

Edge Virtual Bridge Proposal

Version 0, Rev 0.1

Abstract: This document is a proposal for the development of IEEE Edge Virtual Bridging (EVB) technologies. These proposals cover a suite of mechanisms that may be used to construct an EVB-based solution including architectural overview, discovery, management objects, and state machines.

Keywords: VEPA, VEB, S-Channel, ECP, VDP, CDCP, EVB

This document is a proposal for Virtual Ethernet Port Aggregator, Edge Virtual Bridging TLV, Edge Control Protocol, S-channels, and VSI Discovery Protocol.

Editors: Hewlett-Packard Corp., IBM

Contributing Authors

Company	Contacts
BNT	Daya Kamath
BNT	Jay Kidambi
BNT	Vijoy Pandey
Broadcom	Pat Thaler
Brocade	Anoop Ghanwani
Chelsio	Asgeir Eiriksson
Emulex	Chait Tumuluri
HP	Paul Bottroff
HP	Paul Congdon
HP	Chuck Hudson
HP	Michael Krause
IBM	Vivek Kashyap
IBM	Renato Recio
IBM	Rakesh Sharma
Juniper	Srikanth Kilaru
QLogic	Manoj Wadekar
IEEE Member	Uri Elzur

TABLE OF CONTENTS

<i>Contributing Authors</i>	2
<i>TABLE OF CONTENTS</i>	3
<i>LIST OF FIGURES</i>	4
<i>Related Documents</i>	5
<i>Change History</i>	5
1. Document Scope	6
1.1 Purpose	6
2. Introduction	7
3. Architecture and Operational Overview	17
3.1 VEPA Address Table Management.....	17
3.2 Processing from VSI Instance to bridge	18
3.3 Processing from the bridge to the VSI Instance	19
3.4 S-Channel Operation	21
3.5 Edge Control Protocol Operation.....	25
3.6 VSI Discovery and Configuration Protocol (VDP) Operation	26
3.6.1 VDP Type Configuration and Automation	28
3.6.2 VSI Type Definition and Management	30
3.6.3 VSI Manager ID	31
4. Ethernet Virtual Bridging TLV Semantics	33
5. S-Channel TLV Semantics and State Machine	36
5.1 S-Channel Bridge Components and Operation.....	36
5.1.1 Introduction	36
5.1.2 S-Component	37
5.2 CDCP Discovery and Configuration	38
5.2.1 CDCP TLV	38
5.2.2 CDCP Configuration Procedures.....	40
5.2.3 CDCP Configuration Variables.....	40
5.2.4 CDCP Configuration Procedures.....	42
5.2.5 CDCP Configuration State Machines.....	43
6. Edge Control Protocol (ECP) and State Machine	45
6.1 Requirements	45
6.2 Edge Control Protocol Data Unit	46
6.3 ECP Procedures.....	47

6.4	ECP State Machines.....	47
6.4.1	ECP Transmit State Machine	48
6.4.2	ECP Receive State Machine	49
7.	<i>Virtual Station Interface (VSI) Discovery and Configuration (VDP) and State Machine.....</i>	50
7.1.1	VSI Discovery and Configuration TLV	50
7.1.2	VDP Requirements and Assumptions	55
7.1.3	VDP – Local Variables and Procedures	56
7.1.4	Station VSI State Machine	58
7.1.5	Edge Bridge VSI State Machine	59
8.	<i>Glossary.....</i>	60
8.1	CDCP Basic Success Scenario (Sparse)	64
8.2	Station Adds an S-channel	64
8.3	Station Removes a S-channel	65
8.4	Drop #S-channels supported.....	66
8.5	Insufficient S-channels on bridge	67
8.6	VSI PreAssociate, Associate and DeAssociate	69
8.7	VSI Transport Error Case	70
8.8	VSI PreAssociate Resource Lease Refresh Exchange.....	71
8.9	VSI Associate Resource Lease Exchange.....	72

LIST OF FIGURES

Figure 1.	Example Hypervisor with Multiple VM, Multiple NIC, attached to an Adjacent Bridge	7
Figure 2.	Example Physical End Station with Multiple VM and Two Software VEB	8
Figure 3.	Example Physical End Station with Multiple Hardware VEB.....	9
Figure 4.	VEB Frame Relay Support	10
Figure 5.	Example Physical End Station with Multiple VM Communicating through a VEPA.....	11
Figure 6.	Example Physical End Station with Multiple Hardware VEPA	12
Figure 7.	VEPA Frame Relay Support.....	13
Figure 8.	S-channel Ethernet Components.....	15
Figure 9.	Multiple S-channels with on CAP at each end	16
Figure 10.	Example VEPA with associated conceptual VEPA Address Table.....	17
Figure 11.	VEPA Egress Processing	18
Figure 12.	VEPA Unicast Ingress Processing from Source A to Destination C	19
Figure 13.	VEPA Multicast Ingress Processing from Source A to Mulicate Group C	20
Figure 14.	Frame fowarding from a directly accessible VSI over a S-channel LAN.....	21
Figure 15.	Frame fowarding when S-channels are configured underneath a VEB	22
Figure 16.	Frame fowarding when S-channel are configured underneath a VEPA	23
Figure 17.	Frame fowarding over S-channel between a VEPA and adirectly attached VSI	24
Figure 18.	Example ECP Exchange	26
Figure 19.	VSI Type Architectural and Operational Overview	29

Figure 20.	VSI Manager ID	31
Figure 21.	VSI Type , VSI Version and VSI Instance ID	32
Figure 22.	VSI Manager Database Lookup.....	32
Figure 23.	EVB TLV Format	33
Figure 24.	Example EVB TLV Exchange.....	35
Figure 25.	Example S-channel Block Diagram.....	36
Figure 26.	Station and Bridge VEB/VEPA and S-comp Diagram	37
Figure 27.	CDCP TLV	39
Figure 28.	Example CDCP TLV Exchange	40
Figure 29.	CDCP State Machine.....	44
Figure 30.	ECP Data Unit	46
Figure 31.	ECP Transmit State Machine.....	48
Figure 32.	ECP Receive State Machine	49
Figure 33.	VDP TLV	50
Figure 34.	VDP Format = 1 Schema.....	51
Figure 35.	MAC-VLAN Information Format 1	54
Figure 36.	Station's VSI State Machine.....	58
Figure 37.	Edge Bridge's VSI State Machine	59
Figure 38.	Sparse S-channel Request.....	64
Figure 39.	Add an S-channel.....	65
Figure 40.	Remove an S-channel	66
Figure 41.	Drop an S-channel	67
Figure 42.	Insufficient Resources	68
Figure 43.	VSI PreAssociate, Associate and DeAssociate Exchange	69
Figure 44.	VSI Transport Error.....	70
Figure 45.	PreAssociate Resource Lease Exchange.....	71
Figure 46.	Associate Resource Lease Exchange.....	72

Related Documents

Specifications	Company

Change History

By	Date	Details
	4/2/2010	Version 0 Proposal

1. Document Scope

This document details Virtual Ethernet Port Aggregator (VEPA) theory of operation and the discovery and capability exchange protocol used to support a VEPA solution. VEPA relies upon LLDP to provide discovery and capability exchange. These exchanges occur between a physical end station and an adjacent bridge.

1.1 Purpose

The purpose of this proposal is to define a proposal for **Discovery and Configuration of Ethernet Virtual Bridging (EVB) capabilities, using: S-channel, Edge Control Protocol (ECP) and Virtual Station Interface (VSI) Discovery and Configuration Protocol (VDP)**. These protocols are used to determine EVB, S-channel, ECP and VDP capability presence within a physical end station and an adjacent bridge.

2. Introduction

Evolving standards combined with the growing size of enterprise and cloud-based networking deployments has led to a significant increase in the complexity of Ethernet networking in the data center. The advent of virtualization technology has compounded this complexity due to the significant increase in the number of Ethernet switches and the change in the solution deployment scenario.

A hypervisor is a software entity that enables multiple Virtual Machines (VM) to share common hardware as illustrated in figure 1. Each VM contains at least one virtual NIC (vNIC) that is associated through the hypervisor with a physical NIC. To create this association, hypervisors have incorporated Virtual Ethernet Bridges (VEB) into the physical end station effectively adding one or more Ethernet switches per end node. A VEB is a frame relay service that supports local bridging between multiple virtual end stations (an internal private virtual network) and (optionally) the external bridging environment. A VEB may be implemented in software as a virtual switch (vSwitch) as illustrated in Figure 2 or as embedded hardware within a Network Interface Controller (NIC) as illustrated in Figure 3.

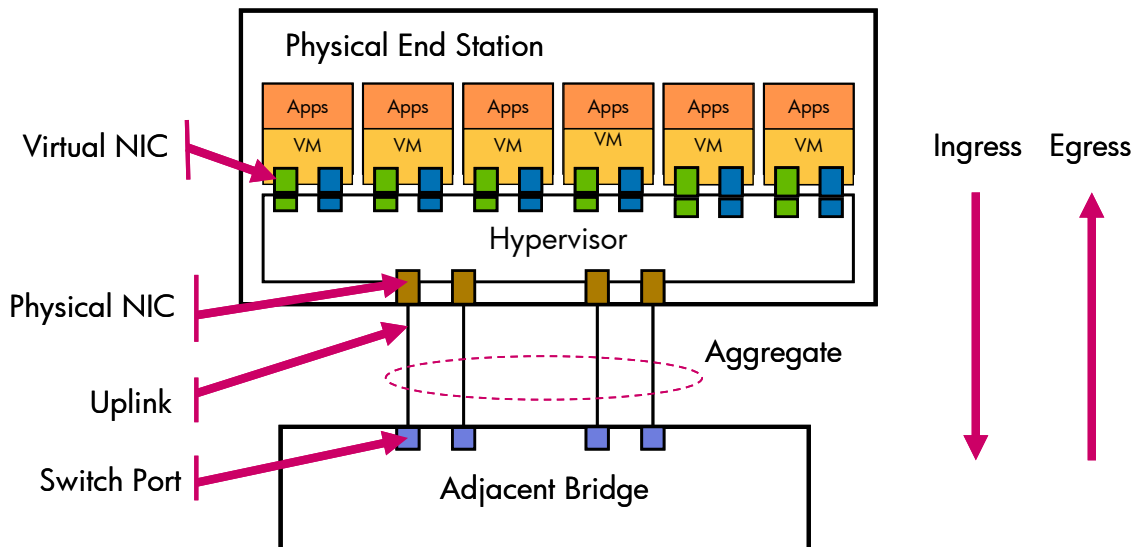


Figure 1. Example Hypervisor with Multiple VM, Multiple NIC, attached to an Adjacent Bridge

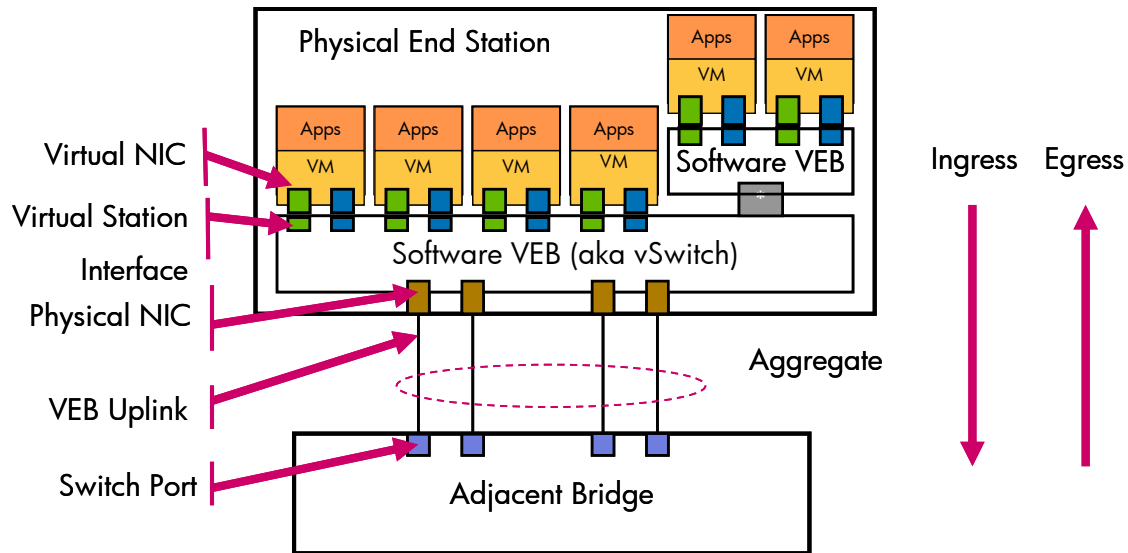


Figure 2. Example Physical End Station with Multiple VM and Two Software VEB

Figure 2 illustrates an example of the data paths for two software VEBs:

- Each VM may support one or more Virtual NICs.
 - Typically, a VM will support a virtual NIC (vNIC) that emulates a physical NIC. Each vNIC is associated with a Virtual Station Interface (VSI) which is connected to a VEB.
- A VEB supports a single logical uplink to the external adjacent bridge. Multiple uplinks can be aggregated via 802.3ad or other techniques.
- A software VEB (vSwitch) is typically implemented within a hypervisor requiring each VM I/O operation to trap to the hypervisor for processing.
 - Hypervisor traps consume system resources and can lead to varying performance loss depending upon the number of I/O operations per second and the amount of rich network functionality performed per operation.
 - Being in the hypervisor allows a software VEB to support one or more physical NICs.
- VEB may be cascaded to provide modularity or additional fan-out.
- Not shown but important to note is a VEB does not require any modifications to the Ethernet frame to operate.

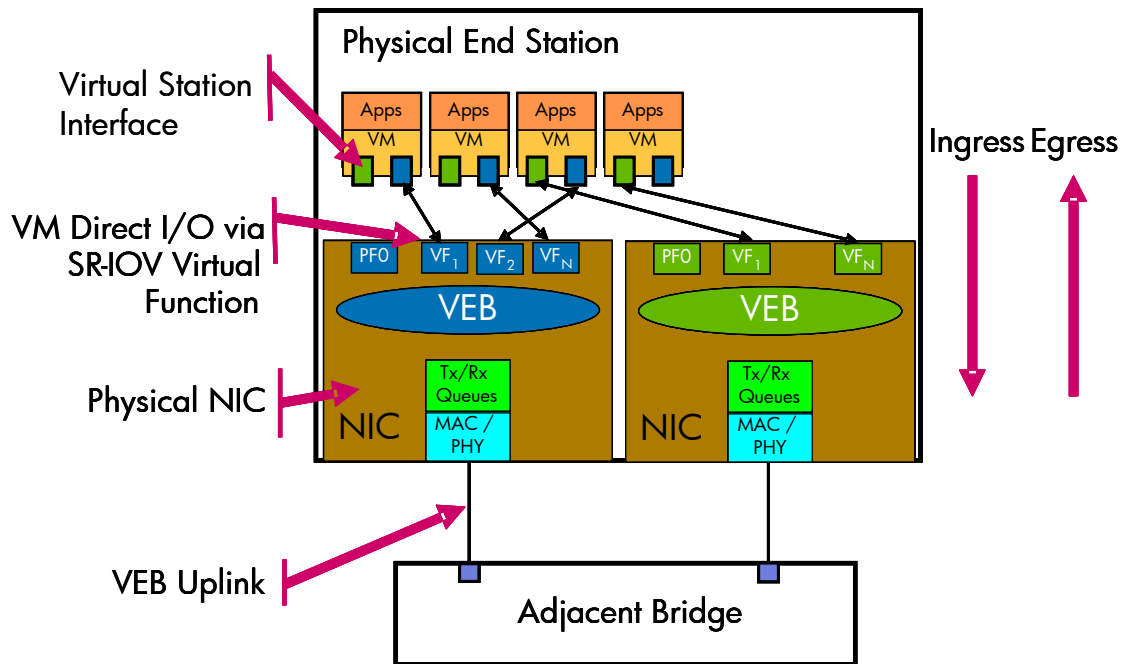


Figure 3. Example Physical End Station with Multiple Hardware VEB

Figure 2 illustrates an example of the data paths for multiple hardware VEBs:

- Each physical NIC supports
 - One (or more) physical ports attached to an adjacent bridge
 - Each physical port represents a single VEB uplink.
 - One or more hardware-embedded VEB. An embedded VEB cannot span multiple physical NIC.
 - Direct I/O support via SR-IOV Virtual Functions (VF).
 - Direct I/O allows a VM to bypass the hypervisor and directly access the NIC to send / receive packets. Bypassing the hypervisor reduces system resource consumption allowing higher performance solutions than traditional software VEB.
- Each VM may support one or more Virtual NICs.
 - In this example, each VM supports two vNIC – one per physical NIC.
 - Each vNIC is associated with a VSI which is associated with a SR-IOV VF to provide Direct I/O support.

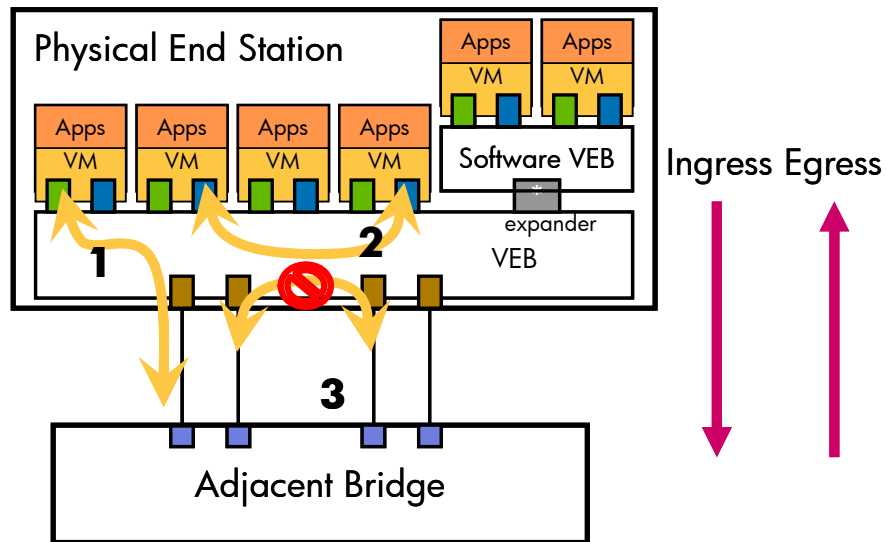


Figure 4. VEB Frame Relay Support

VEB packet forwarding supports both traditional end station-to-adjacent bridge as well as local VSI-to-VSI packet forwarding. As illustrated in Figure 4, a VEB forwards packets as follows:

- VEB forwards packets based on the MAC address and optionally via a port group or VLAN identifier.
- VEB forwards packets from a VSI to the uplink from an adjacent bridge (path 1) or between co-located VSI (path 2)
 - A NIC-embedded VEB can only forward packets between VSI attached to a common NIC. As shown in Figure 3, only VM that share the blue VEB can forward packets via the blue VEB. Similarly, only VM that share the green VEB can forward packets via the green VEB. A VM on the blue VEB cannot forward packets to a VM on the green VEB directly; a software VEB or external bridge would be required to bridge the two NIC-embedded VEBs.
- VEB supports only a single active logical uplink
 - Multiple uplinks can be teamed via 802.3ad or other techniques
 - Uplink-to-uplink packet forwarding is not allowed (path 3)
- VEB does not participate in or affect spanning tree operation.

VEB solutions have been shipping for a number of years and are available from multiple suppliers. Though the functional robustness of solutions will vary, local bridging via a VEB provides a number of common benefits and allows hypervisors to:

- Operate without external bridges attached
- Operate with a broad range of Ethernet environments
- Maximize local bandwidth – bandwidth is limited by end station memory and local I/O bandwidth and not by the Ethernet link bandwidth
- Minimize local latency – no incremental latency due to interaction with the external network
- Minimize packet loss, i.e. no packet loss due to external network events – external bridge or link failure, CRC error detection, congestion-based packet loss, etc.

By definition traffic between VMs connected to a VEB stay within the server. Some clients prefer the traffic to be sent through an external switch, so the external network's access

and security policies can be applied to the traffic. To address this type of requirement a Virtual Ethernet Port Aggregator (VEPA) is documented.

A Virtual Ethernet Port Aggregator (VEPA) is a capability within a physical end-station that collaborates with an adjacent bridge to provide frame relay services between multiple co-located virtual machines (VMs) and the external network. A VEPA collaborates by:

- Forwarding all station-originated frames to the adjacent bridge for frame processing and frame relay.
- Steering all frames and replicating Multicast and Broadcast frames received from the adjacent bridge to the appropriate VM destinations.
- A VEPA takes advantage of a special reflective relay forwarding mode (i.e. allow forwarding back out the port a frame was received) on the adjacent bridge to support inter-VM communication within the same physical host.
 - Clause 8.6.1 of Standard IEEE 802.1Q-2005 [11] states that when a switch reception port is in the forwarding state, each switch port in the forwarding state, *other than the reception port itself*, is a potential transmission port. A VEPA requires an exception to this rule in order to allow inter-VM traffic on the adjacent host over the single uplink. This exception distinguishes the port attached to a VEPA uplink as a VEPA-enabled port which supports forwarding in reflective relay mode.
- Similar to a VEB, a VEPA may be implemented in software or in conjunction with embedded hardware within a NIC.

The VEPA is connected to the adjacent bridge only by a single uplink connection. The connection is attached to a VEPA-enabled port on the adjacent bridge. A conceptual VEPA is shown in Figure 5.

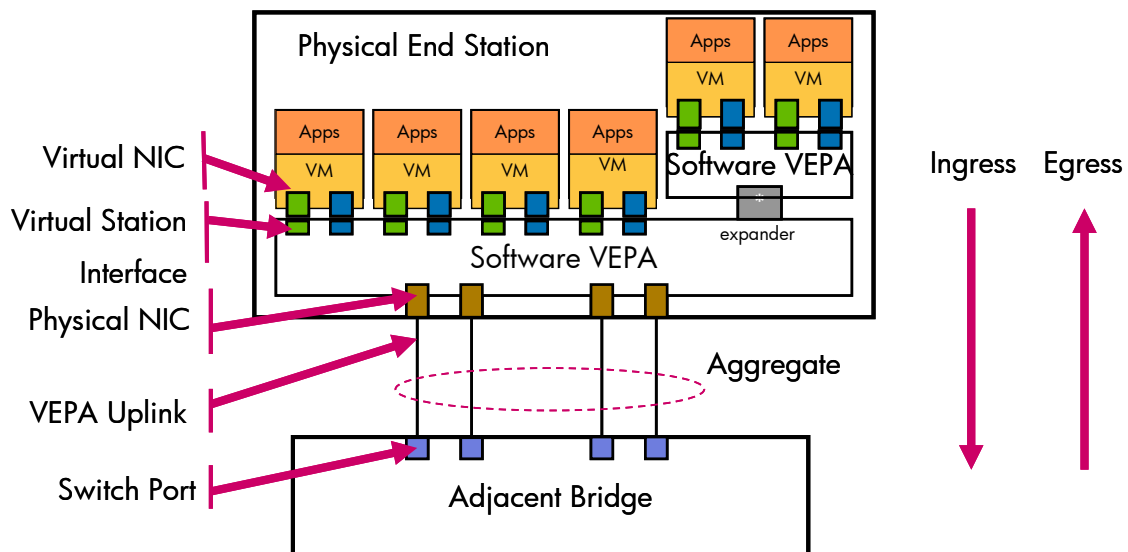


Figure 5. Example Physical End Station with Multiple VM Communicating through a VEPA

Figure 5 illustrates an example of the data paths for multiple VM communicating through a VEPA:

- Each VM may support one or more Virtual NICs..
 - A VM supports a virtual NIC (vNIC) which emulates a physical NIC. A vNIC is attached to a VEPA via a Virtual Station Interface (VSI). A VSI is a physical or software emulated end station connected to a VEB or a VEPA.
- A VEPA supports a single logical uplink.
- Software VEPA (vSwitch) may support one or more physical NICs.
- The total number of VSI made available may be scaled by cascading VEPA in a tree as shown in Figure 5. The port on a root VEPA connected to a leaf VEPA higher in the topology is known as an expander port. A root VEPA will forward all frames with an unknown destination address to the expander port. This eliminates the need for the root VEPA to comprehend all of the MAC addresses of every VM in the physical end station.
- Not shown but important to note is VEPA does not require any modifications to the Ethernet frame to operate.

