

# IETF TRILL UPDATE

なんで独自実装ばっかで IETF TRILL はさっぱりこないのか

2013年2月22日 ENOG19 Meeting

ブロケードコミュニケーションズシステムズ 株式会社

高嶋隆一〈rtakashi @ brocade.com〉



# WHAT IS "TRIL"?

# IRANSPARENT INTERCONNECTION OF LOTS OF LINKS

# コンセプト

Algorhyme V2, by Ray Perlner from RFC6325

I hope that we shall one day see A graph more lovely than a tree.

A graph to boost efficiency While still configuration—free.

> ツリーもうやめよう

➤ でも複雑な設定はNO!



# コンセプト cont.

Algorhyme V2, by Ray Perlner from RFC6325

A network where RBridges can Route packets to their target LAN.

The paths they find, to our elation, \Are least cost paths to destination!

→ 最短経路で必要な所にだけ流したい

# コンセプト cont.

Algorhyme V2, by Ray Perlner from RFC6325

With packet hop counts we now see The network need not be loop-free!

RBridges work transparently, Without a common spanning tree.

- ➤ TTLが欲しい
- ▶ ループ構成でもいい じゃない
- ケ エンドを変えずスパツリなしで動かすんだ!

# 要はアンチSTP!



# というかやり直し?

#### RFC6325 authors

R. Perlman<sup>o</sup>
Intel Labs

D. Eastlake 3rd Huawei

D. Dutt S. Gai Cisco Systems

> A. Ghanwani Brocade

➤ 802.1D Spanning-Treeの生 みの親



From

http://en.wikipedia.org/wiki/Radia\_Perlman

# どんな感じで?

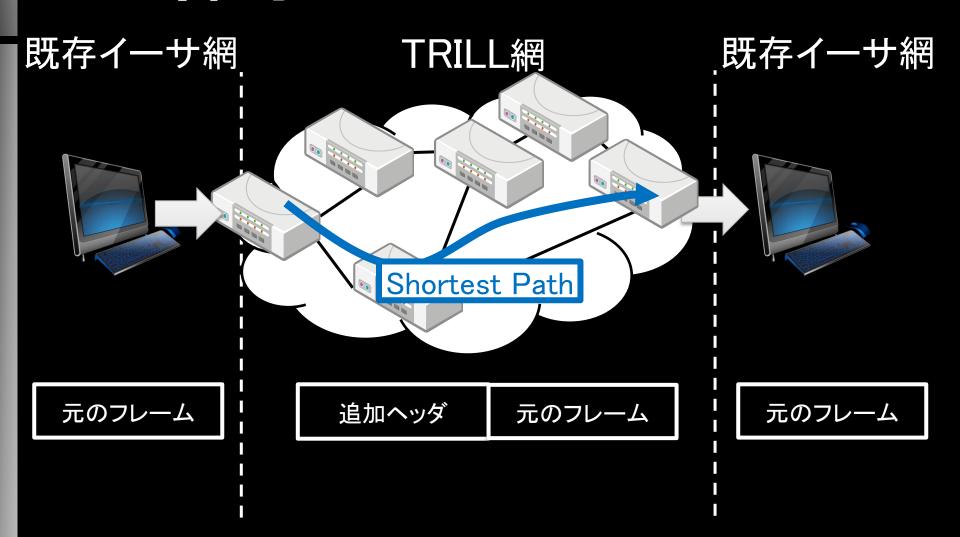
Control Plane

IS-ISによる ルーティング

Data Plane

追加へッダで Encapsulation

# 大雑把なイメージ



#### 実はRFC6325のタイトルは "TRILL" ではない

Routing Bridges (RBridges): Base Protocol Specification

Router (L3) の様な動きをする Bridge (L2) なので "Routing Bridge"

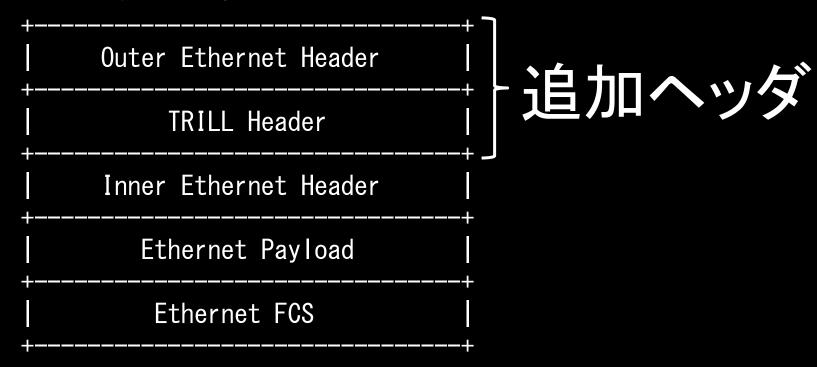
※以下"Rbridge" "RB"



