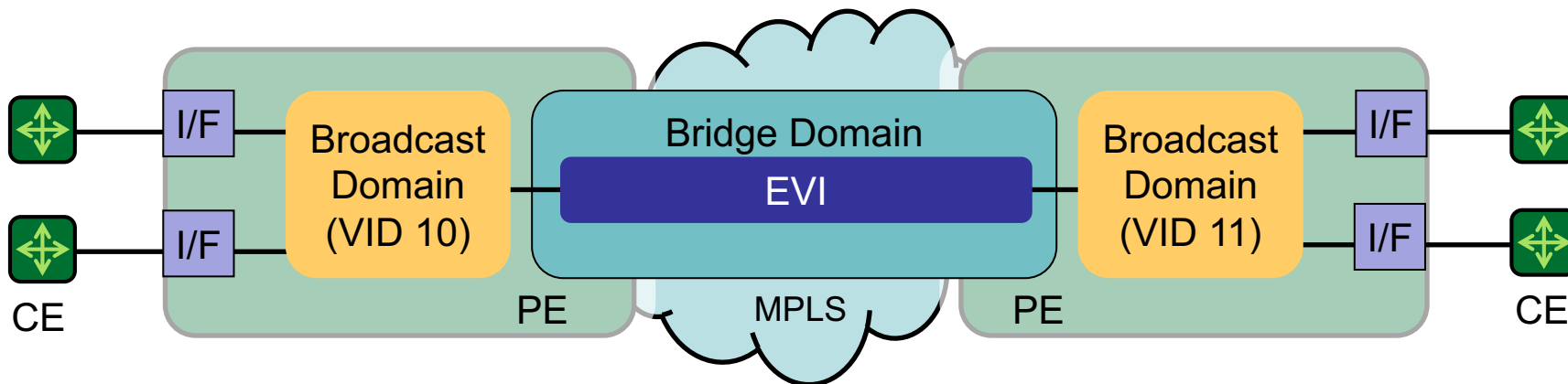


VLAN Bundle Service Interface

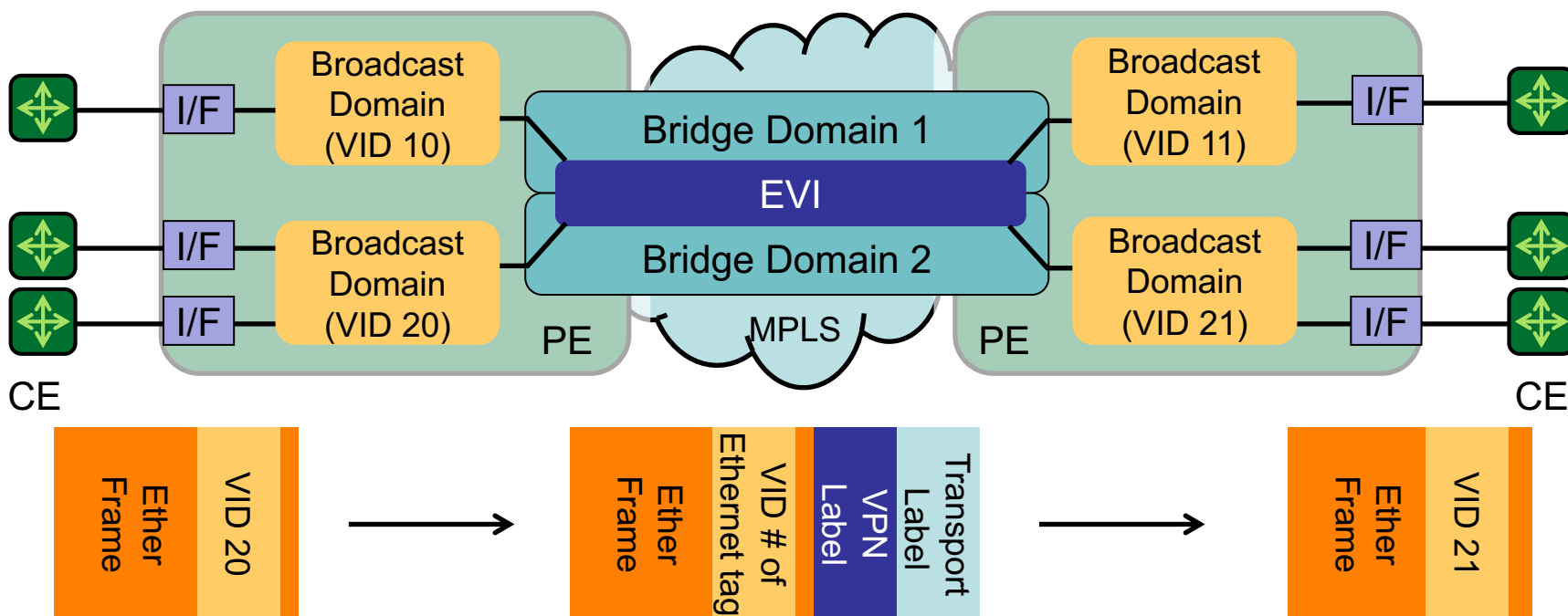


VLAN-Aware Bundle Service Interface





- 1 EVI につき 1 Broadcast Domain (VLAN) が対応
- 1 EVI につき 1 Bridge Domain が対応
- Egress PE にて VID 変換が可能
- VLAN-Based の場合、すべての EVPN Route の Ethernet tag が 0
 - Ethernet tag は Broadcast Domain を識別するために利用されるが、1 EVI につき 1 Broadcast Domain のため、Ethernet tag 値は不要



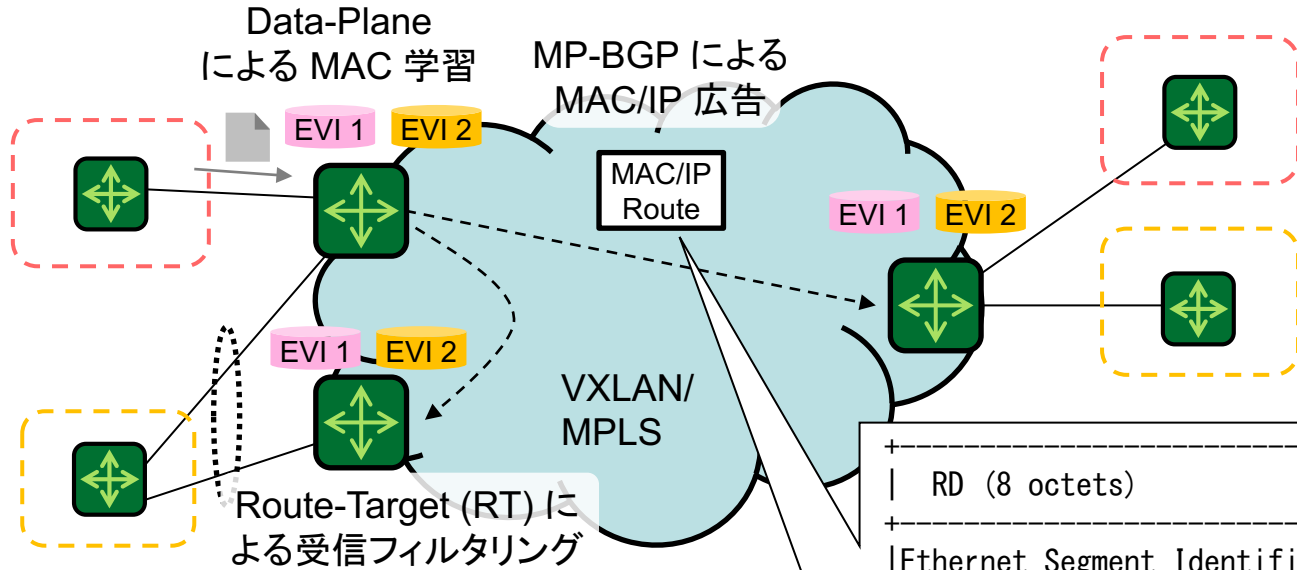
- 1 EVI につき複数 Broadcast Domain (VLAN) が対応
- 1 EVI につき複数 Bridge Domain が対応
- Egress PE にて VID 変換が可能
- EVPN Route の Ethernet tag を Broadcast Domain を識別する値として利用



EVPN 機能紹介

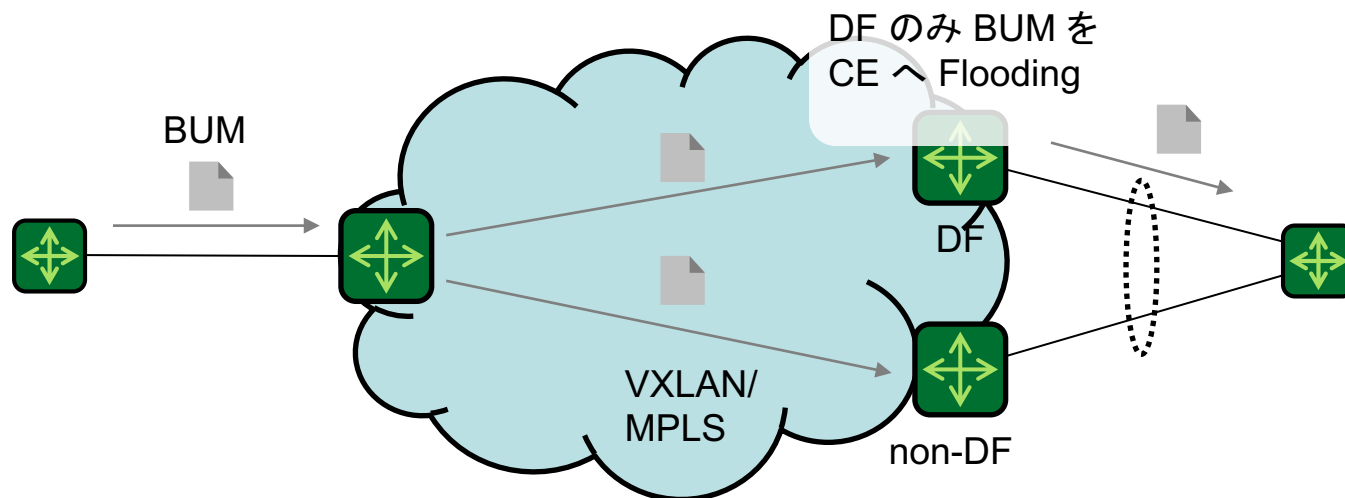
- MP-BGP による MAC 学習
- All-active Multihoming と DF Election
- Aliasing
- Mass Withdraw による Fast Convergence
- ARP/ND proxy
- MAC Mobility (Appendix を参照)
- Split Horizon (※Appendix を参照)
- Inter-subnet Forwarding (※Appendix を参照)

MP-BGP による MAC 学習



- PE は Local-CE から Ethernet frame を受信することで、Local-CE の MAC address を Data-Plane で学習し、Remote-PE へ MAC/IP Route を広告する
- PE は Remote-PE から MAC/IP Route を受信することで、MAC address を学習することが出来る

