

Pseudo-Wire Edge-to-Edge (PWE3) Working Group
Internet Draft
Document: <[draft-ietf-pwe3-arch-06.txt](#)>
Expires: April 2004

Stewart Bryant
Cisco Systems

Prayson Pate
Overture Networks, Inc.

Editors

October 2003

PWE3 Architecture

Status of this Memo

This document is an Internet-Draft and is in full conformance with all provisions of [section 10 of RFC2026](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress".

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/1id-abstracts.txt> The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>.

Abstract

This document describes an architecture for Pseudo Wire Emulation Edge-to-Edge (PWE3). It discusses the emulation of services (such as Frame Relay, ATM, Ethernet, TDM and SONET/SDH) over packet switched networks (PSNs) using IP or MPLS. It presents the architectural framework for pseudo wires (PWs), defines terminology, specifies the various protocol elements and their functions.

Co-Authors

The following are co-authors of this document:

Thomas K. Johnson	Litchfield Communications
Kireeti Kompella	Juniper Networks, Inc.
Andrew G. Malis	Tellabs
Thomas D. Nadeau	Cisco Systems
Tricci So	Caspian Networks
W. Mark Townsley	Cisco Systems
Craig White	Level 3 Communications, LLC.
Lloyd Wood	Cisco Systems
XiPeng Xiao	Riverstone Networks

Conventions used in this document

The key words "MUST", "MUST NOT", "REQUIRED", "SHALL", "SHALL NOT", "SHOULD", "SHOULD NOT", "RECOMMENDED", "MAY", and "OPTIONAL" in this document are to be interpreted as described in [[RFC2119](#)].

Table of Contents

1.	Introduction.....	5
1.1	Pseudo Wire Definition.....	5
1.2	PW Service Functionality.....	6
1.3	Non-Goals of this document.....	6
1.4	Terminology.....	6
2.	PWE3 Applicability.....	9
3.	Protocol Layering Model.....	9
3.1	Protocol Layers.....	9
3.2	Domain of PWE3.....	11
3.3	Payload Types.....	11
4.	Architecture of Pseudo-wires.....	14
4.1	Network Reference Model.....	14
4.2	PWE3 Pre-processing.....	15
4.3	Maintenance Reference Model.....	19
4.4	Protocol Stack Reference Model.....	19
4.5	Pre-processing Extension to Protocol Stack Reference Model.....	20
5.	PW Encapsulation.....	21
5.1	Payload Convergence Layer.....	22
5.2	Payload-independent PW Encapsulation Layers.....	24
5.3	Fragmentation.....	27
5.4	Instantiation of the Protocol Layers.....	27
6.	PW Demultiplexer Layer and PSN Requirements.....	32
6.1	Multiplexing.....	32
6.2	Fragmentation.....	33
6.3	Length and Delivery.....	33
6.4	PW-PDU Validation.....	33
6.5	Congestion Considerations.....	33
7.	Control Plane.....	34
7.1	Set-up or Teardown of Pseudo-Wires.....	34
7.2	Status Monitoring.....	35
7.3	Notification of Pseudo-wire Status Changes.....	35
7.4	Keep-alive.....	37
7.5	Handling Control Messages of the Native Services.....	37
8.	Management and Monitoring.....	37
8.1	Status and Statistics.....	37
8.2	PW SNMP MIB Architecture.....	38
8.3	Connection Verification and Traceroute.....	41

9.	IANA considerations.....	41
10.	Security Considerations.....	41

1. Introduction

This document describes an architecture for Pseudo Wire Emulation Edge-to-Edge (PWE3) in support of [[XIAO](#)]. It discusses the emulation of services (such as Frame Relay, ATM, Ethernet, TDM and SONET/SDH) over packet switched networks (PSNs) using IP or MPLS. It presents the architectural framework for pseudo wires (PWs), defines terminology, specifies the various protocol elements and their functions.

1.1 Pseudo Wire Definition

PWE3 is a mechanism that emulates the essential attributes of a telecommunications service (such as a T1 leased line or Frame Relay) over a PSN. PWE3 is intended to provide only the minimum necessary functionality to emulate the wire with the required degree of faithfulness for the given service definition. Any required switching functionality is the responsibility of a forwarder function (FWRD). Any translation or other operation needing knowledge of the payload semantics is carried out by native service processing (NSP) elements. The functional definition of any FWRD or NSP elements is outside the scope of PWE3.

The required functions of PWs include encapsulating service-specific bit-streams, cells or PDUs arriving at an ingress port, and carrying them across a IP path or MPLS tunnel. In some cases it is necessary to perform other operation such as managing their timing and order, to emulate the behavior and characteristics of the service to the required degree of faithfulness.

From the perspective of a Customer Edge Equipment (CE), the PW is characterised as an unshared link or circuit of the chosen service. In some cases, there may be deficiencies in the PW emulation that impact the traffic carried over a PW, and hence limit the applicability of this technology. These limitations must be fully described in the appropriate service-specific documentation.

For each service type, there will be one default mode of operation that all PEs offering that service type MUST support. However, OPTIONAL modes MAY be defined to improve the faithfulness of the emulated service, if it can be clearly demonstrated that the additional complexity associated with the OPTIONAL mode is offset by the value it offers to PW users.

1.2 PW Service Functionality

PWs provide the following functions in order to emulate the behavior and characteristics of the native service.

- o Encapsulation of service-specific PDUs or circuit data arriving at the PE-bound port (logical or physical).
- o Carriage of the encapsulated data across a PSN tunnel.
- o Establishment of the PW including the exchange and/or distribution of the PW identifiers used by the PSN tunnel endpoints.
- o Managing the signaling, timing, order or other aspects of the service at the boundaries of the PW.
- o Service-specific status and alarm management.

1.3 Non-Goals of this document

The following are non-goals for this document:

- o The on-the-wire specification of PW encapsulations.
- o The detailed definition of the protocols involved in PW set-up and maintenance.

