



Cisco ACI 転送

この章は、次の内容で構成されています。

- [ACI ファブリックは現代のデータセンタートラフィックフローを最適化する \(1 ページ\)](#)
- [ACI で VXLAN \(2 ページ\)](#)
- [サブネット間のテナントトラフィックの転送を促進するレイヤ 3 VNID \(4 ページ\)](#)
- [スパンニングツリープロトコル BPDU の送信 \(6 ページ\)](#)

ACI ファブリックは現代のデータセンタートラフィックフローを最適化する

Cisco ACI アーキテクチャは、従来のデータセンター設計から来る制限を解放して、最新のデータセンターで増大する East-West トラフィックの需要に対応します。

今日のアプリケーション設計は、データセンターのアクセスレイヤを通る、サーバ間の East-West トラフィックを増大させています。このシフトを促進しているアプリケーションには、Hadoop のようなビッグデータの分散処理の設計、VMware vMotion のようなライブの仮想マシンまたはワークロードの移行、サーバのクラスタリング、および多層アプリケーションなどが含まれます。

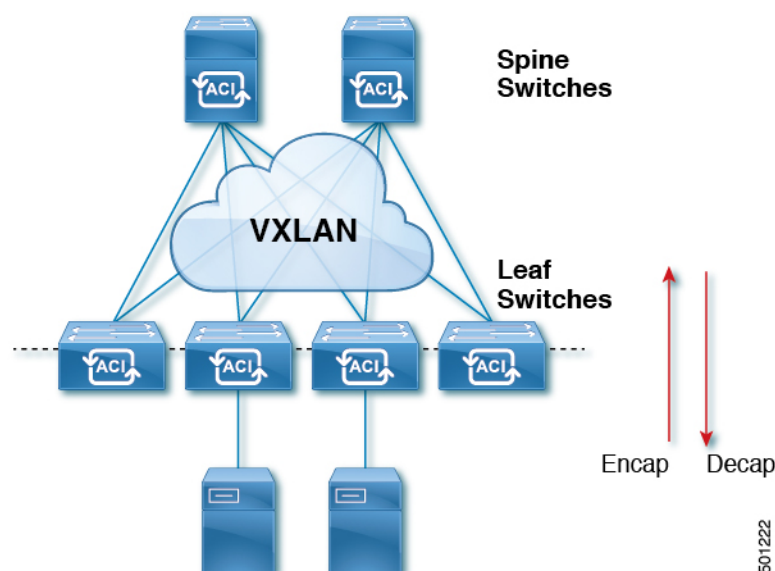
North-South トラフィックは、コア、集約、およびアクセスレイヤ、またはコラプストコアとアクセスレイヤが重要となる、従来型のデータセンター設計を推進します。クライアントデータは WAN またはインターネットで受信され、サーバの処理を受けた後、データセンターを出ます。このような方式のため、WAN またはインターネットの帯域幅の制限により、データセンターのハードウェアは過剰設備になりがちです。ただし、スパンニングツリープロトコルが、ループをブロックするために要求されます。これは、ブロックされたリンクにより利用可能な帯域幅を制限し、トラフィックが準最適なパスを通るように強制する可能性があります。

従来のデータセンター設計においては、IEEE 802.1Q VLAN がレイヤ 2 境界の論理セグメンテーションまたはブロードキャストドメインを提供します。ただし、ネットワークリンクの VLAN の使用は効率的ではありません。データセンターネットワークでデバイスの配置要件は柔軟性に欠け、VLAN の最大値である 4094 の VLAN が制限となり得ます。IT 部門と

クラウドプロバイダが大規模なマルチテナントデータセンターを構築するようになるにつれ、VLAN の制限は問題となりつつあります。

スパインリーフアーキテクチャは、これらの制限に対処します。ACI ファブリックは、外界からは、ブリッジングとルーティングが可能な単一のスイッチに見えます。レイヤ3 のルーティングをアクセスレイヤに移動すると、最新のアプリケーションが必要としている、レイヤ2 の到達可能性が制限されます。仮想マシンワークロードモビリティや一部のクラスタリングのソフトウェアのようなアプリケーションは、送信元と宛先のサーバ間がレイヤ2 で隣接していることを必要とします。アクセスレイヤでルーティングを行えば、トランクダウンされた同じVLANの同じアクセススイッチに接続したサーバだけが、レイヤ2 で隣接します。ACI では、VXLAN が、基盤となるレイヤ3 ネットワークインフラストラクチャからレイヤ2 のドメインを切り離すことにより、このジレンマを解決します。