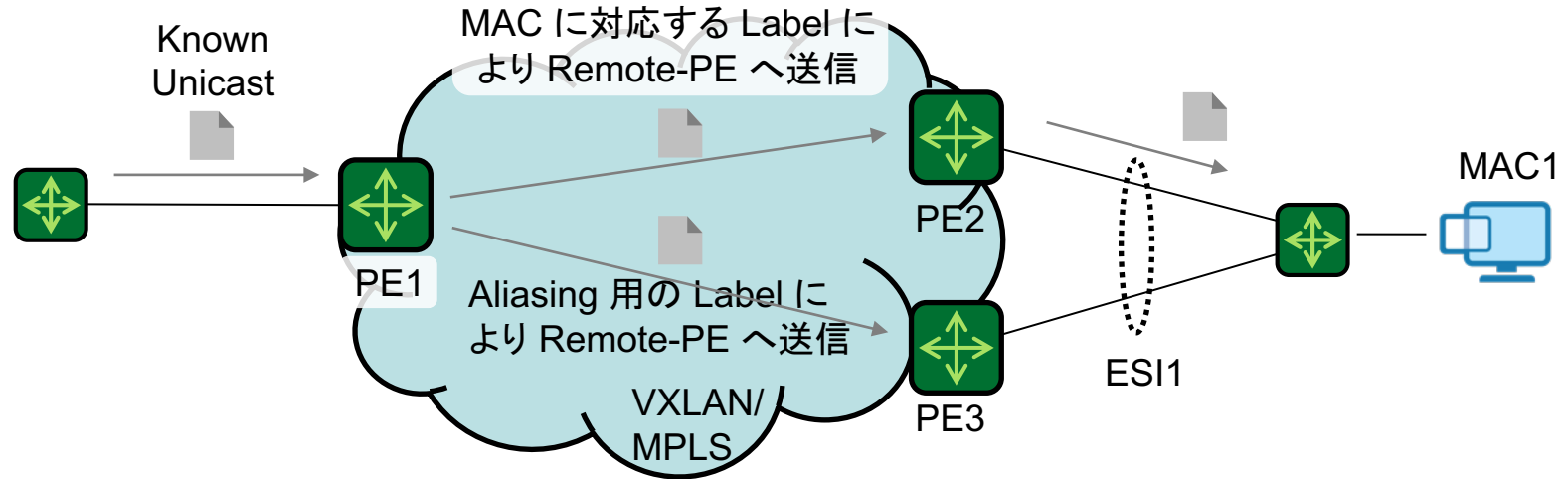
RD (8 octets)
Ethernet Segment Identifier (10 octets)
Ethernet Tag ID (4 octets)
MAC Address Length (1 octet)
MAC Address (6 octets)
IP Address Length (1 octet)
IP Address (0, 4, or 16 octets)
MPLS Label1 (3 octets)
MPLS Label2 (0 or 3 octets)



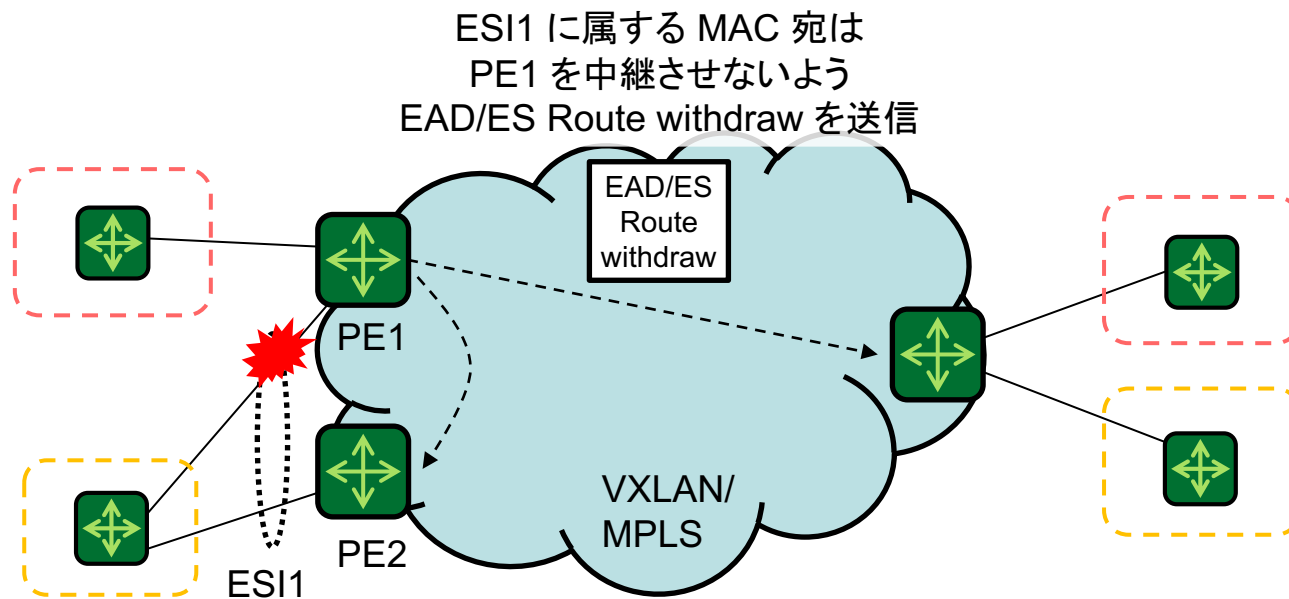
- EVPN は CE が 2 つ以上の PE と接続することで All-active Multihoming を実現可能
- 同一 ES に接続する PEs は、BUM を CE へ中継する PE (DF: Designated Forwarder) を決定し、CE への BUM パケットが複数送信されることを防止する

※BUM: Broadcast, Unknown unicast, Multicast

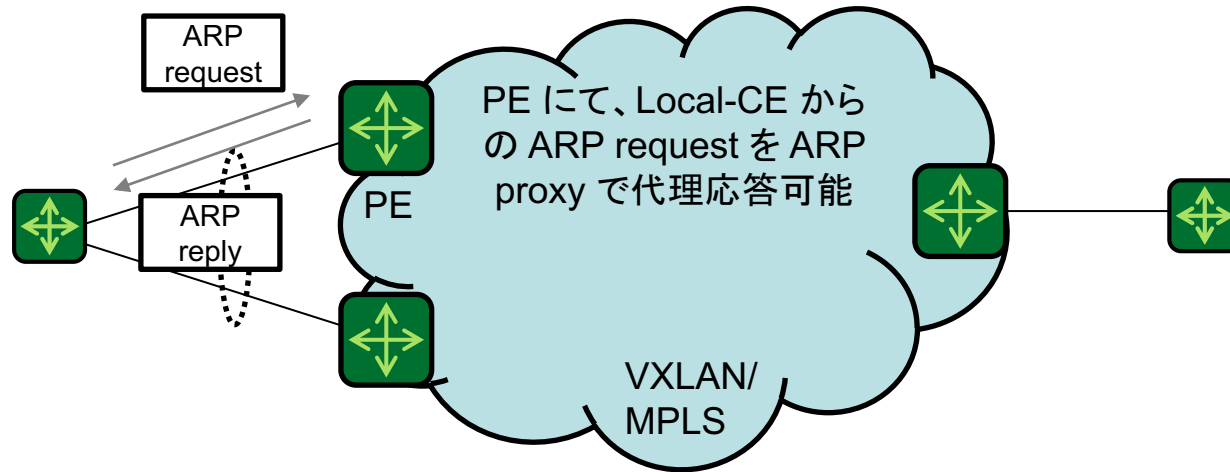
Aliasing



- All-active 環境の場合、同一 ES に接続する PEs への Load Balancing を実現可能
 - 同一 ES に接続する PEs (PE2, PE3) は、同一 ES への到達性があることを Ethernet Auto Discovery per EVI (EAD/EVI) Route により、EVPN 環境構築時に Remote-PE (PE1) へ通知している
 - PE1 が PE2 からのみ MAC1 の MAC/IP Route を受信している場合
 - ✓ PE1 が MAC1 宛のパケットを PE2 へ中継する際、MAC/IP Route で学習した Label を利用する
 - ✓ PE1 が MAC1 宛のパケットを PE3 へ中継する際、EAD/EVI Route で学習した Aliasing 用の Label を利用する



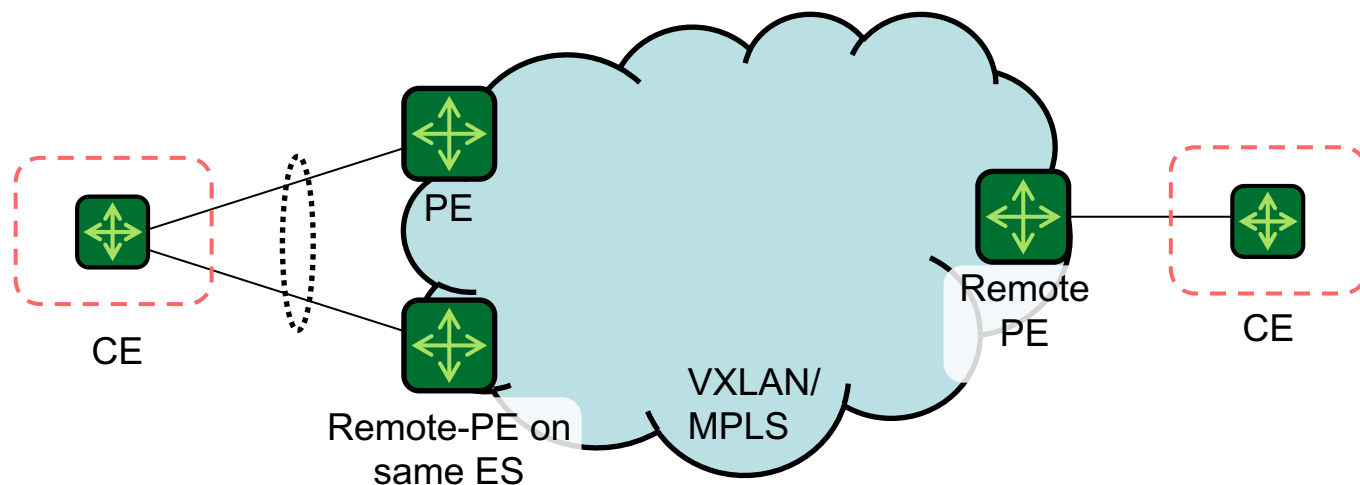
- Multihoming 環境の場合、PE と CE 間で回線障害が発生した時、高速な Convergence を実現可能
 - PE が接続している ES 情報を Ethernet Auto Discovery per ES (EAD/ES) Route により、EVPN 環境構築時に Remote-PE へ通知している
 - PE は CE 間で回線障害が発生した時、Remote-PE へ EAD/ES Route の withdraw を送信することで、Remote-PE は障害が起こっている PE への冗長パスを削除し、通信を早急に復旧する



- PE は Local-CE からの ARP request, ND solicitation が送信されたとき、ARP/ND Proxy により代理応答することで、コア網への BUM を抑制することが可能
- PE は Local-CE からのパケットを Snooping することで、MAC Address 及び IP address を学習し、ARP/ND Table を更新する
- PE は Remote-PE へ MAC/IP Route により、MAC/IP binding 情報を広告し、ARP/ND Table を更新させ、Remote-PE でも ARP/ND proxy を実現することが可能