# まずは Data Plane

# TRILL Frame ×

MAC-in-MAC Style Encapsulation



RFC6325 Figure 2: An Ethernet Encapsulated TRILL Frame

※Ethernet以外にPPPもサポートするが今回は省略



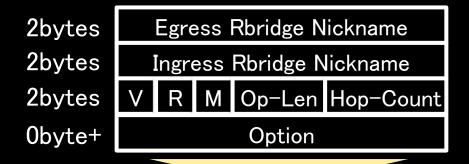
# TRILL Frame

#### もうちょっと詳しく

8bytes	Preamble/SFD		
6bytes	Dst MAC		
6bytes	Src MAC	Outer MAC	
4bytes	802.1Q	ヘッダ	
2bytes	Ether Type (=TRILL) / Size		
6bytes+	TRILL Header	}T	RILL ヘッダ
6bytes	Dst MAC	ר ו	
6bytes	Src MAC	Inner MAC	TRILL
4bytes	802.1Q	ヘッダ	ペイロード
2bytes	Ether Type / Size		
	Payload		
4bytes	CRC/FCS		

### TRILL Frame

#### TRILL Header



Egress Rbridge Nickname 宛先のRbridge
Ingress Rbridge Nickname 送信元のRbridge

M Multi Destination

Hop-Count TTL

その他

V: Version, R: Reserved, Op-Len: Option Length

16

# 次に Control Plane

# L2-IS-IS

#### IS-IS

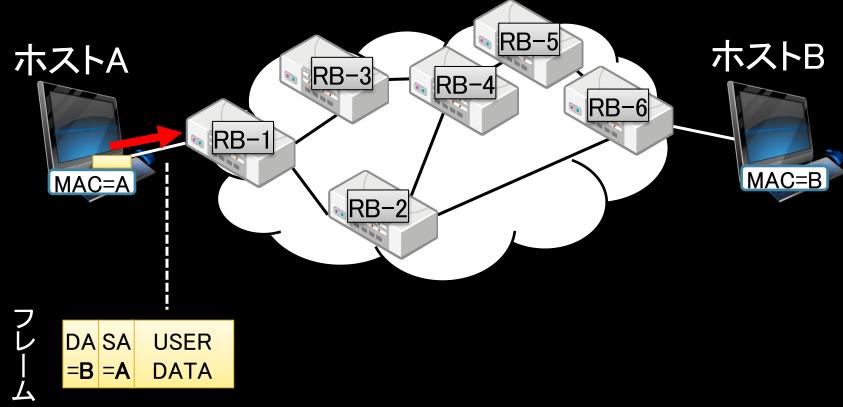
- √ ISIS
  - リンクステート型のルーティングプロトコル
  - > L2上で直接動作し、L3設定不要
  - ➤ TLV形式のサポートにより、マルチプロトコ ル拡張が容易
  - > SPF (Dijkstra法)による最短経路検索
  - ➤ ECMPのサポート

#### L2-IS-IS

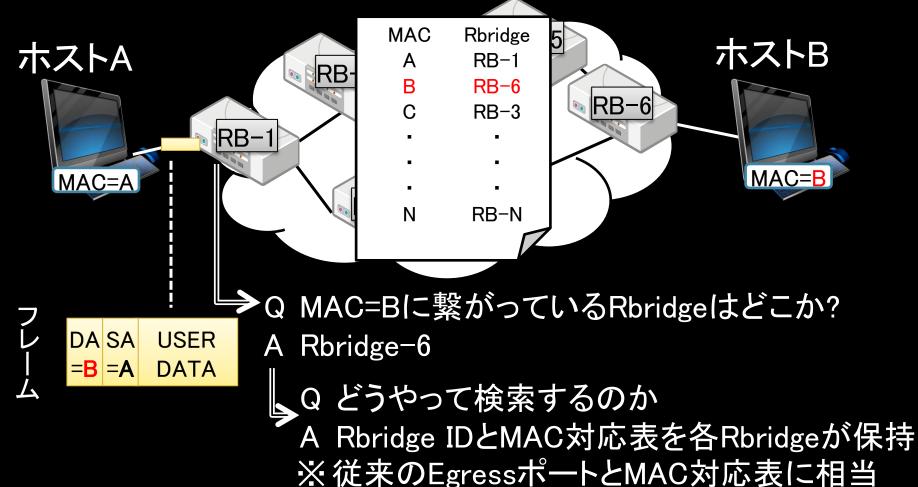
- ✓ L2-IS-IS
  - ➤ L3 IS-ISとは別の Ether Type を持ち、インス タンスも独立
  - ➤ Level-1 のみで構成されるシングルエリアで 動作
- ✓ Hello
  - ▶ リンク上のNeighborの確立
  - ➤ リンク上のDesignated Rbridge (DRB)の決定
  - ➤ DRBはTRILLに使用するVLAN ID※を決定
    - ※ Outer MACのVLAN ID, Ethernetでの動作の場合