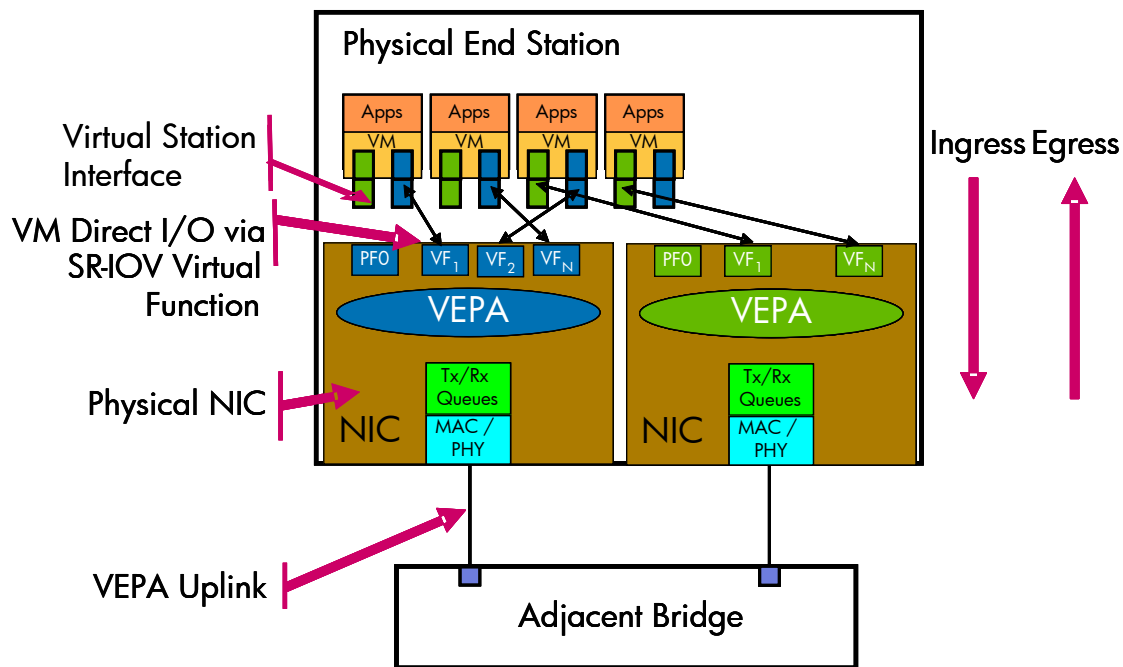


Figure 6.

Example Physical End Station with Multiple Hardware VEPA

Figure 6 illustrates an example of the data paths for multiple hardware VEPAs:

- The end station contains two independent physical NICs. Each NIC supports
 - One or more hardware-embedded VEPA. While this figure illustrates only one VEPA (blue or green) per NIC, an implementation may support multiple VEPA, each on a separate S-channel.
 - In this example, the blue and the green VEPA are completely separate, independent entities, i.e. they do not share any resources and cannot directly communicate with one another.

- One or more physical ports attached to an adjacent bridge. A VEPA has only one logical uplink.
- In this example, each NIC supports direct I/O via PCI SR-IOV technology.
 - Direct I/O allows a VM to bypass the hypervisor and directly access the NIC to send / receive packets. Direct I/O is achieved by using a light-weight PCI Function called a Virtual Function (VF) to act as a conduit between the VM and the NIC hardware. This is analogous to a NIC supporting multiple traditional PCI Functions but is less hardware-intensive. Each VF is associated with a Physical Function (PF) which can be used by the hypervisor as a management conduit to provide overall control of the device or the port depending upon the implementation.
- Each VM may support one or more Virtual NICs.
 - In this example, each VM supports two vNIC – one per physical NIC.
 - Each vNIC contains a VSI which is associated with a SR-IOV VF to provide the direct I/O conduit.

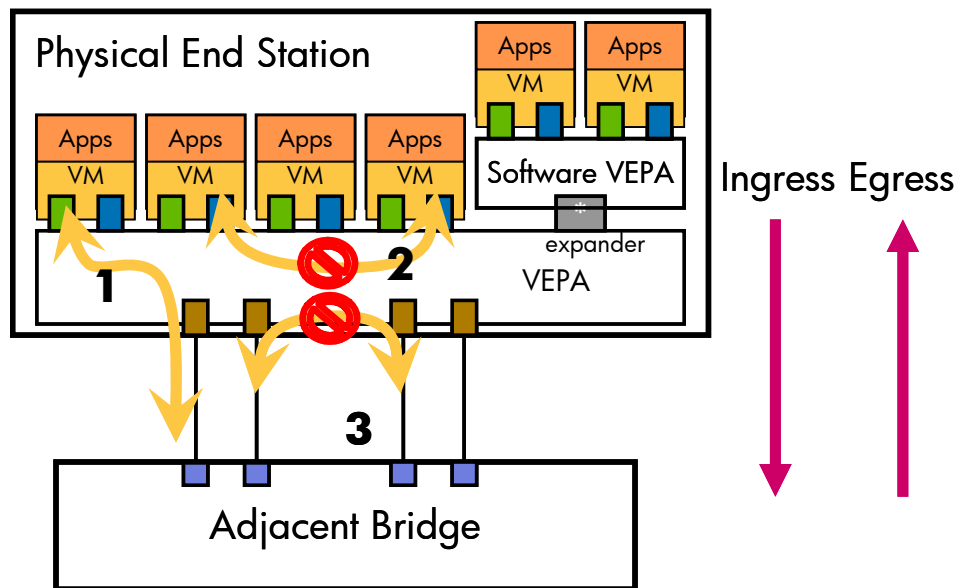


Figure 7. VEPA Frame Relay Support

As illustrated in Figure 7, a VEPA forwards packets as follows:

- VEPA forwards packets based on the MAC address and optionally via a port group or VLAN identifier.
- All VEPA traffic must be forwarded from the VSI to the uplink of an adjacent bridge (path 1).
 - VSI-to-VSI packet forwarding is not allowed (path 2).
- A VEPA supports only a single active logical uplink
 - Uplink-to-uplink packet forwarding is not allowed (path 3)
 - A VEPA may be partitioned into multiple logical VEPA each associated with its own independent uplink.
- A VEPA does not participate in or affect spanning tree operation, i.e. VEPA internal topology is not visible to the adjacent bridge, except for management associated TLVs (e.g. EVB TLV in this document).

Based on the prior materials, the reader should note the significant overlap in functionality and potential implementation between a VEB and a VEPA with the primary difference occurring in frame relay support. Further, this difference determines where and how network features are surfaced and their associated impact on system functional robustness and performance. This difference allows VEPA solutions to provide the following benefits:

1. Reduces complexity and potentially enables higher performance by off-loading advanced network functions from the VM or hypervisor to the adjacent bridge.
2. Allows NICs to maintain low cost circuitry by leveraging advanced functions on the adjacent bridge.
3. Enables a consistent level of network policy enforcement by routing all network traffic through the adjacent bridge with its more complete policy-enforcement capabilities.
4. Provides visibility of inter-VM traffic to network management tools designed for adjacent bridge.
5. Reduces the amount of network configuration required by server administrators, and as a consequence, reduces the complexity for the network administrator.
6. Can increase solution performance by off-loading advanced network functionality that may be computationally intensive to implement within a hypervisor or VM to the adjacent bridge.

As it can be seen, a VEPA provides a number of benefits but it too has limitations:

- Promiscuous support – To support a promiscuous VSI, a VEPA address table must be configured with all VM source MAC addresses. This requires either adding MAC address learning support or provisioning large address tables. Either option adds implementation cost and complexity.
- Support for simultaneous VEB, VEPA, and directly accessible ports on the same physical link – The adjacent bridge can only process a frame based on its contents and therefore lacks sufficient information to delineate between these three operating modes.
- Hierarchy of unrestricted physical ports – Normal bridge learning and flooding is not possible due to the lack of information within a frame.

To address these limitations, IEEE Std. 802.1Qbc-2011 is applied. This standard enables multiple virtual channels to be multiplexed on a single physical LAN – referred to as S-channel functionality. Individual S-channels are delineated by a tag which is added to the frame and processed by S-VLAN Components (a bridge component) which are logically inserted into the adjacent bridge and the physical end station below the virtual bridge layer as illustrated in the following figure.

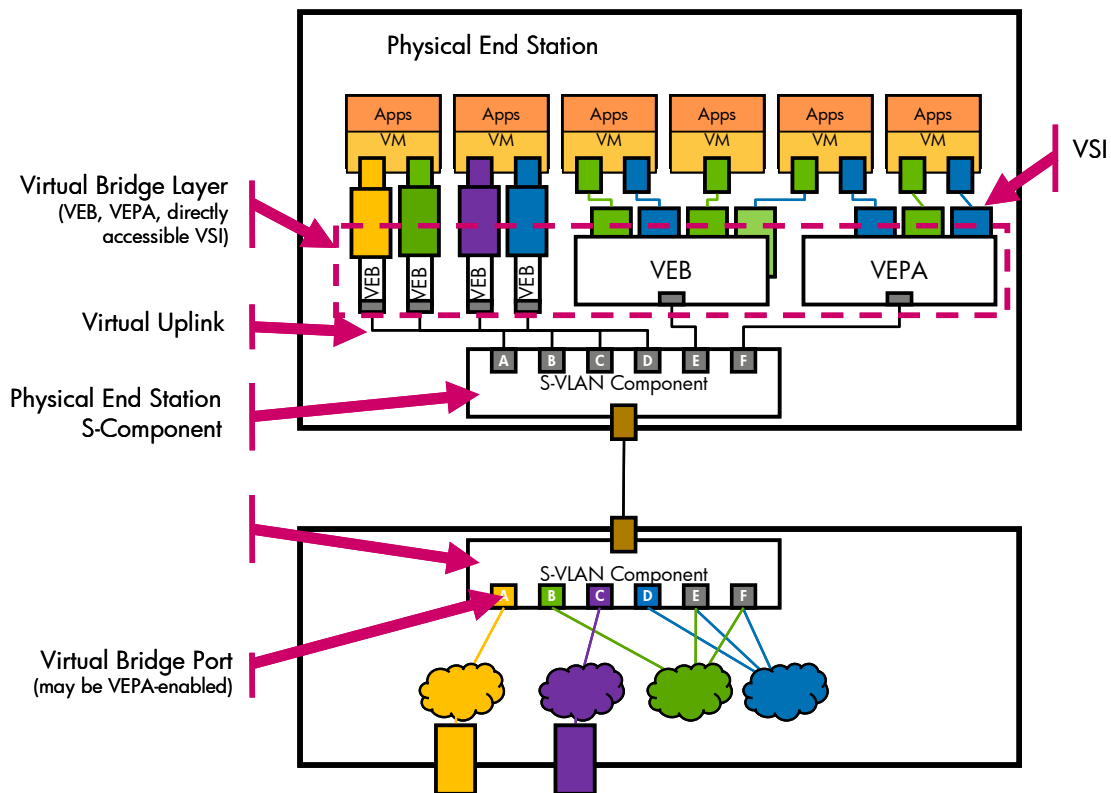


Figure 8. S-channel Ethernet Components

The S-VLAN component recognizes, inserts and removes service VLAN tags (S-Tag) to enable multiple S-channels in the bridged network. Adding an S-VLAN component to an end-station allows VEPA, VEB, and individual VSI to operate independently and simultaneously. Each VEPA, VEB, or individual VSI operates over its own virtual uplink instantiated by a pair of S-VLAN components - one in the adjacent bridge and one on the end-station.

The virtual uplinks created by the end-station's S-VLAN component are effectively connected over an S-channel uplink to virtual ports (S-channel Access Ports or CAPs) created by the S-VLAN component on the adjacent bridge as illustrated in the following figure.

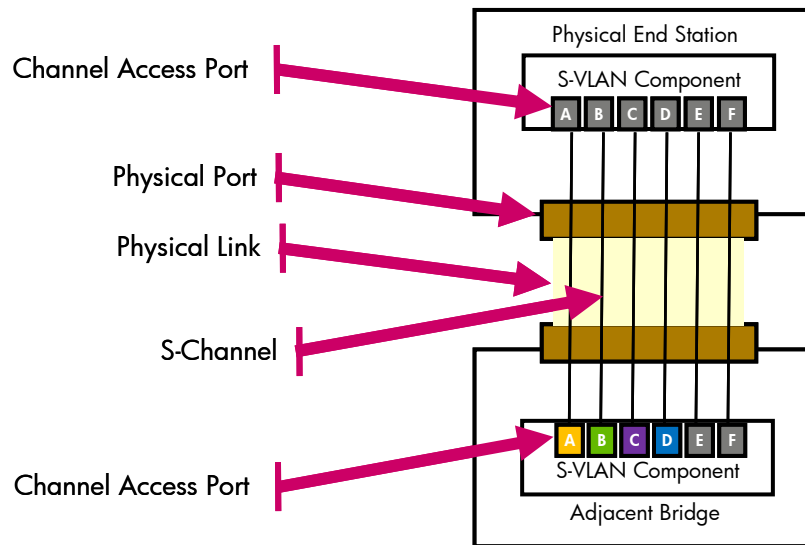


Figure 9. Multiple S-channels with on CAP at each end

Each frame traversing the physical S-channel LAN will all have an S-Tag inserted by the first S-VLAN component it encounters and removed by the second S-VLAN component as it reaches the far side of the S-channel LAN. The S-Tag inserted by the end-station identifies the particular source virtual uplink and the S-Tag inserted by the adjacent bridge identifies the destination virtual uplink. Any frames that must be broadcast, multicast or flooded to more than one VSI are replicated by the adjacent bridge and delivered across the S-channel LAN as many times as needed, each with the proper S-Tag inserted.

Adding the S-channel capability to the end-station solves the problem of supporting virtual machines needing promiscuous ports by isolating such VSI in a separate S-channel. By doing so, normal learning and forwarding behavior is pushed to the adjacent bridge, isolating it from the simple forwarding of the VEPA. It also allows the end station administrator to choose how virtual VM are connected to the network. A group of VM that require direct connectivity between each other for high performance and low latency can be attached to a VEB. Another group that requires traffic visibility, firewall inspection or other services on the adjacent bridge can be attached to a VEPA. Finally any individual VM that requires an isolated promiscuous VSI can be attached directly to a virtual uplink.

The subsequent chapters within this proposal provide additional details on VEPA and S-channel operation, discovery, and configuration.

3. Architecture and Operational Overview

This chapter will describe VEPA and S-Channel architectural components and illustrate how these components are used via example operations.

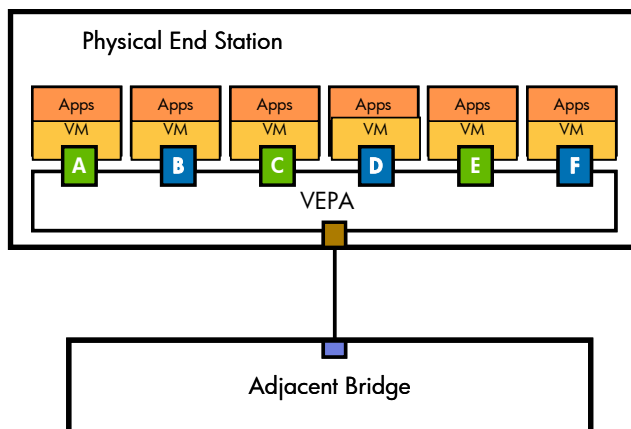
3.1 VEPA Address Table Management

As a network edge end station, a VEPA is not required to support address learning. Instead, the VEPA address table is populated through a registration process. As an address, filter, or VLAN identifier is registered, the server virtualization infrastructure (e.g. the Hypervisor) updates the corresponding VEPA address table entry. This applies to:

- VSI default MAC address
- Locally Administered Address (LAA)
- Multicast addresses
- Promiscuous address mode support

Note: A MAC address may appear on multiple VLANs. Hence, a <MAC, VLAN Identifier> pair is required to identify a unique VSI.

The following figure illustrates an example physical end station, VEPA, and the associated VEPA address table.



Destination MAC	VLAN	Copy To (ABCDEF)
A	1	100000
B	2	010000
C	1	001000
D	2	000100
E	1	000010
F	2	000001
Broadcast	1	101010
Broadcast	2	010101
Multicast C	1	101010
Unknown Multicast	1	100010
Unknown Multicast	2	010101
Unknown Unicast	1	000000
Unknown Unicast	2	000000

Figure 10. Example VEPA with associated conceptual VEPA Address Table

In this example, the VEPA address table holds the following:

- A unicast MAC address (and VSI Instance Identifier) per VSI
- Per VLAN broadcast address – frames are forwarded to the VSIs indicated by the corresponding bit mask
- Specified multicast addresses – VSI A, C, and E are the only participants in the specified multicast group.
- Per VLAN unknown unicast address – unknown frames are discarded on ingress (i.e. from the bridge to the VSI Instance), they are sent to the adjacent bridge on egress (i.e. from VSI Instance to bridge).
- Per VLAN unknown multicast address – unknown frames are forwarded to the VSI indicated by the corresponding bit mask

The address table is configured by the server virtualization infrastructure (e.g. the Hypervisor) simplifying the VEPA implementation by eliminating the need to support dynamic address learning. Further, the hypervisor can configure additional address table fields (not shown in the figure examples) such as QoS settings, VLAN configuration, promiscuous listening support, and so forth to provide additional functional capabilities.

3.2 Processing from VSI Instance to bridge

VEPA processing from VSI Instance to bridge is defined as the set of operations required to transfer a frame from a VSI to the VEPA uplink. Since all frames are required to be forwarded to the uplink, the frame is moved from the VSI to the physical uplink for transmission and the S-VLAN-Tag is added.

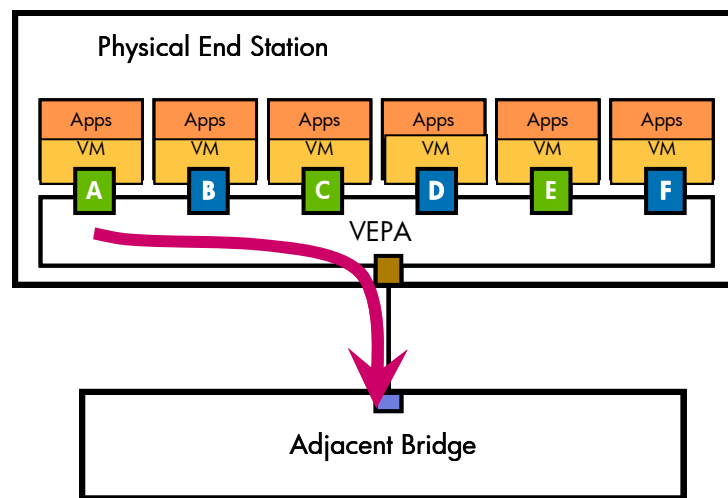


Figure 11. VEPA Egress Processing

While the main objective of VEPA is to forward frames to the adjacent bridge for advanced processing, VEPA implementations may provide additional processing on the frame since the VEPA comprehends the source VSI Instance. For example, source VSI MAC address validation to prevent spoofing, application of QoS and bandwidth management policies, VLAN formatting validation (tagged or untagged), and so forth. The adjacent bridge can only operate on the frame's contents (MAC address, VLAN ID,

etc.) and is unaware of the source VSI Instance therefore eliminating such functional possibilities.

3.3 Processing from the bridge to the VSI Instance

VEPA processing from the bridge to the VSI Instance is defined as the set of operations required to steer and transfer a frame received on the uplink to the appropriate VSI or set of VSIs. The VEPA must make use of the address table to perform this operation correctly. Address table access is illustrated in the following figure as VSI A transmits a unicast frame to VSI C.

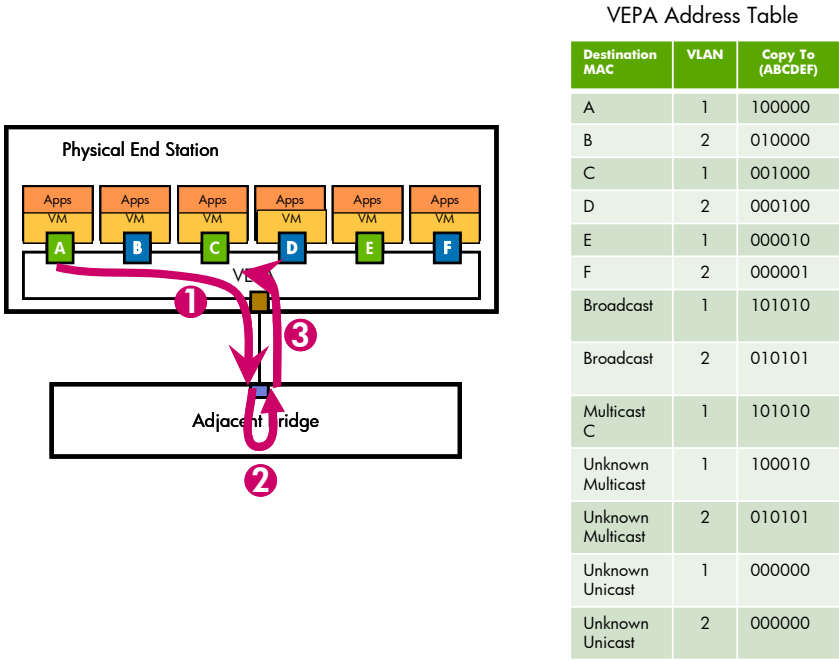


Figure 12. VEPA Unicast Ingress Processing from Source A to Destination C

In this example, the transmission and reception processing is:

1. VSI A performs egress processing and performs any additional functionality prior to the frame being transmitted out the egress port (step 1)
2. The adjacent bridge has enabled VEPA communication. The bridge applies the appropriate network processing to the frame and reflects the frame to the VEPA uplink (step 2).
3. Upon frame receipt and validation, the VEPA searches the address table to locate the destination VSI Instance based on the contents of the unicast frame (minimally the Destination MAC address and the VLAN Identifier are required to uniquely identify a destination). In this example, the "Copy to" mask indicates VSI C is the destination and the VEPA delivers the frame (step 3).
 - a. If the unicast address is unknown, then the "Unknown Unicast" for the associated VLAN identifier would determine the appropriate "Copy to" mask, which is x000000, and the frame is discarded.

The address table access and associated processing is similar for multicast and broadcast with one exception. The originator of a multicast or broadcast frame may have been one of the VSI before the adjacent bridge reflected the frame back. In this case, the VEPA must perform additional filtering to avoid delivering the frame to its originator. This is illustrated in the following figure.

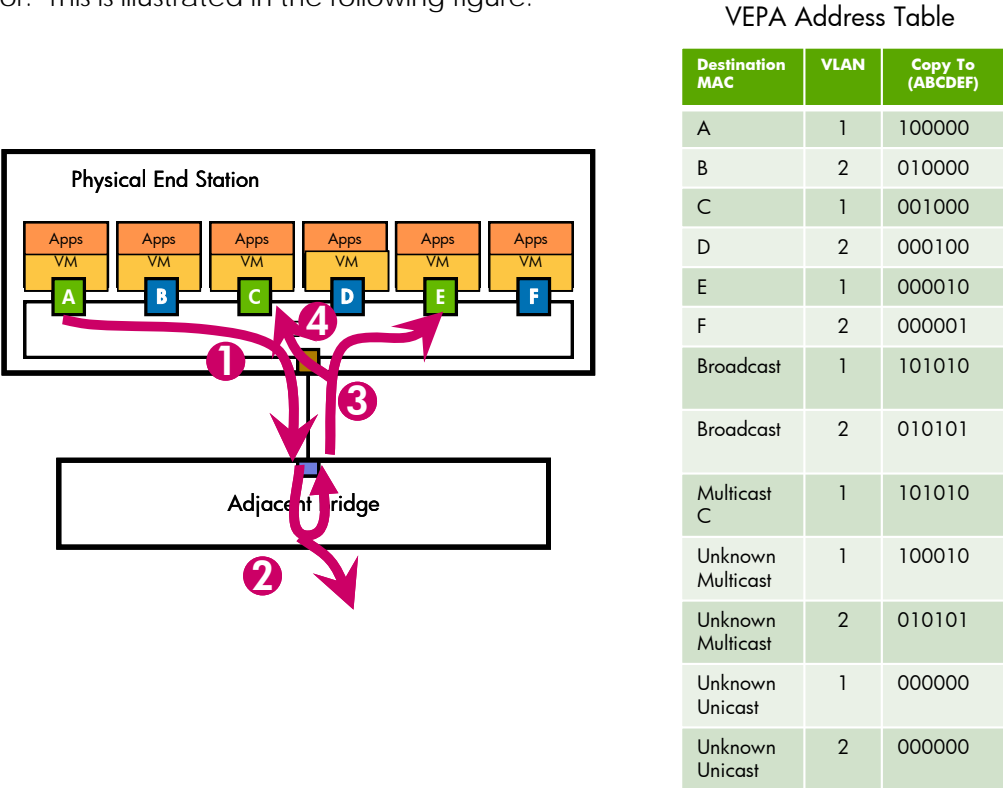


Figure 13. VEPA Multicast Ingress Processing from Source A to Mulicate Group C

In this example, the transmission and receive processing is:

1. VSI A performs egress processing and performs any additional functionality prior to the frame being transmitted out the egress port (step 1)
2. The adjacent bridge has enabled VEPA communication. The bridge applies the appropriate network processing to the frame and reflects the frame to the VEPA uplink (step 2).
3. Upon frame receipt and validation, the VEPA searches the address table to locate the destination VSI based on the contents of the multicast frame (step 3). To prevent the frame from being delivered to its originator, the VEPA performs a source address lookup and filters out the VSI associated with the source address from the original "Copy To" mask associated with the destination address (step 3). The delivery mask is constructed via $(\text{Copy To} = (\text{Destination Copy To}) \text{ AND } \neg(\text{Source Copy To}))$. In this example,

Destination Copy To

Source Copy To

Delivery Mask

= 101010

= 100000

= 001010
4. The frame is replicated using the delivery mask (step 4).

