

The ARISTA logo is displayed in a bold, dark blue, sans-serif font in the top left corner. The background of the slide features a light blue grid with a network of nodes and connecting lines, some of which are highlighted in a darker blue.

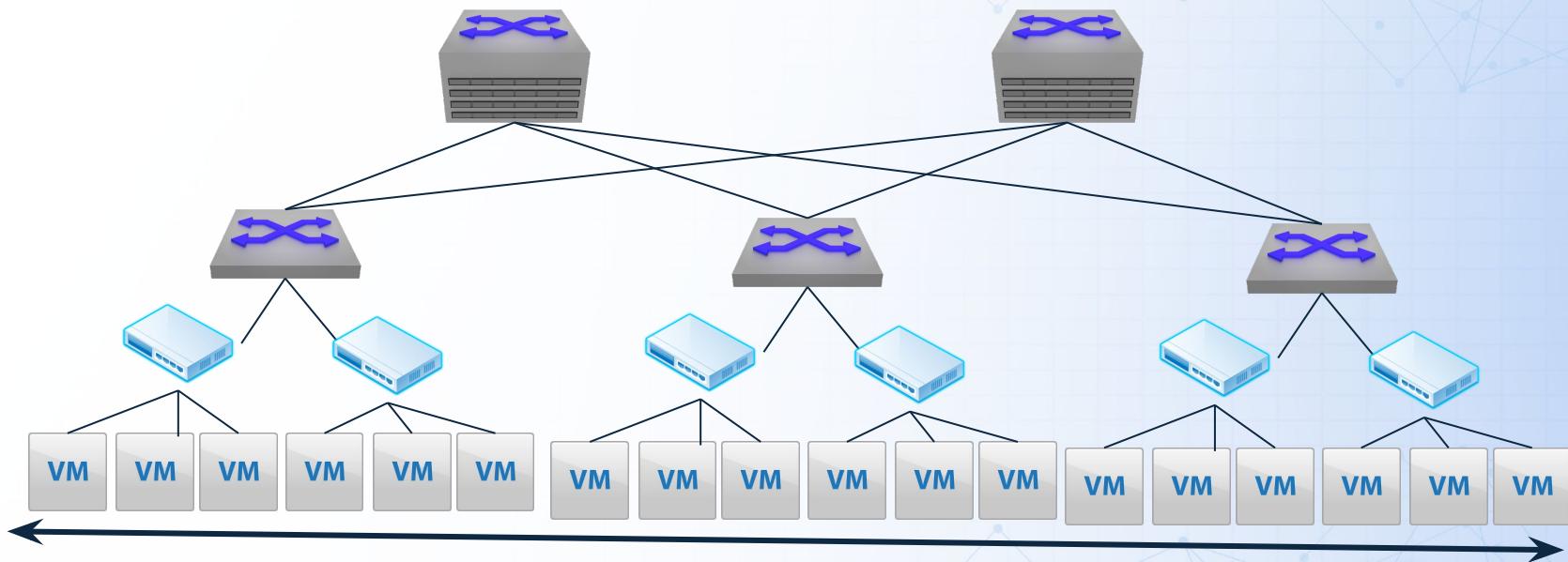
ARISTA

# 爆発するMACアドレスとの戦い

EVPN VESPA/仮想イーサネットセグメントとProxy ARPの活用

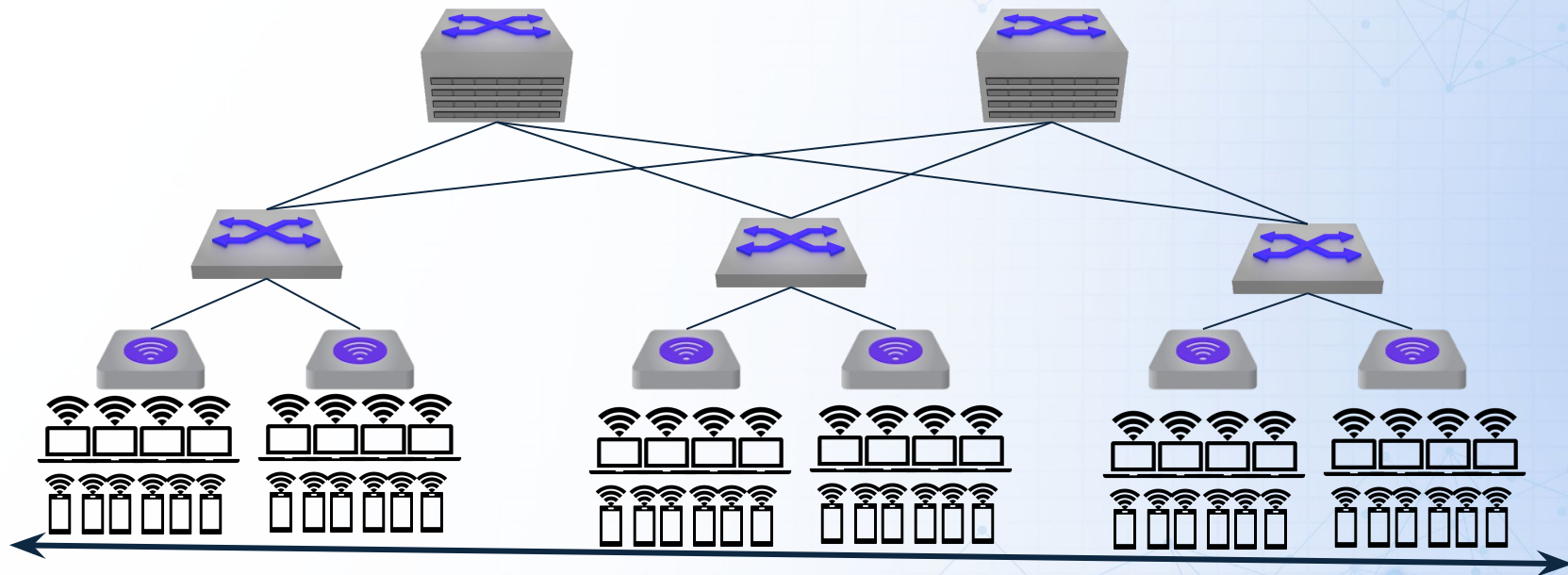
Shishio Tsuchiya  
[shtsuchi@arista.com](mailto:shtsuchi@arista.com)

# データセンターでの MACアドレスの増加要因



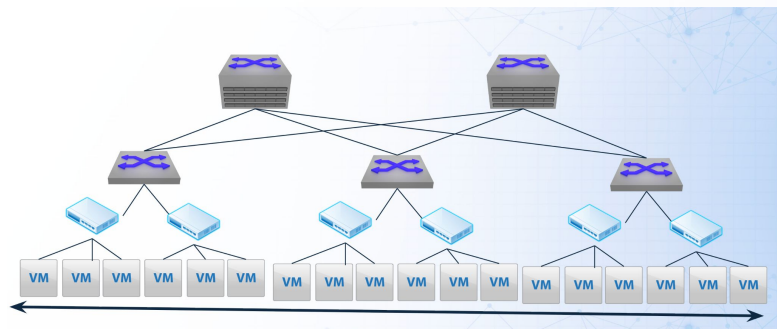
- ムーアの法則でCPU/メモリのスペックは上がっていく1台のサーバーの集積密度は上がっていく
- 一方仮想マシンのサービス要求は変わらない
- コンテナなど仮想化の技術も増えてきた

# 企業ネットワークでの MACアドレスの増加要因

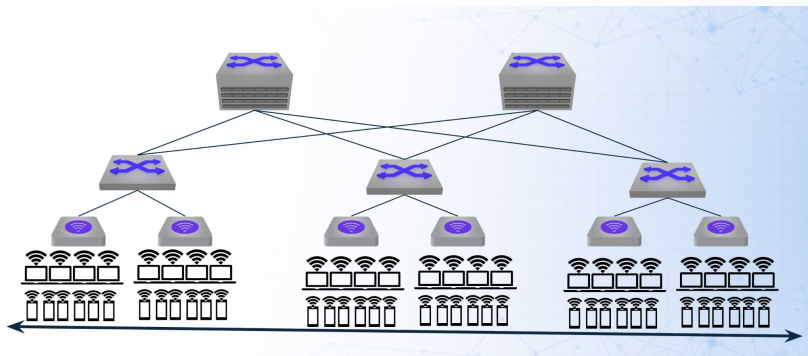


- BYOD端末やノートPCのモビリティの確保
- MACアドレスのランダム化も要因の1つに!?