図 1: ACI ファブリック



トラフィックがファブリックに入ると、ACI がカプセル化してポリシーを適用し、必要に応じてスパインスイッチ (最大 2 ホップ) によってファブリックを通過させ、ファブリックを出るときにカプセル化を解除します。ファブリック内では、ACI はエンドポイント間通信でのすべての転送について、Intermediate System-to-Intermediate System プロトコル (IS-IS) および Council of Oracle Protocol (COOP) を使用します。これにより、すべての ACI リンクがアクティブで、ファブリック内での等コストマルチパス (ECMP) 転送と高速再コンバージョンが可能になります。ファブリック内と、ファブリックの外部のルータ内でのソフトウェア定義ネットワーク間のルーティング情報を伝播するために、ACI はマルチプロトコル Border Gateway Protocol (MP-BGP) を使用します。

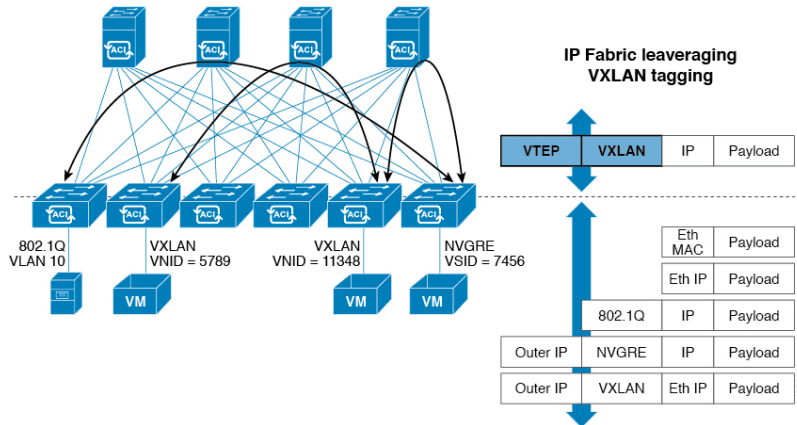
ACI で VXLAN

VXLAN は、レイヤ2 オーバーレイの論理ネットワークを構築するレイヤ3 のインフラストラクチャ上でレイヤ2 のセグメントを拡張する業界標準プロトコルです。ACI インフラストラク

チャレイヤ2 ドメインが隔離ブロードキャストと障害ブリッジ ドメインをオーバーレイ内に存在します。このアプローチは大きすぎる、障害ドメインの作成のリスクなしで大きくなるデータセンター ネットワークを使用できます。

すべてのトラフィック、ACIファブリックはVXLAN パケットとして正規化されます。入力でACI VXLAN パケットで外部 VLAN、VXLAN、および NVGRE パケットをカプセル化します。次の図は、ACIカプセル化の正規化を示します。

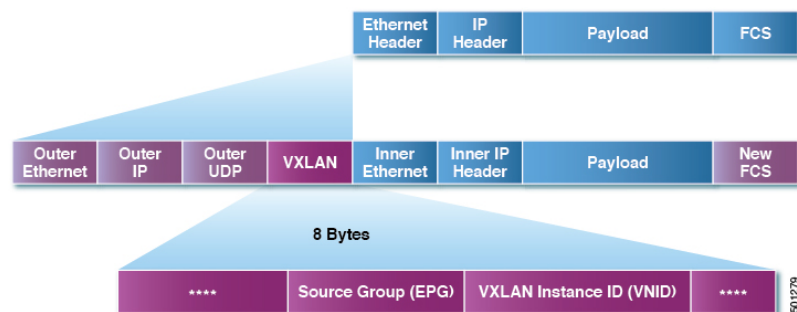
図 2: ACI カプセル化の正規化



ACI ファブリックでの転送は、カプセル化のタイプまたはカプセル化のオーバーレイ ネットワークによって制限または制約されません。ACI ブリッジ ドメインのフォワーディング ポリシーは、必要な場合に標準の VLAN 動作を提供するために定義できます。

ファブリック内のすべてのパケットにACIポリシー属性が含まれているため、ACIは完全に分散された方法でポリシーを一貫して適用できます。ACIにより、アプリケーションポリシーのEPG IDが転送から分離されます。次の図に示すように、ACI VXLANヘッダーは、ファブリック内のアプリケーションポリシーを特定します。