# Unicast

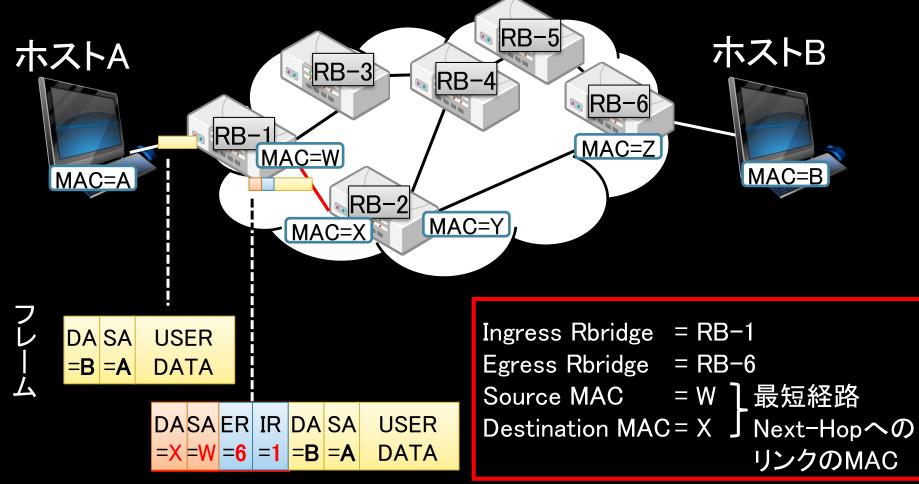
① ホストAはMACアドレスB宛てのフレームを送信



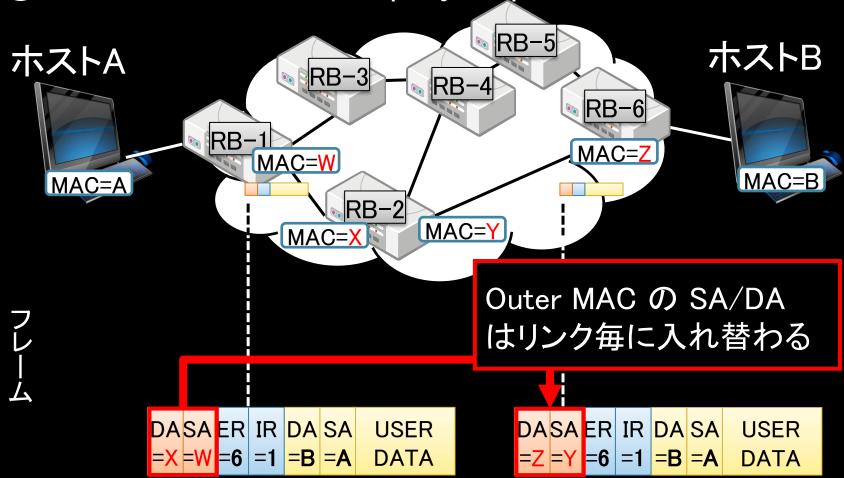
② Ingress Rbridge は Egress Rbridgeを検索



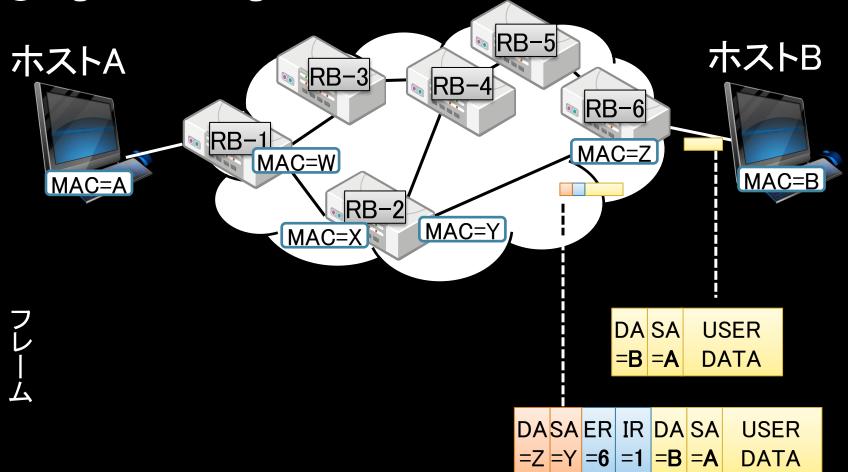
③ 検索結果を元にTRILLへッダ、Outer MACを付与



④ 最短経路に沿って Hop by Hop で転送



⑤ Egress Rbridge は TRILLヘッダ, Outer MAC を除去



# Multi-Destination

Broadcast

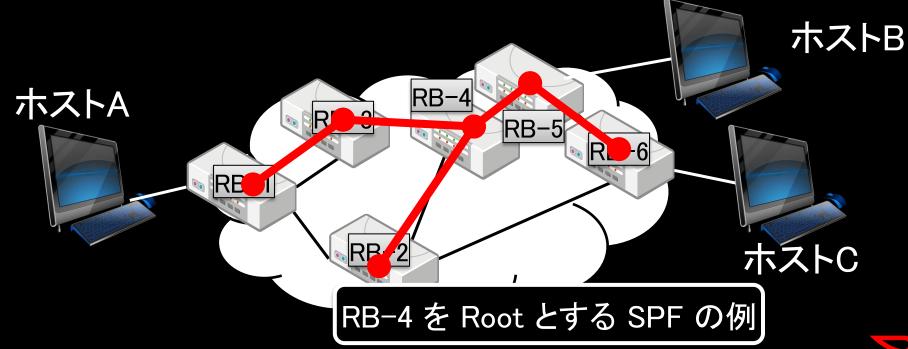
BUM Unknown-Unicast