# 解決しなければいけないスケール



- クライアント数(VMやWi-Fi端末)
  - 500,000
- アクセスデバイス(vSwitchやAP)
  - 30,000
- レイヤー2ドメイン
  - 制限では無く自分で決める事が出来る



# こんな事がありました

🦁 いうてDCと比べて少ない？

🦁 あ——そういう事



👤 JANOGの会場ネットワークもMACアドレスがやばい

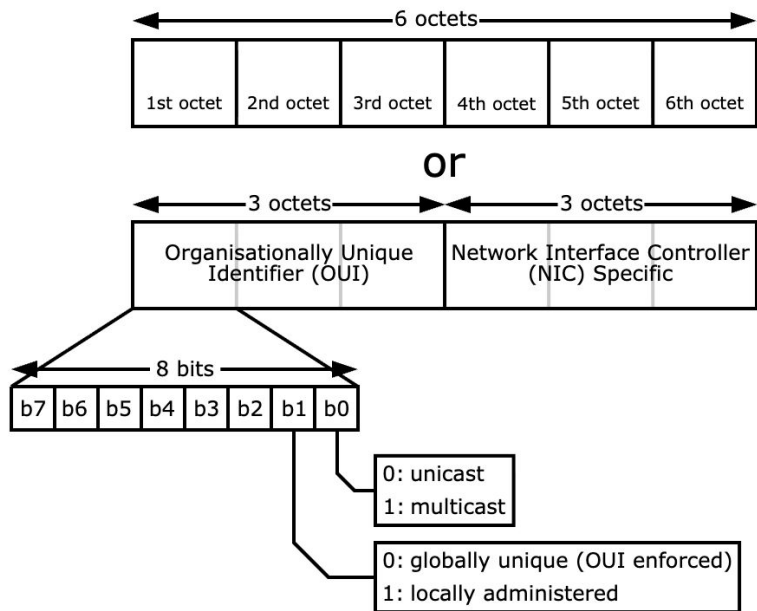
👤 MACアドレスのランダム化って知ってる？

# MACアドレスランダム化の背景



- カフェやショッピングモールなどWi-Fi環境においてMACアドレスは収集され下記の様な目的に使用される
  - 現在の位置情報
  - 滞在時間
  - 滞在時間
- プライバシーの観点から各OSベンダーが積極的に採用
- 接続要求時のみでは無く、接続中もランダム化
- 通常運用にも影響が!?

# どのようにランダム化されるか



- 48ビットのMACアドレスの7bit目はグローバルにユニークか、ローカルに管理されているかを示す
- 1が立っているとローカルに管理されたアドレスであることを示す。

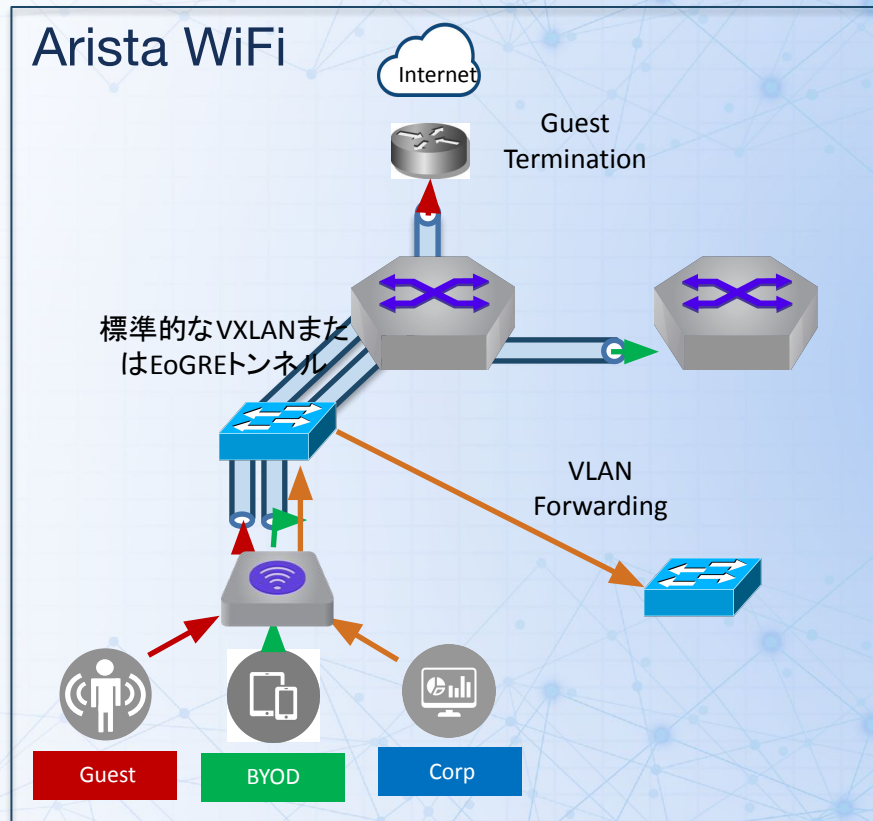
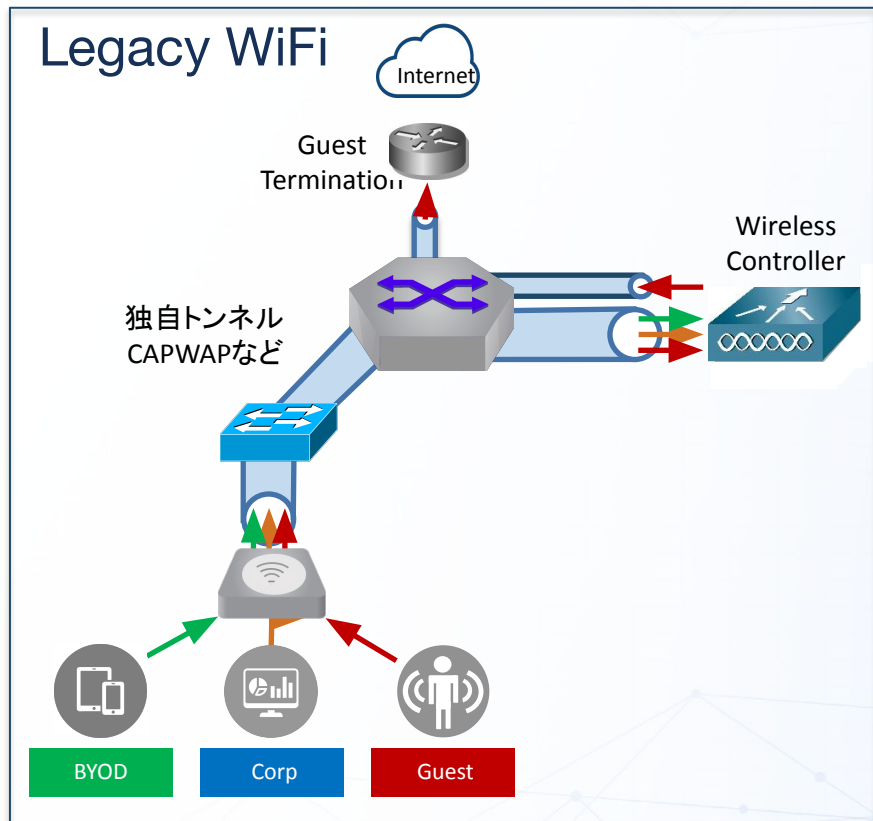
# OSベンダー毎のMACアドレスランダム化サポート状況