図 3: ACI VXLAN のパケット形式



ACI VXLAN パケットには、レイヤ2のMACアドレスとレイヤ3 IPアドレスの送信元と宛先フィールド、ファブリック内の効率的な拡張性の転送を有効にします。ACI VXLAN パケットヘッダーの送信元グループフィールドは、パケットが属するアプリケーションポリシーエンドポイントグループ (EPG) を特定します。VXLAN インスタンス ID (VNID) は、テナントの仮想ルーティングおよび転送 (VRF) ドメインファブリック内で、パケットの転送を有効にし

まず、VXLAN ヘッダーで 24 ビット VNID フィールドでは、同じネットワークで一貫レイヤ2 のセグメントを最大 16 個の拡張アドレス空間を提供します。この拡張アドレス空間は、大規模なマルチテナントデータセンターを構築する柔軟性 IT 部門とクラウドプロバイダーを提供します。

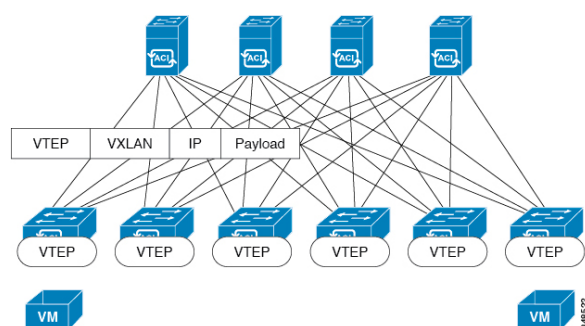
VXLAN を有効に ACI ファブリック全体にわたってスケールでの仮想ネットワーク インフラストラクチャのレイヤ3 のアンダーレイ レイヤ2 を展開します。アプリケーションエンドポイント ホスト柔軟に配置できます、アンダーレイ インフラストラクチャのレイヤ3 バウンダリのリスクなしでデータセンターネットワーク間をオーバーレイ ネットワーク、VXLAN でレイヤ2 の隣接関係を維持します。

サブネット間のテナントトラフィックの転送を促進するレイヤ3 VNID

ACI ファブリックは、ACI ファブリック VXLAN ネットワーク間のルーティングを実行するテナントのデフォルトゲートウェイ機能を備えています。各テナントに対して、ファブリックはテナントに割り当てられたすべてのリーフ スイッチにまたがる仮想デフォルトゲートウェイを提供します。これは、エンドポイントに接続された最初のリーフ スイッチの入力インターフェイスで提供されます。各入力インターフェイスはデフォルトゲートウェイ インターフェイスをサポートします。ファブリック全体のすべての入力インターフェイスは、特定のテナントサブネットに対して同一のルータの IP アドレスと MAC アドレスを共有します。

ACI ファブリックは、エンドポイントのロケータまたは VXLAN トンネルエンドポイント (VTEP) アドレスで定義された場所から、テナントエンドポイントアドレスとその識別子を切り離します。ファブリック内の転送は VTEP 間で行われます。次の図は、ACI で切り離された ID と場所を示します。

図 4: ACI によって切り離された ID と場所



VXLAN は VTEP デバイスを使用してテナントのエンドデバイスを VXLAN セグメントにマッピングし、VXLAN のカプセル化およびカプセル化解除を実行します。各 VTEP 機能には、次の 2 つのインターフェイスがあります。

- ブリッジングを介したローカルエンドポイント通信をサポートするローカル LAN セグメントのスイッチ インターフェイス
- 転送 IP ネットワークへの IP インターフェイス

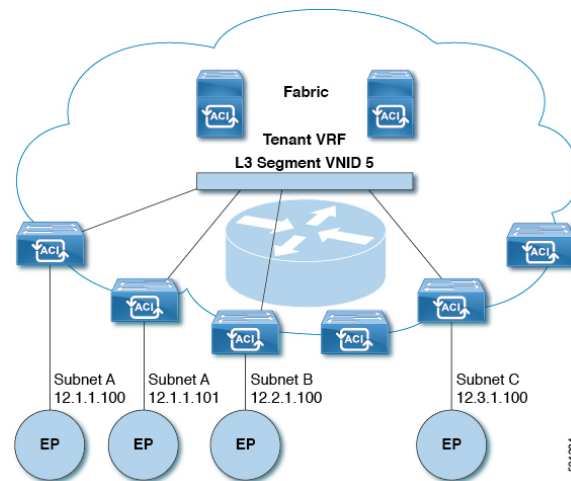
IP インターフェイスには一意の IP アドレスがあります。これは、インフラストラクチャ VLAN として知られる、転送 IP ネットワーク上の VTEP を識別します。VTEP デバイスはこの IP アドレスを使用してイーサネット フレームをカプセル化し、カプセル化されたパケットを、IP インターフェイスを介して転送ネットワークへ送信します。また、VTEP デバイスはリモート VTEP で VXLAN セグメントを検出し、IP インターフェイスを介してリモートの MAC Address-to-VTEP マッピングについて学習します。