Multicast

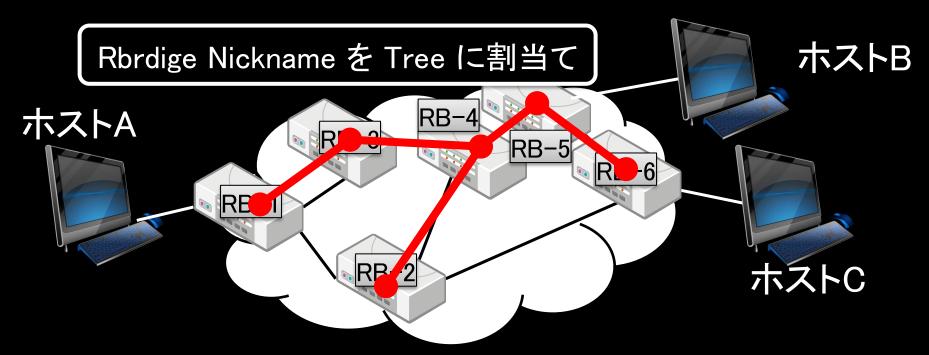
#### Multi-Destination 概要

Distribution Tree による配布

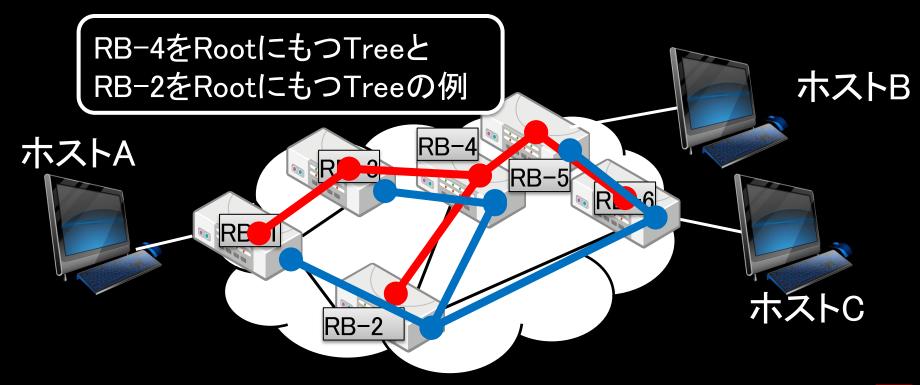
✓ TRILL網内では IS-IS SPF で作成された 双方向 Distribution Tree に沿って BUM 通信が行われる



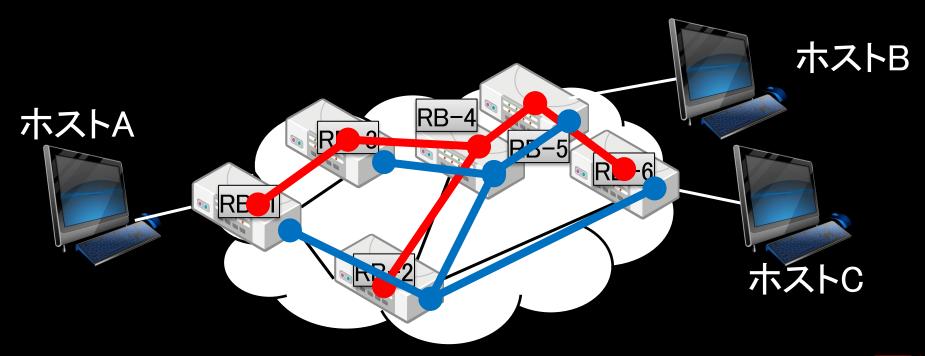
- ✓ Distribution TreeにもRbrdige Nicknameが振られる
  - ➤ Egress RbridgeにDistribution Treeを指定
  - ➤ L2/L3 の Broadcast/Multicast Address と同じ