	OS毎の挙動		
挙動	iOS14	Android 11	Windows 10
デフォルト	ON/ユーザはデフォルトの動作を変更出来ない	ON/ユーザはデフォルトの動作を変更出来ない	OFF/ユーザーはデフォルト設定変更可能
初めてのSSIDへの接続	SSIDに初回接続時、接続用に新しいランダムMACを生成	SSIDに初回接続時、接続用に新しいランダムMACを生成	OFF:ハードウェアMACアドレスで接続 ON:SSIDに初回接続時、接続用に新しいランダムMACを生成
既存のSSIDへの接続	同じSSIDから切断し再接続する際、接続には同じランダムMACアドレスを使用	同じSSIDから切断し再接続する際、接続には同じランダムMACアドレスを使用	ONの場合 同じSSIDから切断し再接続する際、接続には同じランダムMACアドレスを使用
特定のSSIDに対してのランダムMACの無効化	ハードウェアMACで接続	ハードウェアMACで接続	ハードウェアMACで接続
すべてのSSIDでのMACアドレスランダム化の無効	N/A	N/A	ハードウェアMACアドレスを使用
SSID プロファイル削除と再接続	同じランダムMACアドレスを使用	同じランダムMACアドレスを使用	OFF/ハードウェアMACアドレスを使用 ON/新たに生成されたランダムMACが使用される

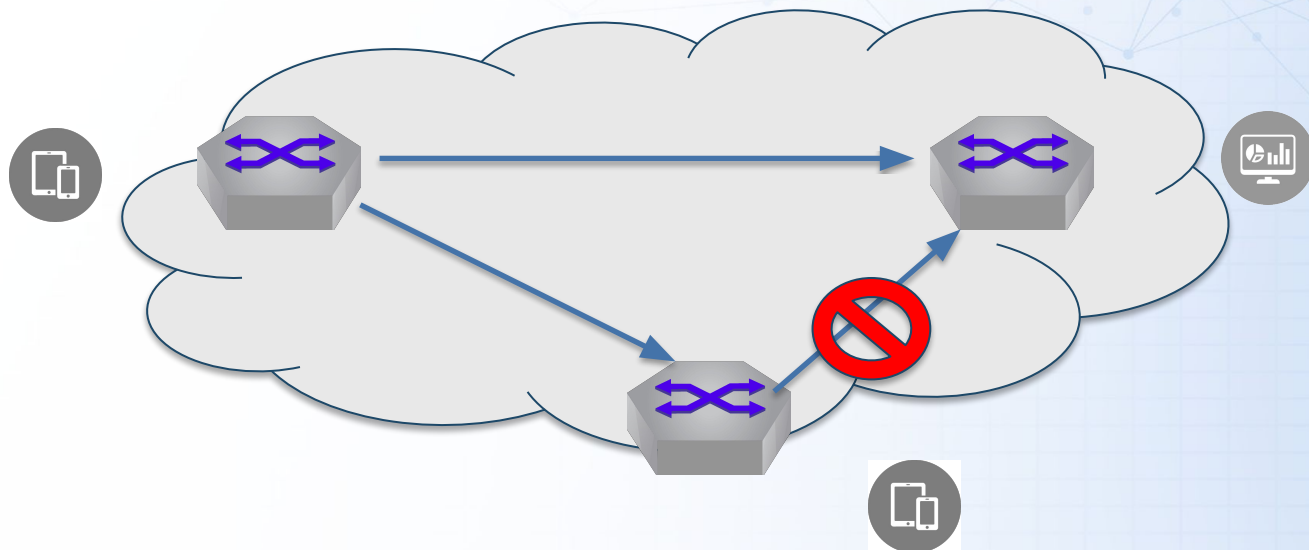
# インパクト

- iOS14にアップデート後に最初にSSIDにログインするとゲストWiFiなどの情報はリセットされる
- 更に最新のOS(Windows11/iOS18/macOS15では接続時間やネットワークの安全性など特定条件において同じSSIDでもMACアドレスのランダム化が行われる(Rotating Wi-Fi Address)
- 下記の様なネットワークの機能に影響がある
  - MACアドレス認証やPortセキュリティやMAC ACLなどのセキュリティ
  - ユーザ情報分析
  - DHCP IP アドレスの消費とサーバー負荷
  - スケーラビリティなどにインパクトが起こり得る

# コントローラーレス WiFiアーキテクチャ

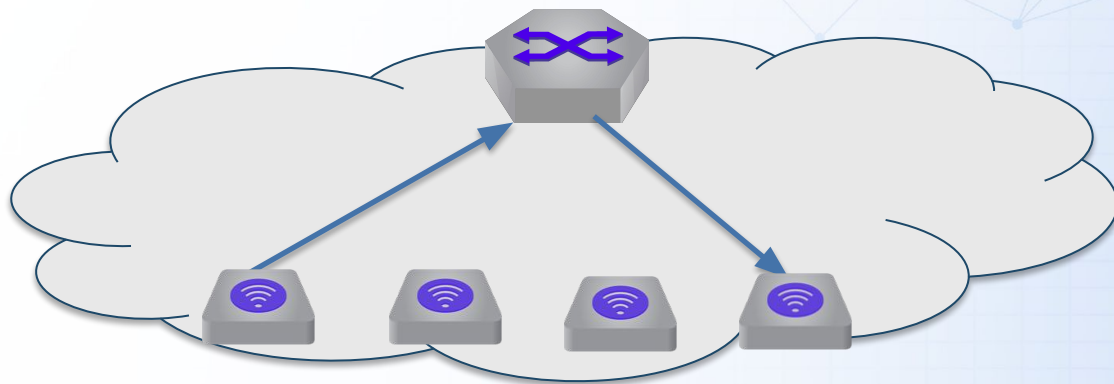


# 通常のvxlanの場合



- すべてのvxlanのエンドポイントはフルメッシュでつながっている
- 他VTEPからフラッドされたトラフィックを別なVTEPにフラッドする事は無い
- これでルーティンググループを防ぐ

# VXLANをHub & Spokeとして動作させる



- APとの組み合わせでは従来のVXLANとは違う動作が求められる
- “vxlan bridging vtep-to-vtep” をVTEPに設定する

# AP側の設定

<https://arista.my.site.com/AristaCommunity/s/article/VXLAN-Configuration>

- SSIDとトンネルが紐づく
- StaticにVTEPを設定

The screenshot shows the 'Tunnel Interface' configuration page in Arista CloudVision. The left sidebar contains navigation options: CloudVision, ダッシュボード, モニター, 構成, トラブルシューティング, フロアプラン, レポート, システム, and サービス. The main content area has tabs for WiFi, SSID, RADIUS, Tunnel Interface (selected), Role Profile, Wireless Settings, and Device Settings. The 'Tunnel Interface Name' field is populated with 'トンネル・インターフェイス名を入力'. A dropdown menu for 'Tunnel Type' is open, showing options: L2 Tunnel, EoGRE, EoGRE over IPsec, VXLAN, and VPN Tunnel (IPsec using VPN). Below the dropdown are radio buttons for 'Primary' and 'Secondary'. The 'Remote Endpoint (IP/Host Name)' field contains 'IPアドレス/ホスト名を入力'. The 'Local Endpoint VLAN' is set to 0. The 'MSS Clamp' checkbox is checked. The 'Tunnel MTU Discovery' is set to 'Manual'.

The screenshot shows the 'SSID' configuration page in Arista CloudVision. The left sidebar is the same as the previous screenshot. The main content area has tabs for WiFi (selected), SSID (selected), RADIUS, Tunnel Interface, Role Profile, and Wireless Settings. The 'SSID Name' field is empty. The 'VLAN ID' is set to 0. The 'Mode' is set to 'Bridge' with radio buttons for Bridge, NAT, L2 Tunnel, and VPN Tunnel. The 'Layer 2 Traffic Forwarding and Filtering' checkbox is unchecked.