- a. If the unicast address is unknown, then the “Unknown Unicast” for the associated VLAN identifier would determine the appropriate “Copy to” mask, which is x000000, and the frame is discarded.

3.4 S-Channel Operation

This section illustrates S-channel operation through several example configurations. In these examples, an S-VLAN component is logically inserted into the adjacent bridge and the physical end station. Further, between these S-VLAN Components, six S-channels (A-F) have been established and associated with a directly accessible VSI, a VEB, or a VEPA.

The first example illustrates how a directly accessible VSI operates over an S-channel configuration when communicating to a VSI accessible through the adjacent bridge. A directly accessible VSI includes a piece of VEB/VEPA function between the VSI and S-channel. This function is called a 2-port VEB/VEPA. The purpose of this function is to manipulate the C-tagging and untagging of frames passed between the S-channel and the VSI. The tagging configuration of the 2-port VEB/VEPA may be configured by the hypervisor in the same manner as it configures a standard VEB or VEPA.

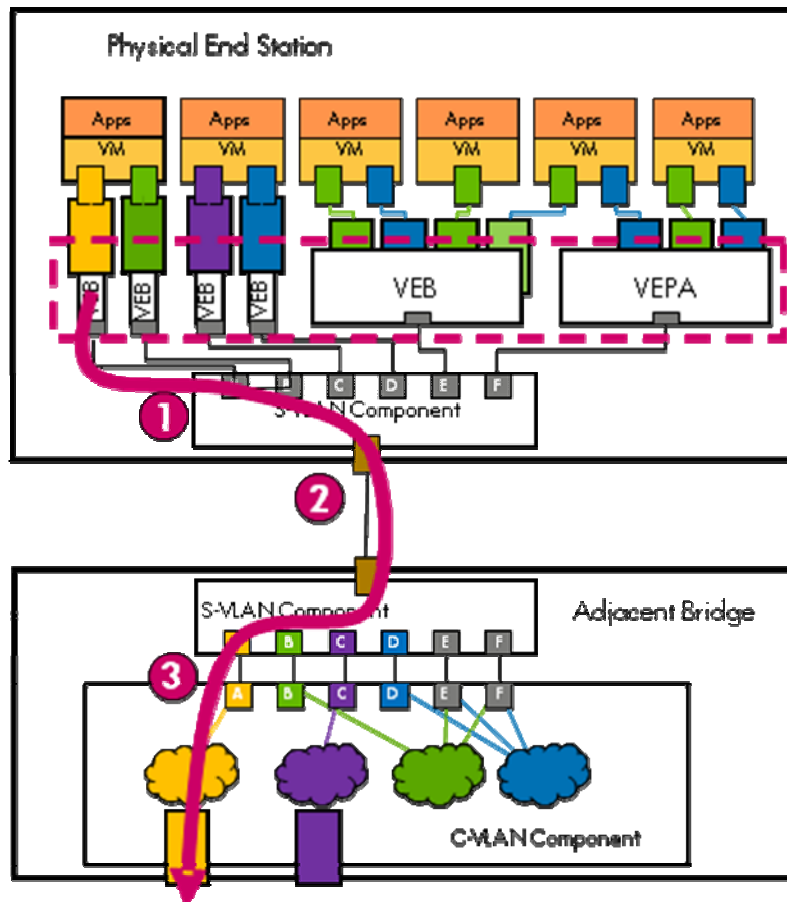


Figure 14. Frame forwarding from a directly accessible VSI over a S-channel LAN

1. VSI performs egress processing and performs any additional functionality prior to the frame being forwarded to the S-VLAN Component within the physical end station. (step 1)
2. If the S-channel A is the default S-channel the S-VLAN Component passes it to the adjacent Bridge un-S-tagged. If the S-channel is not the default then the S-VLAN Component inserts an S-Tag associated with S-channel A into the frame and forwards the frame to the adjacent bridge. (step 2)
3. Within the adjacent bridge, the S-VLAN Component removes the S-Tag and forwards the frame to external destination. (step 3)

This example illustrates how a VEB operates.

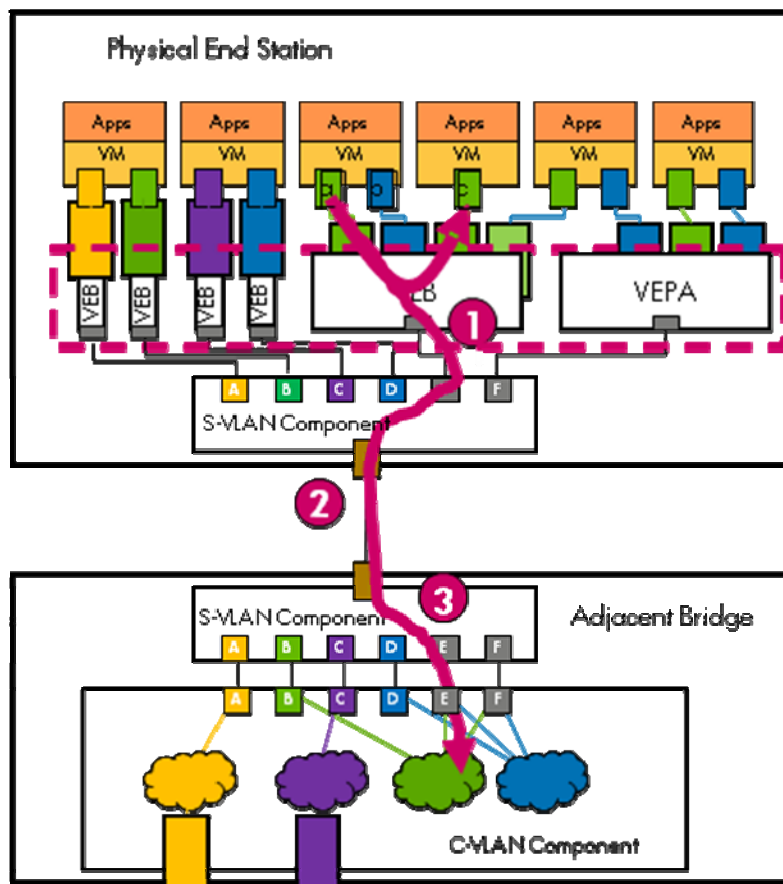


Figure 15. Frame forwarding when S-channels are configured underneath a VEB

1. VM-to-VM communication (a transmitting to c) across a shared VEB does not involve the multi-channel link.
2. The frame forwarding steps to communicate to a VSI not attached to the VEB are identical to the communication used for a directly accessible VSI running over an S-channel.

- a. VSI (a) performs egress processing and performs any additional functionality prior to the frame being forwarded to the S-VLAN Component within the physical end station. (step 1)
- b. The S-VLAN Component inserts an S-Tag associated with S-channel E (or if the S-channel E is the default S-channel no S-tag is inserted) into the frame and forwards the frame to the adjacent bridge. (step 2)
- c. Within the adjacent bridge, the S-VLAN Component removes the S-Tag and forwards the frame to external destination. (step 3)

This example illustrates VM-to-VM unicast communication through a VEPA when S-channels are configured.

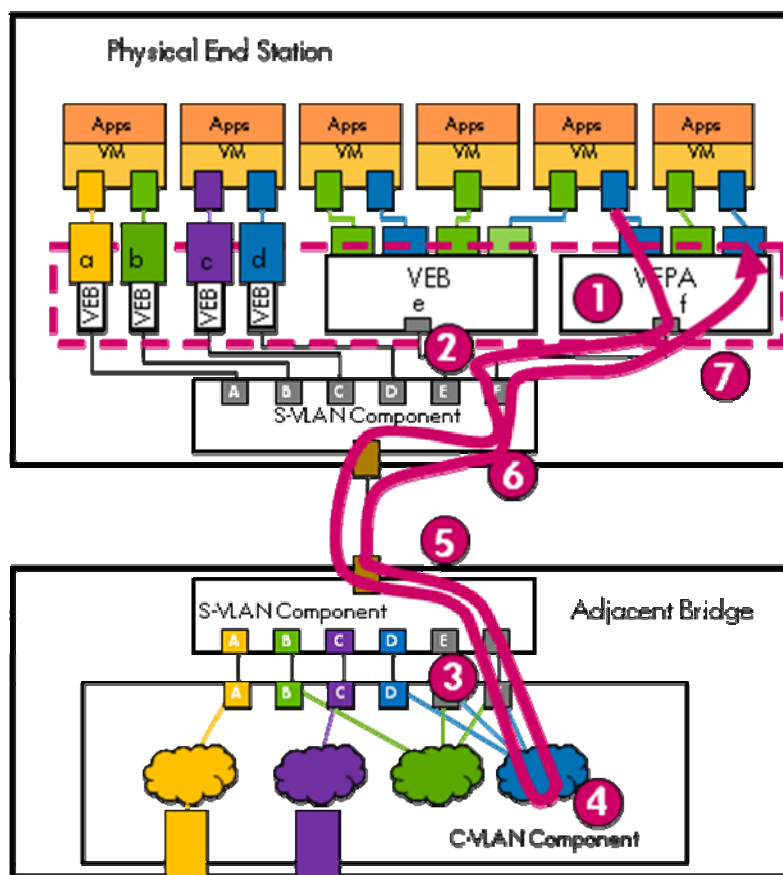


Figure 16. Frame forwarding when S-channel are configured underneath a VEPA

1. VSI performs egress processing and performs any additional functionality prior to the frame being forwarded to the S-VLAN Component within the physical end station. (step 1)
2. The S-VLAN Component inserts an S-Tag associated with S-channel F (or if the S-channel F is the default S-channel no S-tag is inserted) into the frame and forwards the frame to the adjacent bridge. (step 2)
3. Within the adjacent bridge, the S-VLAN Component removes the S-Tag and forwards the frame (step 3).

4. The adjacent bridge determines that the CAP is configured for VEPA mode so it forwards the frame based on the bridge forwarding table (step 4).
5. Within the adjacent bridge, the S-VLAN Component adds the S-Tag associated with S-channel F (or if the S-channel F is the default S-channel no S-tag is inserted) and forwards the frame to the S-VLAN Component within the physical end station (step 5).
6. The S-VLAN Component within the physical end station removes the S-Tag and forwards the frame to the associated VEPA (step 6).
7. The VEPA forwards the frame based on its VEPA address table to the associated VSI (step 7).

The following example illustrates how a VEPA-attached VM communicates to a directly attached VSI through a common physical end station.

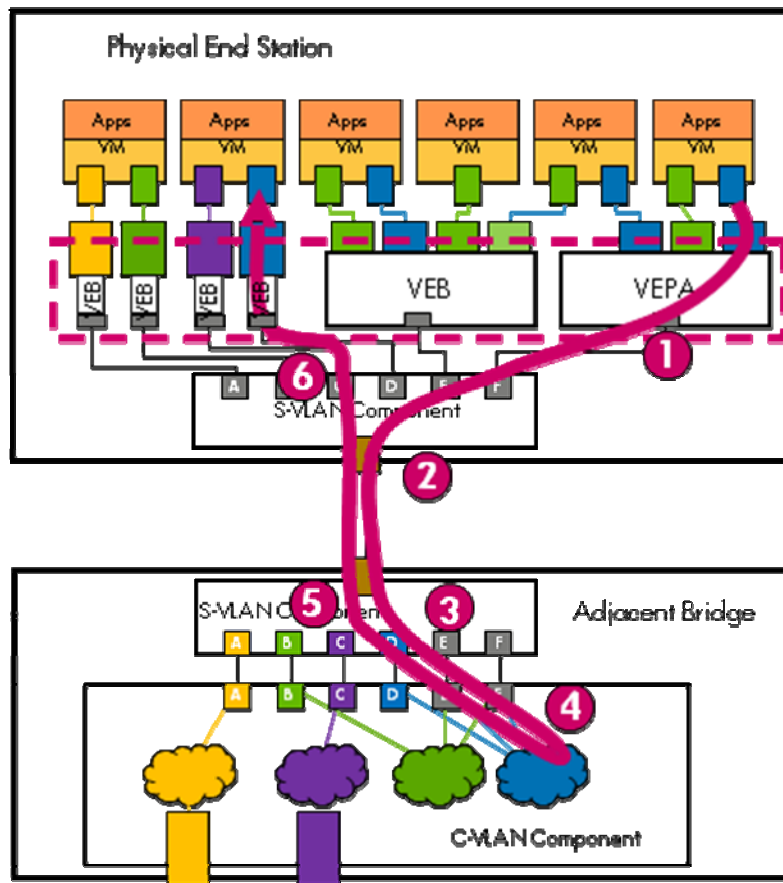


Figure 17. Frame forwarding over S-channel between a VEPA and a directly attached VSI

1. VSI performs egress processing and performs any additional functionality prior to the frame being forwarded to the S-VLAN Component within the physical end station. (step 1)
2. The S-VLAN Component inserts an S-Tag associated with S-channel F (or if the S-channel F is the default S-channel no S-tag is inserted) into the frame and forwards the frame to the adjacent bridge. (step 2)

3. Within the adjacent bridge, the S-VLAN Component removes the S-Tag and forwards the frame (step 3).
4. The adjacent bridge determines that the frame's next hop is associated with S-channel D and forwards the frame to the S-VLAN component.
5. Within the adjacent bridge, the S-VLAN Component adds the S-Tag associated with S-channel D (or if the S-channel D is the default S-channel no S-tag is inserted) and forwards the frame to the S-VLAN Component within the physical end station (step 5).
6. The S-VLAN Component within the physical end station removes the S-Tag and forwards the frame to the directly attached VSI.

3.5 Edge Control Protocol Operation

Today, IEEE control plane discovery operations are performed over unacknowledged protocols, such as LLDP and DCBX. For VSI Discovery (and potentially other ULPs), the Edge Control (ECP) provides acknowledgements, which signal to the sender that the receiver is able to receive an additional ECP Data Unit. ECP enables the sender to transmit discovery operations more frequently than would be the case with timer based approaches. The intent is to have the server's virtualization infrastructure (e.g. Hypervisor) implement ECP, versus having the NIC implement ECP.

The following diagram depicts, at a high level, ECP semantics. In step 1, the ULP passes an outgoing ULP Data Unit to ECP by invoking a transmit request procedure. In step 2, the ULP Data Unit, which for some ULPs (e.g. VSI) may contain a set of ULP TLVs, is transmitted and a ECP low level Acknowledgement (L-ACK in the diagram) timer is set, but the frame is not yet deleted from the transmit buffer until a ECP is received for that ECPDU. In step 3, the arriving ECP frame is received into a receive 'buffer', where it is held until it is removed by an ECP indication procedure that passes the ULP Data Unit to the associated upper level protocol. In step 4, when the receive buffer is emptied, a low-level acknowledge (L-ACK) is sent to the sender. In step 5, if the L-ACK is received before the L-ACK timer expires, then the transmit buffer is cleared and ECP can process another ULP PDU through the ECP procedure. However, if the L-ACK timer expires before the L-ACK is received, then the frame in the transmit buffer is resent (some preset number of times).

ECP is used to support VSI Discovery and Configuration Protocol (VDP) described in section 3.6. To support VDP an instance of ECP is created in the station on each LAN facing Bridge Port of the VEB, VEPA or 2-port VEB. This peers with an ECP peer located on the adjacent Bridge at a Bridge Port of the C-component. The ECP dialog proceeds using the Nearest Customer Bridge LLDP address and is configured using the EVB TLV described in chapter 4.

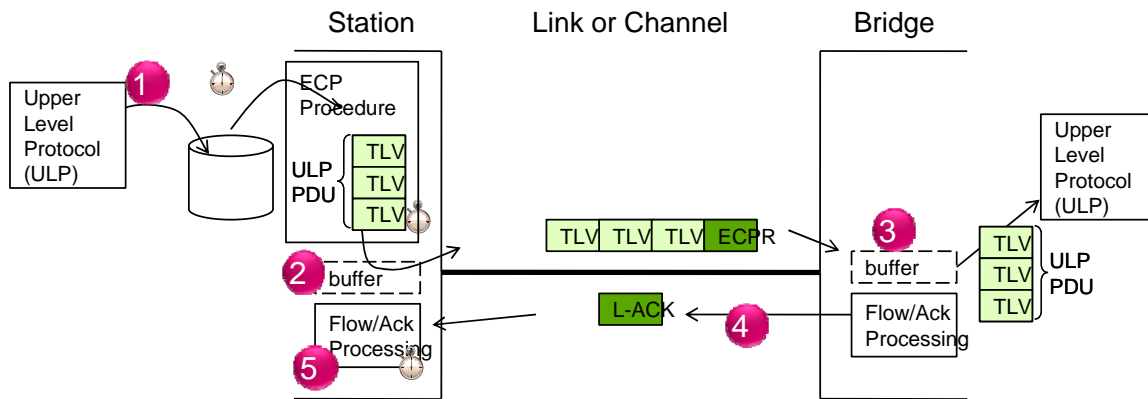


Figure 18. Example ECP Exchange

3.6 VSI Discovery and Configuration Protocol (VDP) Operation

Enterprise and cloud-based networking deployments have been rapidly growing in size leading to a significant increase in the complexity of Ethernet networking in data centers. The advent of virtualization technology brings unprecedented network configuration complexity due to the significant increase in the number of Ethernet switches and very large number of Virtual Station Interfaces (VSIs). The problem is made more complex by advent of Virtual Machine (VM) mobility and solutions requiring external network state to move with the VM, when the VM moves.

Virtual Ethernet Bridges (VEBs) are embedded in hypervisors or NICs have been around for years. VEBs provide efficient VM to VM communications between co-located VM. However, today's virtual switch management is too manual and scale-out server sprawl and virtualization magnifies this complexity. Two of the major challenges associated with today's virtualization approaches is the ability to automate the association of a VSI Instance with it's network state and automate VM migration, including all the network state associated with the VM.

Today, when a VM migrates from one server to another, the internal VEB's VSI Type that is associated with the VM must be migrated. The VSI Type consists of the network state associated with the VM and may include access and QoS Controls. In today's implementations, the external switch's port profiles do not move with the VM. The client has three options for dealing with this issue. Option 1 is to use the same VSI Type for all VMs - the problem with this approach is that it limits virtualization's value, because all VMs in the network must have identical network capabilities and attributes. As a result, if a group of servers doing the same type of work gets over utilized, the VMs from those servers cannot be moved to a group of servers doing another type of work (e.g. file/print).

Option 2 is to move the VSI Type after the VM moves. This can be done by having the external switch look up the VSI Type when the VM starts sending messages on the new server. For example, when the VM starts sending messages, the external switch uses the VM's MAC Address to look-up the port profile. This approach suffers from two problems:

The external switch cannot tell if the MAC Address used by the VM is a migrated MAC address (i.e. from a migrated VM) or a re-incarnated MAC address (i.e. from a new VM that is using a previously destroyed VM's MAC address). The second problem is that there is a VSI Type exposure window between the VM's first message and the time it takes the external switch to obtain the VSI Type from the switch's fabric manager.

Option 3 is to simply configure the link between the server and the edge switch as a trunk port. The issue with this approach is all physical servers must be in the same security domain, which has the similar VM movement limitations as option 1. For example, a physical server cannot be managed by tenant A in the same fabric as a physical server that is managed by tenant B.

The VDP Protocol specified in this document enables the association of a VSI Type with a VSI instance (e.g. a VM virtual port) and the de-association of a VSI Type with a VSI instance (e.g. a VM virtual port). VDP simplifies and automate Virtual Server (VS) network configuration by enabling the movement of the VSI Type when the VSI Instance moves.

3.6.1 VDP Type Configuration and Automation

A virtualized server hosts a set of VMs. Each VM may support one or more Virtual Station Interface (VSI) Instances. Typically, a VM will support a virtual NIC (vNIC) that emulates a physical NIC. Each vNIC will contain a VSI which is connected to a VEB or VEPA. The server's virtualization infrastructure (e.g. a Hypervisor) assigns one or more VSIs to a VM to access the network. The VM is able to communicate with other VMs on the same physical server through the VSI Instance. Similarly, the VM is able to communicate to external stations through the VSI Instance.

Each VSI Instance is assigned VSI Type ID (VTID). The contents of a VSI Type are outside the scope of this proposal. For information context purposes only, a VSI Type definition may include port access or rate limiting controls. Prior to the activation of a VM, VDP exchanges are used to associate a VSI Instance with a VLAN Identifier, a MAC Address and a VTID in the adjacent bridge and, if VEB is used, VEB. Similarly, a VDP exchange is used to de-associate a VSI Instance with a VLAN Identifier, a MAC Address and a VTID in the adjacent bridge and, if VEB is used, VEB, when a VM is either destroyed or moved.

The following sections provide an operational overview of how VDP can be used to automate the configuration of network state (e.g. VSI Type) and the association of network state to a VSI Instance. It will then describe the management elements required to support such an example.

3.6.1.1 VDP – Operational Example

An example of the steps associated with VDP is depicted in the following figure.

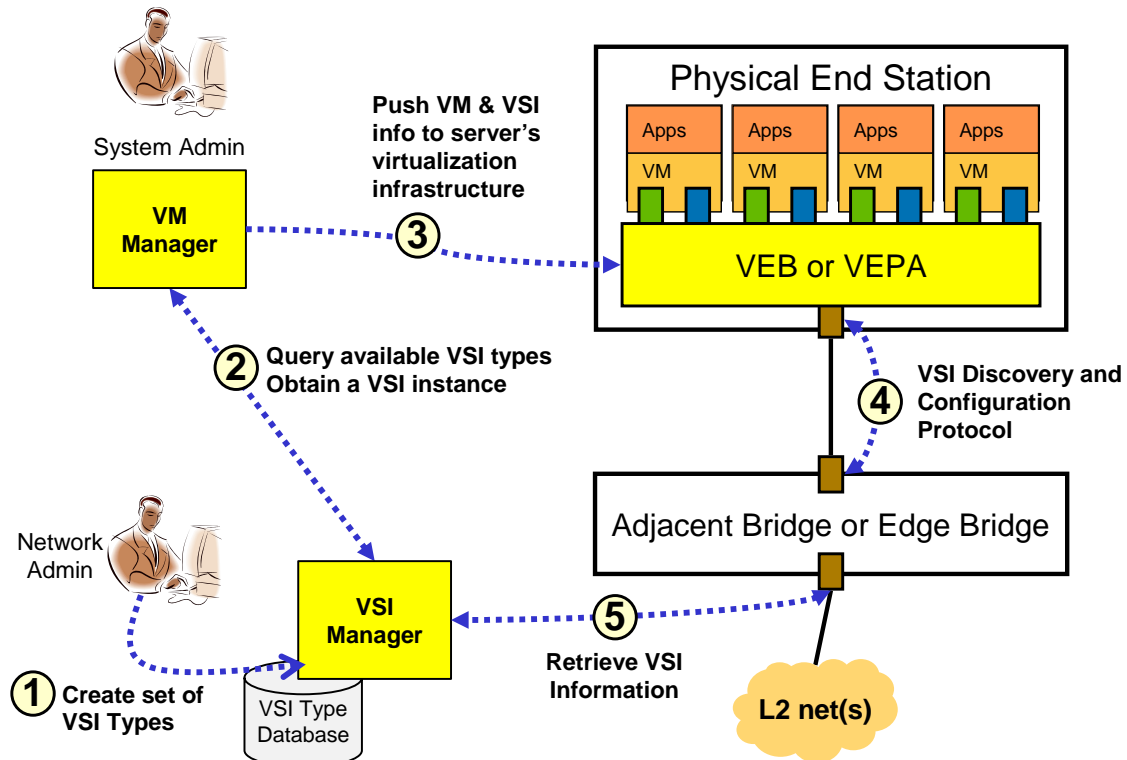


Figure 19. VSI Type Architectural and Operational Overview

Following are the steps depicted in the figure above:

Step 1: Network Manager creates a set of VSI Types. Each VSI Type has a VSI Version and the Network Manager can deploy one or more VSI Versions at any given time.