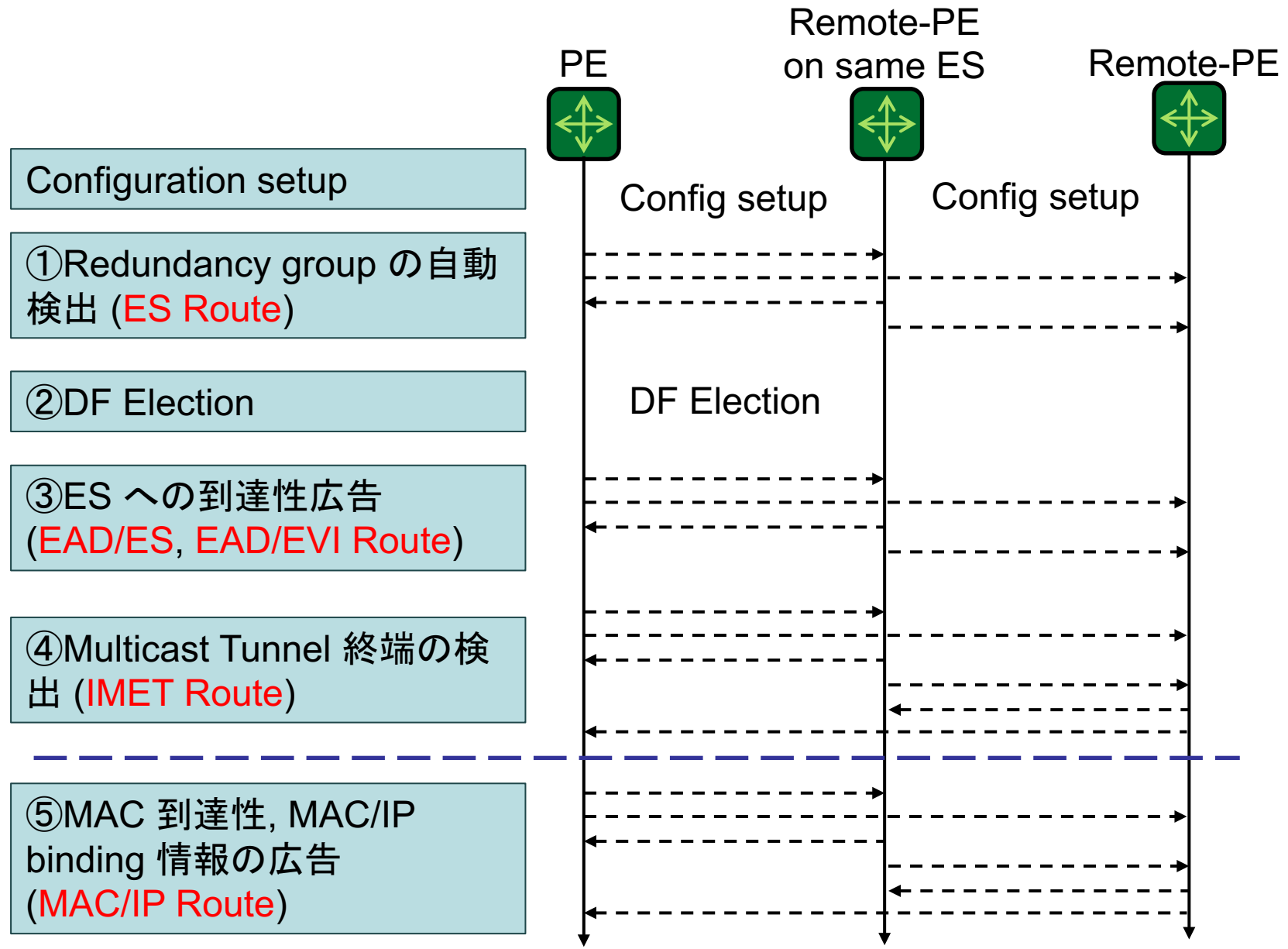


EVPN プロトコル紹介

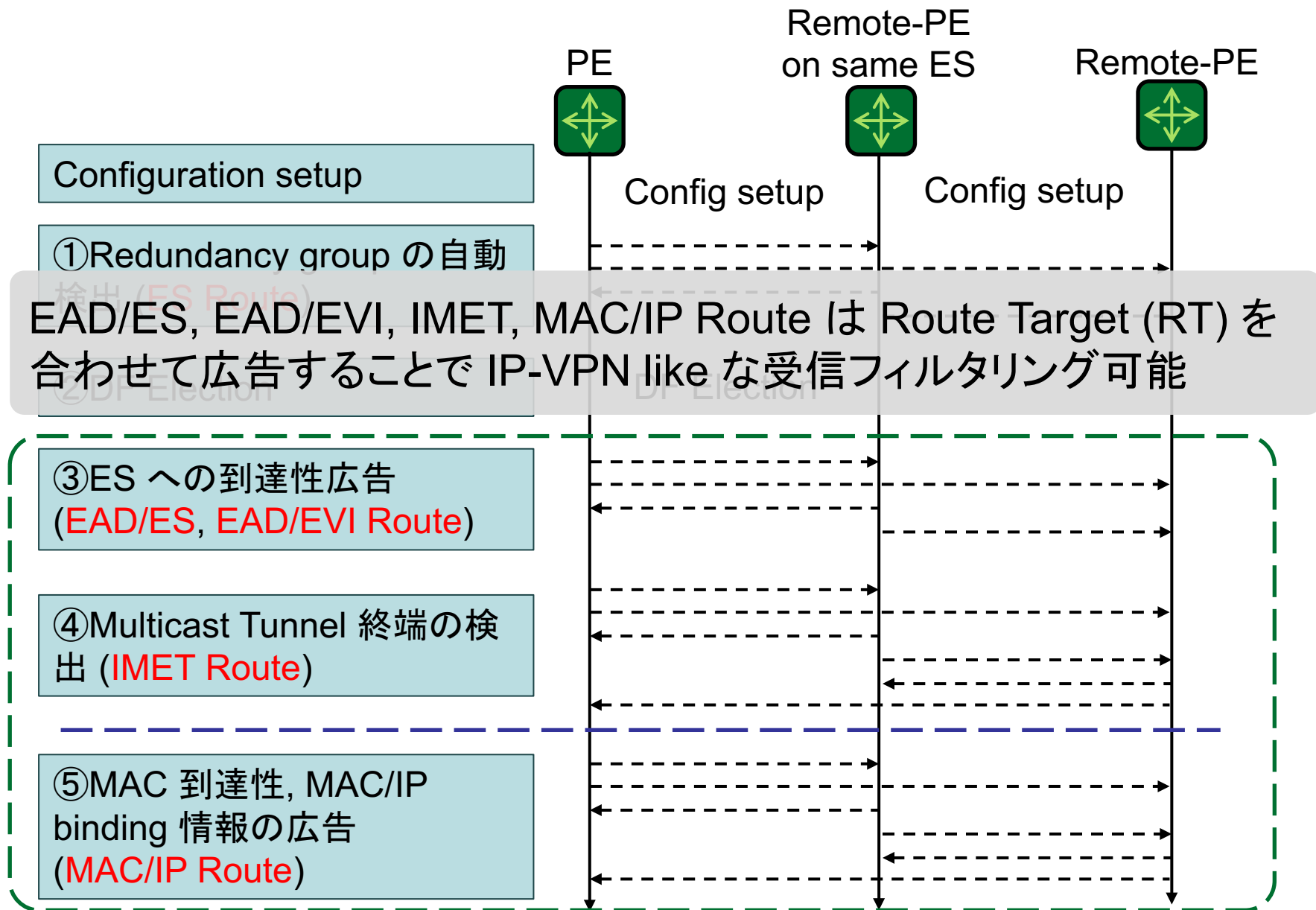


- Multihoming が設定されている PE, Remote-PE on same ES と、Singlehoming が設定されている Remote PE で構成される EVPN 網における EVPN startup scenario 例を紹介