# VXLANコントロールプレーンの選択

- VXLANコントロールプレーンはMAC学習とパケットフラッディングに使用される
  - リモートVTEPの背後にいるホストを発見するメカニズム
  - どのようにVTEPとVNIのメンバーを発見するのか?
  - レイヤー2セグメント(VNI)内でブロードキャスト/マルチキャストを転送する為のメカニズム

## IPマルチキャストコントロールプレーン

- VTEPはVNIの所属するマルチキャストグループにjoin
- VNI内のUnknownユニキャストはVTEPにマルチキャストで転送
  - サードパーティVTEPをサポート
- 転送と学習はIPマルチキャストを要求-限られた展開例

## HeadEnd Replication (HER)

- BUMトラフィックはVNIの中のリモートVTEPに複製される
  - ingress VTEPで複製される
- サードパーティのVTEPをサポート
- MAC学習は転送され学習されるしかし、マルチキャストは必要が無い

## HER with CloudVision eXchange (CVX)

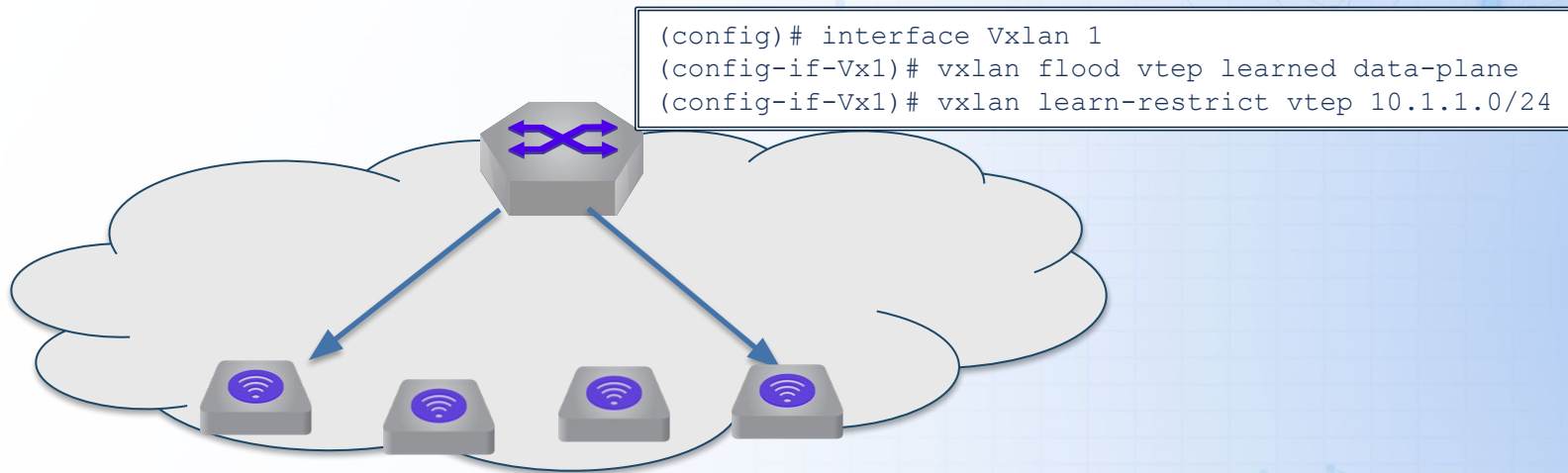
- ローカルに学習されたMACとVNIバインドはCVXに公開される
- CVXはダイナミックにリモートVTEPにステートを分配する
  - サードパーティVTEPをサポート
- MAC学習とfloodlistの設定は自動にされる
- 信頼性の為にHAクラスタサポート

## EVPNモデル Model

- VTEP間のローカルに学習されたMACとIPの紐付けの為にBGPを使用
- ブロードキャストトラフィックはIPマルチキャストやHERによりハンドルされる
- 設定されたBGPIによりMACアドレスとVNIはダイナミックに分配される
- サードパーティVTEPをサポート

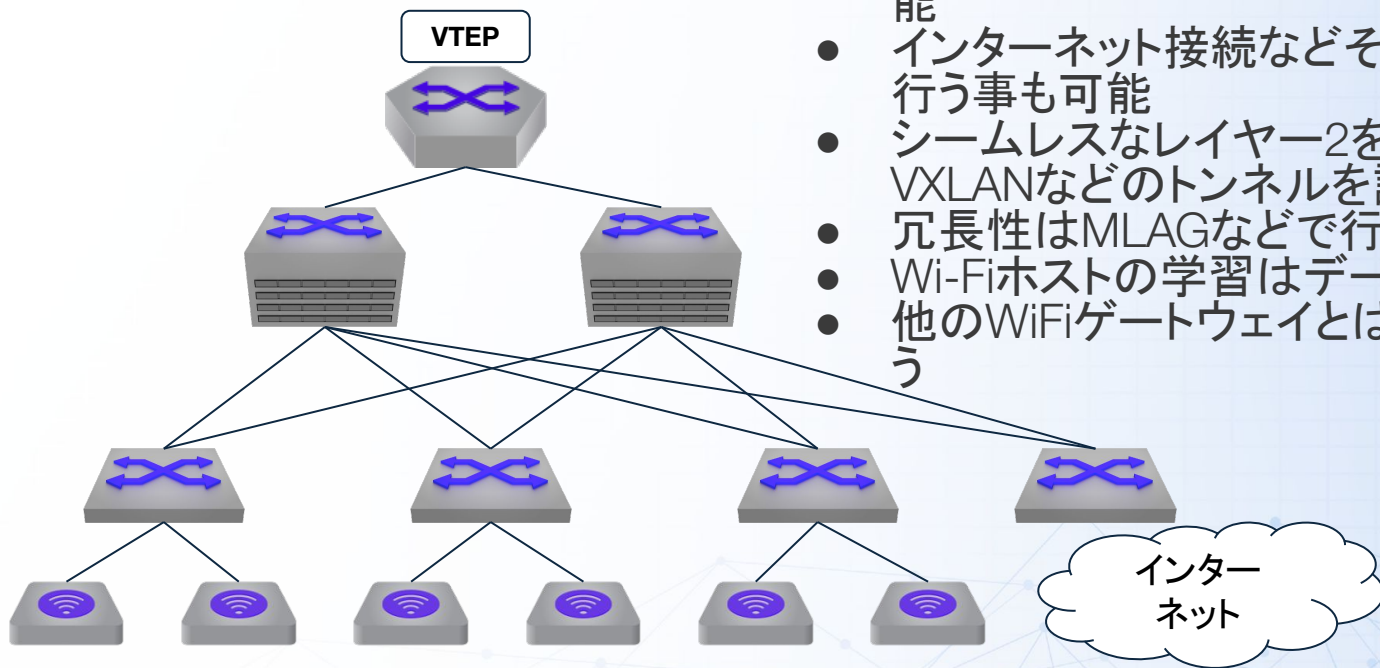
# 従来の方法の Flood-listではAPは管理出来ない

<https://www.arista.com/en/support/toi/eos-4-23-0f/14368-vxlan-auto-flood-list-construction>



- 実際にパケットを送ってきた(データプレーン)VTEPをflood-listに登録する
- アドレス範囲を制限する事も可能

# 現在のコントローラーレス WiFiアーキテクチャー



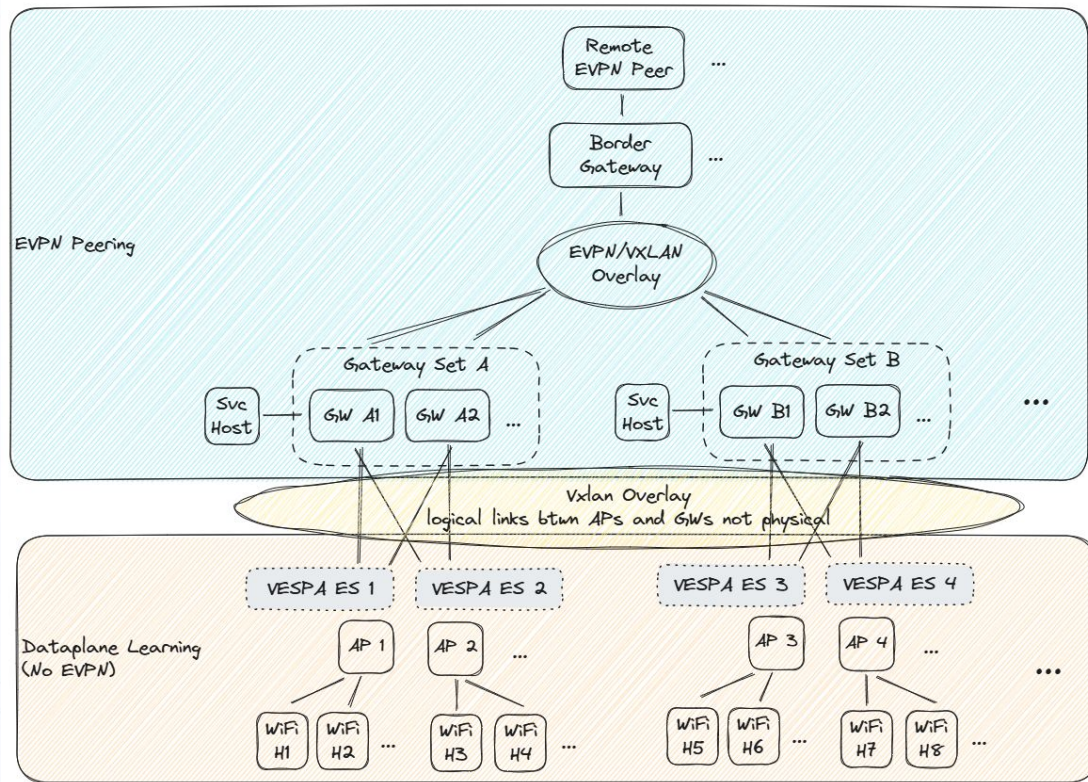
- 各SSIDでそれぞれのデータプレーンを選択可能
- インターネット接続などそのままルーティングで行う事も可能
- シームレスなレイヤー2を構築する為にはVXLANなどのトンネルを設定する
- 冗長性はMLAGなどで行う
- Wi-Fiホストの学習はデータプレーンで行う
- 他のWiFiゲートウェイとはEVPNでやり取りを行う