Step 2: VM Manager queries available Version Type IDs (VTIDs) and creates a VSI Instance consisting of VSI Instance ID and the chosen VTID. The VTDB server may create and track VSI Instance.

Step 3: VM Manager configures VSI with VTID and VSI Instance ID obtained from VSI Manager's VTDB.

Step 4: Before VSI Instance (VM) activation, the VDP Module performs VSI Discovery and Configuration protocol exchanges to associate the VSI instance with a VTID, MAC Address and VLAN Identifier. The VDP Module is intended to be implemented as part of the server's virtualization infrastructure (e.g. in the Hypervisor or a service VM guest running on top of the Hypervisor). The VDP Module is also implemented in the adjacent bridge.

Step 5: As part of the VDP exchange the adjacent Bridge retrieves the VSI Type from the VTDB by using the VTID and possibly the VSI Type Version and VSI Instance ID. The adjacent Bridge stores the association of VLAN ID, VSI Type, VSI Type Version and MAC Address in its local memory. This association is then applied to the traffic flow from/to the VSI Instance. Note the VTDB access protocol is not part of this document.

3.6.1.2 VSI Type Database (VTDB)

The VSI Type Database described above is used to store detailed definition of VSI types. Again these definitions are outside the scope of this document. For information purposes only, a VSI Type may contain access and traffic controls. Also for information purposes only, a VSI Type Database is expected to be part of the database used by the edge switch's Network Change and Configuration Manager and the server's virtualization infrastructure.

VSI Type Definitions within a VTDB are identified by VSI Type ID (VTID) and VTID version. Optionally, VSI instance specific definitions are possible.

The mechanisms used to create VSI Types in a VTDB are outside the scope of this document. For information purposes only, each VSI Type may refer to different use models, such as a server type, where each server type (e.g. web, file/print, e-mail) has a unique VSI Type. Many other use models are possible, for example the Server Virtualization Partitioning and Clustering from the DMTF.

3.6.2 VSI Type Definition and Management

VSI Type Definition and Management is outside the scope of this document. In other words, the content of a VSI Type entry in the VTDB and how that content is managed are outside the scope of this document.

Similarly, VSI Type management and access protocols are outside the scope this proposal. This is not a hindrance to deployment of VDP because current Data Center Network (DCN) infrastructure includes mature tools for management and configuration and can be easily deployed to manage VSI Types. Further, VSI Type Management approach proposed in this document matches well with currently deployed DCN management practices. It is achieved by aligning management and configuration responsibility with current organization structure e.g. VSI Types can be managed by Network Administrator and deployed on servers by server administrators.

3.6.3 VSI Manager ID

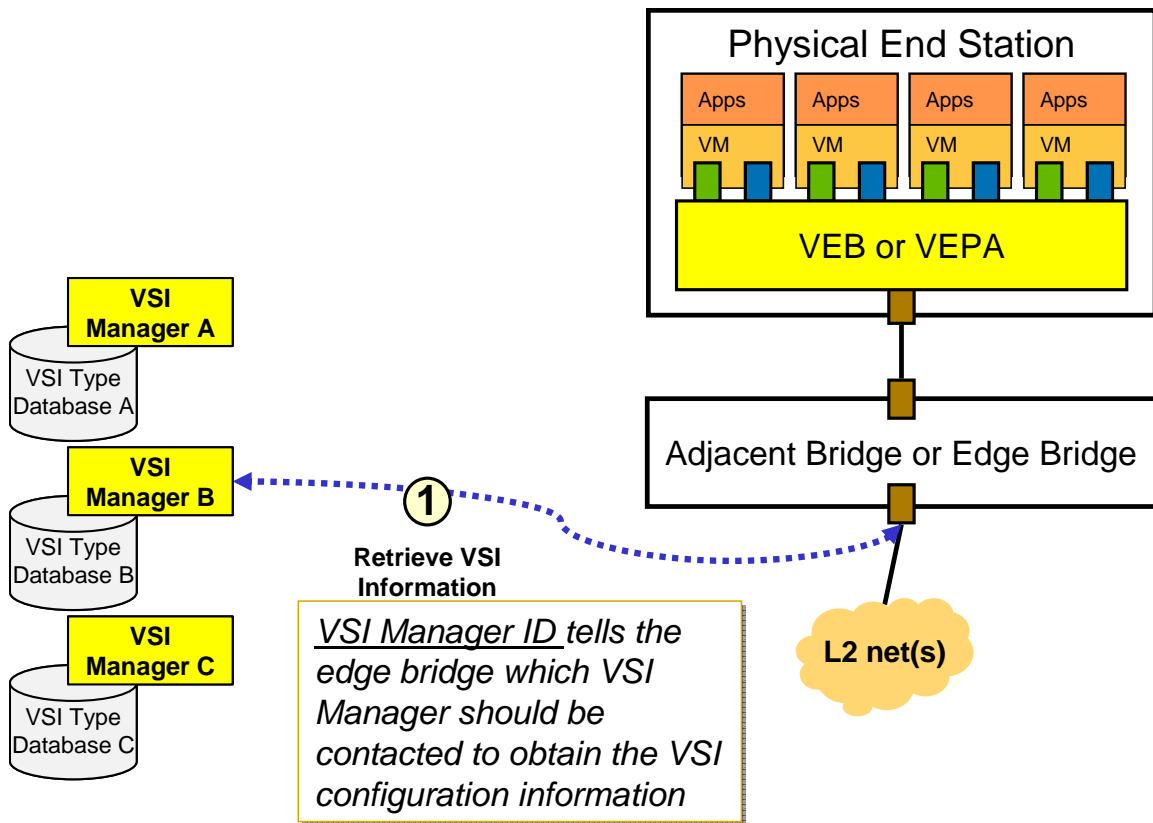


Figure 20. VSI Manager ID

VSI Manager ID tells the edge bridge which VSI Type Manager should be contacted to obtain the VSI configuration information. The VSI Manager ID is part of VDP exchange between Station and the Edge Bridge.

Note, the VSI Type ID or VSI Instance ID can also be used as index to look up VSI Type configuration in VSI Type Database, see the following figure:

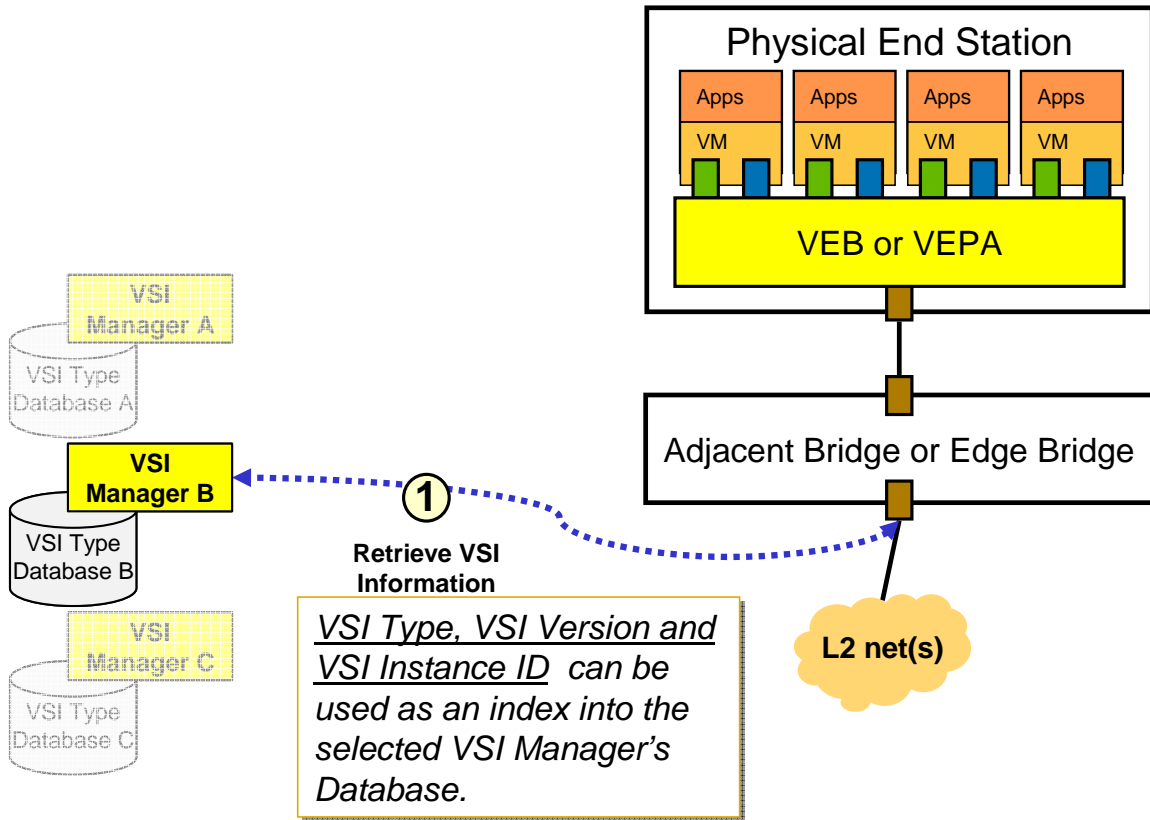


Figure 21. VSI Type , VSI Version and VSI Instance ID

3.6.3.1 VSI Manager ID Usage Example

The VSI Manager ID Identifies the VSI Manager with the Database that holds the detailed VSI Type and/or VSI Instance Identifier definitions. The contents of the VSI Manager Database are outside the scope of this proposal. The VSI Manager Database may use a combination of the following fields to index into the VSI Manager Database:

- VSI Type Identifier
- VSI Type Version
- VSI Instance Identifier

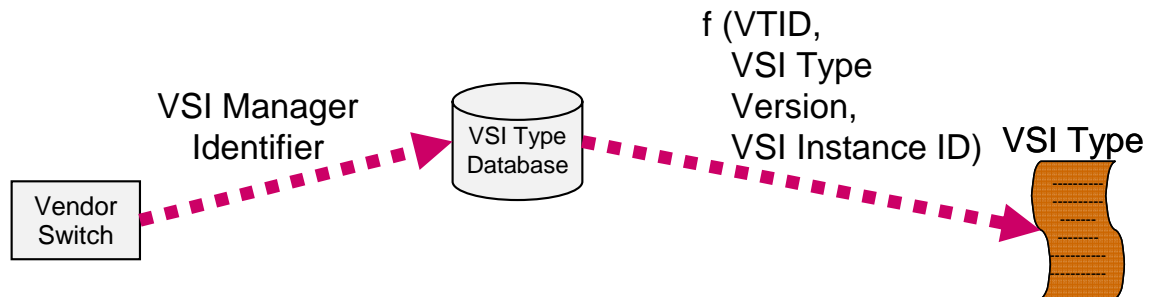


Figure 22. VSI Manager Database Lookup

4. Ethernet Virtual Bridging TLV Semantics

The EVB TLV is used to:

- Advertise a station or bridge's EVB functional and resource capabilities
- Activate common functional capabilities
- Negotiates resource capabilities to a maximum common value, for example the number of VSIs supported.

The EVB TLV is exchanged via LLDP and conforms to the LLDP TLV specification. The LLDP database carrying the EVB TLV is addressed using the Nearest Customer Bridge LLDP address. One LLDP database is built at the Bridge Port of the VEB/VEPA/2-port VEB attached to an S-channel by an internal LAN. When no S-VLAN components are available in the station the LLDP database is built at the Bridge Port of the VEB/VEPA/2-port VEB attached to the external LAN. In this case the VEB/VEPA/2-port VEB is attached to the default S-channel.

The EVB TLV allows setting the adjacent Bridge's C-component Port in reflective relay. Reflective relay is implemented by changing the active topology enforcement rules described in IEEE 802.1Q 8.6.1 to allow forwarding on the reception Bridge Port. When reflective relay is enabled unicast and multicast frames are reflected.

The EVB TLV also allows turning off SA learning and enabling VDP based learning on the adjacent Bridge's C-component Bridge Port. When VDP learning is enabled addresses are not aged. Instead, they exist in the filtering database as long as the VSI is pre-associated or associated.

The EVB TLV is illustrated in the following figure:

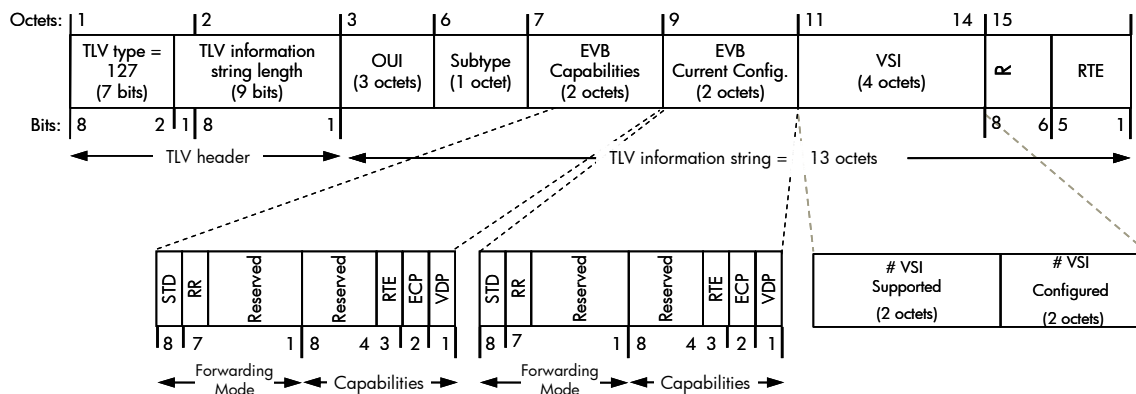


Figure 23. EVB TLV Format

The EVB TLV fields are:

The OUI used to identify the EVB protocols (EVB, CDCP, ECP and VDP) is XX-XX-XX.

The Subtype 0x0000 is used to identify the EVB TLV.

EVB Capabilities - The TLV describes EVB capabilities that are supported by the sender. The capabilities are:

- Forwarding Mode:

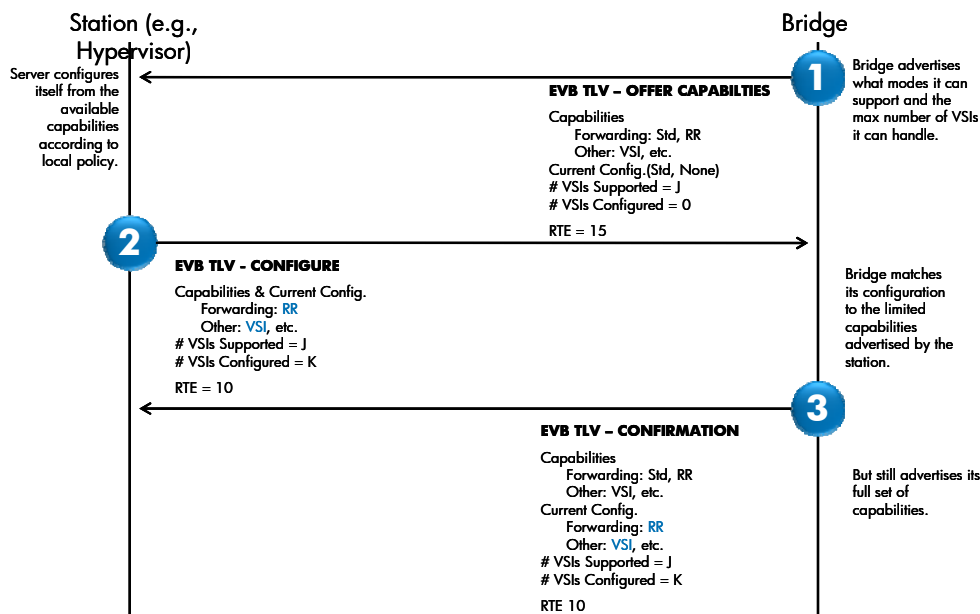
- Standard 802.1Q forwarding
- Reflective Relay – enables frames to be reflected back through the ingress port. For example, in a VEPA solution, frames exchanged between co-located VM must flow through the adjacent bridge. Reflective relay allows these exchanges to flow through a common uplink between the station and the adjacent bridge.
 - From the station, RR = TRUE indicates the station requests reflective relay support.
 - From the adjacent bridge, RR = TRUE indicates the bridge supports reflective relay support.
 - If the station and the adjacent bridge set RR = TRUE, then reflective relay can be enabled. The EVB TLV Current Configuration RR bit is set to TRUE.
 - If either side does not set RR = TRUE, the reflective relay cannot be enabled. The EVB TLV Current Configuration RR bit is set to FALSE.
- Retransmission Timer Exponent (RTE) – Indicates the current RTE value is present
- Edge Control Protocol (ECP) – Indicates the sender supports ECP
 - From the station, ECP = TRUE indicates the station supports ECP.
 - From the adjacent bridge, ECP = TRUE indicates the bridge supports ECP.
 - If the station and the adjacent bridge set ECP = TRUE, then ECP can be enabled. The EVB TLV Current Configuration ECP bit is set to TRUE.
- If either side does not set ECP = TRUE, then ECP cannot be enabled. The EVB TLV Current Configuration ECP bit is set to FALSE.
- VSI Discovery Protocol (VDP) – Indicates the sender supports VDP. VDP is dependent upon ECP being enabled.
 - From the station, VDP = TRUE indicates the station supports VDP.
 - From the adjacent bridge, VDP = TRUE indicates the bridge supports VDP.
 - If the station and the adjacent bridge set VDP = TRUE and ECP == TRUE, then VDP can be enabled. The EVB TLV Current Configuration VDP bit is set to TRUE.
 - If either side sets VDP = FALSE or ECP == FALSE, then VDP cannot be enabled. The EVB TLV Current Configuration VDP bit is set to FALSE.

EVB Current Configuration – The TLV describes the EVB capabilities that are currently configured at the sender. Current configuration represents the intersection of the capabilities and resources between the two senders on a link.

- Number of VSI Supported – The maximum number of VSI that can be supported by the sender.
 - From the station, it indicates the number of VSIs the station is able to support.
 - From the adjacent bridge, it indicates the number of VSI the bridge is able to support.

- Number of VSI configured – The number of VSI that has been configured by the sender.
 - From the station, it indicates the number of VSI the station is requesting the bridge to reserve.
 - From the adjacent bridge, it indicates the number of VSI the bridge has reserved for use by the station.
- Retransmission Exponent (RTE) – RTE is an EVB link or S-channel attribute used to calculate the minimum ULP PDU retransmission time. The ULP PDU retransmission time is calculated as follows:
 - The Retransmission Granularity (RTG) is set to 10 micro-seconds.
 - The Retransmission Multiplier (RTM) is set to 2^{RTE}
 - The sender's ULP transmission timer is set to $RTM * RTG$
 - Both sides agree to the largest common value

The following illustrates an example EVB TLV exchange between a station (e.g. a hypervisor) and the adjacent bridge. This exchange is accomplished using LLDP. In this example, both the station and the bridge support Reflective Relay, VDP, and a set of VSI resources.



Station or Bridge may initiate EVB TLV Exchange
Figure 24. Example EVB TLV Exchange

5. S-Channel TLV Semantics and State Machine

This chapter provides an overview, detailed semantics, and state machines for the S-Channel Discovery and Configuration Protocol (CDCP).

5.1 S-Channel Bridge Components and Operation

5.1.1 Introduction

The purpose of CDCP is to configure S-channels, which are used by a station to simplify the internal configuration and operation of Virtual Station Interfaces (VSIs), Virtual Ethernet Bridges (VEBs) and Virtual Ethernet Port Aggregators (VEPA). S-channels must be used when more than one VEB, VEPA or 2-Port VEB is used simultaneously within the station. S-channels are implemented in stations and bridges using an IEEE 802.1Qbc Port-mapping S-VLAN component.

When no Port-mapping S-Component is in the system all traffic is un-S-tagged and is considered part of the default S-channel (SCID=1). The SVID 1 is reserved for the default S-channel and must not be assigned by the Bridge when a new S-channel is requested. When the Port-mapping S-components used to create S-channels exist they may start exchanging un-S-tagged frames which are assigned to the default S-VID 1 and are considered the default S-channel 1. The default S-channel is always un-S-tagged even when S-channels are enabled.

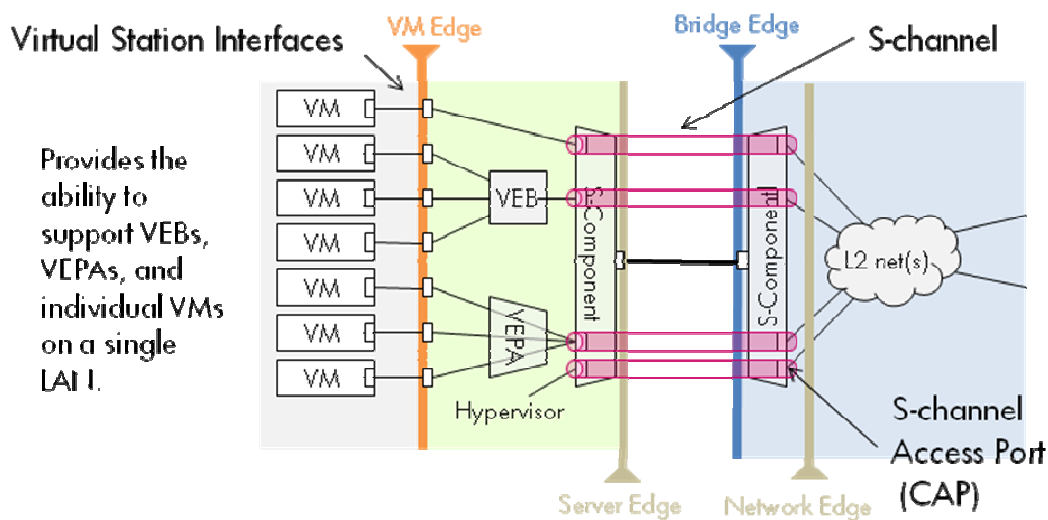


Figure 25. Example S-channel Block Diagram

The S-channel configuration is determined by the bridge's capabilities and by requests made using CDCP described in this chapter. The hypervisor determines what S-channels are needed and then requests S-channels using CDCP. CDCP in turn uses an LLDP TLV exchange to co-ordinate the creation and deletion of S-channels. The LLDP database used by CDCP is addressed using the Nearest non-TPMR Bridge LLDP address. The port-mapping S-VLAN component filters both the Nearest Bridge and the Nearest non-TPMR Bridge addresses on all ports and passes the Nearest Customer Bridge LLDP address. Only

the exterior facing Station-Bridge Access Port will build an LLDP database. This database is built from frames exchanged on the Nearest non-TPMR Bridge address.

The SCID=1 and SVID=1 are always reserved for the exclusive use as the un-S-tagged default S-channel. The default S-channel number 1 with S-VID 1 is created unilaterally by the Station and the Bridge, without the use of CDCP. CDCP reports the default S-channel in the CDCP TLV as the first SCID,SVID pair (i.e. <1,1>). The Bridge must not assign this S-VID except to the default S-channel. If the Bridge deletes the default S-channel from its configured S-channels list it means the Bridge will not process the un-S-tagged set as an S-channel (though control information may still be passed as un-S-tagged traffic).

5.1.2 S-Component

The figure below is a "baggy pants" Bridge relay architecture model for the station and Bridge. The S-Component conforms to the Port-mapping S-VLAN component specified in 802.1Qbc. The S-Component is used to create S-channels. The C-Component of the Bridge is a standard Bridge C-Component relay with the exception of additions for the reflective relay feature and support for EVB and VSI discovery and configuration protocols.

Not all the represented components need to be present in an implementation. If the S-Components depicted are present and CDCP is present the system will be able to create and delete S-channels. If no CDCP is present then the S-Components may be present, however disabled or may not be present at all. It is also possible that one or both of the S-Components will be absent. It is possible to have a VEB or VEPA without an S-Component.