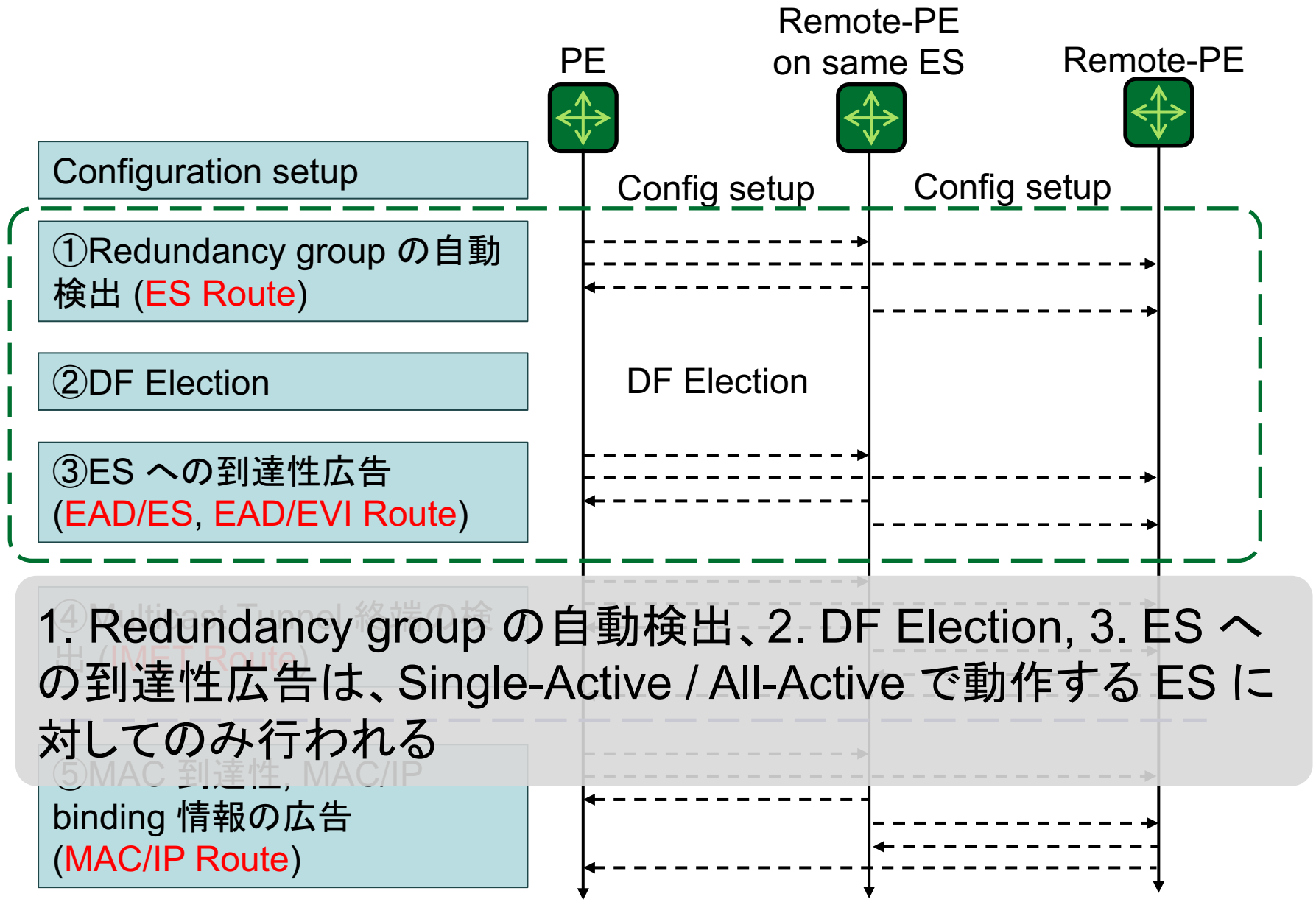
EVPN startup scenario



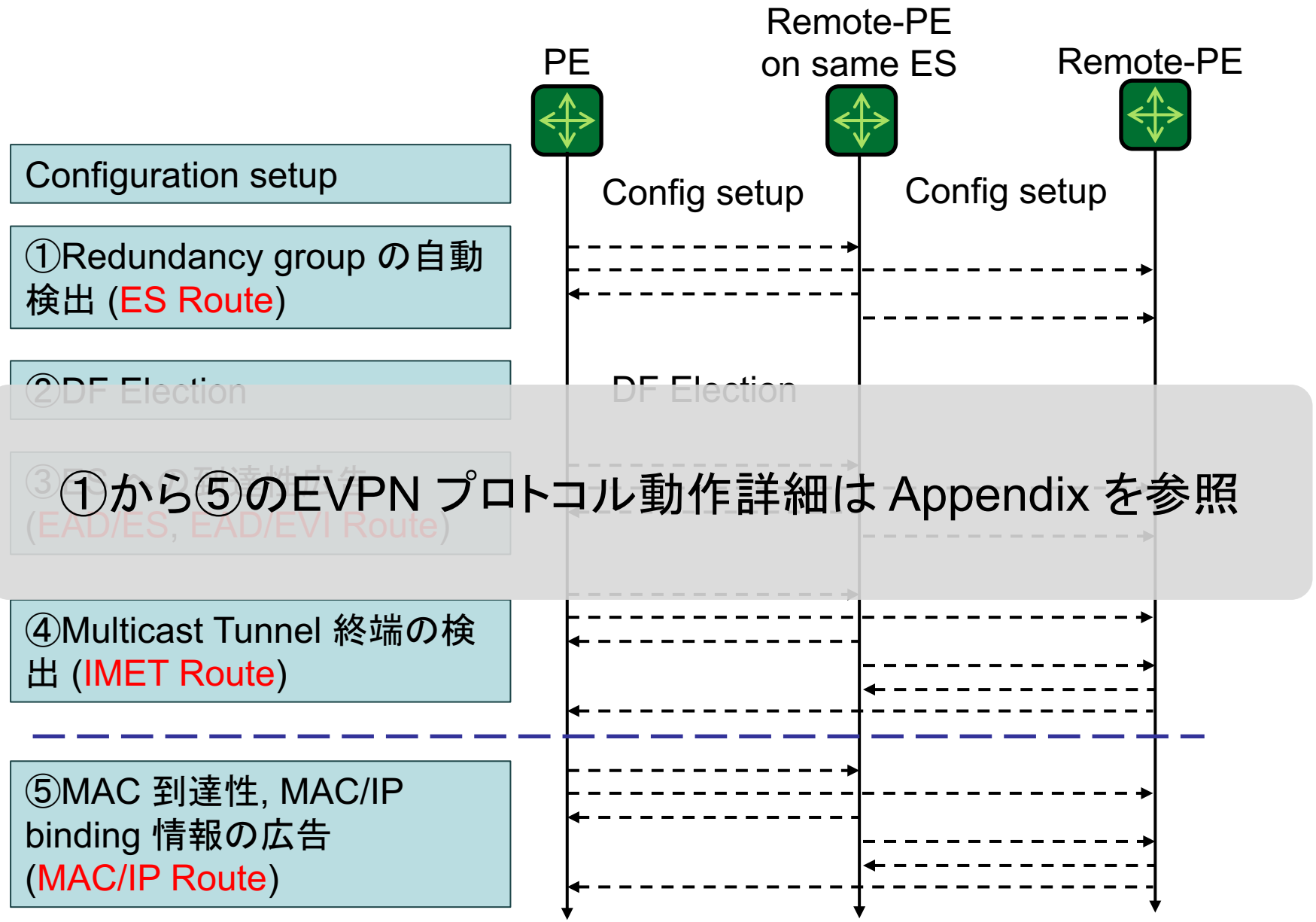
EVPN startup scenario



EVPN startup scenario



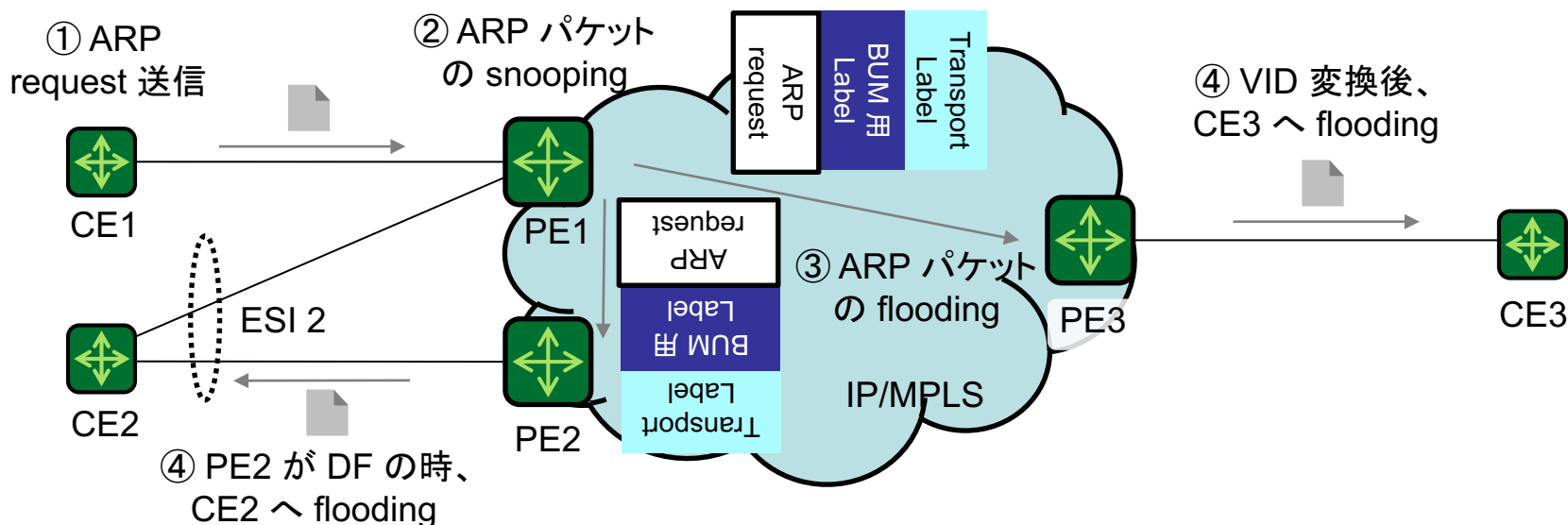
EVPN startup scenario





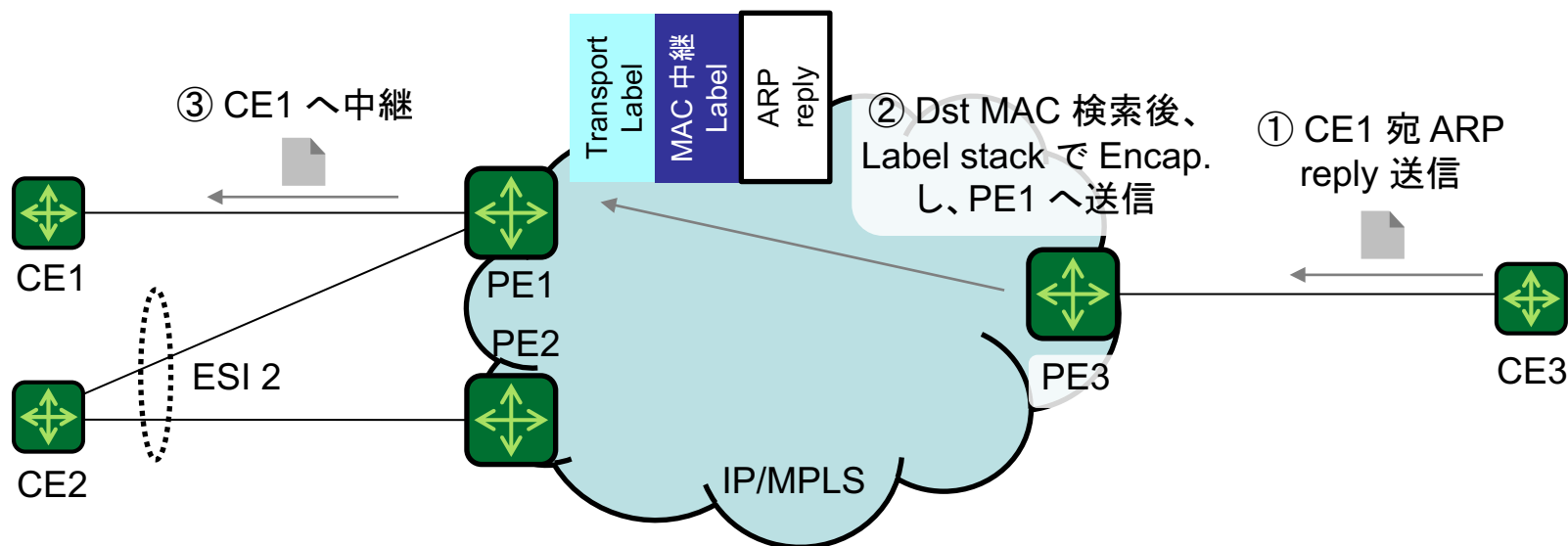
EVPN 動作例

Packet Walk (BUM from CE1)



- ① CE-VID = 1 の ARP request (source-MAC CE1, destination-IP CE3) が PE1 に送信される
- ② PE1 では CE-VID に従って、フレームは EVI1 のものと判別される。PE1 は ARP Snooping により、EVI-1 の MAC Table, Proxy-ARP table の更新を行なう
- ③ PE1 は PE2, PE3 へ ARP を Flooding する
- ④ PE2 と PE3 に MPLS-encapsulated frame が到達する。PE2 は ESI=2 の DF であれば、CE2 へ Broadcast を送信する。PE3 は必要に応じて VID translation を行なった後、CE3 へ Flooding する

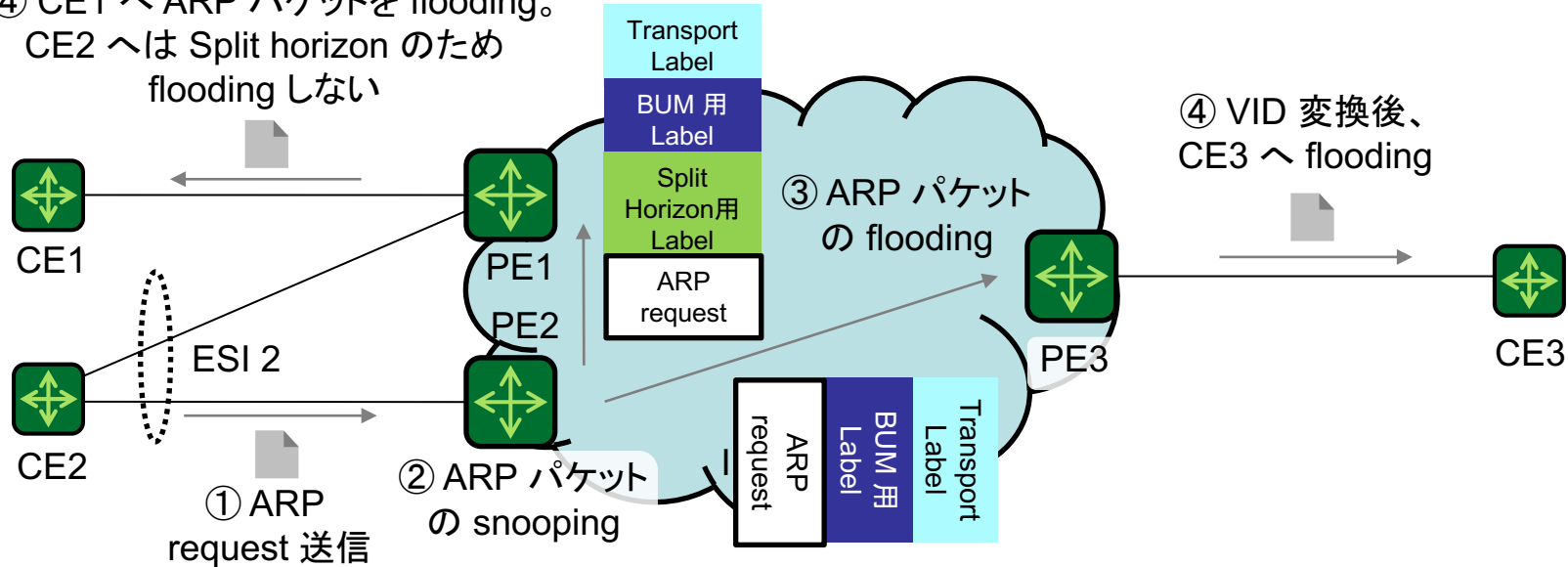
Packet Walk (Unicast from CE3 to CE1)