The following are outside the scope of PWE3:

- o Any multicast service not native to the emulated medium.
Thus, Ethernet transmission to a "multicast" IEEE-48 address is in scope, while multicast services like MARS [[RFC2022](#)] that are implemented on top of the medium are out of scope.
- o Methods to signal or control the underlying PSN.

1.4 Terminology

This document uses the following definitions of terms. These terms are illustrated in context in Figure 2.

Attachment Circuit (AC)	The physical or virtual circuit attaching a CE to a PE. An attachment Circuit may be for example a Frame Relay DLCI, an ATM VPI/VCI, an Ethernet port, a VLAN, a PPP connection on a physical interface, a PPP session from an L2TP tunnel, an MPLS LSP, etc. If both physical and virtual ACs are of the same technology (e.g., both ATM, both Ethernet, both Frame Relay) the PW is said to provide "homogeneous transport"; otherwise it is said to provide "heterogeneous transport".
-------------------------	---

CE-bound	The traffic direction where PW-PDUs are received on a PW via the PSN, processed and then sent to the destination CE.
CE Signaling	Messages sent and received by the CEs control plane. It may be desirable or even necessary for the PE to participate in or monitor this signaling in order to effectively emulate the service.
Control Word (CW)	A four octet header used in some encapsulations to carry per packet information when the PSN is MPLS.
Customer Edge (CE)	A device where one end of a service originates and/or terminates. The CE is not aware that it is using an emulated service rather than a native service.
Forwarder (FWRD)	A PE subsystem that selects the PW to use to transmit a payload received on an AC.
Fragmentation	The action of dividing a single PDU into multiple PDUs before transmission with the intent of the original PDU being reassembled elsewhere in the network. Fragmentation MAY be performed in order to allow sending of packets of a larger size than the network MTU which they will traverse.
Maximum transmission unit (MTU)	The packet size (excluding data link header) that an interface can transmit without needing to fragment.
Native Service Processing (NSP)	Processing of the data received by the PE from the CE before presentation to the PW for transmission across the core, or processing of the data received from a PW by a PE before it is output on the AC. NSP functionality is defined by standards bodies other than the IETF, such as ITU-T, ANSI, ATMF, etc.)
Packet Switched Network (PSN)	Within the context of PWE3, this is a network using IP or MPLS as the mechanism for packet forwarding.
PE-bound	The traffic direction where information

	from a CE is adapted to a PW, and PW-PDUs are sent into the PSN.
PE/PW Maintenance	Used by the PEs to set up, maintain and tear down the PW. It may be coupled with CE Signaling in order to effectively manage the PW.
Protocol Data Unit (PDU)	The unit of data output to, or received from, the network by a protocol layer.
Provider Edge (PE)	A device that provides PWE3 to a CE.
Pseudo Wire (PW)	A mechanism that carries the essential elements of an emulated service from one PE to one or more other PEs over a PSN.
Pseudo Wire Emulation Edge to Edge (PWE3)	A mechanism that emulates the essential attributes of service (such as a T1 leased line or frame relay) over a PSN.
Pseudo Wire PDU (PW-PDU)	A PDU sent on the PW that contains all of the data and control information necessary to emulate the desired service.
PSN Tunnel	A tunnel across a PSN inside which one or more PWs can be carried.
PSN Tunnel Signaling	Used to set up, maintain and tear down the underlying PSN tunnel.
PW Demultiplexer	Data-plane method of identifying a PW terminating at a PE.
PWE3 Payload Type Identifier (PWE3-PID)	A identifier used to distinguish between an MPLS IP payload and a CW that is not ECMP safe.
Time Domain Multiplexing (TDM)	Time Division Multiplexing. Frequently used to refer to the synchronous bit-streams at rates defined by G.702.
Tunnel	A method of transparently carrying information over a network.

2. PWE3 Applicability

The PSN carrying a PW will subject payload packets to loss, delay, delay variation, and re-ordering. During a network transient there may be a sustained period of impaired service. The applicability of PWE3 to a particular service depends on the sensitivity of that service (or the CE implementation) to these effects, and the ability of the adaptation layer to mask them. Some services, such as IP over FR over PWE3, may prove quite resilient to IP and MPLS PSN characteristics. Other services, such as the interconnection of PBX systems via PWE3, will require more careful consideration of the PSN and adaptation layer characteristics. In some instances, traffic engineering of the underlying PSN will be required, and in some cases, the constraints may be such that it is not possible to provide the required service guarantees.

3. Protocol Layering Model

The PWE3 protocol-layering model is intended to minimise the differences between PWs operating over different PSN types. The design of the protocol-layering model has the goals of making each PW definition independent of the underlying PSN, and maximizing the reuse of IETF protocol definitions and their implementations.

3.1 Protocol Layers

The logical protocol-layering model required to support a PW is shown in Figure 1.

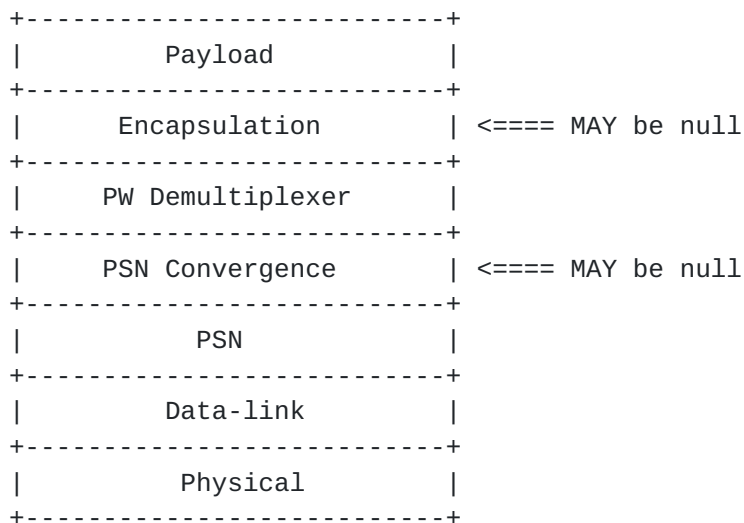


Figure 1: Logical Protocol Layering Model

The payload is transported over the Encapsulation Layer. The Encapsulation Layer carries any information, not already present within the payload itself, that is needed by the PW CE-bound PE interface to send the payload to the CE via the physical interface. If no information is needed beyond that in the payload itself, this layer is empty.