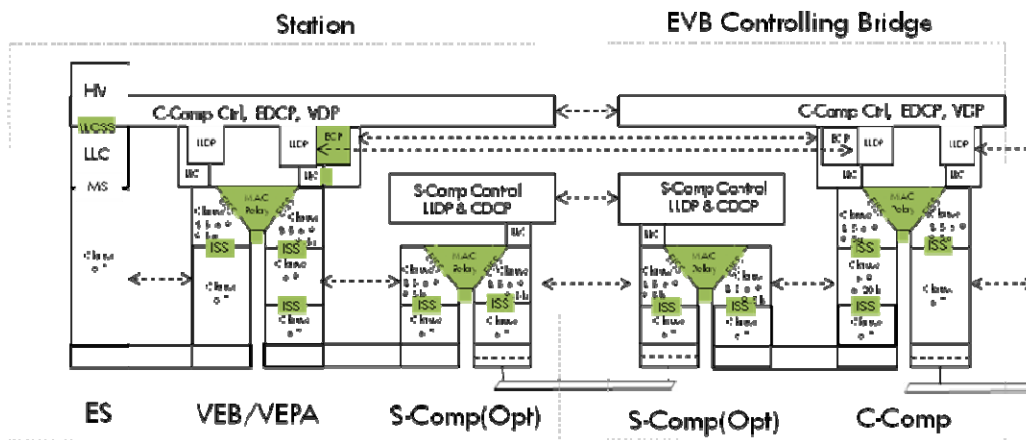


Figure 26. Station and Bridge VEB/VEPA and S-comp Diagram

S-channels are implemented using the S-Components. The C-VLANs carried over each S-channel are determined by configuration of the C-Component within the Bridge (under the control of EVB TLV and VDP). Each S-channel terminates at an internal S-Comp Bridge Port called an S-channel Access Port (CAP). Internal LANs, within the Bridge and station, each span between one CAP and one internal C-Comp Bridge Port (Bridge) or VEB/VEPA/2-port VEB (station).

The C-VLAN configuration and reflective relay configuration of the Bridge is determined by the configuration of the C-VLAN aware component of the Bridge.

Each S-channel terminates at an internal S-Comp Bridge Port called an S-channel Access Port (CAP). Internal LANs, within the Bridge and station, each span between one CAP and one internal C-Comp Bridge Port (Bridge) or VEB/VEPA/2-port VEB (station). Each CAP must be configured as follows:

1. Admit Only Untagged and Priority-tagged frames (802.1Q 6.9)
2. a PVID parameter equal to the S-VID associated with its S-channel (802.1Q 6.9)
3. must be a member of the untagged set (802.1Q 8.8.2) for the VLAN identified by the S-channel's S-VID
4. must not be a member of any S-VLAN member set except the one identified by the SVID associated with the CAP's S-channel (each CAP supports a single S-channel)
5. each CAP filters the Nearest Bridge and Nearest non-TPMR Bridge LLDP addresses

Each S-Comp has external facing Bridge Ports called Station-Bridge Access Ports (SBAP). All S-channels are multiplexed over these ports. Each SBAP must be configured as follows:

1. Admit All frames (802.1Q 6.9)
2. a PVID parameter equal to the default S-channel SVID of 1 (802.1Q 6.9)
3. the SBAP must be a member of the untagged set (802.1Q 8.8.2) for the S-VLAN identified by the default S-channel SVID 1
4. the SBAP must be a member of the tagged set (802.1Q 8.8.2) for all S-VLANs identified by the SVID of an active S-channels
5. each SBAP filters the Nearest Bridge and Nearest non-TPMR Bridge LLDP addresses

5.2 CDCP Discovery and Configuration

S-channels are configured by the exchange of LLDP TLV using the Nearest non-TPMR bridge address. The exchange begins when the system is initialized. The configuration protocol begins with the station, which makes a request for S-channel resources from the Bridge. In response the Bridge provides the best matching set of S-VLANs it is capable of providing. It is possible the Bridge does not have all the resources requested in which case the Bridge response will provide a subset of the requested S-VLANs.

After initialization it is possible for the station to change it's S-channel configuration. The Bridge seeing a change in the stations request will alter it's configuration to match the needs of the station.

5.2.1 CDCP TLV

The station and Bridge both use the same LLDP TLV to configure S-channels. This TLV is in LLDP OUI format (802.1AB sub-clause 8.6). The S-channel's capabilities, requests and running configuration is encoded in the info field of this TLV as follows:

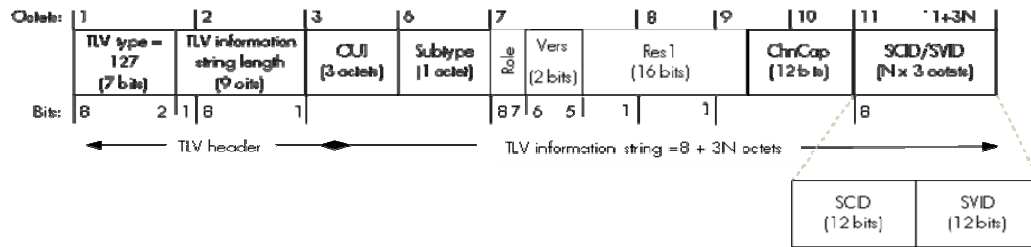


Figure 27. CDCP TLV

- The OUI used to identify the EVB protocols (EVB, CDCP, ECP and VDP) is XX-XX-XX.
- The Subtype 0x0001 is used to identify the CDCP TLV.
- **Role Bits** (see note 1 regarding ties)
 - S(01b) – Indicates the sender assigns channels numbers and a default SVID for the default channel 1 and requests SVID assignments from the neighboring 'B'.
 - B(10b) – Indicates the sender accepts S-channel configuration requests from its neighboring 'S' and that the sender will do the best it can to fill the SVID assignment requests from the neighboring 'S'.
- **Vers** – 2 bits: 10b identifies this version, 00b disables Sch (incompatibility)
- **Res1**- 16 bits: must be set to zero, ignored on receipt (for new features)
- **ChnCap**- S-channels supported identifies the number of S-channels that are supportable by the sender.
- **SCID/SVID Pairs**
 - **SCID** - indicates the index number of the S-channels. The station assigns S-channel numbers in the range 0-167. Zero is reserved. The S-channel index should be between 1 and the maximum number of S-channels supported by the port.
 - **SVID** – The VLAN ID assigned to the S-channel. The Bridge assigns SVIDs to channels in the range 1-0xfe. A station uses the 0 SVID to request an SVID assignment from the Bridge.

Note1: The first entry in the list of SCID/SVID pairs must contain the default S-channel. (i.e. if the default S-channel is being used then the first channel pair must be <1,1>).

Note2: A maximum of 167 S-channels can be supported. Other formats (assuming sequential SVIDs) could be defined to allow support for 4K+ S-channels.

Note3: This listing could be sparse (in order to indicate arrival and removal of S-channels). The S-channel going away is indicated by removing the SCID/SVID pair.

Note4: The order of the list will determine the priority of SVID assignments. If the Bridge does not have resources for all channels it will assign the first channels in the list.

5.2.2 CDCP Configuration Procedures

The CDCP protocol used to discover and configure S-channels begins by announcing the presence of CDCP along with the station and Bridge capabilities (1). After the initial announcement the Bridge will look for a request from the station (1). Once the Bridge sees a station request it will configure itself with and provide the best matching configuration to the station (2). The station seeing that the Bridge is now configured goes operative using the Bridge's configuration (3).

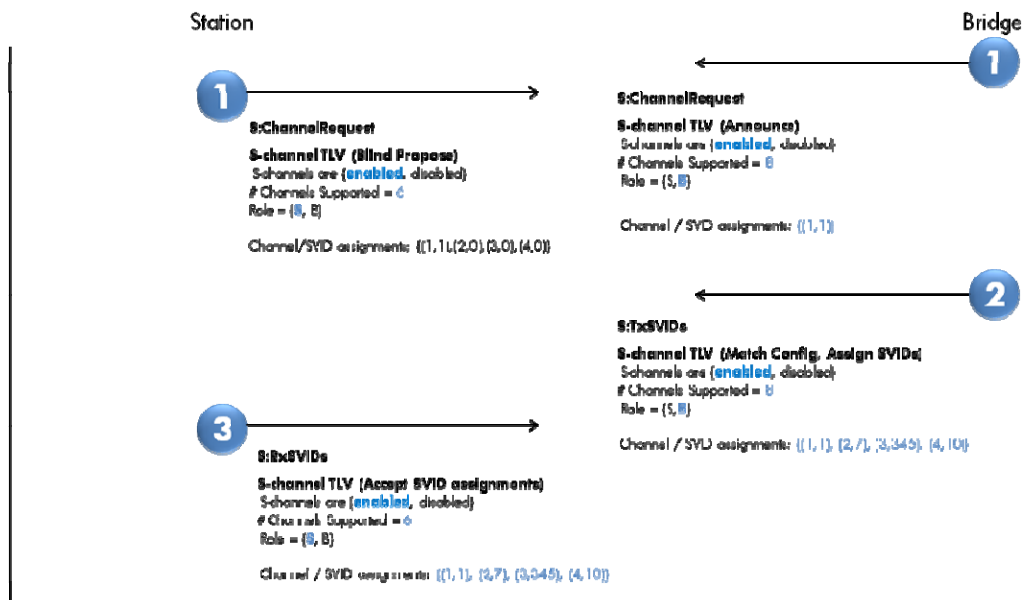


Figure 28. Example CDCP TLV Exchange

5.2.3 CDCP Configuration Variables

The following variables are used by the CDCP state machine to perform S-channel configuration. The CDCP requires each side of the configuration be assigned a role as a Bridge or a Station. This is done by setting the AdminRole variable. In most pieces of equipment the station or bridge role will not be settable, though the protocol allows for equipment which can take either role. For CDCP to configure an S-channel one side must take the station role and one side must take the Bridge role. If both sides of the LAN have equipment configured as stations or as bridges the protocol will not configure S-channels.

- **AdminRole:** Is the administratively configured value for the local port's role parameter. The value of AdminRole is not reflected in the S-channel TLV. The AdminRole may take the value S or B. S indicates the sender is unwilling to accept S-channels configuration (mode, # channels supported, channel index) from its neighbor and that the sender is willing to accept SVID assignments from the neighbor. Stations usually take the S role. B indicates the sender is willing to accept S-channels configuration (mode, # channels supported, channel index)

- from its neighbor and that the sender is willing do the best it can to fill the SVID assignments from the neighbor. Bridges usually take the B role.
- **OperRole:** The current operational value of the Role parameter in the local port. This value is included as the Role parameter in the S-channel TLV and may take values S or B as described for AdminRole.
 - **RemoteRole:** Indicates the value in the remote S-channel TLV role field. `rwNull` indicates either the TLV was not present in the last LLDP PDU or that no LLDP PDUs have been received. `rwS` and `rwB` indicate that the Role field was set in the S-channel TLV received and that it had a value of S or B respectively as described for the AdminRole variable.
 - **schState:** The current running state of CDCP. The values for this variable are `NOTRUNNING` or `RUNNING`.
 - **AdminVersion:** The administratively configured value for the Vers parameter. This value is included as the Vers parameter in the S-channel TLV. If the value is `DISABLE = 00b` it means S-channels are disabled. If the value is `VER0 = 10b` it means this version.
 - **AdminChnCap:** The administratively configured value for the Number of Channels supported parameter. This value is included as the ChnCap parameter in the S-channel TLV.
 - **AdminSVIDWants:** The administratively configured value for (SCID, SVID) pairs wanted by an S. Not used by a B. The first value is always the pair (1, 1) for the default S-channel assignment. The S-channel numbers may be any valid number from 0-0xffe. A 0 S-channel number indicates reserved space in the TLV. If the SVID value is 0 it means the S is requesting any available SVID. SVID value 1 is reserved for exclusive use for the default S-channel SVID. . The AdminSVIDWants parameter is used to form the (SCID, SVID) pairs in the S-channel TLV.
 - **LastSVIDWants:** A local temporary copy of the AdminSVIDWants.
 - **LocalSVIDPool:** The set of SVIDs and bridge ports available for S-channel assignment. These are determined by both administrative resource assignments and by resource availability. The OperSVIDList for a B role must be drawn from the LocalSVIDPool.
 - **LastLocalSVIDPool:** A temporary copy of the LocalSVIDPool.
 - **OperVersion:** The current value for the Vers parameter. This value is included as the Vers parameter in the local S-channel TLV. The value `VER0 = 10b` means this version. The value `DISABLE = 00b` mean don't run CDCP.
 - **OperChnCap:** The current value for the ChnCap parameter. This value is included as the ChnCap parameter in the local S-channel TLV. The range for this variable is 1-0xffe.
 - **OperSVIDList:** The current value for (SCID, SVID) assignments. This is the list of (SCID, SVID) pairs included in the local S-channel TLV. The total size of the list may not exceed 167 pairs. The list must always include the default S-channel pair (1,1).The valid range for each S-channel of this list is from 1-0xffe. The valid range for each SVID in the list is from 0 to 0xffe. For the S role a SVID of 0 indicates a request for a channel. For the B role an SVID of 0 indicates a non-configured channel.
 - **RemoteVersion:** The current value for the remote S-channel Vers parameter. This value is included as the Vers parameter in the remote S-channel TLV. `NULL` means no remote S-channel TLV exists in the local LLDP database. The value for this

variable may be VER0=100b setting any other value will result in stopping S-channel operation.

- RemoteChnCap: The current value for the ChnCap parameter. This value is included as the ChnCap parameter in the remote S-channel TLV. NULL means no remote S-channel TLV exists in the local LLDP database. The range for this variable is 1-0xffe.
- RemoteSVIDList: The current value for (SCID, SVID) assignments. This is the list of (SCID, SVID) pairs included in the remote S-channel TLV. NULL means no remote S-channel TLV exists in the local LLDP database. If the list is empty but the S-channel TLV is present its value is NONE. The total size of the list may not exceed 167 pairs. The valid range for each S-channel of this list is from 1-0xffe. The valid range for each SVID in the list is from 0 to 0xffff. When the SVID is value is 0 the SVID is not configured. For the S role a SVID of 0 indicates a request for a channel. For the B role an SVID of 0 indicates a non-configured channel.
- LastRemoteSVIDList: Temporary local copy of the RemoteSVIDList. This variable is not included in the S-channel TLV. The LastRemoteSVIDList has the same syntax as RemoteSVIDList.

5.2.4 CDCP Configuration Procedures

The CDCP state machine uses three procedures. The SetSVIDRequest() procedure places a new request from the station or sets the initial TLV for a Bridge. The RxSVIDConfig() procedure is used by the station to configure a new set of S-channel and SVID assignments. The TxSVIDConfig() is used by the Bridge to respond to the station's request for S-channels.

- SetSVIDRequest(OperRole, AdminSVIDWants, OperSVIDList)
 - This function creates the OperSVIDList placed in the Local TLV database.
 - If the OperRole for the equipment is B then the OperSVIDList remains unchanged.
 - If the OperRole for the equipment is S two possible cases exist. In the first case we don't have any configured S-channels, indicated by OperSVIDList being equal to NONE. In this case the function places the AdminSVIDWants in OperSVIDList. In the second case we already have a running configuration indicated by the OperSVIDList not equal to NONE. In this case the function compares the AdminSVIDWants with the OperSVIDList. All active S-channels in the OperSVIDList which are in the AdminSVIDWants are kept active and in addition any channels not currently in the OperSVIDList are requested by including them in the OperSVIDList along with a 0 SVID number. The OperSVIDList S-channel order is set to match the AdminSVIDWants.
- RxSVIDConfig (OperSVIDs, LastRemoteVIDList)
 - This function creates the OperSVIDList placed in the Local TLV database for an S role port.
- The function compares the AdminSVIDWants with the LastRemoteSVIDList. For each AdminSVIDWants S-channel with an SVID assignment in the LastRemoteSVIDList a (SCID, SVID) pair is generated in the OperSVIDList. For each AdminSVIDWants S-channel without an SVID assignment in the LastRemoteSVIDList a (SCID,0) pair is generated in the OperSVIDList. The OperSVIDList S-channel order

is set to match the AdminSVIDWants.TxSVIDConfigB(OperChnCap, RemoteChnCap, LastLocalSVIDPool, RemoteSVIDList, OperSVIDList)

- This function creates the OperSVIDList placed in the Local TLV database for an S role equipment
- First the function takes the smaller of the OperChnCap and RemoteChnCap and truncates the RemoteSVIDList to the smaller of the two.
- A new OperSVIDList is created as follows:
 - For each S-channel in the RemoteSVIDList with a (SCID, SVID) pair in the OperSVIDList the (SCID, SVID) remains unchanged unless the SVID is no longer part of the LastLocalSVIDPool. If the SVID is no long in the pool a new one is selected if available. If no SVID is available the (SCID, SVID) pair will be deleted from the OperSVIDList.
 - For each S-channel in the RemoteSVIDList without a (SCID, SVID) pair in the OperSVIDList an SVID is obtained from the LastLocalSVIDPool (the pool for Bridge resources) if available. If no SVID is available the (SCID, SVID) pair will be deleted from the OperSVIDList.

5.2.5 CDCP Configuration State Machines

The CDCP state machine operates on TLV exchanged using LLDP operating on the nearest non_TMPR bridge address (figure 29).

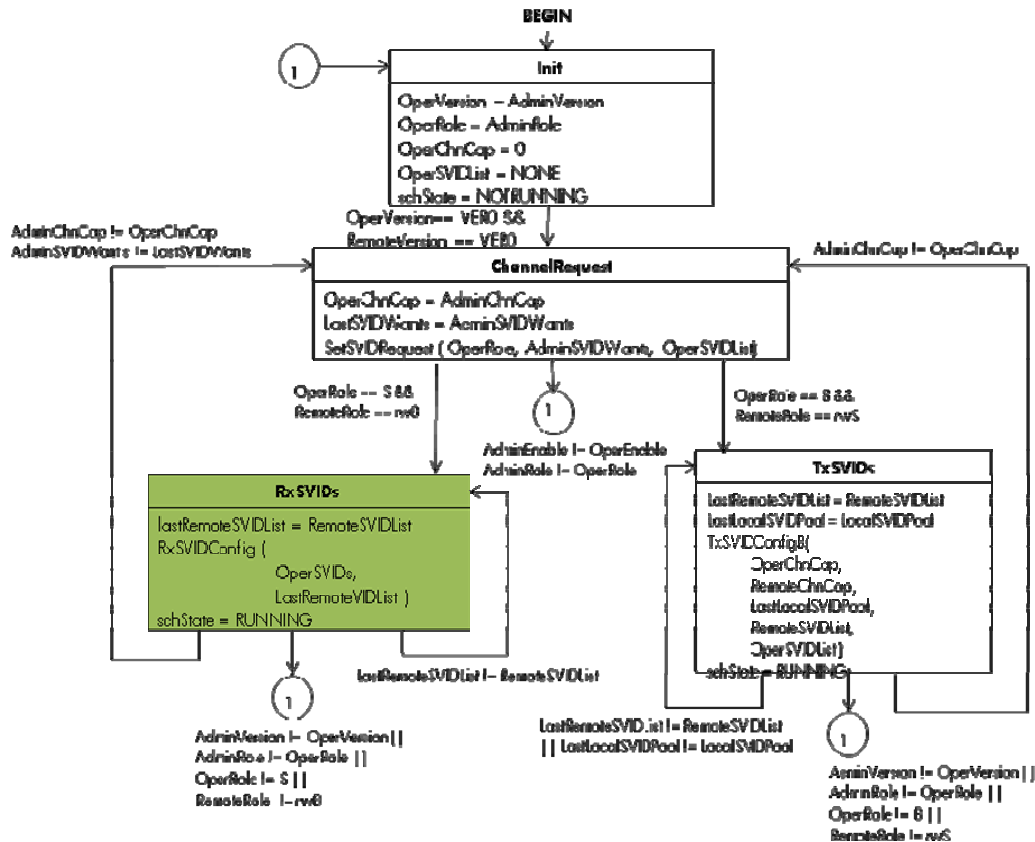


Figure 29. CDCP State Machine

This LLDP instance is one per physical LAN associated with the Station/Bridge Access Port of the S-Component, which faces the LAN connecting the station to the bridge. If either the station or Bridge is not capable of S-channel operation no S-channel-TLV will be inserted in the LLDP database. The absence of a S-channel-TLV therefore indicates that the station or Bridge is only capable of creating the default S-channel. For CDCP to progress both sides must indicate they are capable of dynamic S-channels operation, have the same version number and one side must have the 'B' role while the other side must have the 'S' role as indicated by the role bits of the S-channel TLV.

If both sides are dynamic S-channel capable, exactly one side has the 'S' role and one side has the 'B' role, and the 'B' has at least some of the resources requested by the 'S' side, the state machine will configure S-channels. The configuration proceeds by the 'B' providing the best match it can to the 'S's requested channels and configuration. The 'S' makes the resource request, the 'B' responds with its best matching resources, the 'S' then goes operational and reports its running configuration to the 'B', and finally the 'B' goes operational with the running configuration of the 'S'.

In the event the 'S' wishes to change its configuration it alters the request in its S-channel TLV and then follows the same process as above. If the 'B' loses its ability to support the current configuration it can alter the current configuration in its S-channel TLV at which time the 'S' must drop down to the resources supplied by the 'B'.

In the event of a change of the administered parameters the current operating S-VLANs must be terminated the configuration machine re-initialized.

6. Edge Control Protocol (ECP) and State Machine

This chapter will describe an architectural overview of the Edge Control Protocol (ECP) protocol, followed by the semantics and associated state machines.

6.1 Requirements

ECP was designed with the following protocol requirements:

- Semantics associated with the <ULP, ECP> interface:
 - A single LAN multiplexed using S-channels has one ECP per S-channel.
 - For VSI, there is one VSI agent per S-channel ECP and that agent may have multiple VSI sharing a single S-channel.
 - The <ULP, ECP> interface is based on a complete ULDPDU (i.e. the group of TLVs that are handed to ECP for transmission).
 - The number of octets in the ULDPDU may be less than the maximum number of octets that can fit into a ECP frame.
 - The number of ULP TLVs may be less than the maximum number that can fit into a ECP frame.
 - Procedures are used to describe how the ULP hands off ULP PDUs to ECP and how ECP hands off ULP PDUs to the ULP.
 - Given the <ULP, ECP> interface is based on full ULPDUs, no immediate processing is needed at the ECP level.
 - Outside the scope of this proposal are the semantics for handling: link down; and how multiple ULPs arbitrate when sharing the same ECP.
- ECP acknowledgement and sequencing semantics.
 - A ECP acknowledge means the ECPDU was received and there is a free buffer available to enable another send.
 - It doesn't mean the ECPDU were delivered to the ULP.
 - At the transmit side, if an ECP Acknowledge is not received within an Ack timer period, ECP will retransmit the ECPDU up to 2 times.
 - The Station must set the EVB Maximum Retry Count to 3.
 - Once the receive side ECP delivers the ULP PDU to the ULP through the receive side hand-off procedure, the ECP buffer becomes available for another send.
 - The Acknowledgement must be sent in a separate ECPDU (vs piggy backing onto a Transmit message in the opposite direction).
 - The receive side will issue a ECP Acknowledgement after the completing the receive side hand-off procedure.
 - If the receive side hand-off procedure takes too long, the receive side ECP may toss the ULP PDU and send back an ACK to indicate the ECP buffer is free on the receive side.
 - Note: The length of time the send side waits, before tossing the ULDPDU should not be less than the retransmission period times the maximum number of retries.
 - Semantics associated with slow ULP Data Unit reception (e.g. raising a flag) are outside the scope of this proposal.

- Sequence numbering must be used to detect duplicate vs new ECPDUs.
- ECP will not provide a keep-alive mechanism. Instead each ULP must do so.
- ECP will not provide a digest at the ECP level and any ULP Data Unit (or TLV) database synchronization is left up to the ULP.

Note, the intent is to have the server's virtualization infrastructure (e.g. Hypervisor) implement ECP, versus having the NIC implement ECP.

6.2 Edge Control Protocol Data Unit

This section specifies the format of a ECP Data Unit, along with the header that is added to and removed from ECP frames by the ECP function. The ECP header allows each ECP Data Unit from the sender to be identified through a sequence number, which the receiver acknowledges by sending a ECP Acknowledgement frame.

Ethertype = TBD	Sub-type	Mode	Sequence Number	ULPDU
← 2 Octets →	← 2 Octets →	← 1 Octet →	← 2 Octets →	Optional

Figure 30. ECP Data Unit

The destination address of the Ethernet frame that contains a ECPDU has the following semantics:

- Nearest bridge (01-80-C2-00-00-0E) for ECP running at the link layer.
- Nearest Customer Bridge (01-80-C2-00-00-00) for ECP running over a S-channel.
- Note, ECP should also be allowed using a Uni-cast address.