# 既存の構成の問題点

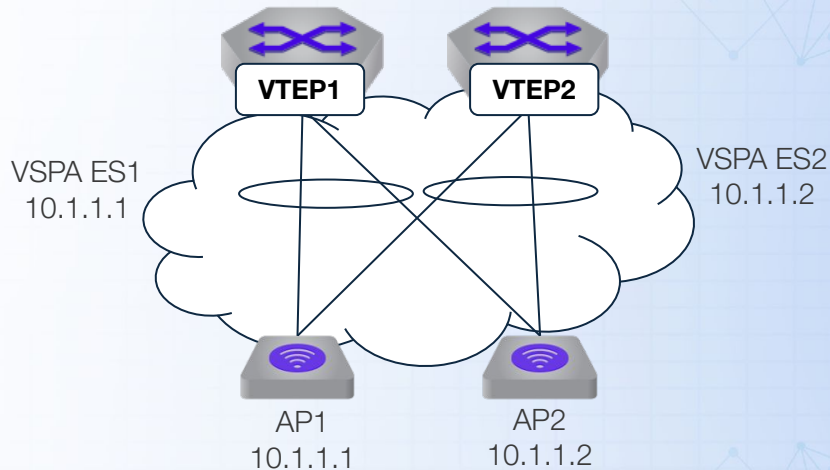
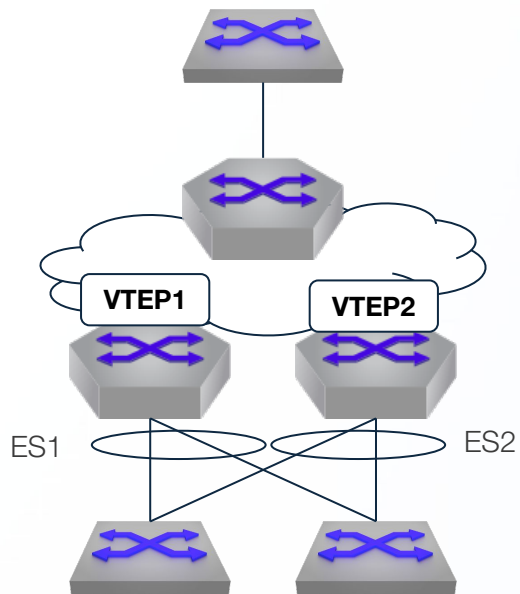
- 冗長性の欠如
  - WiFiゲートウェイが単一障害ポイントとなり得る
  - プライマリーとセカンダリーの設定
    - ICMPを定期的にポーリングし、プライマリーが切れた時はセカンダリーに切り替わる。
- 解決策
  - MLAG
    - 独自プロトコルとなる
    - 地理的制限がある
    - 2台のみの冗長となる
  - EVPN All Active Multihoming
    - EVPN マルチホーミングのイーサネットセグメントはStaticに設定される。APの数が多いほど、Static設定の量も増える
    - APはゲートウェイに直接接続されておらず、アンダーレイネットワークを経由したVXLANTンネルを介して接続されている

# VESPA(Virtual Ethernet Segment with Proxy Arp)

- VESPAは、本質的に通常のEVPNマルチホーミングの拡張機能
- ローカルインターフェースに加えてトンネルにも対応し、イーサネットセグメントの動的な構築を可能に
- MACアドレスのスケールビリティはAPがL2 Proxyとして動作する事により実現する



# EVPN VESPA(Virtual Ethernet Segment with Proxy Arp)



- 通常のEVPN Ethernet Segment
  - セグメントは物理インターフェースに紐づく
- EVPN VESPA(Virtual Ethernet Segment with Proxy Arp)
  - セグメントはトンネルに紐づく(VXLAN VTEP)

# L2 Proxy(MRO:MAC Rewrite Offload)を有効にする

**Network Profiles** ▾ **Tunnel**

*Changes to this tunnel interface will affect all SSIDs and LAN Ports that use this tunnel interface*

← ASU Demo

Tunnel Interface Name \*  Tunnel Type

---

L2 Proxy

**Primary** **Secondary** *Secondary tunnel settings are disabled because L2 Proxy is enabled.*

Remote Endpoint (Enter IP Address/Hostname)

Local Endpoint VLAN \*  [0 - 4094] VXLAN VNI Offset \*  [0 - 16777215]