This layer also provides support for real-time processing, and sequencing, if needed.

The PW Demultiplexer Layer provides the ability to deliver multiple PWs over a single PSN tunnel. The PW demultiplexer value used to identify the PW in the data-plane may be unique per PE, but this is not a PWE3 requirement. It MUST, however, be unique per tunnel endpoint. If it is necessary to identify a particular tunnel, then that is the responsibility of the PSN layer.

The PSN Convergence Layer provides the enhancements needed to make the PSN conform to the assumed PSN service requirement. This layer therefore provides a consistent interface to the PW, making the PW independent of the PSN type. If the PSN already meets the service requirements, this layer is empty.

The PSN header, MAC/Data-link and Physical Layer definitions are outside the scope of this document. The PSN can be IPv4, IPv6 or MPLS.

3.2 Domain of PWE3

PWE3 defines the Encapsulation Layer, the method of carrying various payload types, and the interface to the PW Demultiplexer Layer. It is expected that the other layers will be provided by tunneling methods such as L2TP or MPLS over the PSN.

3.3 Payload Types

The payload is classified into the following generic types of native data unit:

- o Packet
- o Cell
- o Bit-stream
- o Structured bit-stream

Within these generic types there are specific service types. For example:

Generic Payload Type	PW Service
-----	-----
Packet	Ethernet (all types), HDLC framing, frame-relay, ATM AAL5 PDU.
Cell	ATM.
Bit-stream	Unstructured E1, T1, E3, T3.
Structured bit-stream	SONET/SDH (e.g. SPE, VT, NxDS0).

3.3.1. Packet Payload

A packet payload is a variable-size data unit delivered to the PE via the AC. A packet payload may be large compared to the PSN MTU. The delineation of the packet boundaries is encapsulation-specific. HDLC or Ethernet PDUs can be considered as examples of packet payloads. Typically a packet will be stripped of transmission overhead such as HDLC flags and stuffing bits before transmission over the PW.

A packet payload would normally be relayed across the PW as a single unit. However, there will be cases where the combined size of the packet payload and its associated PWE3 and PSN headers exceeds the PSN path MTU. In these cases, some fragmentation methodology needs to be applied. This may, for example, be the case when a user is providing the service and attaching to the service provider via Ethernet, or where nested pseudo-wires are involved. Fragmentation is

discussed in more detail in [Section 5.3](#)

A packet payload may need sequencing and real-time support.

In some situations, the packet payload MAY be selected from the packets presented on the emulated wire on the basis of some sub-multiplexing technique. For example, one or more frame-relay PDUs may be selected for transport over a particular pseudo-wire based on the frame-relay Data-Link Connection Identifier (DLCI), or, in the case of Ethernet payloads, using a suitable MAC bridge filter. This is a forwarder function, and this selection would therefore be made before the packet was presented to the PW Encapsulation Layer.

[3.3.2. Cell Payload](#)

A cell payload is created by capturing, transporting and replaying groups of octets presented on the wire in a fixed-size format. The delineation of the group of bits that comprise the cell is specific to the encapsulation type. Two common examples of cell payloads are ATM 53-octet cells, and the larger 188-octet MPEG Transport Stream packets [[DVB](#)].

To reduce per-PSN packet overhead, multiple cells MAY be concatenated into a single payload. The Encapsulation Layer MAY consider the payload complete on the expiry of a timer, after a fixed number of cells have been received or when a significant cell (e.g. an ATM OAM cell) has been received. The benefit of concatenating multiple PDUs should be weighed against a possible increase in packet delay variation and the larger penalty incurred by packet loss. In some cases, it may be appropriate for the Encapsulation Layer to perform some type of compression, such as silence suppression or voice compression.

The generic cell payload service will normally need sequence number support, and may also need real-time support. The generic cell payload service would not normally require fragmentation.

The Encapsulation Layer MAY apply some form of compression to some of these sub-types (e.g. idle cells MAY be suppressed).

In some instances, the cells to be incorporated in the payload MAY be selected by filtering them from the stream of cells presented on the wire. For example, an ATM PWE3 service may select cells based on their VCI or VPI fields. This is a forwarder function, and the selection would therefore be made before the packet was presented to the PW Encapsulation Layer.

3.3.3. Bit-stream

A bit-stream payload is created by capturing, transporting and replaying the bit pattern on the emulated wire, without taking advantage of any structure that, on inspection, may be visible within the relayed traffic (i.e. the internal structure has no effect on the fragmentation into packets).

In some instances it is possible to apply suppression to bit-streams. For example, E1 and T1 send "all-ones" to indicate failure. This condition can be detected without any knowledge of the structure of the bit-stream, and transmission of packetized data suppressed.

This service will require sequencing and real-time support.

3.3.4. Structured bit-stream

A structured bit-stream payload is created by using some knowledge of the underlying structure of the bit-stream to capture, transport and replay the bit pattern on the emulated wire.

Two important points distinguish structured and unstructured bit-streams:

- o Some parts of the original bit-stream MAY be stripped in the PSN-bound direction by NSP block. For example, in Structured SONET the section and line overhead (and, possibly more) may be stripped. A framer is required to enable such stripping. It is also required for frame/payload alignment for fractional T1/E1 applications.
- o The PW MUST preserve the structure across the PSN so that the CE-bound NSP block can insert it correctly into the reconstructed unstructured bit-stream. The stripped information (such as SONET pointer justifications) may appear in the encapsulation layer to facilitate this reconstitution.