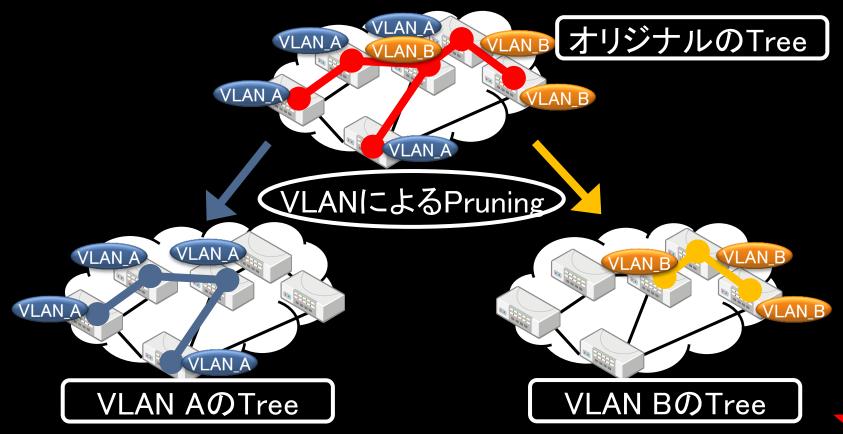
- ✓ 複数の Tree があってもよい
  - ▶ Ingress Bridge は最も近いRootを持つTreeを選択



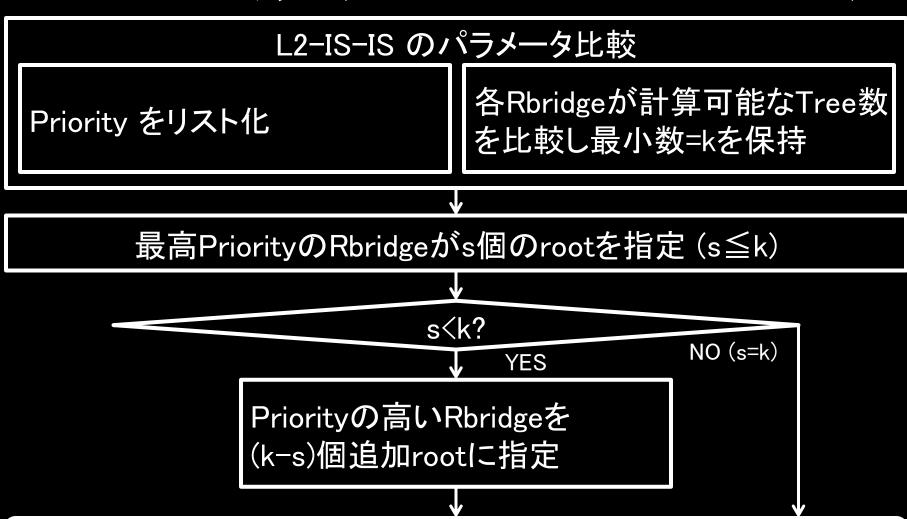
- ✓ 複数の Tree は同じ Root をもっていてもよい
  - Multicast でもロードバランス可能



- ✓ 各Distribution Tree は用途に応じてPruningされる
  - ➤ VLAN, IP Multicast由来のL2 Multicastアドレス



#### Tree生成(Root Selection)



各RbridgeはRootのリストを受けとり、SPFに従いTreeを生成

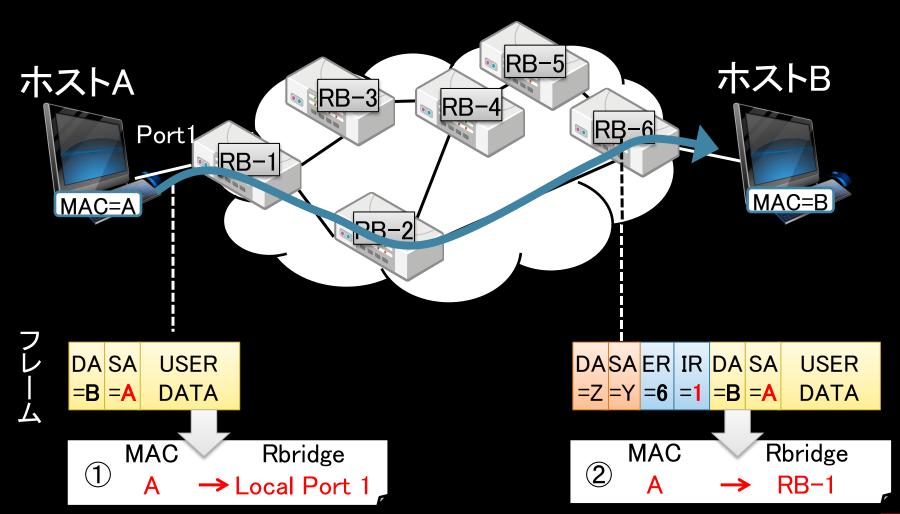
# MAC学習

#### 5種類のMAC学習方法

- ① Ingress Rbridgeで学習
  - ▶ ローカルの"Source Port→MAC"を学習
- ②Egress Rbridgeで学習
  - ▶ Decap時に"Ingress Rbridge↔MAC"を学習
- ③L2 Registration プロトコル (IEEE802.11 Association等)
- ④TRILL ESADI プロトコルによる学習
  - ➤ Rbridge間で"Rbridge↔MAC"情報を同期
- 5手動設定