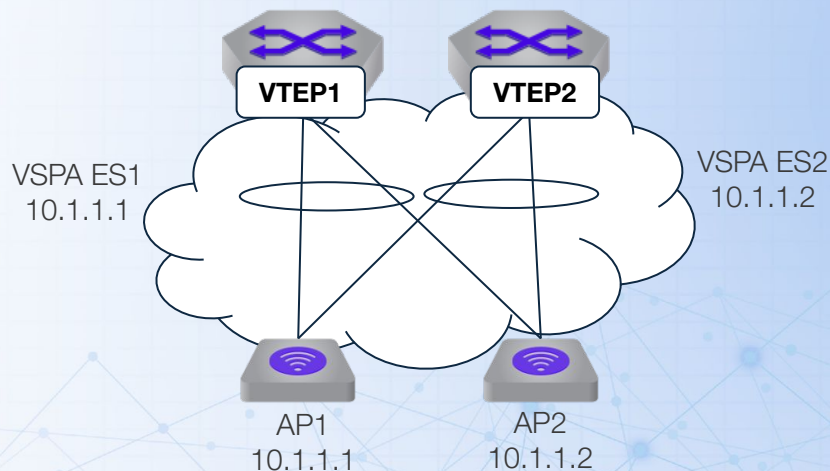
---

**Tunnel MTU**

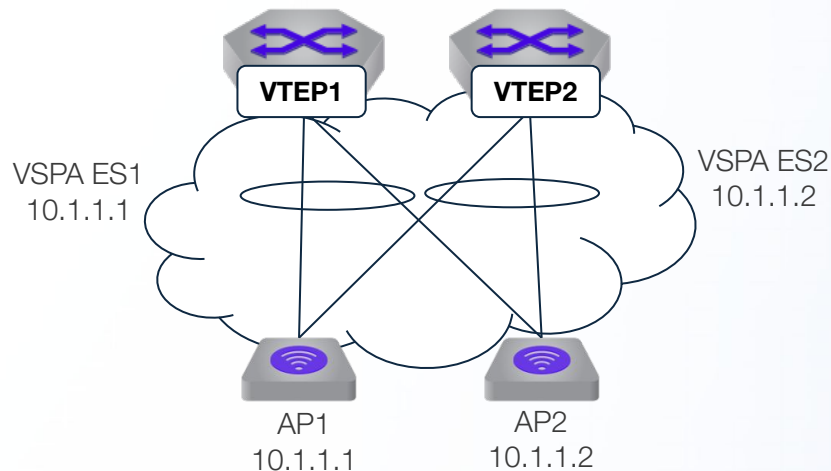
Tunnel MTU Discovery  Auto  Manual

Enforce Fragmentation

- Network Profiles → VxLAN Tunnel TypeでL2 Proxyにチェックする



# VESPA GW Set IDの設定

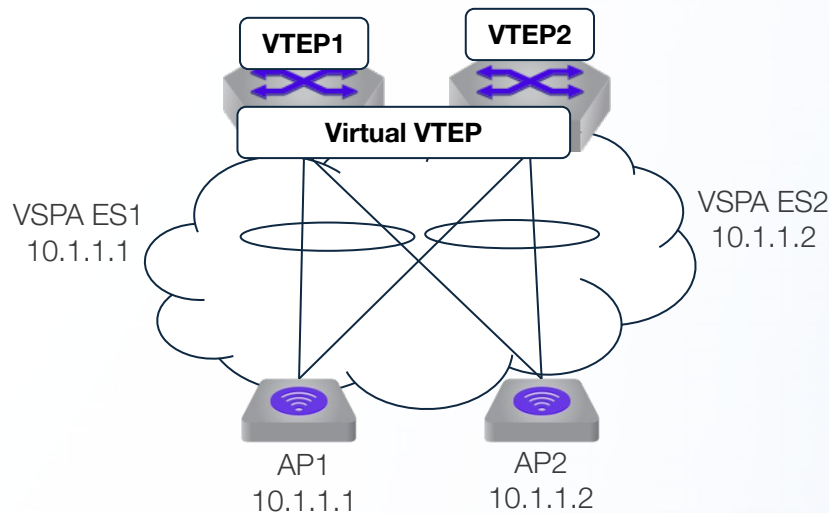


- VESPAゲートウェイにはSet IDを設定する必要がある。
- 接続されているAPに関して、互いに同じイーサネットセグメントに属する必要があるすべてのVESPAゲートウェイには、同じセットIDを設定する必要がある

```
router bgp 65000
  address-family evpn
    evpn ethernet-segment dynamic-tunnel mh-set-id
```

10

# VESPA GW での Virtual VTEP の設定

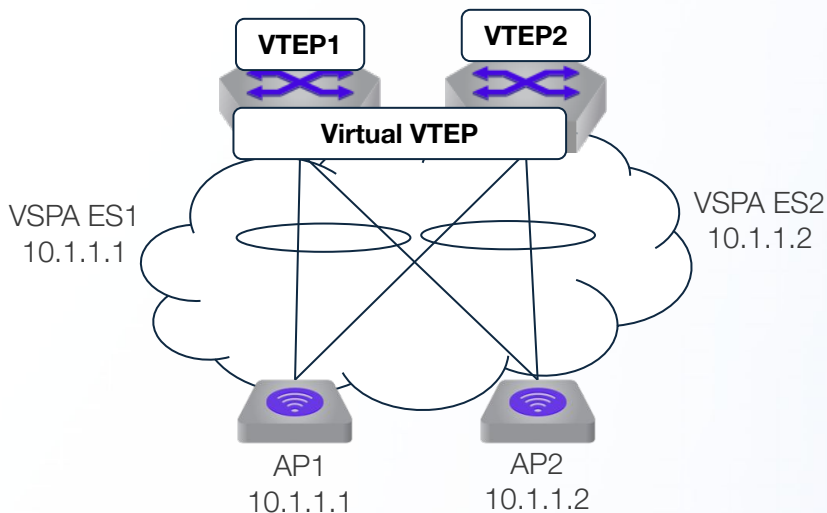


```
VESPA-GW1 (config)#interface loopback 1
VESPA-GW1 (config-if-Lo1)# description VXLAN VTEP
VESPA-GW1 (config-if-Lo1)# ip address 10.255.255.1/32
VESPA-GW1 (config-if-Lo1)#interface loopback 2
VESPA-GW1 (config-if-Lo2)# description VXLAN VIRTUAL VTEP
VESPA-GW1 (config-if-Lo2)# ip address 10.255.255.250/32
VESPA-GW1 (config-if-Lo2)#interface vxlan 1
VESPA-GW1 (config-if-Vx1)# vxlan source-interface Loopback1
VESPA-GW1 (config-if-Vx1)# vxlan virtual-vtep local-interface
Loopback2
```

```
VESPA-GW2 (config)#interface loopback 1
VESPA-GW2 (config-if-Lo1)# description VXLAN VTEP
VESPA-GW2 (config-if-Lo1)# ip address 10.255.255.2/32
VESPA-GW2 (config-if-Lo1)#interface loopback 2
VESPA-GW2 (config-if-Lo2)# description VXLAN VIRTUAL VTEP
VESPA-GW2 (config-if-Lo2)# ip address 10.255.255.250/32
VESPA-GW2 (config-if-Lo2)#interface vxlan 1
VESPA-GW2 (config-if-Vx1)# vxlan source-interface Loopback1
VESPA-GW2 (config-if-Vx1)# vxlan virtual-vtep local-interface
Loopback2
```

- GWでは2つのVTEP IPを設定する必要がある。
- VXLANの送信元インターフェースとして使用するアドレス(loopback 1)
- virtual-vtep interface アクセスポイントがトンネルエンドポイントとして設定するVESPA Anycast GWアドレス(loopback 2)

# SVIの設定

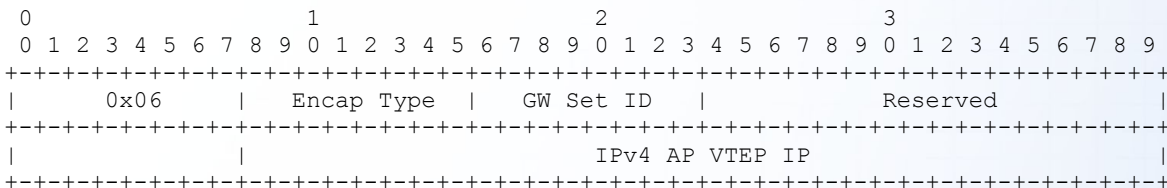


```
interface Vlan3525
  description USER_VRF_A
  mtu 9090
  no autostate
  vrf USER-3520
  arp aging timeout 3600
  arp refresh retries unicast 5 broadcast
  0
  ip helper-address 192.168.21.2
  ip address virtual 10.1.1.250/16
```

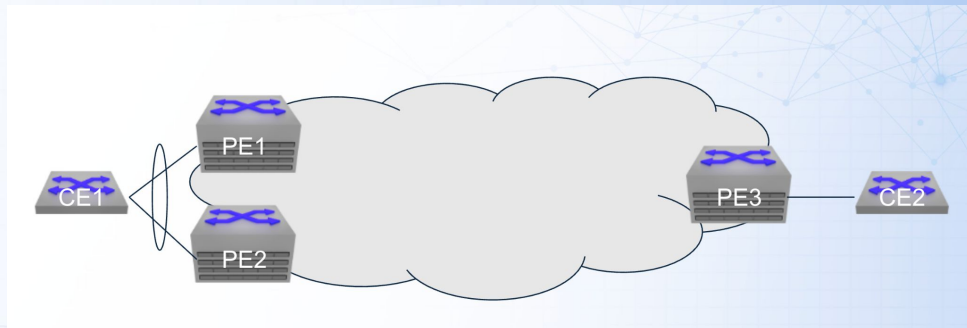
- アクセスポイントは、そこに接続されているすべてのワイヤレスクライアントに対して、レイヤ2-MACプロキシとして機能
- VESPA GWはこれらすべてのワイヤレスクライアントのARPエントリを維持する必要がある
- ワイヤレスクライアントのARPエントリは、APが能動的にGARPパケットを送信することに基づいて作成
- GWは、エージェントベースの定期的なユニキャストリフレッシュを使用して、これらのARPエントリを維持する

# ESIタイプフォーマット

- RFC7432では5つのESIフォーマットが定義
- VESPAではType6を使用



Type	名前	フォーマット	モード
0	User	手動設定	手動設定
1	LACP	LACP System MAC + port key + 00	自動生成
2	Bridged	Root bridge MAC + priority key + 00	自動生成
3	MAC	System MAC + Local Discriminator	手動設定/自動生成
4	Router-ID	RouterID + Local Discriminator + 00	手動設定/自動生成
5	ASN	AS番号 + Local Discriminator + 00	手動設定/自動生成