As an option, the Encapsulation Layer MAY also perform silence/idle suppression or similar compression on a structured bit-stream.

Structured bit-streams are distinguished from cells in that the structures may be too long to be carried in a single packet. Note that "short" structures are indistinguishable from cells and may benefit from the use of methods described in [section 3.3.2](#).

This service REQUIRES sequencing and real-time support.

3.3.5. Principle of Minimum Intervention

To minimise the scope of information, and to improve the efficiency of data flow through the Encapsulation Layer, the payload SHOULD be transported as received with as few modifications as possible [[RFC1958](#)].

This minimum intervention approach decouples payload development from PW development and requires fewer translations at the NSP in a system with similar CE interfaces at each end. It also prevents unwanted side-effects due to subtle misrepresentation of the payload in the intermediate format.

An approach which does intervene can be more wire-efficient in some cases and may result in fewer translations at the NSP where the CE interfaces are of different types. Any intermediate format effectively becomes a new framing type, requiring documentation and assured interoperability. This increases the amount of work for handling the protocol the intermediate format carries, and is undesirable.

4. Architecture of Pseudo-wires

This section describes the PWE3 architectural model.

4.1 Network Reference Model

Figure 2 illustrates the network reference model for point-to-point PWs.

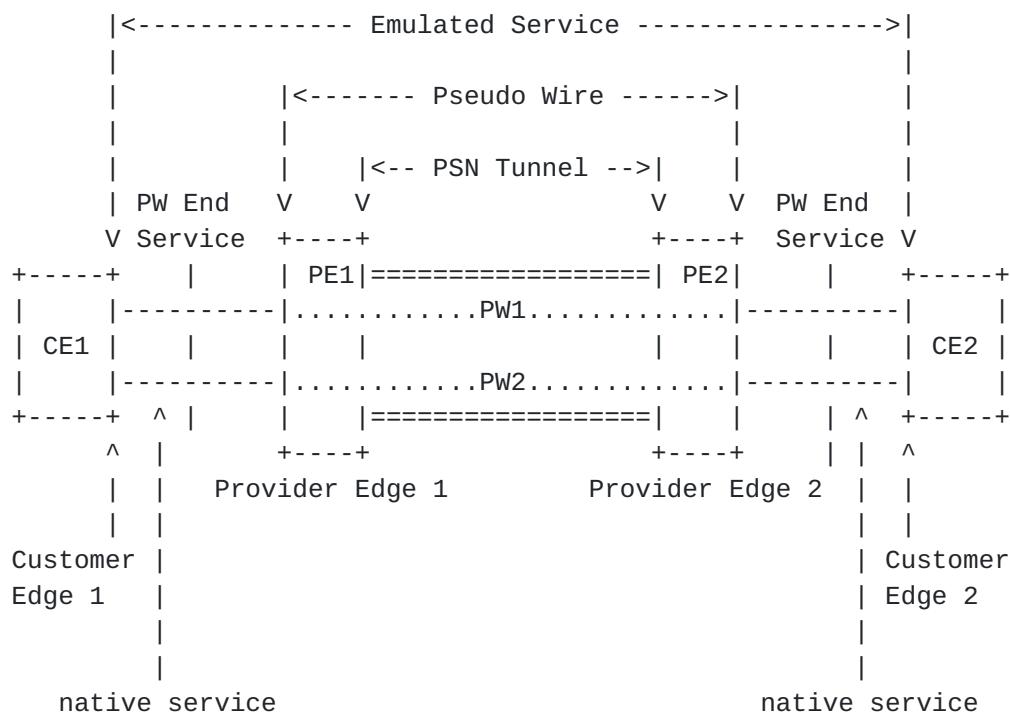


Figure 2: PWE3 Network Reference Model

The two PEs (PE1 and PE2) need to provide one or more PWs on behalf of their client CEs (CE1 and CE2) to enable the client CEs to communicate over the PSN. A PSN tunnel is established to provide a data path for the PW. The PW traffic is invisible to the core network, and the core network is transparent to the CEs. Native data units (bits, cells or packets) arrive via the AC, are encapsulated in a PW-PDU and are carried across the underlying network via the PSN tunnel. The PEs perform the necessary encapsulation and decapsulation of PW-PDUs, as well as handling any other functions required by the PW service, such as sequencing or timing.

4.2 PWE3 Pre-processing

In some applications, there is a need to perform operations on the native data units received from the CE (including both payload and signaling traffic) before they are transmitted across the PW by the PE. Examples include Ethernet bridging, SONET cross-connect, translation of locally-significant identifiers such as VCI/VPI, or translation to another service type. These operations could be carried out in external equipment, and the processed data sent to the PE over one or more physical interfaces. In most cases, there are cost and operational benefits in undertaking these operations within the PE. This processed data is then presented to the PW via a virtual interface within the PE.

- o Forwarder (FWRD)
- o Native Service Processing (NSP)

4.2.1. Forwarders

In some applications there is the need to selectively forward payload elements from one or more ACs to one or more PWs. In such cases there will also be the need to perform the inverse function on PWE3-PDUs received by a PE from the PSN. This is the function of the forwarder.

The forwarder selects the PW based on, for example: the incoming AC, the contents of the payload, or some statically and/or dynamically configured forwarding information.