The source address shall be the sending station or port individual MAC address.

A new Ethertype will be needed for ECP. A ECP exchange will run at the link if the link is not configured for S-channels. If ECP is performed over an S-channel, then the STAG for the S-channel shall precede the ECPDU.

The ECPDU contains:

- Sub-type - a 2 octet field that defines the ULP type included in the PDU. Note for Ack's the sub-type is ignored at the station.
 - The version of VDP described in this document uses a sub-type of 0x0000.
 - The sub-type 0xffff is reserved for organizationally unique use of ECP. In this case the first 3 octets of the ULPDU shall contain the OUI defining the organization using ECP.
- Mode – Identifies whether the operation is a:
 - ECP request (0x00)
 - ECP acknowledgement (0x01).
- Sequence number – identifies the sequential order of the PDU, with respect to other ECPDUs. The starting sequence number may start anywhere for the first ECPDU, but the sequence number for each subsequent new ECPDU is incremented by 1.

6.3 ECP Procedures

Two procedures are used to hand-off Data Units between the ULP and ECP: ECP_UNITDATA.request and ECP_UNITDATA.indication. The implementation of these two procedures is outside the scope of this proposal. These <ULP, ECP> interface procedures may be implemented in many ways, including a queue. Also, the system must have a way of associating the ulptype with a specific ULP.

The ECP_UNITDATA.request is invoked by the ULP at the sender to notify ECP that a ULDPDU is ready to be transmitted. The ulpdu parameter is a unit of work from the ULP. For example, for VSI it consists of a set of VSI TLVs passed from the VSI ULP to ECP for transmission, where the set of TLVs must be less than or equal to the maximum allowed ECPDU. Following is the format for the ECP_UNITDATA.request procedure:

ECP_UNITDATA.request (ulptype, ulpdu)

The ECP_UNITDATA.indication is invoked by ECP at the receiver to indicate a ULDPDU has been successfully received and is available ULP processing. The ulpdu parameter is unit of work from the ULP. For example, for VSI it consists of a set of VSI TLVs passed from the ULP to ECP for transmission, where the set of TLVs must be less than or equal to the maximum allowed ECPDU. Following is the format for the ECP_UNITDATA.indication procedure:

ECP_UNITDATA.indication (ulptype, ulpdu)

6.4 ECP State Machines

There are two state machines used by each ECP instance: transmit and receive. The transmit state machine is invoked through the ECP_UNITDATA.request procedure. The receive state machine is invoked upon reception of an ECP Data Unit and it invokes the ECP_UNITDATA.indication procedure.

6.4.1 ECP Transmit State Machine

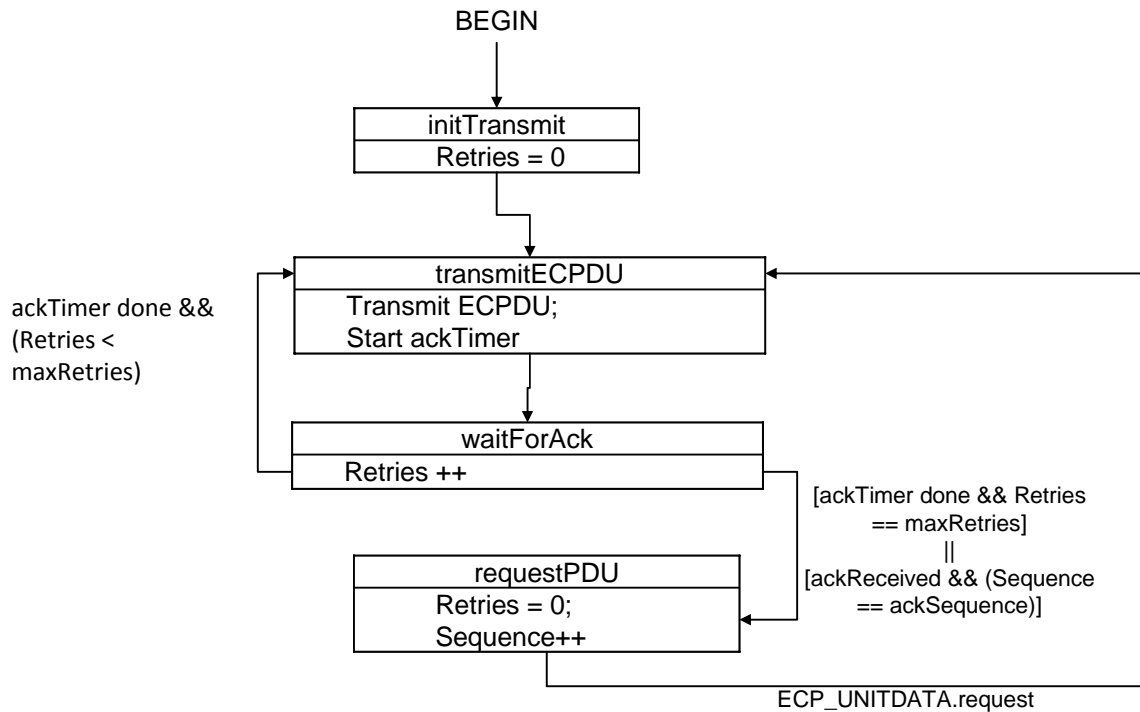


Figure 31. ECP Transmit State Machine

The first entrance into transmitECPDU is used to initiate the sequence counting on the receive side. That is, an ECP Frame that simply contains the ECP header is sent and an ackTimer is started. The waitForAck state waits for the L-ACK to be received that matches the last transmitted ECP sequence number. If an L-ACK is received that matches the last transmitted ECP sequence number or the number of retries exceeds the maximum number of retries, the sender will stop transmitting the ECP Frame and proceed to requestPDU. If an L-ACK is not received within an ackTimer period and the number of retries is less than the maximum number of retries, the sender will retransmit the ECP Frame. The requestPDU state increments the sequence count and waits for the ECP_UNITDATA.request procedure to be invoked.

Note, the starting sequence number may start anywhere for the first ECPDU. Also a Link Down event may restart the sequence number at the same point every time or not.

6.4.2 ECP Receive State Machine

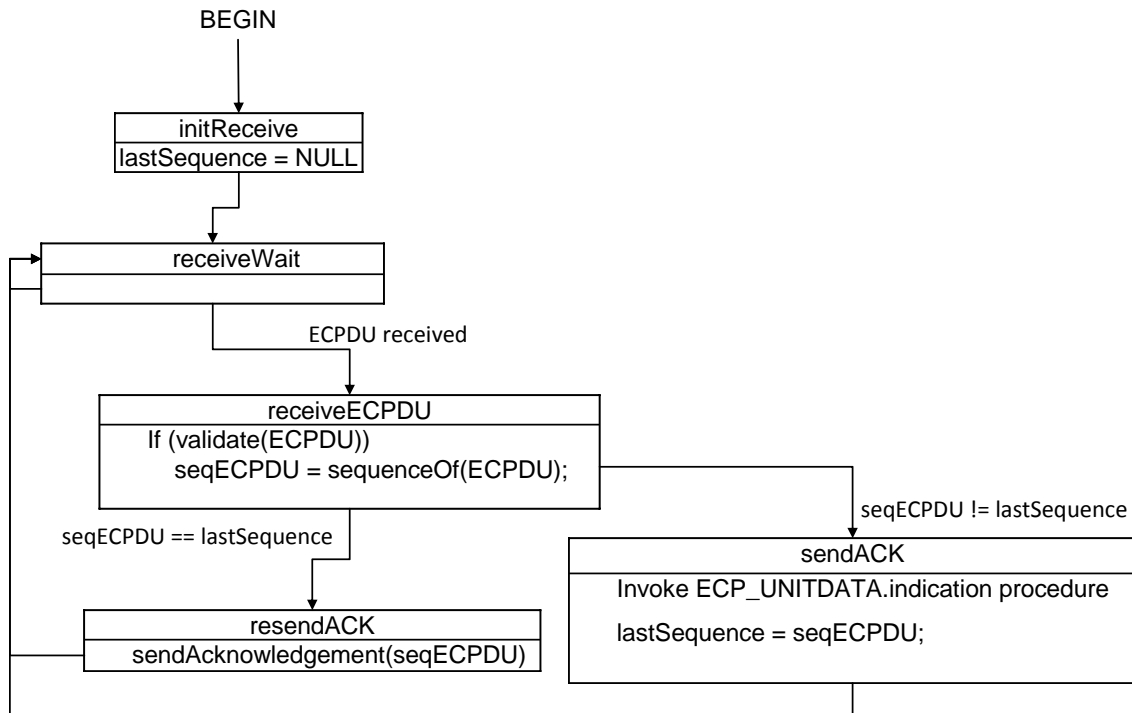


Figure 32. ECP Receive State Machine

The first entrance into InitReceive is used to set the sequence counting to NULL and then proceed to receiveWait, which waits for an ECP Data Unit to be received. The receiveECPDU validates the ECPDU and sets the current sequence number to the sequence number of the transmitted ECPDU. If the current sequence number doesn't match the last transmitted ECP sequence number, then in sendACK the ECP Data Unit is delivered to the ULP and the lastSequence number is set to the current sequence number. If the current sequence number doesn't match the last transmitted ECP sequence number, then in resentACK an L-ACK is sent.

7. Virtual Station Interface (VSI) Discovery and Configuration (VDP) and State Machine

This section covers the Virtual Station Interface (VSI) Discovery and Configuration Protocol (VDP) and State Machine. VDP uses ECP (Edge Control Protocol) for VDP exchanges.

7.1.1 VSI Discovery and Configuration TLV

VSI TLV is used for discovery and configuration and is exchanged between the Station and Bridge. One or more VSI TLVs are transported in an ECP Data Unit. Following is the format and semantics for a VSI TLV:

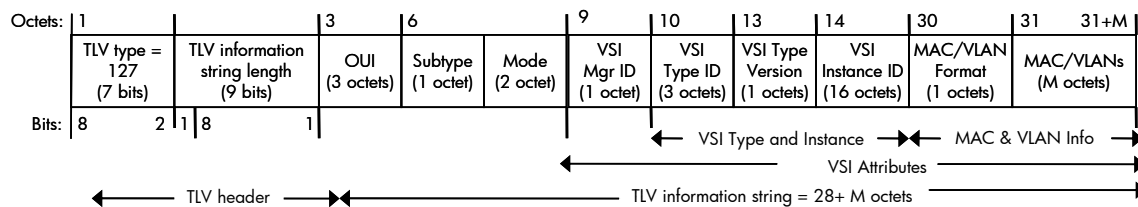


Figure 33. VDP TLV

The OUI used to identify the EVB protocols (EVB, CDCP, ECP and VDP) is XX-XX-XX.

The Subtype 0x0002 is used to identify the VDP TLV.

Mode – Indicates VSI TLV Mode

- First octet identifies a pre-associate, associate, de-associate, or the corresponding confirmation or rejection for each.
- Second octet is used during a rejection to indicate the reason for the pre-assoc or assoc rejection.

VSI Manager ID – Identifies the VSI Manager with the Database that holds the detailed VSI type and or instance definitions. VSI Manager ID can be used to obtain IP address and/or other connectivity and access information for the manager.

VSI Type ID (VTID) – The integer identifier of the VSI Type.

VSI Type ID Version – The integer identifier designating the expected/desired version of the VTID

VSI Instance ID – A globally unique ID for the connection instance. The ID shall be done consistent with IETF RFC 4122.

Format – identifies the format of the MAC and VLAN information that follows in the TLV. Note, the VSI TLV allows multiple formats, which makes possible extensions in the future.

MAC/VLANs – Listing of the MAC/VLANs associated with the Virtual Station Instance (VSI).

Following is the format for Format = 1

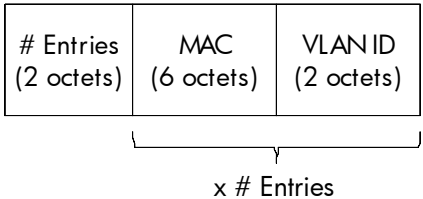


Figure 34. VDP Format = 1 Schema

Note, the station and switch environments and their common understanding of the VTID meaning is outside the scope of this TLV. Also, the contents of a VSI Type are outside the scope of this proposal.

7.1.1.1 VSI TLV – Mode and Mode Response

The purpose of the Mode field is to identify the type of VSI TLV. It is defined as follows.

VSI TLV Request field: 1st octet

Pre-Associate:	0x00
Pre-Associate with resource reservation:	0x01
Associate:	0x02
De-Associate:	0x03

VSI TLV Response field: 2nd octet

For all the responses, the bridge reflects the same VSI TLV fields as the Requester had sent. On requests, response field is initialized to 0x00 (Success). Following are the possible values of the response field.

Success:	0x00	The VSI Request was successfully completed by the switch
Invalid Format:	0x01	The VSI Format is not supported by the switch
Insufficient Resources:	0x02	The switch does not have enough resources to complete the VSI operation successfully.
Unused VTID:	0x03	The VSI referenced by the VSIID does not exist in the VSI Manager database referenced by the VSI Manager Identifier
VTID Violation:	0x04	The VSI referenced by the VSIID is not allowed to be associated with the VTID.
VTID Version Violation:	0x05	The VSI referenced by the VSIID is not allowed to be associated with the VTID Version.
Out of Sync:	0x06	The VTID or one of the VSI List fields used in the Associate is not the same as the corresponding field used in the Pre-Associate.

Reserved

0x08 – 0xFF

These Responses are reserved for future use.

Mode and Mode Response fields are used under the control of VDP state machines.

7.1.1.2 VSI TLV Mode and Responses Semantics

Following are the semantics association with each VSI TLV Request.

7.1.1.2.1 Pre-Associate

The Pre-Associate is used to pre-associate a VSI Instance Identifier to a VSI Type ID. If required, the bridge should obtain VSI Type Definition from the VSI Manager Database. The bridge must validate the request (see below) and fail it in case of errors (see below for responses). Successful Pre-Association does not enable any traffic from VSI. Note that VSI may still be associated at another station. The Pre-Associate enables faster response to an Associate, by allowing the bridge to obtain VSI Type state, prior to an association.

The second Mode octet is used by the bridge to communicate the results of the Pre-Associate requested for the VSI Instance ID (VSIID).

Following are the mode and responses with their semantics:

- Success - Pre-Associate was successful. The switch shall permit a subsequent Associate or De-Associate by the VSI referenced by the VSI Instance Identifier.
- The following are all unsuccessful Pre-Associate Completions. For each of these, the switch shall not permit a subsequent Associate or De-Associate by the VSI referenced by the VSIID.
 - Invalid Format.
 - Insufficient PT Resources.
 - Unused VTID
 - VTID Violation
 - VTID Version Violation

Pre-Associate requires resource lease timer mechanism to conserve Bridge resources. Pre-Associate does not allow any traffic from VSI which is enabled when the VSI is Associated.

7.1.1.2.2 Pre-Associate with Resource Reservation

Pre-Associate with Resource Reservation has same steps as Pre-Associate but also reserves resources.

Bridge should validate required resources and shall reserve resources for subsequent Associate step. Pre-Associate requires resource lease timer mechanism to conserve Bridge resources. Pre-Associate does not allow any traffic from VSI which is enabled when the VSI is Associated.

Second Mode octet contains the results of the Pre-Associate requested for the VSI Instance ID (VSIID). Following are the mode and responses with their semantics.

- Success - Pre-Associate with Resource Reservation was successful. The switch shall permit a subsequent Associate or De-Associate by the VSI referenced by the VSI Instance Identifier.
- The following are all unsuccessful Pre-Associate with Resource Reservation Completions. For each of these, the switch shall not permit a subsequent Associate or De-Associate by the VSI referenced by the VSIIID.
 - Invalid Format.
 - Insufficient PT Resources.
 - Unused VTID
 - VTID Violation
 - VTID Version Violation

7.1.1.2.3 Associate

Associates the VSI Instance ID with the VSI Type ID (VTID). If VSI Type definition is not already cached in the bridge, the bridge fetches the VSI Type definition from the VSI Type definition Database. Bridge allocates required bridge resources for the referenced VSI. The Bridge binds specific MAC/VLAN pairs with the VSI Type ID which allows classification of L2 traffic to the VSI and enforcing of VSI Type controls. Bridge activates the configuration for the VSI Type ID. This association is then applied to the traffic flow from/to the VSI Instance.

For a given VSI Instance ID, a Station may issue an Associate without having previously issued a Pre-Associate or Pre-Associate with Resource Reservation. During normal operations a VSI Instance is Associated on only one port. During network transitions (e.g. VM migration) a VSI Instance might be Associated with more than one port.

In VSI TLV, second octet in the mode field contains the results of the Associate request performed for the VSI Instance Identifier. These are described below.

- Success - Associate was successful. Prior to issuing this response, for a format 1 VSI TLV, the bridge shall associate the VSI Type referenced by the VSI Type Identifier and VSI Type Version with the MAC Address, VLAN and VSIIID.
- The following are all unsuccessful Associate Completions.
 - Invalid Format
 - Insufficient Resources - If the Associate was preceded by a successful Pre-Associate with Resource Reservation, then the bridge shall not issue this response.
 - VTID Violation
 - VTID Version Violation
 - Out of Sync

7.1.1.2.4 De-Associate

De-associate a VSI Instance Identifier from the associated VTID. Pre-Associated and Associated VSIs can be De-Associated. De-Associate releases resources and de-activates the configuration associated with the VSI instance. A VSI Instance may get De-Associated by bridge due to bridge error situation or management action.

In VSI TLV, second octet in the mode field contains the results of the De-Associate request performed for the VSI Instance Identifier. These are described below.

- Success - De-Associate was successful. Prior to issuing this response, for a format 1 VSI TLV, the bridge shall de-associate the VSI Type referenced by the VSI Type Identifier from the the MAC Address, VLAN and VSI Instance ID.
- The following are all unsuccessful De-Associate Completions.
 - Invalid Format
 - VTID Violation
 - VTID Version Violation

Note: The result of the above semantics is that De-Associate can be issued at any time.

7.1.1.2.5 VSI Type ID (VTID) Semantics

VSI Type ID (VTID) is an integer value field used to identify a pre-configured set of controls/attributes that are to be associated with a set of VSIs.

VTID contents and meaning and the database used to contain the VSI Type are outside the scope of this effort. One VTID may describe the VSI Type configuration of multiple VSIs. The VSI Type content referenced by the same VTID may differ between switches and VEBs. For example: same VTID is used by switches from two different vendors; or same VTID is used by a VEB and vendor switches.

7.1.1.3 VSI Type ID Version Semantics

VTID Version is integer identifier designating the expected/desired VTID version. The VTID Version enables a VSI Manager Database to contain multiple VSI Type versions. It allows smooth migration to newer VSI types.

7.1.1.4 VSI Instance ID

VSI Instance ID is a globally unique ID for the VSI instance. The ID shall be done consistent with IETF RFC 4122. VSI ID is generated when VSI instance is created by VSI Instance Manager at request of VM Manager. VSI Instance creation mechanism is outside scope of this proposal but expected to be created by VM Manager or VSI Manager.

7.1.1.5 MAC – VLAN Information Format

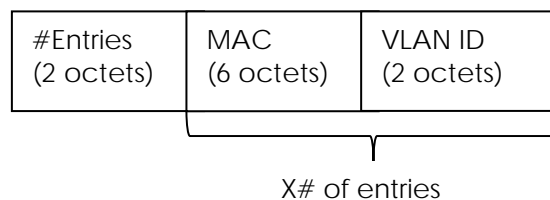


Figure 35. MAC-VLAN Information Format 1

MAC-VLAN Format-1 contains the set of MAC Addresses and VLANs to be associated with the VSI Instance ID. Note the bridge uses MAC+VID to identify traffic from VSI and to steer the frames.

Field:		
	#MAC-VLAN pairs:	2 octets
Per MAC-VLAN Pair Content:		
	MAC address:	48 bits
	VID:	12 bits

7.1.2 VDP Requirements and Assumptions

Following are VDP requirements associated met by VDP state machines described in this section:

1. VDP must support a VSI Pre-Associate (with and without resource reservations), Associate and De-Associate.
2. Associate, Pre-Associate and De-Associate are Idempotent i.e. can be repeated.
3. The bridge must allow for an Associate to be issued without the need for a previous Pre-Associate.
4. VDP may be used in conjunction with both a VEPA and VEB.
5. VDP utilizes ECP as the transport for a VDP Data Unit that contains one or more VDP TLVs. VDP utilizes the following capabilities of ECP:
 1. Transport will be transmitting TLVs in-order and are received in-order.
 2. Flow control
 3. ECP provides best effort delivery of TLV. At the Station, if a VDP Acknowledgement is not received, within an Acknowledgement timeout period, VSI exits the state machine. The Acknowledgement timeout period is defined as $2 \times \text{ECP retransmission period} \times \text{Maximum number of retries}$, plus a locally administered wait that is outside the scope of this document.
6. Health TLV mechanism to ensure:
 1. Bridge resources are not reserved for too long a time period for inactive VSIs (lease semantics)
 2. Allow removing resources from inactive VSIs with the goal of
 1. Conserving bridges resources (Number VSIs being handled by bridge can be large).
 2. Prevent inactive or VMs in error state to continue to hold resources.
 3. For S-channels, timeout out values to be negotiated on a per S-channel basis between station and bridge. One timeout used for all ULPs on the S-channel negotiated using EVB TLV.
 4. If S-channels are not enabled, timeout out values to be negotiated per link basis between station and bridge. One timeout used for all ULPs on the link negotiated using EVB TLV.
7. Ensure VSI state and configuration between the Station and the Bridge remains consistent.
8. Hard errors at the Bridge or the Hypervisor that can impact individual VSI or Hypervisor/Bridge as a whole are handled by removing all VSI configuration.

9. Bridge and Station Errors are detected through one or more of the following mechanisms.
 1. VSI KEEP-ALIVE (periodic transmission of VSI TLV from station and response from Bridge)
 2. ACK Timer
10. The value of the Station's Activity Timer must not be greater than three times ECP Max Retries * EVB's RTM * RTG. That is, Activity Timer $\leq 9 * \text{RTM} * \text{RTG}$.
11. Supports for switch/hypervisor administrator actions that force VSI De-Associate.
12. Should enable statistics and logging capability.

7.1.3 VDP – Local Variables and Procedures

vsiState:	Local variable for current state.
localTLV:	Current local (active) TLV (configuration)
AdminTLV:	TLV from local administration. In addition appropriate localChange variable is set. It allows mode change
RemoteTLV:	TLV received from remote.
TxTLV(<i>vsiTLV</i>):	Transmits AdminTLV using TLV transport service (ECP) service interfaces.
ProcRxAndSetCfg(<i>vsiRemoteTLV</i> , <i>vsiLocalTLV</i> , <i>vsiState</i>):	Processes receive TLV and Sets local TLV variable based on Received Remote TLV and vsiState. In case of error, returns error. This function handles PreAssociate with and without resource reservation case as well as accessing VSI Type definition and fetch, if required.
StartACKtimer():	Resets ACKTimeout local variable to FALSE and Starts ACK timer. Response (ACK or NACK is expected before timer expires.
ACKTimeout:	This local variable is set to true, if ACK timer expires
vsiErrorPerm(<i>vsiRemoteTLV</i>):	Processes the vsiRemoteTLV and returns TRUE if response code is an unrecoverable (permanent) error.

The next sections contain the VSI State Machine. Following are notes regarding those state machines:

1. The purpose of the ACKtimer is to catch the unusual case of a TLV getting lost. The following architectural minimum shall be used: The Acknowledgement timeout period is defined as $2 * \text{ECP retransmission period} * \text{Maximum number of retries}$, plus a locally administered wait that is outside the scope of this document.
2. For any VSI ACK received for a non-active VSI the station shall drop the packet.
3. VSI State is set to NULL on exit.
4. The VSI State Machine does will not implement retry mechanism on NACK. Instead the ULP can process the NACK reasons and retry the VSI operation.

5. VDP state machine will exit on receiving NACK.

7.1.4 Station VSI State Machine

Following is the VSI State Machine for the Station.