# まとめ

- 大規模のレイヤー2ネットワークをサポートするEVPNを使ったVESPA(Virtual Ethernet Segment with Proxy Arp)を紹介
- 大規模なMACアドレス/不安定なMACアドレスに対して安定したネットワークを提供

# 聞きたい事

- MACアドレス爆発してます？
- それはWiFiネットワーク？データセンター？
- MACアドレスのランダム化で困った事ある？
- L2 Proxy議論された事もありますが、やってたりする人います？  
OpenFlow的ななにかで？
- 今どうですかね？

# 参考

- Mega Scale Campus Mobility Powered by Arista VESPA
  - <https://www.youtube.com/watch?v=FUryTCSmUU>
- Arista VESPA User Guide
  - <https://www.arista.com/en/support/toi/eos-4-35-2f/23763-arista-vespa-user-guide>
- RFC7432 BGP MPLS-Based Ethernet VPN 5. Ethernet Segment
  - <https://datatracker.ietf.org/doc/html/rfc7432#section-5>
- The Scalable Address Resolution Protocol (SARP) for Large Data Centers
  - <https://datatracker.ietf.org/doc/html/rfc7586>
- Address Resolution Problems in Large Data Center Networks
  - <https://datatracker.ietf.org/doc/html/rfc6820>
- MAC Randomization: Behavior and Impact
  - <https://www.arista.com/assets/data/pdf/Whitepapers/MAC-Randomization-Behavior-and-Impact.pdf>

ARISTA

Thank You

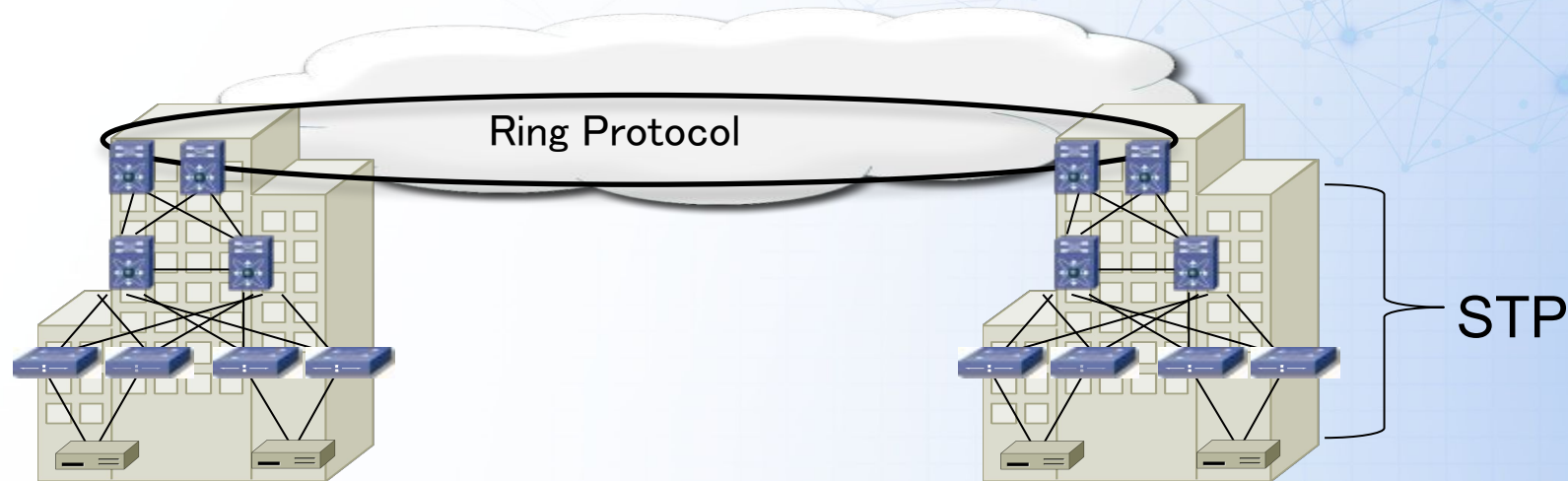
[www.arista.com](http://www.arista.com)



# 技術の成熟と 次代への継承

MACアドレス爆発解決策の  
光と影

# 2011年ごろのデータセンターネットワーク



- フラットなレイヤー2のネットワークデザイン
- データセンター内ではスパンニングツリーが使用
- ベンダー独自のリングプロトコル (またはG.8032)がデータセンター間で使用される
- VMライブマイグレーションは必要 (GARPによって移動を通知)
- VLANがユーザ毎にアサインされる

# 問題点

- データセンター間/データセンター内
- VLANスケーラビリティ > 4K
  - MACアドレステーブルのスケーラビリティ
  - VMライブマイグレーションや簡単に使うためのブロードキャストドメインの拡張
  - East / Westトラフィック帯域の増加
  - 高速収束
  - 自動化/オーケストレーション
- データセンター間
  - 要求に応じた帯域増強
  - ベンダーロックイン技術からの解放
  - 柔軟性のあるトポロジーデザイン
  - トラフィックエンジニアリング
  - BUM(Broadcast/Unknown unicast/Multicast) トラフィックの最適化
- ゲートウェイ
  - ARP/NDPスケーラビリティ
  - IETF ARMD(Address Resolution for Massive numbers of hosts in the Data center) Groupは一つの informational RFCを発行 RFC6820 Address Resolution Problems in Large Data Center Networks

# 問題に関するソリューション

- 仮想オーバーレイプロトコル

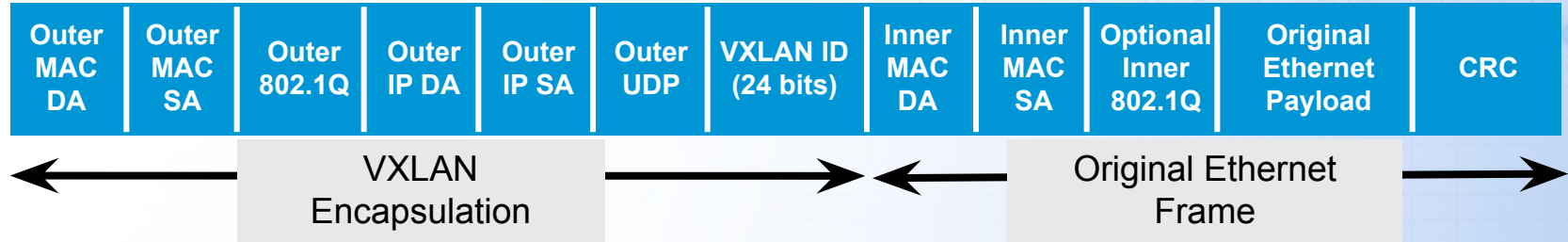
- VXLAN: A Framework for Overlaying Virtualized Layer 2 Networks over Layer 3 Networks
  - draft-mahalingam-dutt-dcops-vxlanは2011年8月に experimental として発行
  - vmware/cisco/arista/broadcom/citrix/Redhat
  - RFC7348は2014年8月に Informationalとして発行
- NVGRE: Network Virtualization using Generic Routing Encapsulation
  - draft-sridharan-virtualization-nvgreは2011年9月に発行
  - microsoft/arista/Intel/Dell/HPE/Broadcom/Emulux
  - RFC7637は2015年にマイクロソフトによって informationalとして発行
- A Stateless Transport Tunneling Protocol for Network Virtualization(STT)
  - draft-davie-sttは2012年2月 nformationとして発行
  - Nicira Networks
  - 2016年にexpire
- Overlay Transport Virtualization(OTV)
  - draft-hasmit-otvは2010年4月にスタンダードとして発行
  - Cisco
  - 2013年にexpire