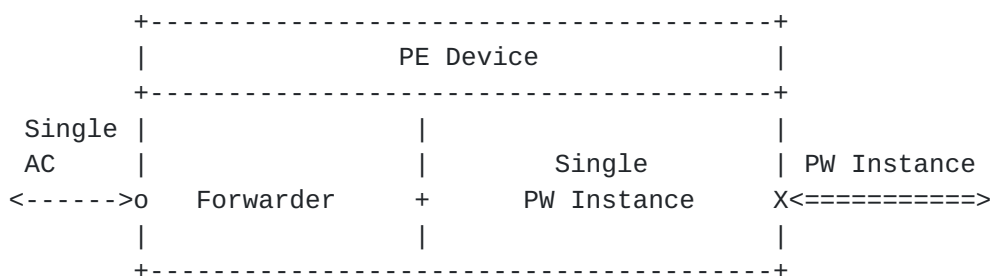


Figure 4a: Simple point-to-point service

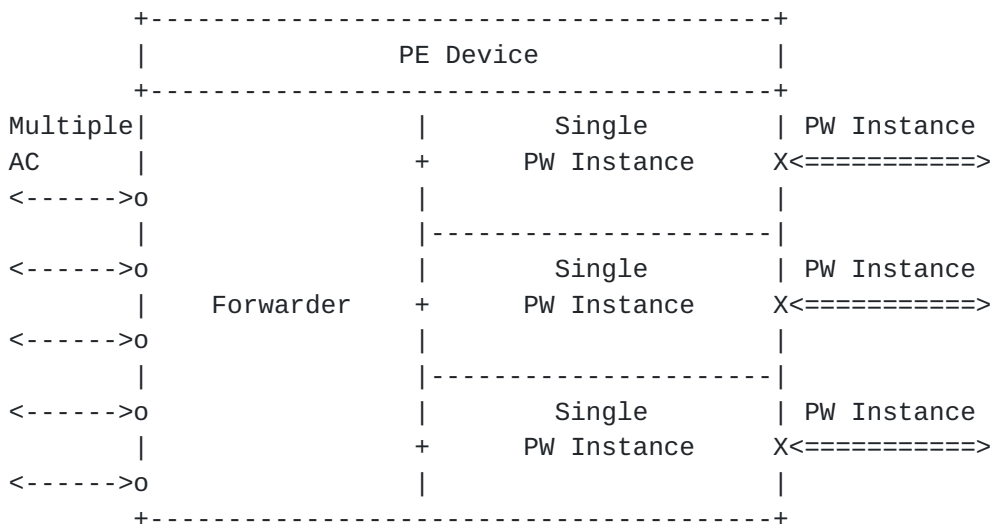


Figure 4b: Multiple AC to Multiple PW Forwarding

Figure 4a shows a simple forwarder that performs some type of filtering operation. Because the forwarder has a single input and a single output interface, filtering is the only type of forwarding operation that applies. Figure 4b shows a more general forwarding situation where payloads are extracted from one or more ACs and directed to one or more PWs. In this case filtering, direction and

Figure 5 illustrates the relationship between NSP, forwarder and PWS in a PE. The NSP function MAY apply any transformation operation (modification, injection, etc.) on the payloads as they pass between the physical interface to the CE and the virtual interface to the forwarder. These transformation operations will of course be limited to those that have been implemented in the data path, and which are enabled by the PE configuration. A PE device MAY contain more than one forwarder.

This model also supports the operation of a system in which the NSP functionality includes terminating the data-link, and applying Network Layer processing to the payload is also supported.

4.3 Maintenance Reference Model

Figure 6 illustrates the maintenance reference model for PWs.

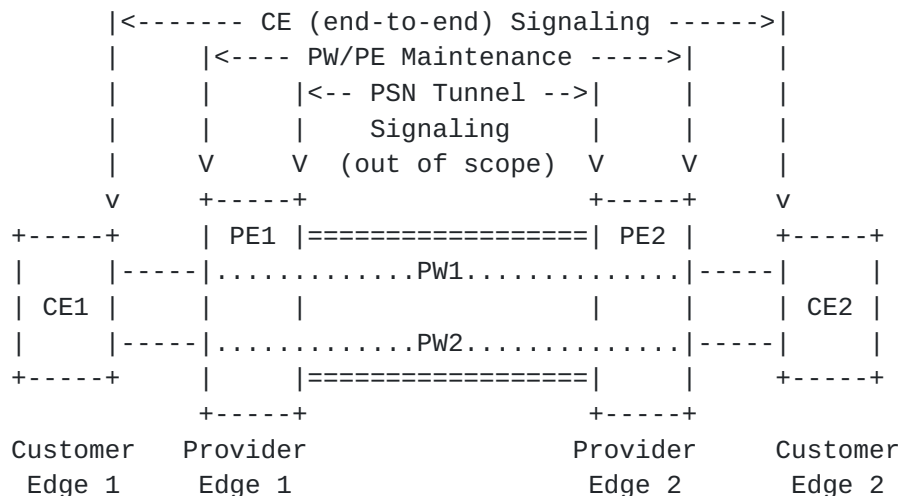


Figure 6: PWE3 Maintenance Reference Model

The following signaling mechanisms are REQUIRED:

- o The CE (end-to-end) signaling is between the CEs. This signaling could be frame relay PVC status signaling, ATM SVC signaling, TDM CAS signaling, etc.
- o The PW/PE Maintenance is used between the PEs (or NSPs) to set up, maintain and tear down PWs, including any required coordination of parameters.
- o The PSN Tunnel signaling controls the PW multiplexing and some elements of the underlying PSN. Examples are L2TP control protocol, MPLS LDP and RSVP-TE. The definition of the information that PWE3 needs to be signaled is within the scope of PWE3, but the signaling protocol itself is not.

4.4 Protocol Stack Reference Model

Figure 7 illustrates the protocol stack reference model for PWs.

Figure 8 illustrates how the protocol stack reference model is extended to include the provision of pre-processing (Forwarding and NSP). This shows the placement of the physical interface relative to the CE.