- ① CE-VID=1 の ARP reply (source-MAC CE3, destination-IP CE1) が PE3 に送信される
- ② PE3 では CE-VID に従って、フレームは EVI1 のものと判別される。MAC Table 検索を行い、CE1-MAC に紐づく Label Stack で Encapsulate して PE1 に送信する
- ③ PE1 はフレームを受信すると CE1 に送信する

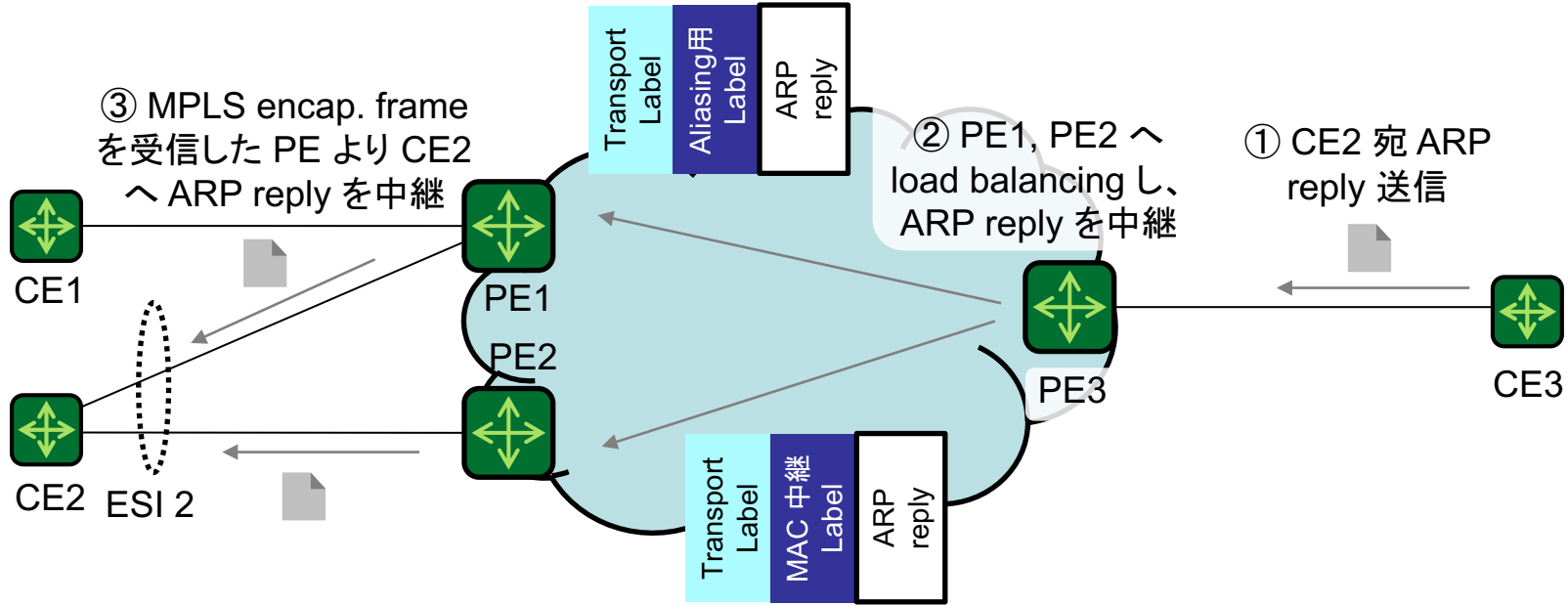
Packet Walk (BUM from CE2)

④ CE1 へ ARP パケットを flooding。
CE2 へは Split horizon のため
flooding しない



- ① CE-VID=1 の ARP request (source-MAC CE2, destination-IP CE3) が PE2 に送信される
- ② PE2 では CE-VID に従って、フレームは EVI1 のものと判別される。PE2 は ARP Snooping により、MAC Table, Proxy-ARP table への更新を行なう
- ③ PE2 は PE1, PE3 へ ARP を Flooding する。PE1 へ送信されるパケットには Split Horizon 用の Label が入る
- ④ PE1 と PE3 に MPLS-encapsulated frame が到達する。PE1 は Split Horizon 用の Label により CE2 への Flooding は行わない。PE3 では、必要に応じて VID translation を行なった後、CE3 へ Flooding する

Packet Walk (Unicast from CE3 to CE2)



- ① CE-VID=1 の ARP reply (source-MAC CE3, destination-IP CE2) が PE3 に送信される
- ② PE3 では CE-VID に従って、フレームは EVI1 のものと判別される。MAC Table 検索を行い、CE2-MAC は ESI2 に属するものであると判断する。CE2-MAC への Next-hop としては PE1 or PE2 のどちらへ中継するか決定し、CE2-MAC に紐づく Label Stack で Encapsulate して PE1 or PE2 に送信する
- ③ PE1 or PE2 は frame を受信すると CE2 に送信する

- EVPN は既存の L2VPN 技術の課題を下記機能により解決
 - All-active Multihoming による PE の Load Balancing
 - ARP/ND proxy による BUM 制御の最適化
 - IP-VPN like なオペレーションによるスケーラビリティ向上
 - Mass Withdraw による Fast Convergence
 - L2/L3VPN 網の統合
- 参考文献
 - RFC 7432 - BGP MPLS-based Ethernet VPN:
<https://tools.ietf.org/html/rfc7432>
 - IETF Draft - A Network Virtualization Overlay Solution using EVPN:
<https://tools.ietf.org/html/draft-ietf-bess-evpn-overlay-05>
 - IETF Draft - Usage and applicability of BGP MPLS based Ethernet VPN:
<https://tools.ietf.org/html/draft-ietf-bess-evpn-usage-03>

- 1. プログラム導入・EVPN 概要紹介**
 - 神谷 尚秀 (古河ネットワークソリューション株式会社)
- 2. EVPN って本当に良いの？ソフトバンクがキャリア EVPN を考えてみた**
 - 三宅 正浩 (ソフトバンク株式会社)
- 3. マルチベンダ環境における EVPN 構築のノウハウ ～Interop Tokyo 2016 ShowNet での相互接続検証を元に～**
 - 大久保 修一 (さくらインターネット株式会社)
- 4. コンテンツ事業者での EVPN 話**
 - 高澤 信宏 (ヤフー株式会社)
- 5. パネルディスカッション**



パネルディスカッション

■ マイグレーション関連

- 現行網からどのように EVPN 網へ移行していくか

■ オペレーション関連

- EVPN 網の運用者に求められるスキルとは
- EVPN を導入することで運用効率・安定性は向上するか

■ EVPN 関連の標準化動向

- EVPN の IETF における標準化動向と EVPN の今後の発展に期待することとは

■ 現行網から EVPN 網へ移行する際のハードルとは？

EVPNのメリット -VPLSとの比較-

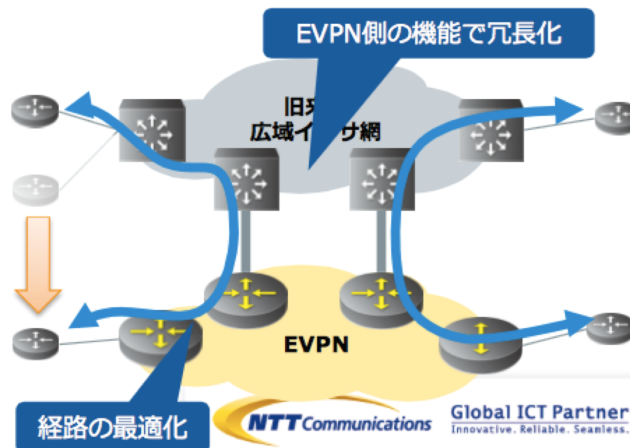
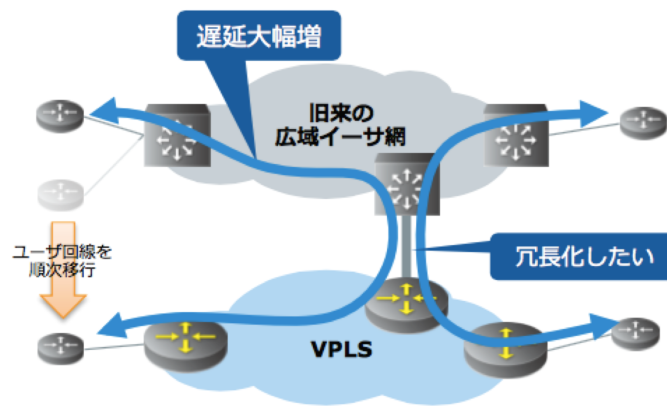
■ マイグレーション

VPLS

- 冗長接続の実現に課題
 - ✓ STP
 - ✓ G.8032
 - E-OAMの実装差分により接続できない事例
 - ✓ VPLS-Multihoming
 - 独自実装に頼る方式、PW数が2倍に
- 結果single接続による遅延の増加

EVPN

- EVPNのMultihoming(single-active)を使用する場合ならNW間のネゴシエーションは不要
- MAC毎にprimary PEを選択できれば、遅延最適化も可能か



- EVPN 網を運用するために必要な知識とは
 - 知識習得のためにどのような取り組みが必要か
- EVPN を導入することで運用効率・安定性は向上するか

- EVPN 関連技術の標準化動向について
- EVPN の今後の発展に期待することとは

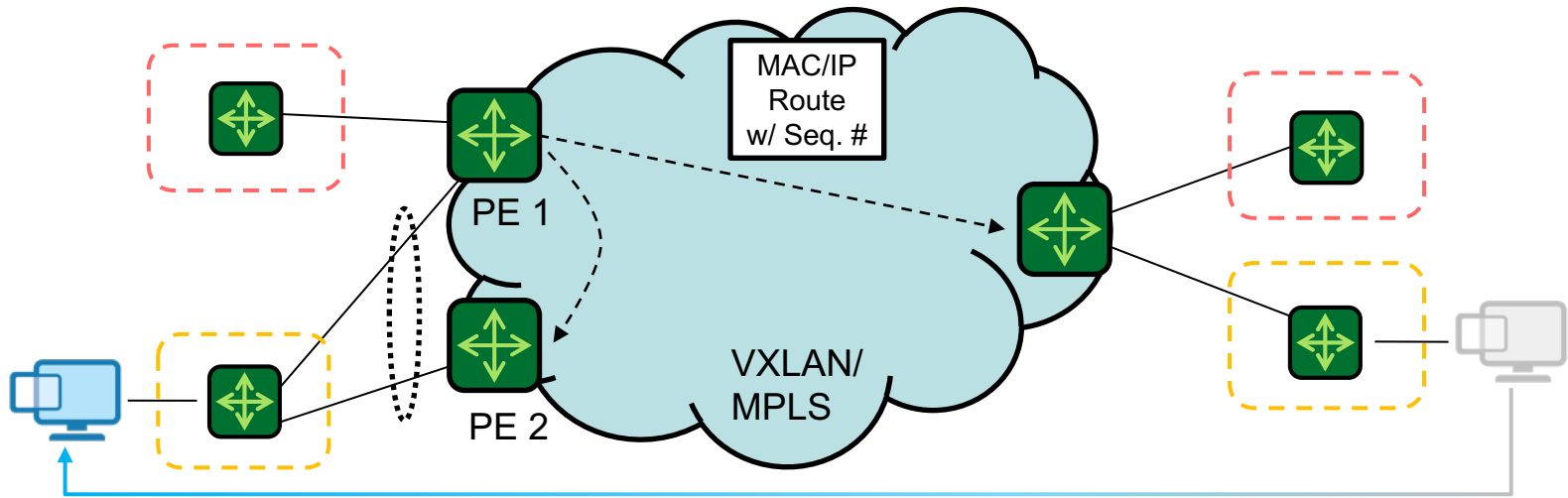


Appendix

■EVPN の各機能 (cont.)