# 問題に関するソリューション

## ● ネットワークデザイン

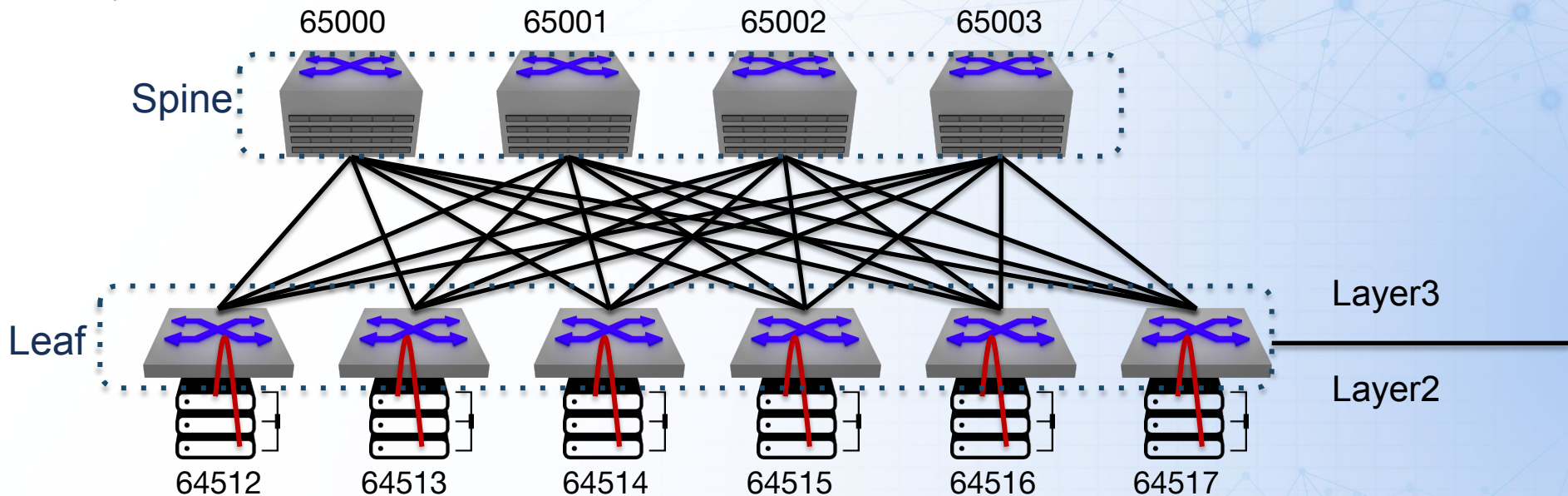
- Use of BGP for routing in large-scale data centers
  - Microsoft Petr Lapukhov がNANOG55 June 3-6, 2012でプレゼン
  - またdraft-lapukhov-bgp-routing-large-dcを2012年7月に発行 Ariff Premji(Arista)が共同著者
  - RFC7938 は2012年8月にInformationalとして発行
- Transparent Interconnection of Lots of Links (trill)
  - draft-perlman-trill-rbridge-protocolは2006年に発行
  - Trill WG はRbridges(RFC6325)を次世代 IEEE802.1Dプロトコルとして定義 July 2011
  - RFCの著者は Intel/Huawei/Cisco/Brocade

# Virtual eXtensible LAN (VXLAN)



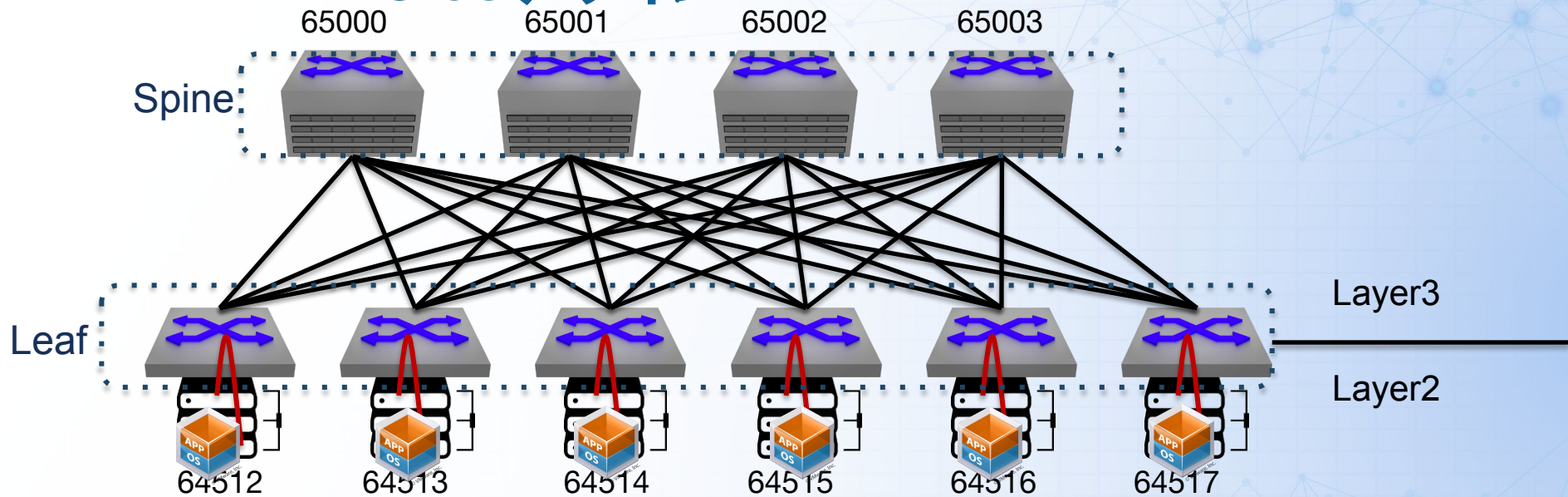
- レイヤー2フレームをIPでカプセルプロトコルを定義
- レイヤー3インフラストラクチャー上でオーバーレイのネットワークを作成する為に使用
- レイヤー2の接続性をユーザに提供

# 大規模データセンタールーティングでの BGPの使用



- スケールアウトするClosデザイン
- 安定した標準的なBGPプロトコルをToR/Leafスイッチに使用
- 安定性にフォーカスし、VMモビリティはラック内に留める

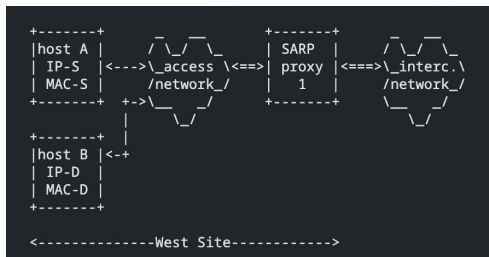
# VXLAN + IP Closデザイン



- データセンターデザインで完全なソリューション
- 標準で安定したBGPプロトコル/広く展開されたVXLANフォーマットを使用
- 特に特別な要求が中間ノードには存在しない

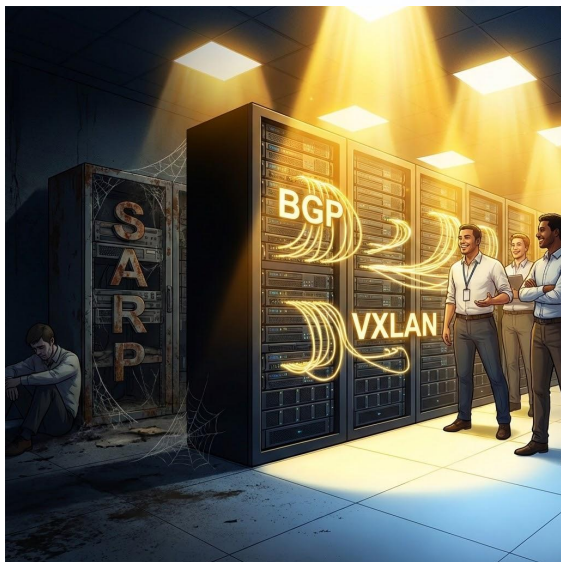
# 2011年のMACアドレス爆発

- BGPのデータセンターネットワークでの活用デザイン+プロトコル標準では無いVXLAN(業界標準)のカプセル化での解決
- つまり標準化IETFは課題定義のみを行いなにも出来なかった
- しかし1つのExperimental RFCを発行していた
- The Scalable Address Resolution Protocol (SARP) for Large Data Centers
  - <https://datatracker.ietf.org/doc/html/rfc7586>



# 2026年のMACアドレス爆発

VXLAN技術の成熟と SARPの継承



- 2011-12年に脚光を浴びたVXLANとBGP
- BGPの拡張EVPNをコントロールプレーンに引っ提げて
- あの時光を浴びなかったL2 Proxyの技術にも注目そして実装