#### 流れるフレームから学習



### TRILL ESADI Protocol

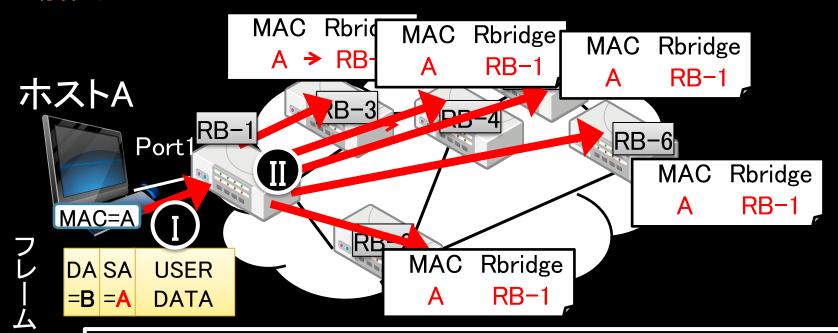
# End Station Address Distribution Information

**RFC6325** 

& draft-ietf-trill-esadi-02 (Expires: 2013-Aug-21)

### TRILL ESADI Protocol

動作イメージ



- I. いずれかのRbridgeがMACを学習もしくは消去
- II. ESADIプロコルにより該当VLANに属する全てのMAC学習テーブルが同期

### **ESADI** Basics

- ✓ ESADIはL2-IS-ISとは独立して動作する
  - ➤ TRILLのデータフレーム扱い
  - ➤ VLAN単位のDistribution Treeを通じて配布
- ✓ ESADIによりRbidge間でのMAC学習テーブルは 同期される
  - > BUMの低減
  - ➤ MAC削除時に直接接続されていない Bridge がMAC Age outまでブラックホールとなる状況を回避

# ESADIに関する注意

TRILL Working Group
INTERNET-DRAFT

Intended status: Proposed Standard

Updates: 6325

Expires: August 21, 2013

Hongjun Zhai Fangwei Hu ZTE Radia Perlman Intel Labs Donald Eastlake Huawei Olen Stokes Extreme Networks February 22, 2013

TRILL (Transparent Interconnection of Lots of Links):
ESADI (End Station Address Distribution Information) Protocol

<draft-ietf-trill-esadi-02.txt>

概要がRFC6325にあるだけで まだRFC化されてません

# 標準化動向

http://datatracker.ietf.org/doc/sear ch/?name=TRILL&rfcs=on&activeDra fts=on&search\_submit=

#### 

# たくさんありますが・・・

# Fine-Grained Labeling

TRILL Working Group

INTERNET-DRAFT

Intended status: Proposed Standard

Updates: 6325

Expires: August 12, 2013

Donald Eastlake
Mingui Zhang
Huawei
Puneet Agarwal
Broadcom
Radia Perlman
Intel Labs
Dinesh Dutt
Cumulus Networks
February 13, 2013

TRILL (Transparent Interconnection of Lots of Links):

Fine-Grained Labeling

<draft-ietf-trill-fine-labeling-05.txt>

### Status: WG Last Call

## Fine-Grained Labeling

VLANスケーラビリティの向上

Preamble/SFD

Outer MAC (DA,SA,802.1Q)

Ether Type (=TRILL) / Size

TRILL Header

**Dst MAC** 

Src MAC

802.1Q Etype 0x8100

Ethe'

/ Size

VLAN 12bit = 4k

O110/100

Original TRILL

Preamble/SFD

Outer MAC (DA,SA,802.1Q)

Ether Type (=TRILL) / Size

TRILL Header

**Dst MAC** 

Src MAC

Inner Label High Etype 0x893b

Inner Label Low Etype 0x893b

Ether

Size

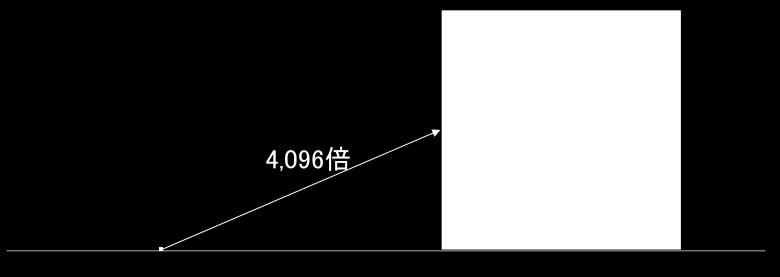
FGL 24bit = 16M

FGL



### FGL Basics

- ✓ VLAN-ID4096個の壁の突破
  - ▶ 4,096\*4,096=16,777,216個のネットワーク
  - ▶ 802.1ad (Q-in-Q)の様に STAG, CTAG の2段構成ではなく16M個をフラットに利用可能