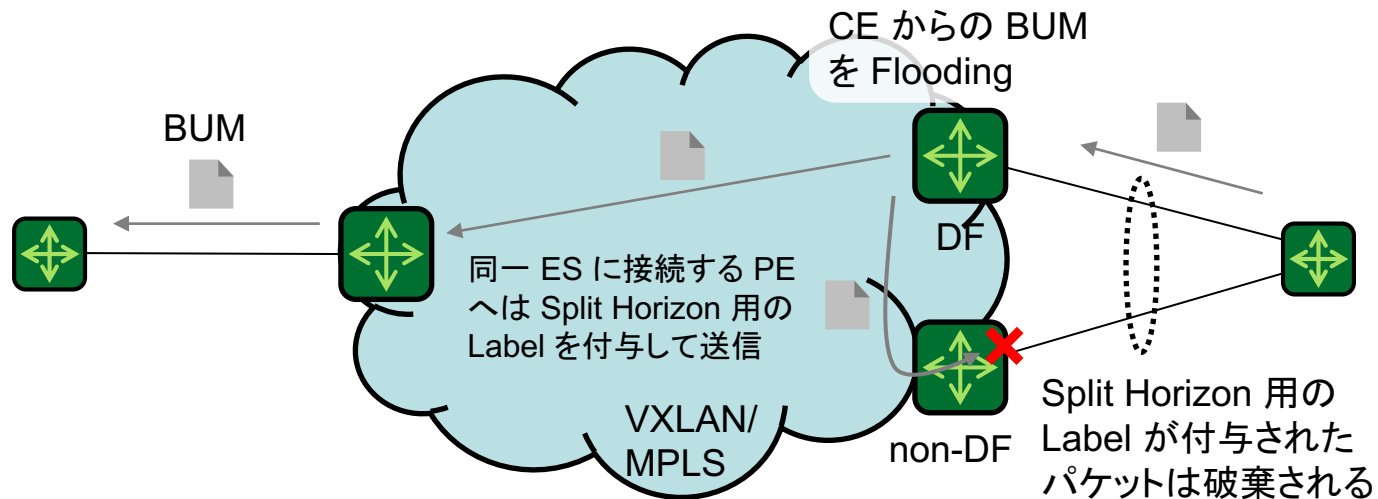
MAC Mobility

- ある ES 接続されている Host がライブマイグレーション等により、同じ EVPN に属している別 ES へ移動 (MAC Mobility) した時、最新の EVPN Route Table へ更新し、MAC address の duplicate を防止可能
- ES 間での MAC Mobility が発生し、Host が移動した先の PE は Data-Plane で MAC address を学習した時、MAC/IP Route と合わせて MAC Mobility Ext-comm を広告することで、Ext-comm 内のシーケンス番号 (MAC Mobility の回数に対応) により、シーケンス番号が一番大きな値の MAC/IP Route を最新の経路情報として各 PE に受信させる



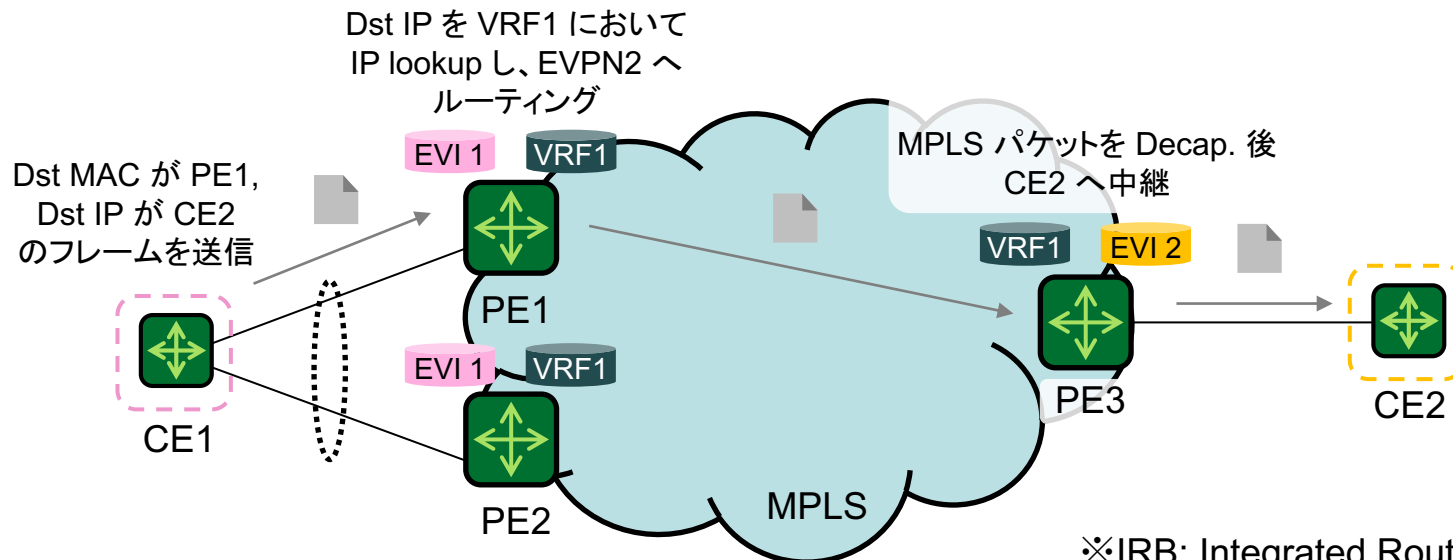
Split Horizon

- Local-CE から送信された BUM は Core 網に Flooding される。
Multihomed 環境では、Split Horizon 機能により、BUM を送信した CE と同じ ES に接続する PE からの BUM 折り返し送信を防止することが可能
- PE は CE から受信した BUM を同一 ES に接続された他の PE に Flooding する場合、Split Horizon 用の Label を付与して送信する
- Split Horizon 用の Label が付与されたパケットを受信した PE はパケットを破棄する



Inter-subnet Forwarding

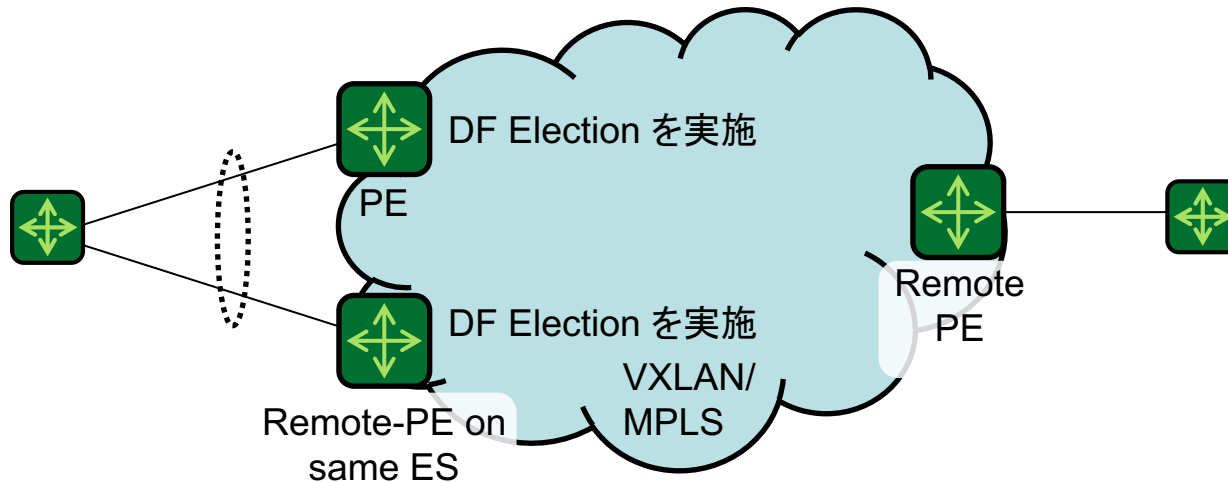
- Inter-subnet forwarding は EVPN で L3 中継を実現する技術
 - PE において IRB IF を利用し、異なる Subnet 間の中継を L3 domain を通じて行う
 - PE は Local 学習や EVPN Route により学習した IP Host address、および、IP prefix をルーティングテーブルへ登録し、必要に応じてパケットをルーティングする
 - ルーティングテーブルは VRF により分離可能



■EVPN プロトコル動作詳細

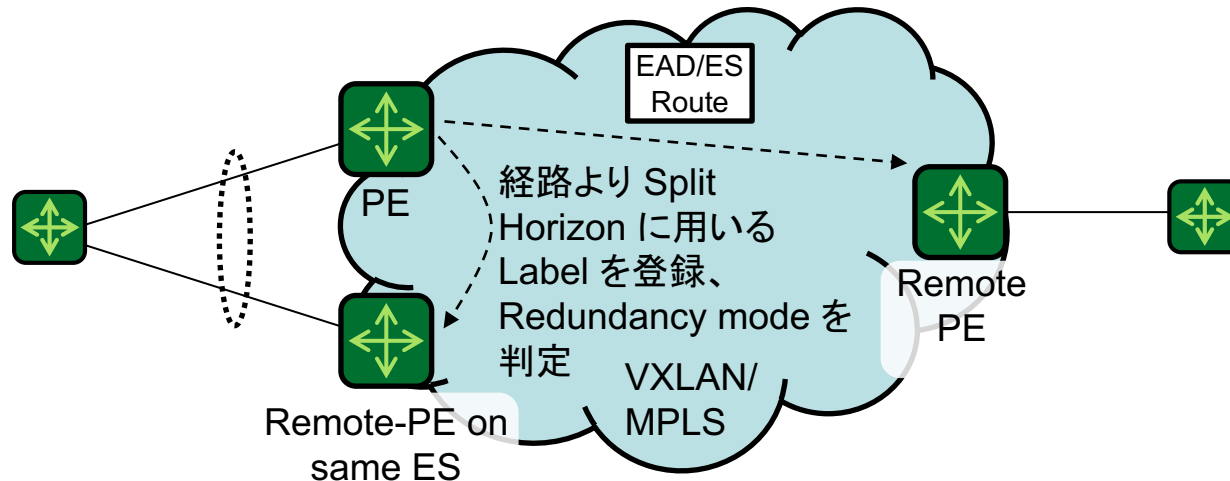
② Designated Forwarder Election

- 同一の Multihomed CE に接続された PEs は ES Route を交換した後、CE へ BUM を送信する PE (DF) の選択を行なう



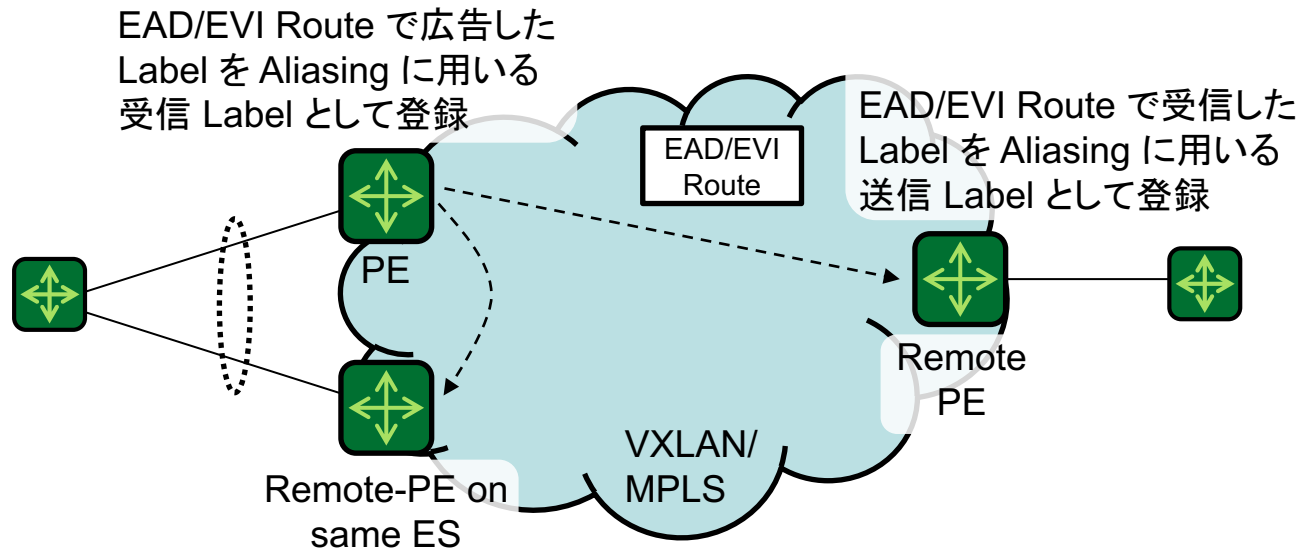
③ ES への到達性広告 (EAD/ES Route)