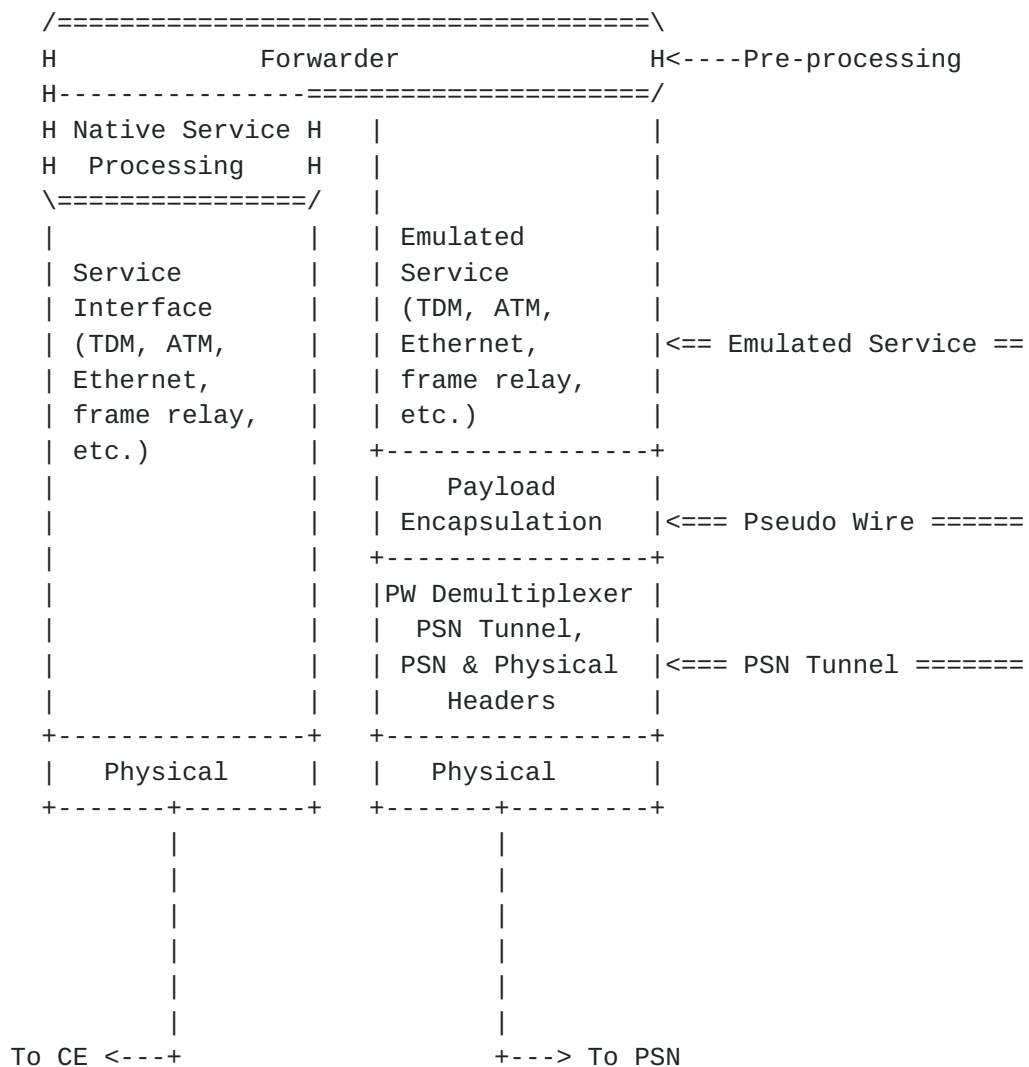


Figure 8: Protocol Stack Reference Model with Pre-processing

5. PW Encapsulation

The PW Encapsulation Layer provides the necessary infrastructure to adapt the specific payload type being transported over the PW to the PW Demultiplexer Layer that is used to carry the PW over the PSN.

The PW Encapsulation Layer consists of three sub-layers:

- o Payload Convergence
- o Timing
- o Sequencing

The PW Encapsulation sub-layering and its context with the protocol stack are shown, in Figure 9.

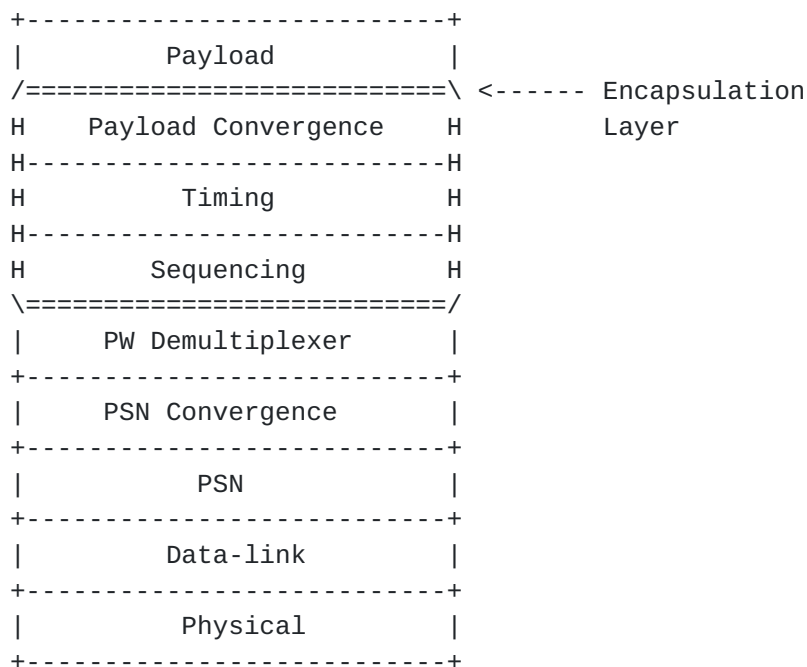


Figure 9: PWE3 Encapsulation Layer in Context

The Payload Convergence Sub-layer is highly tailored to the specific payload type, but, by grouping a number of target payload types into a generic class, and then providing a single convergence sub-layer type common to the group, we achieve a reduction in the number of payload convergence sub-layer types. This decreases implementation complexity. The provision of per-packet signaling and other out-of-band information (other than sequencing or timing) is undertaken by this layer.

The Timing Layer and the Sequencing Layer provide generic services to the Payload Convergence Layer for all payload types that require them.