ACI の VTEP は分散マッピング データベースを使用して、内部テナントの MAC アドレスまたは IP アドレスを特定の場所にマッピングします。VTEP はルックアップの完了後に、宛先リーフ スイッチ上の VTEP を宛先アドレスとして、VXLAN 内でカプセル化された元のデータ パケットを送信します。宛先リーフ スイッチはパケットをカプセル化解除して受信ホストに送信します。このモデルにより、ACI はスパニングツリー プロトコルを使用することなく、フルメッシュでシングル ホップのループフリー トポロジを使用してループを回避します。

VXLAN セグメントは基盤となるネットワーク トポロジに依存しません。逆に、VTEP 間の基盤となる IP ネットワークは、VXLAN オーバーレイに依存しません。これは送信元 IP アドレスとして開始 VTEP を持ち、宛先 IP アドレスとして終端 VTEP を持っており、外部 IP アドレス ヘッダーに基づいてパケットをカプセル化します。

次の図は、テナント内のルーティングがどのように行われるかを示します。

図 5: ACI のサブネット間のテナントトラフィックを転送するレイヤ 3 VNID



ACI はファブリックの各テナント VRF に単一の L3 VNID を割り当てます。ACI は、L3 VNID に従ってファブリック全体にトラフィックを転送します。出力リーフ スイッチでは、ACI によって L3 VNID からのパケットが出力サブネットの VNID にルーティングされます。

ACI のファブリック デフォルト ゲートウェイに送信されてファブリック入力に到達したトラフィックは、レイヤ 3 VNID にルーティングされます。これにより、テナント内でルーティングされるトラフィックはファブリックで非常に効率的に転送されます。このモデルを使用すると、たとえば同じ物理ホスト上の同じテナントに属し、サブネットが異なる 2 つの VM 間では、トラフィックが (最小パス コストを使用して) 正しい宛先にルーティングされる際に経由する必要があるは入力スイッチ インターフェイスのみです。

ACI ルート リフレクタは、ファブリック内での外部ルートの配布にマルチプロトコル BGP (MP-BGP) を使用します。ファブリック管理者は自律システム (AS) 番号を提供し、ルート リフレクタにするスパイン スイッチを指定します。



(注) Cisco ACI は IP フラグメンテーションをサポートしていません。したがって、外部ルータへのレイヤ 3 Outside (L3Out) 接続、または Inter-Pod Network (IPN) を介したマルチポッド接続を設定する場合は、インターフェイス MTU がリンクの両端で適切に設定することを推奨します。

IGP プロトコル パケット (EIGRP、OSPFv3) は、インターフェイス MTU サイズに基づいてコンポーネントによって構築されます。Cisco ACI では、CPU MTU サイズがインターフェイス MTU サイズよりも小さく、構築されたパケットサイズが CPU MTU より大きい場合、パケットはカーネルによってドロップされます (特に IPv6)。このような制御パケットのドロップを回避するには、コントロールプレーンとインターフェイスの両方で常に同じ MTU 値を設定します。

Cisco ACI、Cisco NX-OS、および Cisco IOS などの一部のプラットフォームでは、設定可能な MTU 値はイーサネット ヘッダー (一致する IP MTU、14-18 イーサネット ヘッダー サイズを除く) を考慮していません。また、IOS XR などの他のプラットフォームには、設定された MTU 値にイーサネット ヘッダーが含まれています。設定された値が 9000 の場合、Cisco ACI、Cisco NX-OS および Cisco IOS の最大 IP パケット サイズは 9000 バイトになりますが、IOS-XR のタグなしインターフェイスの最大 IP パケット サイズは 8986 バイトになります。