- PE は EAD/ES Route を広告することで、Split Horizon に用いる Label の登録と、Redundancy mode を設定する
- 1. PE は EAD/ES Route と合わせて ESI Label Ext-comm を広告し、Ext-comm に含まれる Label を Split Horizon に用いる受信 Label として登録する
- 2. Remote-PE は受信した経路情報より、Split horizon に用いる送信 Label として登録し、PE が接続する ES の Redundancy mode (All-active/Single-active) を判定する



③ ES への到達性広告 (EAD/EVI Route)

- PE は EAD/EVI Route を広告することで、Aliasing に用いる Label を設定する
 1. PE は EAD/EVI Route を広告し、EAD/EVI Route に含まれる Label を Aliasing に用いる受信 Label として登録する
 2. Remote-PE は受信した経路情報に含まれる Label を Aliasing に用いる送信 Label として登録する

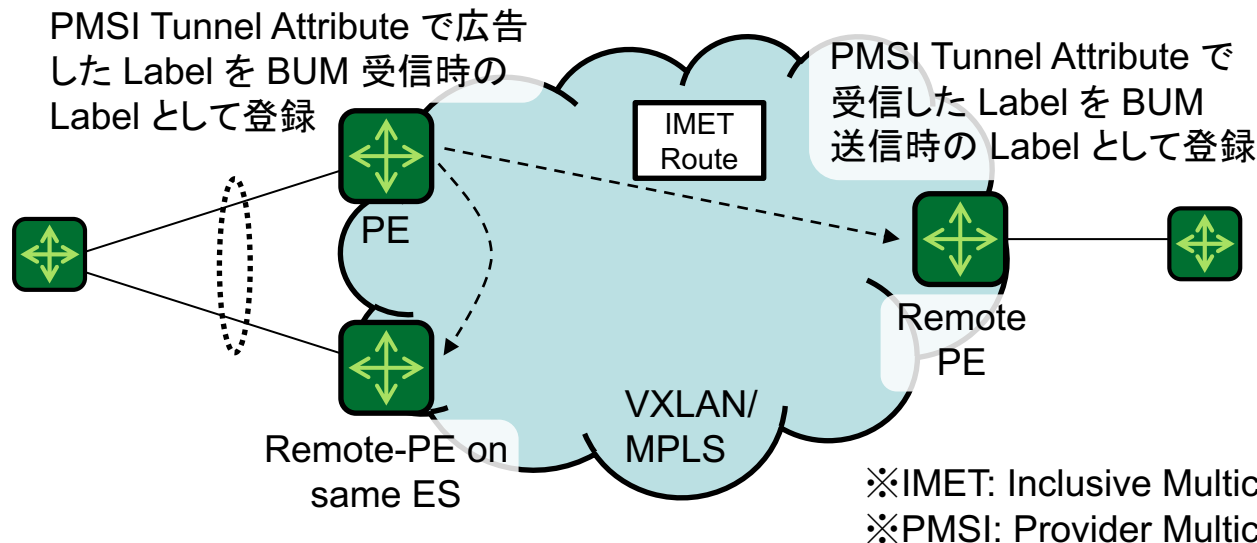


④ Multicast Tunnel 終端の検出

- PE は IMET Route を広告することで、BUM 送受信時に用いる Label の設定を行う

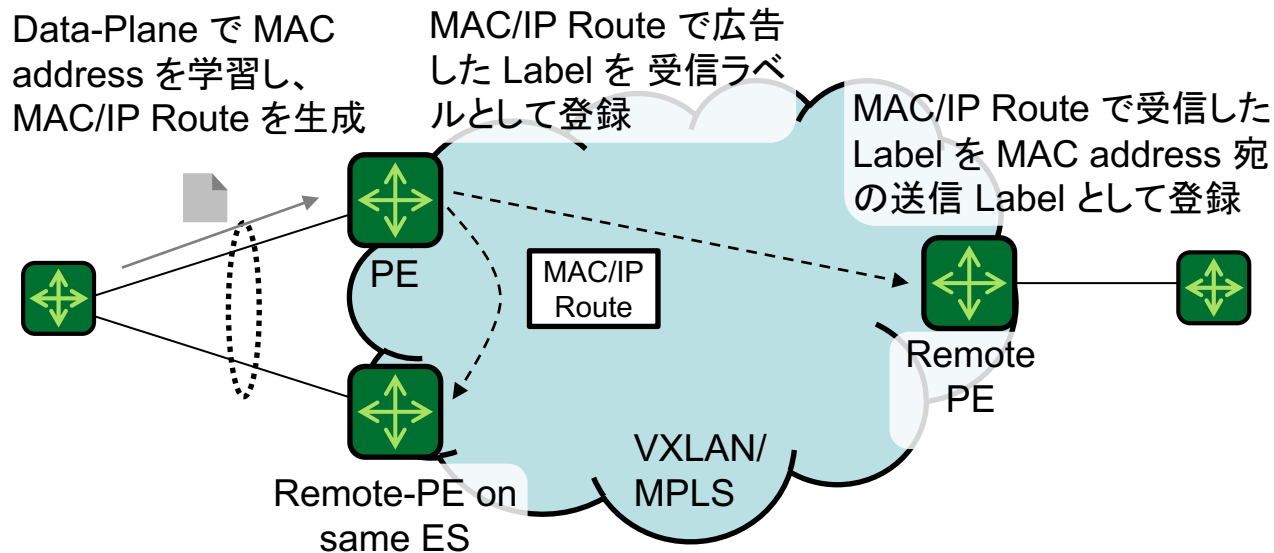
※下記は BUM 処理に Ingress replication を用いる場合の動作例

1. PE は IMET Route と合わせて PMSI Tunnel Attribute を広告し、PMSI Tunnel Attribute で広告した Label を BUM 受信時に付与させる Label として登録する
2. Remote-PE は受信した経路情報に含まれる Label を BUM 送信時に付与する Label として登録する



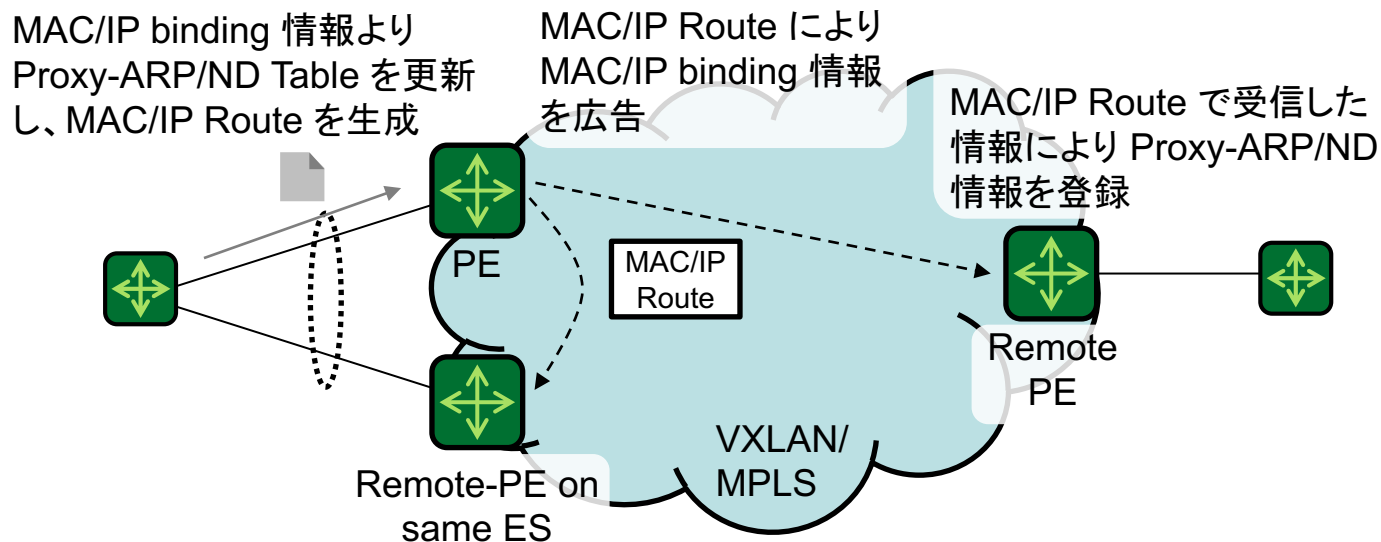
⑤MAC 到達性の広告

- PE は Data-Plane で学習した MAC address 情報を MAC/IP Route として広告することで、L2VPN に用いる Label の設定を行う
- 1. PE は学習した MAC address 情報を保持した MAC/IP Route を広告し、MAC/IP Route で広告した Label を MAC address に対する受信 Label として登録する
- 2. Remote-PE は MAC/IP Route に含まれる Label を MAC address 宛の送信 Label として登録する

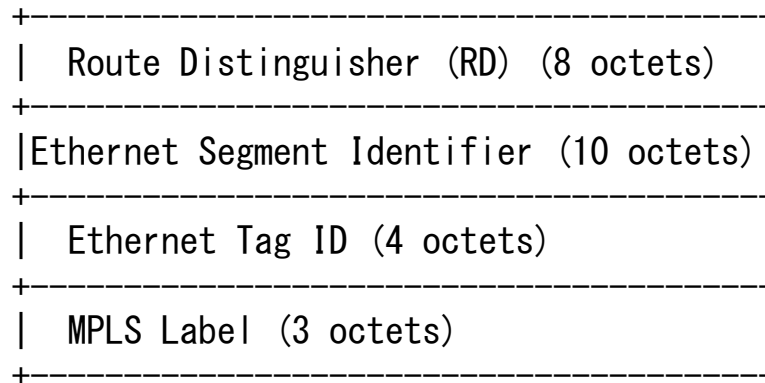


⑤ MAC/IP Binding 情報の広告

- PE は、ARP snooping 等により学習した MAC/IP binding 情報を MAC/IP Route として広告することで、Proxy-ARP/ND Table の設定を行う
- 1. PE は学習した MAC/IP binding 情報を保持した MAC/IP Route を広告する
- 2. Remote-PE は MAC/IP Route で受信した MAC/IP binding 情報より Proxy-ARP/ND Table を更新する



■各種 Route Type フォーマット



- EAD Route には EAD/ES Route および EAD/EVI Route がある
 - Ethernet Tag ID が MAX-ET なら EAD/ES, MAX-ET {0xFFFFFFFF} でなければ EAD/EVI となる
- EAD/ES Route と合わせて ESI Label Ext-comm (Optional, Transitive) を送信する