5.1 Payload Convergence Layer

5.1.1. Encapsulation

The primary task of the Payload Convergence Layer is the encapsulation of the payload in PW-PDUs. The native data units to be encapsulated MAY contain a L2 header or L1 overhead. This is service specific. The Payload Convergence header carries the additional information needed to replay the native data units at the CE-bound physical interface. The PW Demultiplexer header is not considered as

part of the PW header.

Not all the additional information needed to replay the native data units need to be carried in the PW header of the PW PDUs. Some information (e.g. service type of a PW) MAY be stored as state information at the destination PE during PW set-up.

5.1.2. PWE3 Channel Types

The PW Encapsulation Layer and its associated signaling require one or more of the following types of channels from its underlying PW Demultiplexer and PSN Layers:

1. A reliable control channel for signaling line events, status indications, and, in some exceptional cases, CE-CE events that must be translated and sent reliably between PEs.

For example, this capability is needed in [[PPPoL2TP](#)] (PPP negotiation has to be split between the two ends of the tunnel). PWE3 may also need this type of control channel to provide faithful emulation of complex data-link protocols.

plus one or more data channels with the following characteristics:

2. A high-priority, unreliable, sequenced channel. A typical use is for CE-to-CE signaling. "High priority" may simply be indicated via the DSCP bits for IP or the EXP bits for MPLS, giving the packet priority during transit. This channel type could also use a bit in the tunnel header itself to indicate that packets received at the PE SHOULD be processed with higher priority [[RFC2474](#)].
3. A sequenced channel for data traffic that is sensitive to packet reordering (one classification for use could be for any non-IP traffic).
4. An un-sequenced channel for data traffic insensitive to packet order.

The data channels (2, 3 and 4 above) SHOULD be carried "in band" with one another to as much of a degree as is reasonably possible on a PSN.

Where end-to-end connectivity may be disrupted by address translation [[RFC3022](#)], access-control lists, firewalls etc., there exists the possibility that the control channel may be able to pass traffic and set-up the PW, but the PW data traffic is blocked by one or more of these mechanisms. In these cases unless the control channel is also

carried "in band" the signaling to set-up the PW will not confirm the existence of an end-to-end data path.

In some cases there is a need to synchronize CE events with the data carried over a PW. This is especially the case with TDM circuits (e.g., the on-hook/off-hook events in PSTN switches might be carried over a reliable control channel, whilst the associated bit-stream is carried over a sequenced data channel).

PWE3 channel types that are not needed by the supported PWs need not be included in such an implementation.

5.1.3. Quality of Service Considerations

Where possible, it is desirable to employ mechanisms to provide PW Quality of Service (QoS) support over PSNs.

5.2 Payload-independent PW Encapsulation Layers

Two PWE3 Encapsulation Sub-layers provide common services to all payload types: Sequencing and Timing. These services are optional and are only used if needed by a particular PW instance. If the service is not needed, the associated header MAY be omitted in order to conserve processing and network resources.

There will be instances where a specific payload type will be required to be transported with or without sequence and/or real-time support. For example, an invariant of frame relay transport is the preservation of packet order. Some frame-relay applications expect in-order delivery, and may not cope with reordering of the frames. However, where the frame relay service is itself only being used to carry IP, it may be desirable to relax that constraint in return for reduced per-packet processing cost.

The guiding principle is that, where possible, an existing IETF protocol SHOULD be used to provide these services. Where a suitable protocol is not available, the existing protocol should be extended or modified to meet the PWE3 requirements, thereby making that protocol available for other IETF uses. In the particular case of timing, more than one general method may be necessary to provide for the full scope of payload timing requirements.

5.2.1. Sequencing

The sequencing function provides three services: frame ordering, frame duplication detection and frame loss detection. These services allow the emulation of the invariant properties of a physical wire. Support for sequencing depends on the payload type, and MAY be

omitted if not needed.

The size of the sequence-number space depends on the speed of the emulated service, and the maximum time of the transient conditions in the PSN. A sequence number space greater than 2^{16} may therefore be needed to prevent the sequence number space wrapping during the transient.

[5.2.1.1](#) Frame Ordering

When packets carrying the PW-PDUs traverse a PSN, they may arrive out of order at the destination PE. For some services, the frames (control frames, data frames, or both control and data frames) MUST be delivered in order. For such services, some mechanism MUST be provided for ensuring in-order delivery. Providing a sequence number in the sequence sub-layer header for each packet is one possible approach to out-of-sequence detection. Alternatively it can be noted that sequencing is a subset of the problem of delivering timed packets, and that a single combined mechanism such as [[RFC3550](#)] MAY be employed.

There are two possible misordering strategies:

- o Drop misordered PW PDUs.
- o Try to sort PW PDUs into the correct order.

The choice of strategy will depend on:

- o How critical the loss of packets is to the operation of the PW (e.g. the acceptable bit error rate).
- o The speeds of the PW and PSN.
- o The acceptable delay (since delay must be introduced to reorder)
- o The incidence of expected misordering.

[5.2.1.2](#) Frame Duplication Detection

In rare cases, packets traversing a PW may be duplicated by the underlying PSN. For some services frame duplication is not acceptable. For such services, some mechanism MUST be provided to ensure that duplicated frames will not be delivered to the destination CE. The mechanism MAY be the same as the mechanism used to ensure in-order frame delivery.

5.2.1.3 Frame Loss Detection

A destination PE can determine whether a frame has been lost by tracking the sequence numbers of the received PW PDUs.

In some instances, a destination PE will have to presume that a PW PDU is lost if it fails to arrive within a certain time. If a PW-PDU that has been processed as lost subsequently arrives, the destination PE MUST discard it.

5.2.2. Timing

A number of native services have timing expectations based on the characteristics of the networks that they were designed to travel over, and it can be necessary for the emulated service to duplicate these network characteristics as closely as possible, e.g. in delivering native traffic with bit-rate, jitter, wander and delay characteristics similar to those received at the sending PE.

In such cases, it is necessary for the receiving PE to play out the native traffic as it was received at the sending PE. This relies on either timing information sent between the two PEs, or in some case timing information received from an external reference.