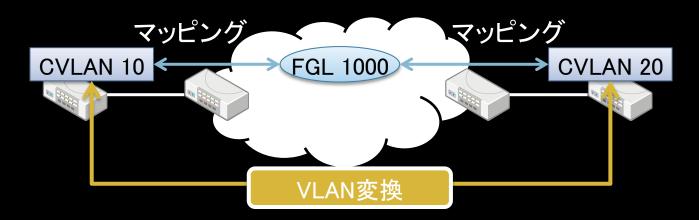


VLAN-ID 4,096個

FGL 16,777,216個

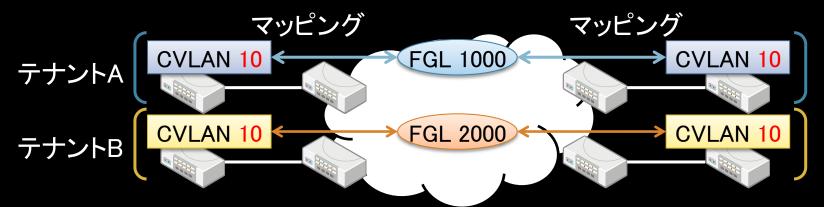
### FGL Basics

- ✓ FGLのエッジポート
  - ▶ Untag/802.1QタグをFGLにマッピングし、同一のフラッディングドメインを構成する
  - ► 結果としてFGLネットワークはVLANタグ変 換機能も持つ事になる
    - → VPLSのインスタンスに類似



## マルチテナントの考慮

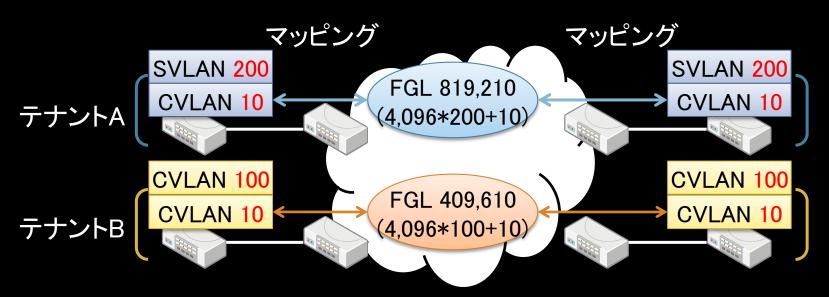
#### ✓ VLAN重複の許容



テナントAとテナントBは重複する CVLAN10を利用しているが、 マッピングされているFGLが異なる為、別のネットワークとして動作

### 802.1adとの共存

- ✓ 相互マッピング
  - ▶ 802.1ad のVLANスペース 4,096x4,096 と同等の空間を持っている為、802.1ad と相互マッピングも可能



単純にCTAGをInner Label Low に、STAGをInner Label Highに割り当てた例

# TRILLの実装

# IETF標準TRILLの実装

ベンダ	実装	Control Plane	Data Plane
Oracle	Solaris 11 (*2)	OSの標準機 能として動作	OSの標準機能として動作
NUST (*1)	Linux (*3)	Quagga/こtrilld を実装	Linux Kernel に拡張

- (\*1) National University of Sciences & Technology (NUST), Pakistan
- (\*2) http://docs.oracle.com/cd/E26924\_01/html/E25834/rbridgesoverview.html
- (\*3) http://www.ietf.org/proceedings/84/slides/slides-84-trill-2.pptx

# TRILL Like な商用実装

### どこが TRILL \*Like\*?

ベンダ	実装	Control Plane	Data Plane
Cisco	Fabric Path (*1)	Pre Standard(*3), L2MP	Pre Standard(*3), 独自フレーム
Brocade	VCS Fabric (*2)	FibreChannel標準 FSPF(*4)	IETF標準 フレーム

- (\*1) http://www.cisco.com/web/JP/product/hs/switches/nexus7000/prodlit/white\_paper\_c11-605488.html
- (\*2) http://www.brocadejapan.com/solutions-technology/technology/vcs-technology/overview
- (\*3) TRILLの標準化に先行してI-Dベースの情報とProprietary技術の組合せで実装, IS-IS ベース
- (\*4) FibreChanel標準 (T11 FC-SW-2) からImport, IS-ISベース

http://www.t11.org/t11/stat.nsf/89dcdcf87b06cf7e85256ebd000a3a96/70d54e897278fd078525660b

なんで独自実装ばっ かで IETF TRILL は さっぱりこないのか

# 標準化の遅れ

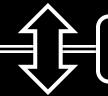
#### CY2010Q2

- ✓ Brocade VCS ファブリックのコンセプトを発表
- ✓ Cisco Fabric Pathのコンセプトを発表

#### CY2010Q4

???