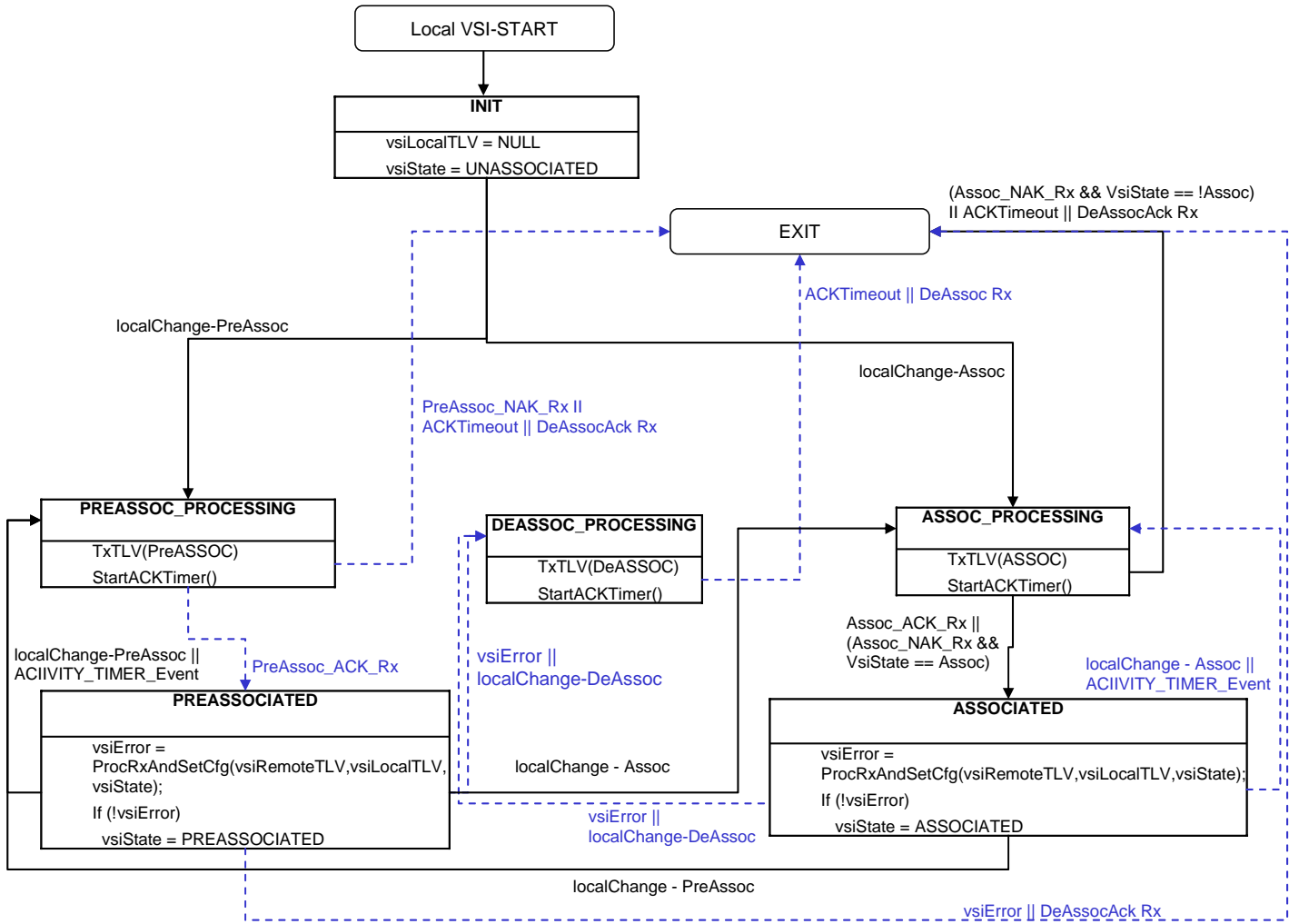


Figure 36. Station's VSI State Machine

7.1.5 Edge Bridge VSI State Machine

Following is the VSI State Machine for the Bridge.

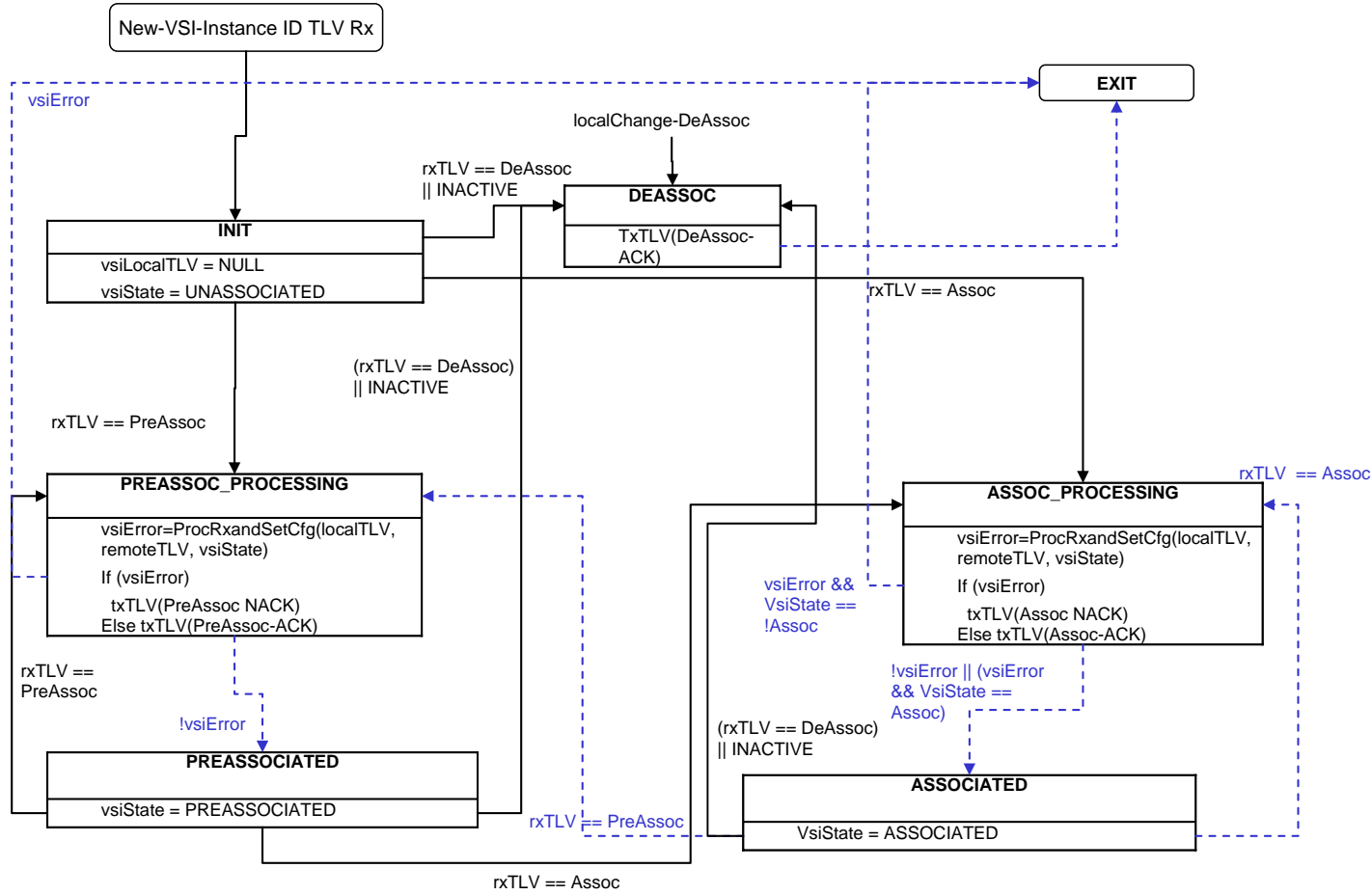


Figure 37. Edge Bridge's VSI State Machine

8. Glossary

Term	Description
S-channel	An S-VLAN used to associate a set of VSI with a physical Ethernet LAN. Traffic within one S-channel is isolated from traffic in another S-channel on the same LAN through the use of a S-Tag.
Chassis	A physical component incorporating one or more IEEE 802 LAN stations and their associated application functionality.
Chassis identifier	An administratively assigned name that identifies the particular chassis within the context of an administrative domain that comprise one or more networks.
CVID	Customer VLAN Identifier
DA	Destination Address
DS	Distribution System
Edge Virtual Bridging (EVB)	The environment where physical end stations, containing multiple VSI, all require the services of adjacent bridges forming a LAN. EVB environments are unique in that virtual NIC configuration information is available to the EVB device that is not normally available to an 802.1Q bridge.
EUI	Extended Unique Identifier
Hypervisor	Computer software and / or hardware platform virtualization software that enables multiple operating systems to operate on top of common, shared hardware.
ID	Identifier
IEEE 802 LAN	Local area network (LAN) technologies that provide a media access control (MAC) Service equivalent to the MAC Service defined in ISO/IEC 15801-1. IEEE 802 LANs include IEEE Std. 802.3, IEEE Std 802.11, IEEE Std 802.16, IEEE Std 802.17, and ISO 9314-2 LANs.
IEEE 802 LAN Station	An IEEE 802-compatible entity that incorporates all the necessary mechanisms to participate in media access control of an IEEE 802 LAN, and that is at least capable of providing the MAC service plus the mandatory capabilities of the LLC.
LLC	Logical Link Control (sub-layer)
Link Layer Discovery Protocol (LLDP)	A media-independent protocol capable of running on all IEEE 802 LAN stations and to allow an LLDP agent to learn the connectivity and management information from adjacent stations.
LLDP agent	The protocol entity that implements LLDP for a particular MSAP associated with a Port.
LLDPDU	Link Layer Discovery Protocol Data Unit
LSAP	Link Service Access Point
MAC	Media Access Control
MAC service access point (MSAP)	The access point for MAC services provided to the LLC sub-layer.
MSAP Identifier	The identifier of a MAC service access point.

Term	Description
Management entity	The protocol entity that implements a particular network management protocol and that provides access support to a MIB associated with the protocol and implemented in a host chassis.
Management Information Base (MIB)	The instantiation of all MIB modules in a managed entity (e.g. system or device)
Management Information Base module (MIB module)	The specification or schema for a data base that can be populated with information required to support a network management information system.
S-Channels	The capability to multiplex multiple virtual channels over a single physical Ethernet LAN.
Network	An interconnected group of systems, each comprising one or more IEEE 802 LAN stations.
Network Interface Controller (NIC)	A device that includes a non-forwarding IEEE 802 LAN station.
Network Management System (NMS)	A management system that is capable of utilizing the information in a MIB.
Object identifier (OID)	An identifier used to name an objective. Structurally, an OID consists of a node in a hierarchically-assigned namespace, formally defined in ISO/IEC 8824-1. Abstract Syntax Notation 1 (ASN.1). OIDs are used in this standard to identify MIB modules and the objects they contain.
OUI	Organizationally Unique Identifier
Physical network topology	The identification of systems, of IEEE 802 LAN stations that compose each system, and of the IEEE 802 LAN stations that attach to the same IEEE 802 LAN.
PCI	Peripheral Component Interface as defined by the PCI-SIG. http://www.pcisig.com . PCI Express (PCIe) represents the latest incarnation of PCI technology within the industry.
PD	Powered Device
Port	The entity in a chassis/system to support an MSAP. A port incorporates one and only one MSAP and identifies the collection of manageable entities that provide the MAC Service at the MSAP.
Port identifier	An administratively assigned name that identifies the particular port within the context of a system, where the identification is convenient, local to the system, and persistent for the system's use and management (whereas the MAC address that globally identifies the MSAP can not be).
PVID	Port VLAN ID
Reflective Relay	Frame relay where the destination port is also the source port
SA	Source Address
Service VLAN	A VLAN identified by a S-VID
Service VLAN ID (S-VID)	A VLAN identifier conveyed in an S-TAG
Service VLAN Tag (S-Tag)	A VLAN tag with a Tag Protocol Identification value allocated for "802.1Q Service Tag Type"

Term	Description
Single-Root I/O Virtualization (SR-IOV)	PCI-SIG specification that enables a PCIe Device to be simultaneously shared by multiple operating systems. A SR-IOV Device supports multiple PCI physical functions (PF) and virtual functions (VF). A PF or a VF is made visible to an operating system by a hypervisor as though it is a single, non-shared PCI Function.
S-VLAN component	A VLAN-aware bridge component with each Port supported by an instance of the IESS that can recognize, insert, and remove Service VLAN tags.
SVID	Service VLAN Identifier
System	A managed collection of hardware and software components incorporating one or more chassis, stations, and ports.
Type, length, value (TLV)	A short, variable length encoding of an information element consisting of sequential type, length, and value fields where the type field identifies the type of information, the length field indicates the length of the information field in octets, and the value field contains the information itself.
VID	VLAN ID
VDP	Virtual Station Interface Discovery and Configuration Protocol. The protocol used to discover and configure a Virtual Station Interface Instance.
Virtual Ethernet Bridge (VEB)	A VEB is a frame relay service that supports local bridging between multiple VSI and (optionally) the external bridging environment. A VEB may be implemented in software as a vSwitch or as embedded hardware within a NIC.
Virtual Ethernet Port Aggregator (VEPA)	<p>A Virtual Ethernet Port Aggregator (VEPA) is a capability within a physical end station that collaborates with an adjacent, external bridge to provide bridging support between multiple virtual end stations and external networks. The VEPA collaborates by forwarding all station-originated frames to the adjacent bridge for frame processing and frame relay (including reflective relay forwarding) and by steering and replicating frames received from the VEPA uplink to the appropriate destinations.</p> <p>May be implemented in software or in conjunction with embedded hardware.</p> <p><i>Note: As with the case of VEBs, VEPAs have access to vNIC configuration information that normally is not available to an 802.1Q bridge.</i></p>
Virtual Machine (VM)	An operating system running on top of a hypervisor.
Virtual NIC (vNIC)	An entity which performs the MAC, LLC, management and control functions needed to attach a VM to an internal LAN.
Channel Access Port	<p>A CAP is a logical Port associated with one end of a S-channel.</p> <p>Within a physical end station, one or more VSI may be multiplexed on top of a CAP.</p> <p>Within an adjacent bridge, a CAP represents a virtual bridge port.</p>
Virtual Station Interface (VSI)	An internal point-to-point Ethernet LAN which connects a Bridge Port of a VEB or VEPA to a vNIC. Each VSI carries a single MAC service instance. Sometimes the term VSI is used to refer to the reference point where the internal LAN attaches to the vNIC.
VSI Instance Identifier	References a specific VSI.
VSI Type Identifier	References a specific VSI Type.

Term	Description
VSI Type	Defines the class of network attributes that can be associated with the VSI.
Virtual Switch (vSwitch)	A software emulated bridge typically implemented within the server virtualization infrastructure (e.g. a Hypervisor). A vSwitch switches network packets between multiple operating systems executing on common, shared hardware. See also VEB.

Appendex CDCP Examples

8.1 CDCP Basic Success Scenario (Sparse)

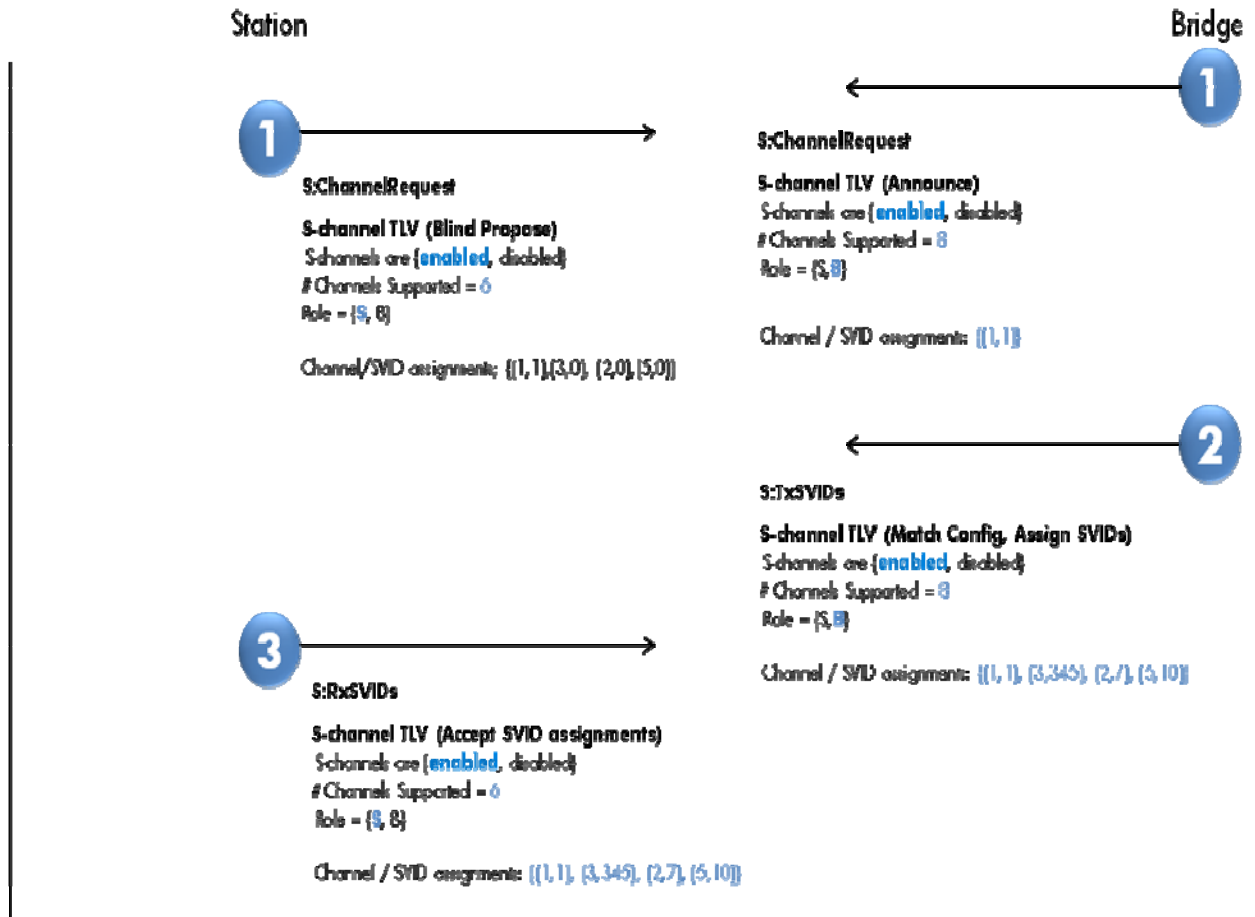


Figure 38. Sparse S-channel Request

8.2 Station Adds an S-channel

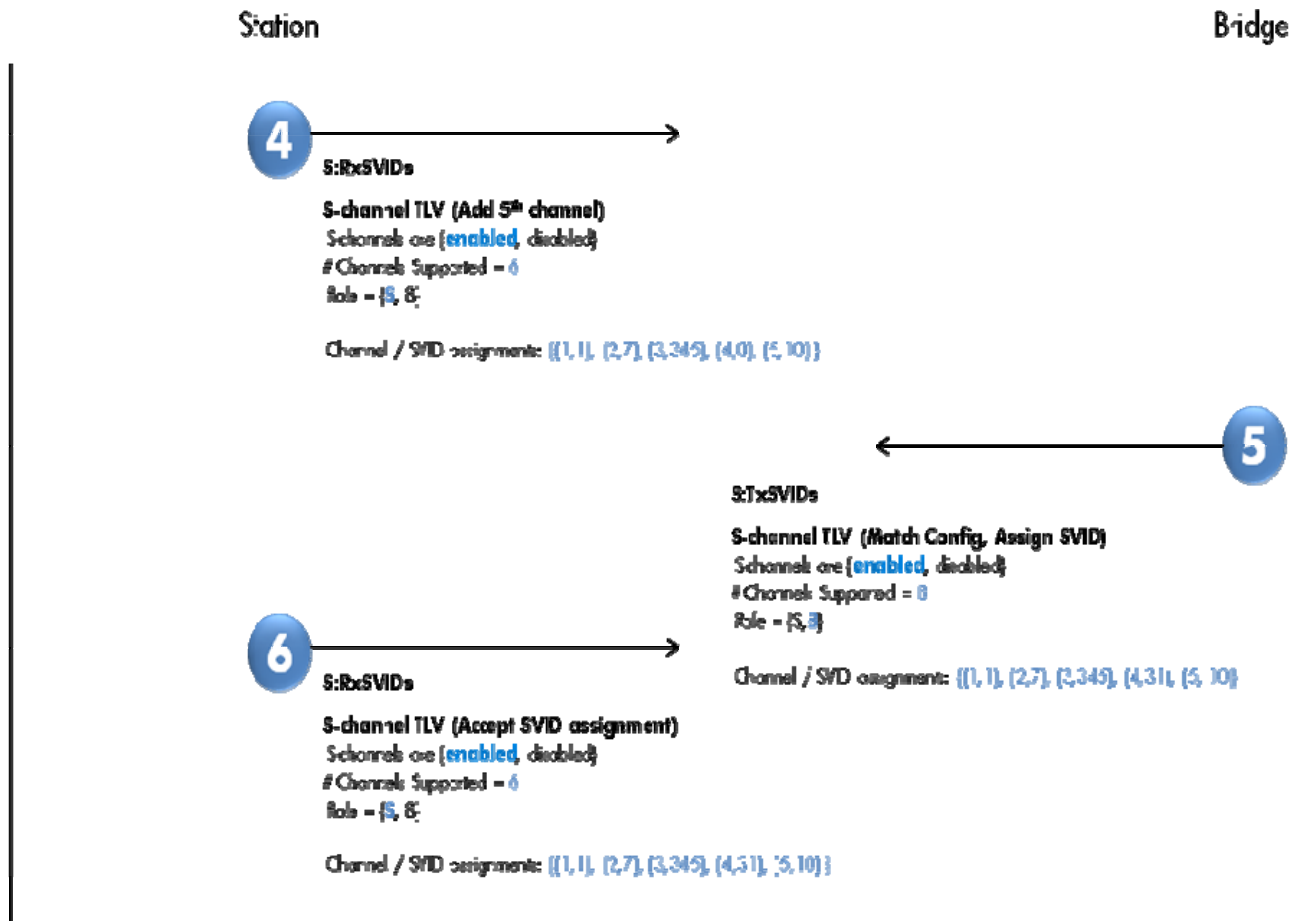


Figure 39. Add an S-channel

8.3 Station Removes a S-channel

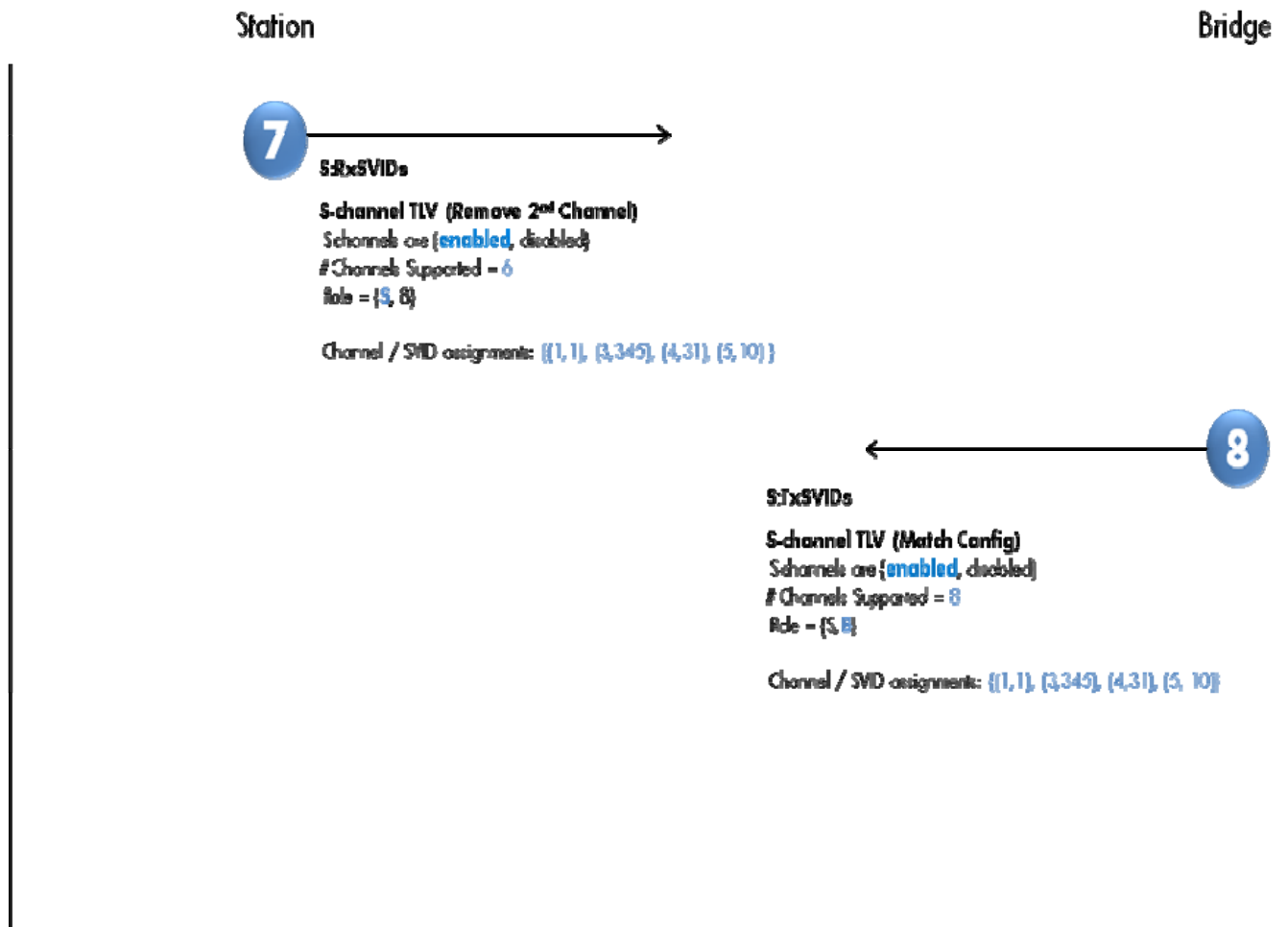


Figure 40.