各プラットフォームの適切な MTU 値については、それぞれの設定ガイドを参照してください。

CLI ベースのコマンドを使用して MTU をテストすることを強く推奨します。たとえば、Cisco NX-OS CLI で、コマンド、`ping 1.1.1.1 df-bit packet-size 9000 source-interface ethernet 1/1` を使用してください。

スパニング ツリー プロトコル BPDU の送信

スパニング ツリー プロトコル (STP) を実行している 2 つ以上のスイッチが EPG の Cisco Application Centric Infrastructure (ACI) に接続されており、スタティック ポートが次のように割り当てられている場合：

- EPG で静的に割り当てられたすべてのポートは、タグなしでアクセスされます。STP ブリッジ プロトコル データ ユニット (BPDU) はタグなしで送受信されます。
- スタティックに割り当てられたトランク ポートとスタティックに割り当てられたアクセス タグなしポートが混在している場合：トランク ポートで受信された STP BPDU は、dot1q タグ付きのアクセス タグなしポートに送信されます。したがって、アクセス ポートは不整合状態になります。
- EPG で静的に割り当てられたトランク ポートと静的に割り当てられたアクセス ポートの組み合わせの場合、Cisco ACI は dot1q タグを使用して STP BPDU を送信し、アクセス ポートは 802.1p アクセスを使用します。

この場合、タグ付き STP パケットを受信して処理するには、レイヤ 2 スイッチで 802.1p アクセスを使用する必要があります。

802.1p がレイヤ 2 スイッチで許可されていない場合は、トランク ポート アクセスを使用します。

- Cisco ACI は全二重ハブとして機能し、BPDU が受信されたカプセル化 VLAN に関連付けられた VxLAN VNID 内でスパニングツリー BPDU をフラッディングします。Cisco ACI は全二重メディアであるため、高速スパニングツリープロトコル (RSTP) または高速 VLAN 単位スパニングツリー (RPVST) のバージョンを実行する外部スイッチは、デフォルトでポイントツーポイントリンクタイプになります。その結果、STP を実行し、同じカプセル化 VLAN および EPG VNID に接続する 2 つ以上の外部スイッチがある場合、コンバージェンスと不安定性の問題を回避するために、外部スイッチインターフェイスでリンクタイプを「共有」に設定する必要があります。これらの問題は、スイッチがこのカプセル化に接続されているすべてのブリッジ (または STP 対応スイッチ) から BPDU を受信するために発生する可能性があります。

スパニングツリー BPDU は、EPG パスで定義された特定の VLAN ID 内でフラッディングされます。この VLAN は、リーフスイッチでは `FD_VLAN` と呼ばれます。リーフスイッチ間で `FD_VLAN` 内のトラフィックを転送するために、Cisco ACI は、`fabric_encap` と呼ばれる VxLAN VNID を割り当てます。`fabric_encap` は、VLAN プールに属する数値ベース識別子を取得し、VLAN プールから割り当てられた VLAN ID のインデックス値を追加することによって取得されます。たとえば、VxLAN VNID 9000 は、VLAN 範囲 10 - 20 を含む VLAN プール A に割り当てられます。VLAN プール A の VLAN 10 には VNID 9000 が割り当てられ、VLAN 11 には VNID 9001 が割り当てられます。

このため、2 つの異なる EPG が同じ VLAN ID を使用しており、同じ VLAN プールからその VLAN ID を割り当てている場合は、異なるファブリックスイッチ上の 2 つの EPG に対して同じ `fabric_encap` VNID を導出できます。これにより、2 つの EPG 間でスパニングツリー BPDU が意図せずフラッディングされる可能性があります。

この動作を回避するには、物理ドメインなどの個別の VLAN プールを持つ異なるドメインを各 EPG に割り当て、特定の VLAN ID を個別の VLAN プールから割り当てます。これにより、ベース ID が異なるようになるので、`fabric_encap` VNID の重複が防止されます。

`fabric_encap` の値は、次のコマンドを使用して確認できます。また、特定の 802.1q VLAN ID のリーフスイッチの出力の「Fabric_enc」列でも確認できます。

```
vsh_lc -c "show system internal eltmc info vlan br"
```


翻訳について

このドキュメントは、米国シスコ発行ドキュメントの参考和訳です。リンク情報につきましては、日本語版掲載時点で、英語版にアップデートがあり、リンク先のページが移動/変更されている場合がありますことをご了承ください。あくまでも参考和訳となりますので、正式な内容については米国サイトのドキュメントを参照ください。