- ✓ Brocade VCS 対応製品を発売
- ✓ Cisco Fabric Path 対応ファームウェアをリリース



#### 商用実装が先行

#### CY2011-July

- ✓ TRILLのコアとなるI-DがRFC化
  - RFC6325: Routing Bridges: Base Protocol Specification
  - RFC6326: TRILL Use of IS-IS
  - RFC6327: Routing Bridges: Adjacency
  - RFC6439: Routing Bridges: Appointed Forwarders



#### CY2013-Feb 現在

- ✓ RFC化されていない基本機能
  - draft-ietf-trill-esadi-02 (MAC学習)
  - draft-ietf-trill-rbridge-channel-08 (コントロールメッセージの取り扱い)
  - draft-ietf-trill-rbridge-vlan-mapping-08 (VLAN PCPとのマッピング)

#### ✓コアとなるRFCのリバイス

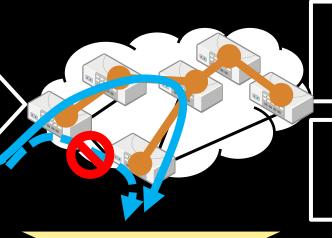
- draft-ietf-isis-rfc6326bis-00 (RFC6326)
- draft-ietf-trill-rbridge-extension-05 (RFC6325)
- draft-ietf-trill-clear-correct-06 (RFC6325,6327,6439)

#### WG Draft に限定し、機能拡張を 除いてもこれだけある

# 実装の難しさ

### Multi-Destination

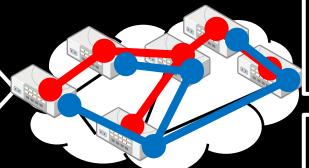
Distribution Tree による配布



Shortest Path とは限らない

RootにBUM が集中する

複数 Tree による Traffic Engineering?



最適Root配置は? 自動? 手動?\_\_\_

障害時の取り扱いは?

# ターゲットの違い

# RFC6325のターゲット

アンチSTP

元々のSTPの領域

キャンパスLAN

# 商用実装のターゲット

FabricPath および Cisco FabricPath Switching System による<mark>データセンターの</mark>

柔軟な拡張

#### 概要

従来のネットワークアーキテクチャは、静的なアプリケーションに対して高い有用性を提供するように設計化、きわめてスケーラブルな分散型アプリケーションなど、業界のトレンドは物理データセンター ゾーン間をAny-to-Any コミュニケーションをサポートする帯域幅のスケーラビリティ向上を必要としています。

http://www.cisco.com/web/JP/prod uct/hs/switches/nexus7000/prodlit/white\_paper\_c11-605488.html

#### Brocade VCS Fabric Technology

http://www.brocade.com/solutionstechnology/technology/vcsBrocade VCS Fabric technology helps organizations create efficient data center networks that just work. Ethernet fabrics built on Brocade VCS Fabric technology provide unmatched VM awareness and automation compared to traditional network architectures and competitive fabric offerings. In addition, they increase flexibility and IT agility, enabling organizations to transition smoothly to elastic, mission-critical networks in highly virtualized data centers.

#### データセンタネットワーク

technology/

### データセンタネットワークの要求と TRILLのカバー範囲

広帯域

耐障害性向上 断時間の短縮

管理工数の削減

仮想化アシスト

ストレージ トラフィック統合 TRILLのカバー範囲

**O**ECMP

○ブロックポートの排除

OL2ループの排除

△中継ノード故障時の高速切替

×マルチシャーシLAG

× ロジカルシャーシ

×増設時の設定工数の削減

× ハイパーバイザ連携

× VM Aware QoS

×低遅延

× DCB, FCoE, Loss-less iSCSI対応

#### では残りの要求の実装は?

#### TRILLでは満たせない要求

<u>実装方法</u>

- △中継ノード故障時の高速切替
- ×マルチシャーシLAG
- ×ロジカルシャーシ
- ×増設時の設定工数の削減
- ×ハイパーバイザ連携
- × VM Aware QoS
- ×低遅延
- × DCB,FCoE\*,Loss-less iSCSI対応

TRILL以外の標準技術

ベンダ固有拡張

ハードウェア実装

TRILL以外の標準技術

※BrocadeがFSPFを利用しているのは multihop-FCoE の為という側面もある

#### では残りの要求の実装は?

# 現状はベンダ固有拡張に頼る部分が大きい

# X TRILLスイッチを作った

O DCスイッチにTRILL の一部機能をいれた



# データセンタ以外 でのTRILLの適用

### DC以外のユースケース

- ■キャンパスLAN
  - ✓元々の要求にはあっているが、コストが見合うか
  - ✓STPとの併用ではあまり意味がない ※
  - ▶エンドユーザからの需要が増えればあり得る
    - ※TRILL自体はSTPとの相互接続を考慮している
- ■メトロエリアイーサネット
  - ✓既存のPB/PBB、MPLS/VPLSに対して優位性があるか
  - ✓OAMやBUM制御等の機能が必要
  - ▶ 機能の大幅追加がないと厳しい

#### どちらも今すぐには難しそう

# まとめ

- ✓ TRILLはSTPにかわるループフリーの L2制御プロトコル
- ✓ 標準化は基本的な機能を含めてまだ まだまだ進行中
- ✓ 商用実装はRFC準拠機能に加え、独 自機能を実装して製品化されている
- ✓ 現状はデータセンタでの利用がメイン





# Questions?