Remove an S-channel

8.4 Drop #S-channels supported

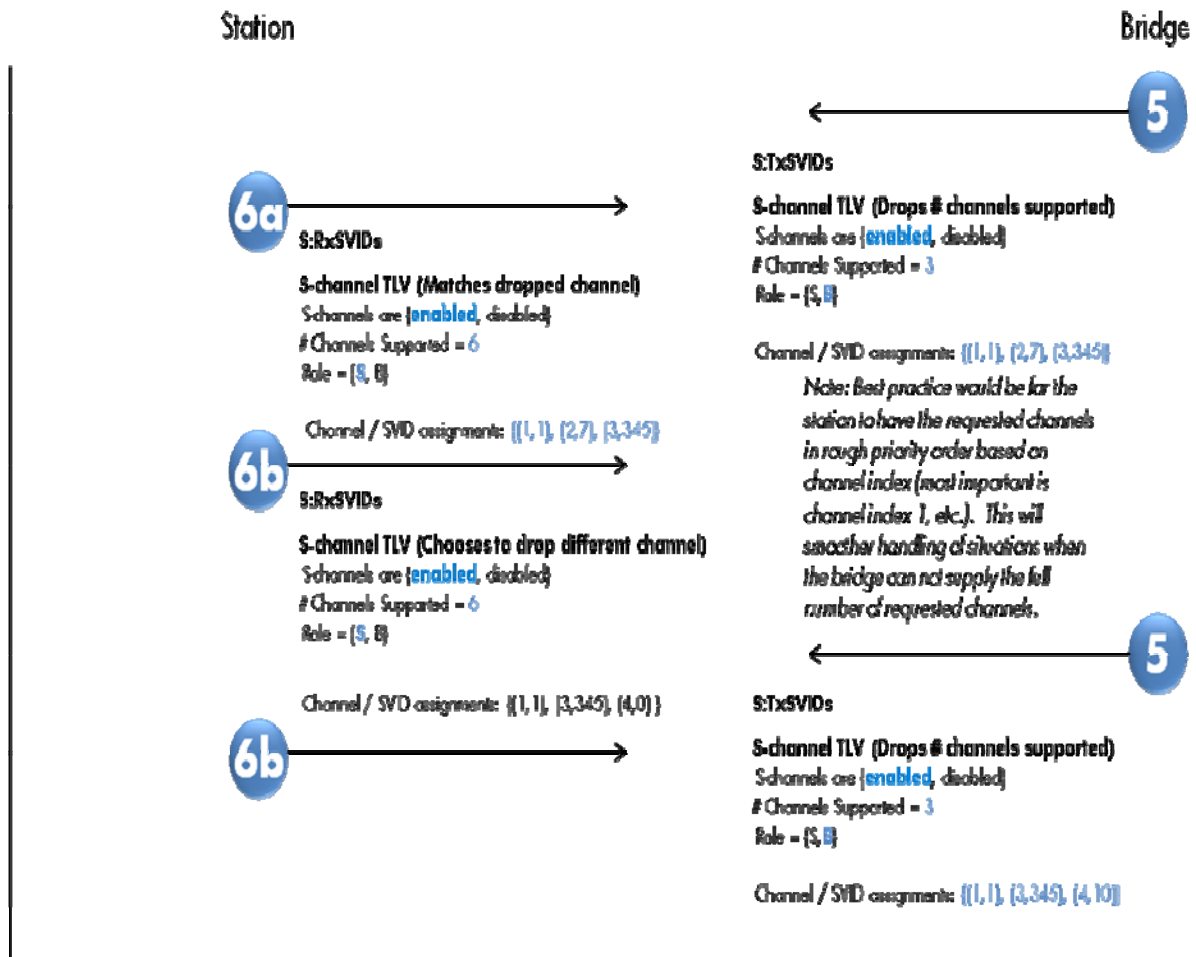


Figure 41. Drop an S-channel

8.5 Insufficient S-channels on bridge

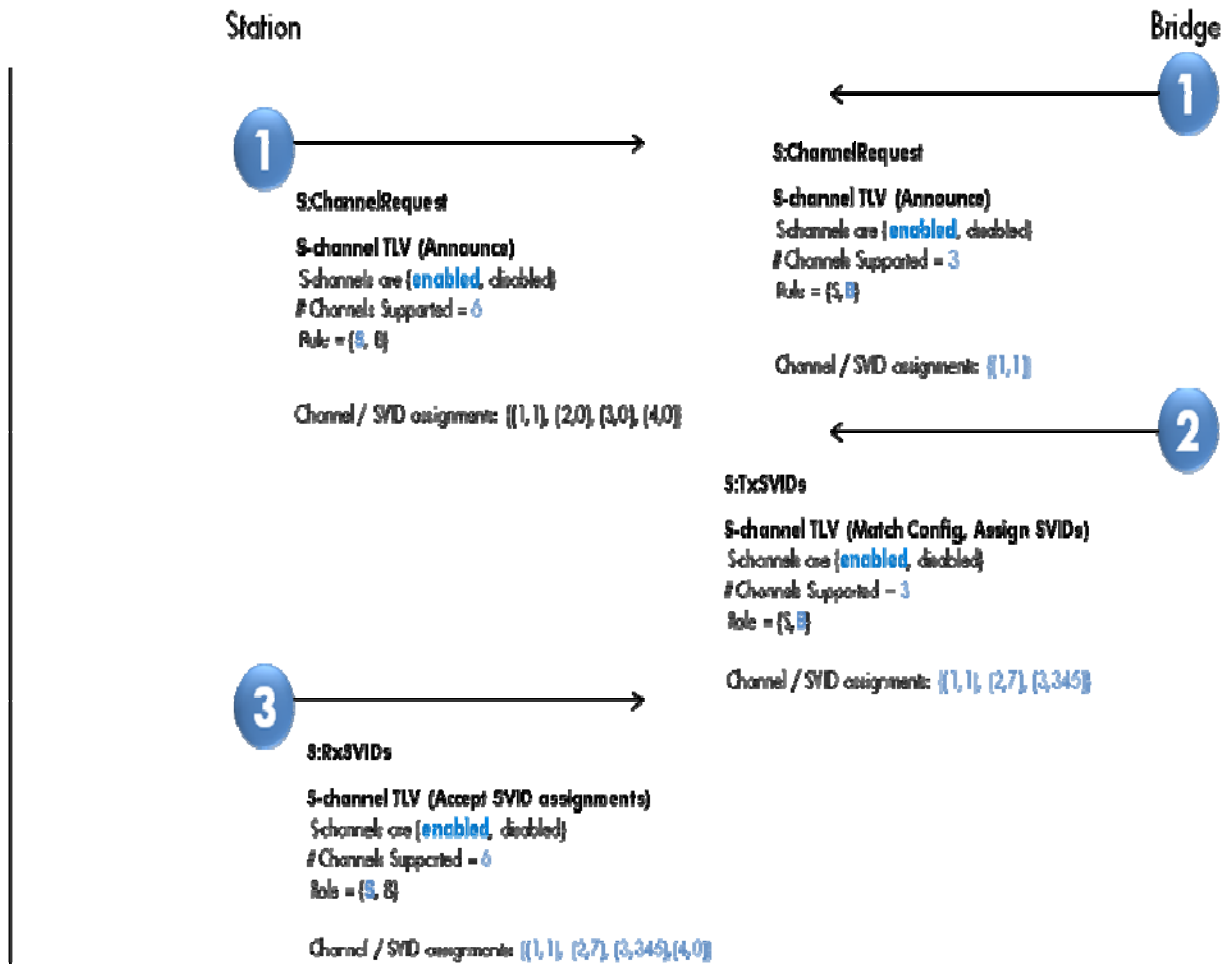


Figure 42. Insufficient Resources

Appendix – VDP Exchange Examples

8.6 VSI PreAssociate, Associate and DeAssociate

The following example depicts the VDP exchanges used to Pre-Associate, Associate and De-Associate a VSI Instance with a VSI Type, VSI Type Version and set of MAC Address and VLAN pairs.

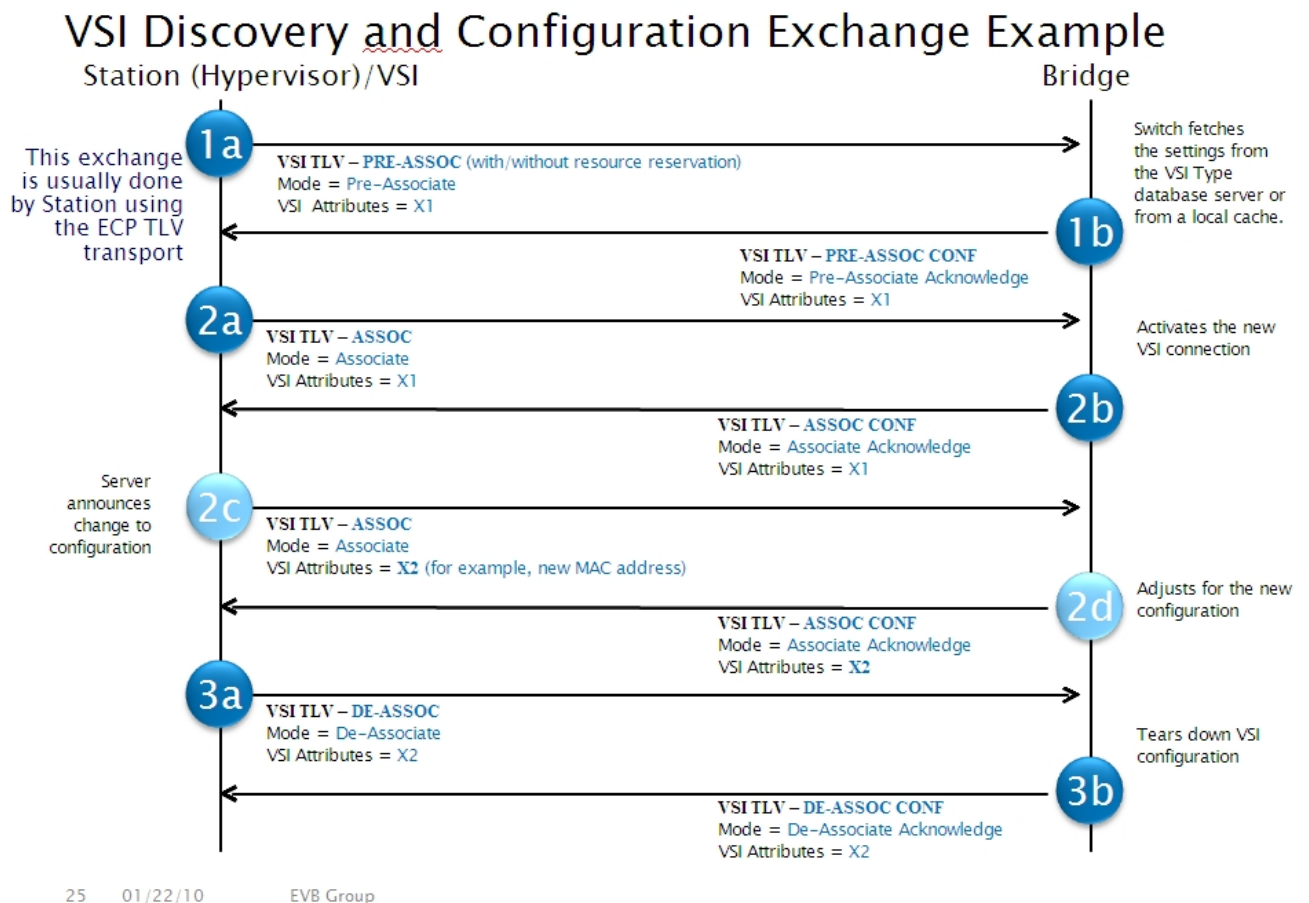


Figure 43. VSI PreAssociate, Associate and DeAssociate Exchange

8.7 VSI Transport Error Case

The following example depicts the VDP exchange associated with a lost EETP transmission of a VSI Associate Request Acknowledgement, showing ECP retrying the transmission.

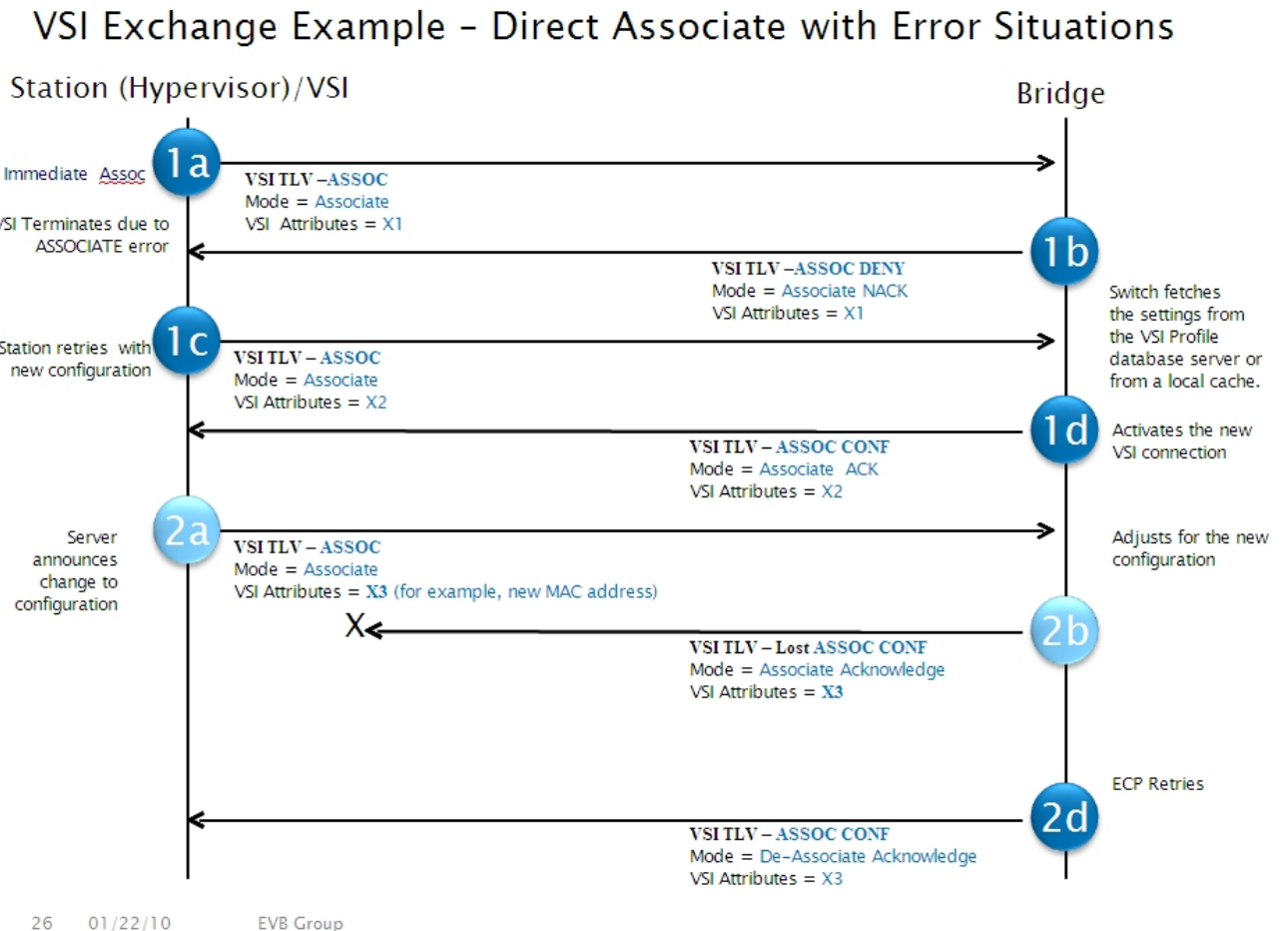


Figure 44. VSI Transport Error

8.8 VSI PreAssociate Resource Lease Refresh Exchange

The following example depicts the VDP exchange associated with an inactive VSI Instance in the Pre-Associated state, where the bridge's VSI State Machine forces a De-Association.

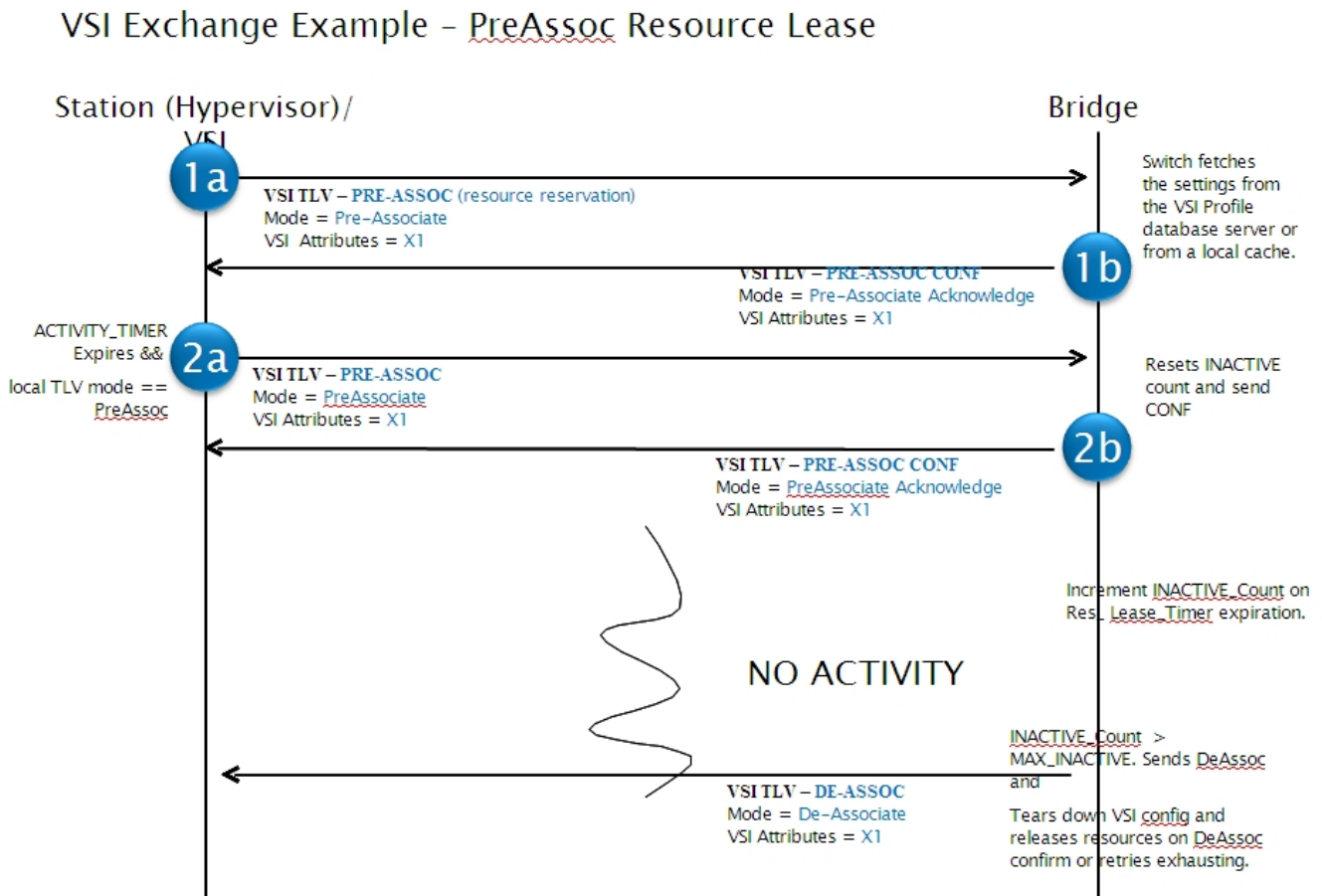


Figure 45. PreAssociate Resource Lease Exchange

8.9 VSI Associate Resource Lease Exchange

The following example depicts the VDP exchange used with an inactive VSI Instance in the Associated state, where the bridge's VSI State Machine forces a De-Association.

VSI Exchange Example - Assoc Resource Lease Refresh

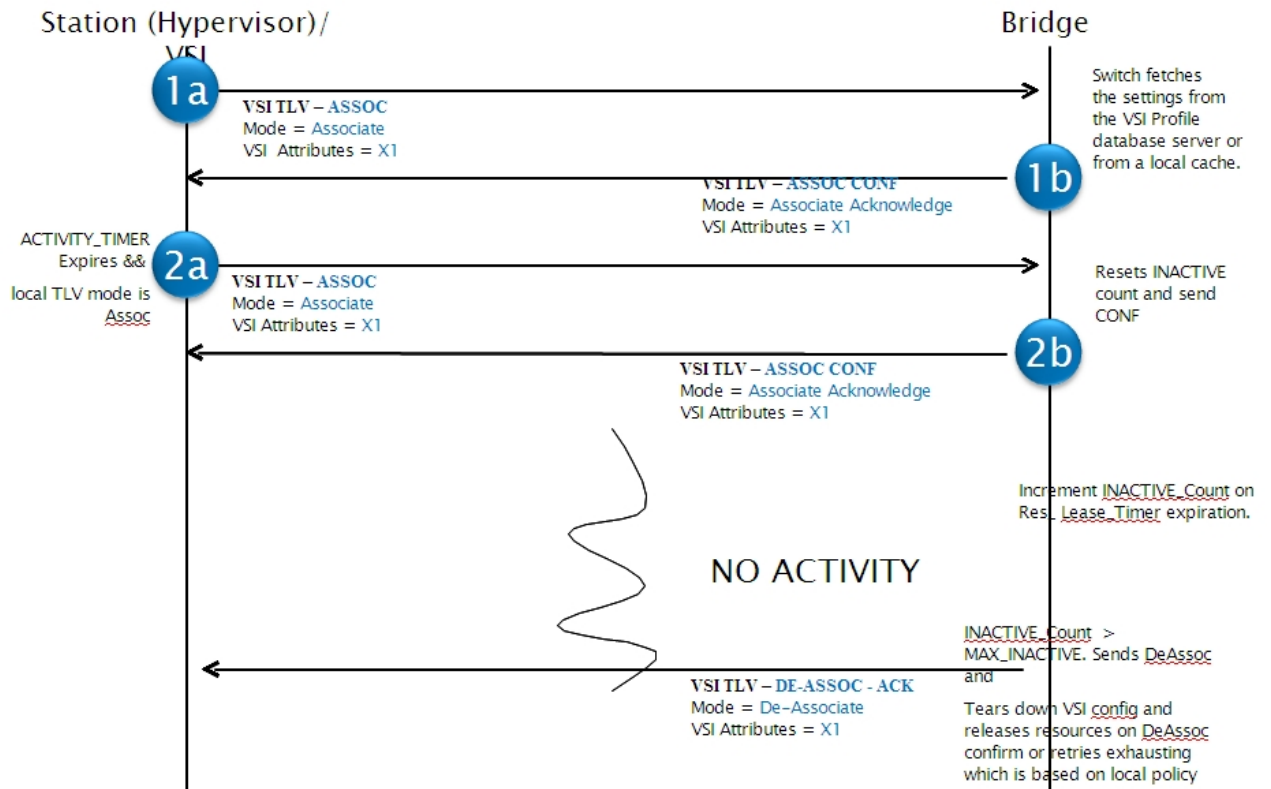


Figure 46. Associate Resource Lease Exchange