フォーマット

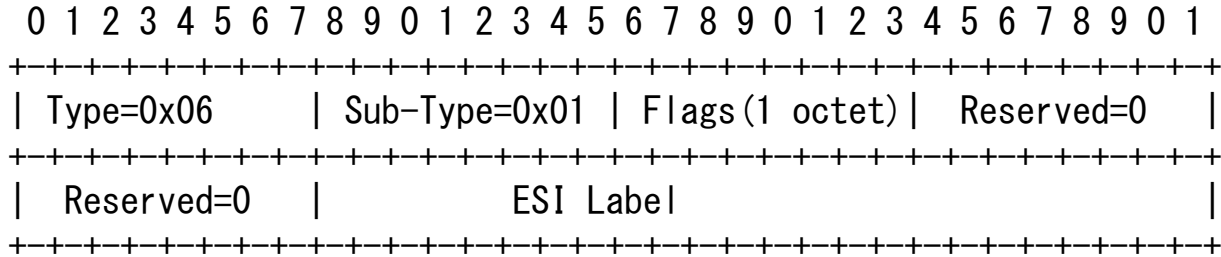
[EAD/ES]

- RD: EVI RD
- ESI: ESI 値
- Ethernet Tag ID: MAX-ET 固定 (不使用)
- MPLS Label: Zero (explicit-null) 固定 (不使用)

[EAD/EVI]

- RD: EVI RD
- ESI: ESI 値
- Ethernet Tag ID: Ethernet Tag ID 値 (VLAN Aware Bundle 以外なら 0)
- MPLS Label: local-PE で割当てたラベル (Aliasing for Load Balancing に利用)

ESI Label Ext-comm (Optional, Transitive)



■ EAD/ES Route と合わせて送信される

フォーマット

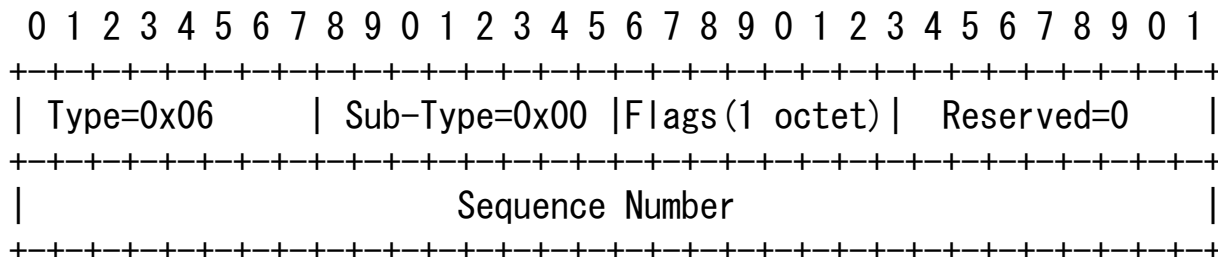
- Flags: $1 \ll 0 \dots$ redundancy mode が “Single-Active” ならばセット
- ESI Label: BUM の送受信に ingress replication を使う場合は受信用 Label、P2MP LSPs を使う場合は送信用 Label (Split Horizon に利用)

RD (8 octets)
Ethernet Segment Identifier (10 octets)
Ethernet Tag ID (4 octets)
MAC Address Length (1 octet)
MAC Address (6 octets)
IP Address Length (1 octet)
IP Address (0, 4, or 16 octets)
MPLS Label1 (3 octets)
MPLS Label2 (0 or 3 octets)

- MAC Mobility Ext-comm (Optional, Transitive)、Default Gateway Ext-comm (Optional, Transitive) を合わせて送信する

フォーマット

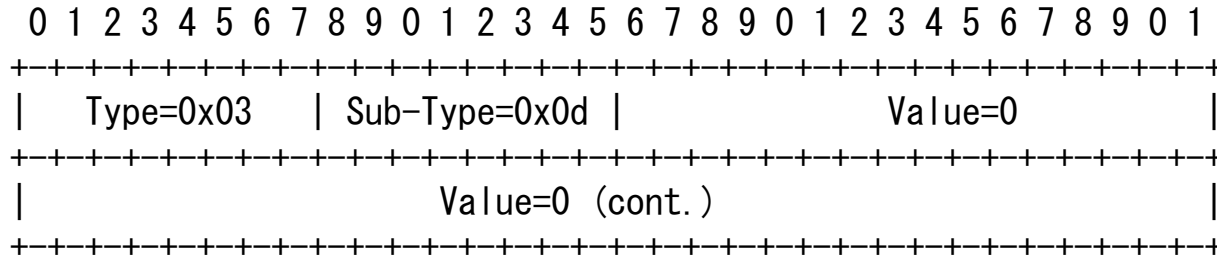
- RD: EVI RD (設定値)
- ESI: ESI値
- Ethernet Tag ID: Ethernet Tag ID 値 (VLAN Aware Bundle 以外なら 0)
- MAC Address Length, MAC Address: MAC 情報
- IP Address Length, IP Address: Optional。MAC と対応する IP address 情報が入る (IP Address は、remote-PE での ARP proxy や Inter-Subnet Routing に利用)
- MPLS Label 1: MAC に割り当てたラベル
- MPLS Label 2: Optional



■ MAC/IP Advertisement Route と合わせて送信される

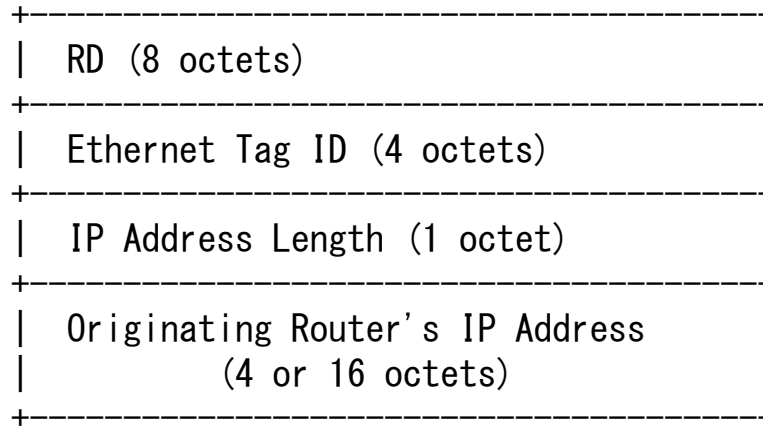
フォーマット

- Flags: $1 \ll 0 \dots$ Flag が set されている場合、MAC address は static であり、MAC move は発生しない
- Sequence Number: 同一の MAC address に関して複数回の update が発生する場合、MAC/IP Route の順番整合性を確保するために付与される



- MAC/IP Advertisement Route と合わせて送信される
- Opaque Type (RFC4360) の Ext-comm を利用する
- IP address フィールドが 0.0.0.0 or 0::0 である MAC/IP Route と共に広告される

Inclusive Multicast Ethernet Tag Route (IMET Route 0x3)



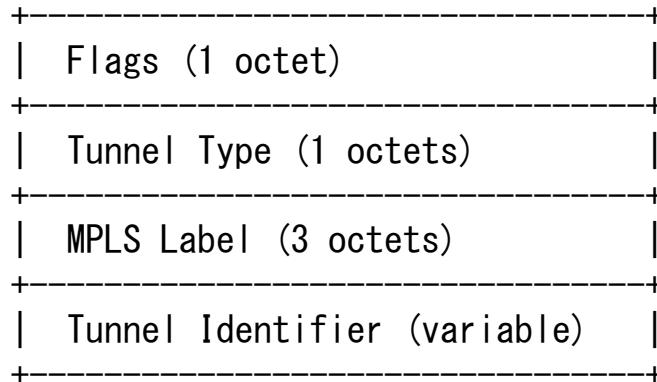
- IMET Route と合わせて PMSI Tunnel Attribute (Optional, Transitive) を送信する

フォーマット

- RD: EVI RD (設定値)
- Ethernet Tag ID: Ethernet Tag ID 値 (VLAN Aware Bundle 以外なら 0)
- IP Address Length, Originating Router's IP address: local PE の IP address (router-id) の length (bit) および address

※PMSI: Provider Multicast Service Interface

PMSI Tunnel Attribute (Optional, Transitive)



■ IMET Route と合わせて送信される

フォーマット

➤ Flags:

0 1 2 3 4 5 6 7

+++++

| reserved |L|

+++++

L: Leaf Information Required

※Leaf Information Label については、送信時には Zero、受信時には無視

- MPLS Label: high-order 20bit を BUM の送受信に用いる Label とする。BUM の送受信に ingress replication を使う場合、受信用 Label として利用

➤ Tunnel Type:

0 - No tunnel information present

1 - RSVP-TE P2MP LSP

2 - mLDP P2MP LSP

3 - PIM-SSM Tree

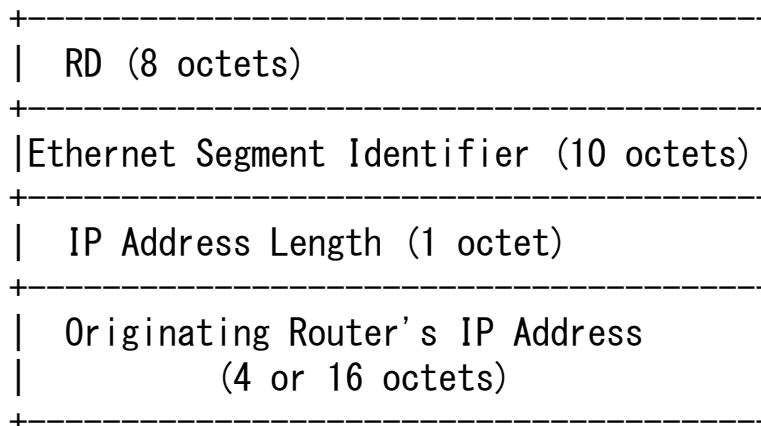
4 - PIM-SM Tree

5 - BIDIR-PIM Tree

6 - Ingress Replication

7 - mLDP MP2MP LSP

- Tunnel Identifier: Tunnel Type が Ingress Replication の場合、local PE の tunnel endpoint IP address (BGP session の local address) をセット

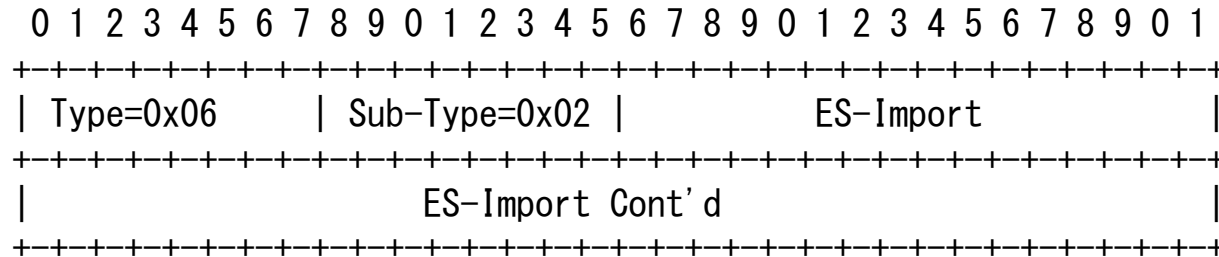


- ES Route と合わせて ES-Import Route-Target (Optional, Transitive) を送信する

フォーマット

- RD: EVI RD (設定値)
- Ethernet Segment Identifier: ESI値
- IP Address Length, Originating Router's IP address: local PE の IP address (router-id) の length (bit) および address

※PMSI: Provider Multicast Service Interface



■ ES Route と合わせて送信される

フォーマット

- ES-Import Value: ESI Type が 0, 1,2,3 ならば ESI の、Type を除いた部分の上位 6 Byte (MAC address に対応)