The Timing Sub-layer must therefore support two timing functions: clock recovery and timed payload delivery. A particular payload type may require either or both of these services.

5.2.2.1 Clock Recovery

Clock recovery is the extraction of output transmission bit timing information from the delivered packet stream, and requires a suitable mechanism. A physical wire carries the timing information natively, but it is a relatively complex task to extract timing from a highly jittered source such as packet stream. It is therefore desirable that an existing real-time protocol such as [[RFC3550](#)] be used for this purpose, unless it can be shown that this is unsuitable or unnecessary for a particular payload type.

5.2.2.2 Timed delivery

Timed delivery is the delivery of non-contiguous PW PDUs to the PW output interface with a constant phase relative to the input interface. The timing of the delivery may be relative to a clock derived from the packet stream received over the PSN clock recovery, or with reference to an external clock.

5.3 Fragmentation

A payload would ideally be relayed across the PW as a single unit. However, there will be cases where the combined size of the payload and its associated PWE3 and PSN headers exceeds the PSN path MTU. When a packet size exceeds the MTU of a given network, fragmentation and reassembly have to be performed in order for the packet to be delivered. Since fragmentation and reassembly generally consume a considerable network resources as compared to simply switching a packet in its entirety, efforts SHOULD be made to reduce or eliminate the need for fragmentation and reassembly throughout a network to the extent possible. Of particular concern for fragmentation and reassembly are aggregation points where large numbers of PWs are processed (e.g. at the PE).

Ideally, the equipment originating the traffic being sent over the PW will be configured to have adaptive measures (e.g. [[RFC1191](#)], [[RFC1981](#)]) in place that ensure that packets that need to be fragmented are not sent. When this fails, the point closest to the sending host with fragmentation and reassembly capabilities SHOULD attempt to reduce the size of packets to satisfy the PSN MTU. Thus, in the reference model for PWE3 [Figure 3] fragmentation SHOULD first be performed at the CE if at all possible. If and only if the CE cannot adhere to an acceptable MTU size for the PW should the PE attempt its own fragmentation method.

In cases where MTU management fails to limit the payload to a size suitable for transmission of the PW, the PE MAY fall back to either a generic PW fragmentation method, or, if available the fragmentation service of the underlying PSN.

It is acceptable for a PE implementation not to support fragmentation. A PE that does not support fragmentation will drop packets that exceed the PSN MTU, and the management plane of the encapsulating PE MAY be notified.

If the length of a L2/L1 frame, restored from a PW PDU, exceeds the MTU of the destination AC, it MUST be dropped. In this case, the management plane of the destination PE MAY be notified.

5.4 Instantiation of the Protocol Layers

This document does not address the detailed mapping of the Protocol Layering model to existing or future IETF standards. The instantiation of the logical Protocol Layering model is shown in Figure 9.

5.4.1. PWE3 over an IP PSN

The protocol definition of PWE3 over an IP PSN SHOULD employ existing IETF protocols where possible.

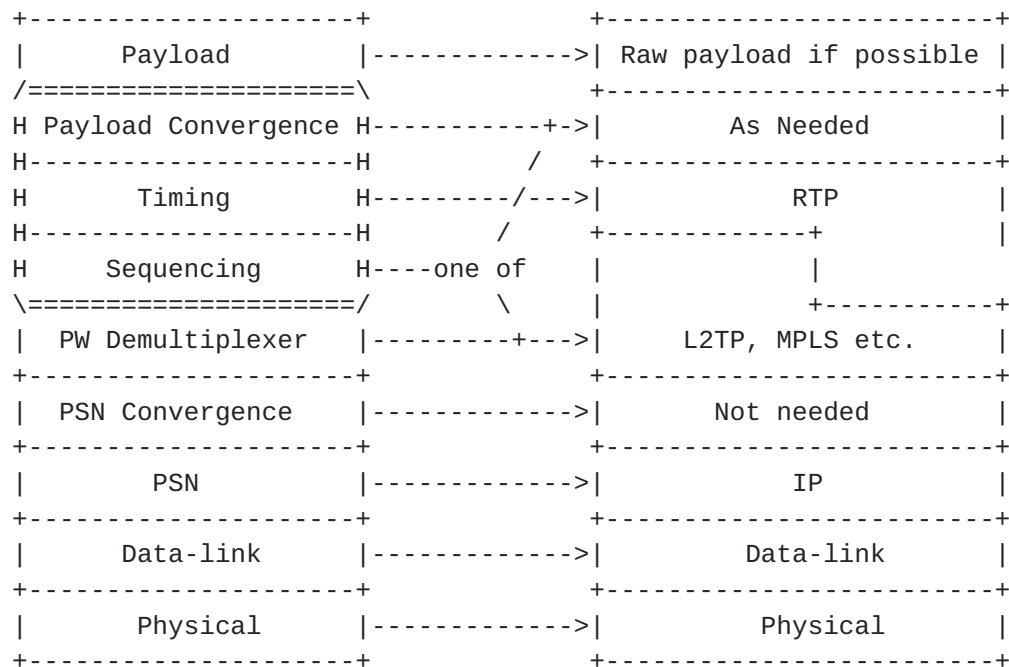


Figure 10: PWE3 over an IP PSN

Figure 10 shows the protocol layering for PWE3 over an IP PSN. As a rule, the payload SHOULD be carried as received from the NSP, with the Payload Convergence Layer provided when needed. (It is accepted that there MAY sometimes be good reason not to follow this rule, but the exceptional circumstances need to be documented in the Encapsulation Layer definition for that payload type).

Where appropriate, timing is provided by RTP [[RFC3550](#)], which when used also provides a sequencing service. PW Demultiplexing may be provided by a number of existing IETF tunnel protocols. Some of these tunnel protocols provide an optional sequencing service. (Sequencing is provided either by RTP, or by the PW Demultiplexer Layer, but not both).

RTP is normally carried over UDP, however the tunnel protocols that are capable of carrying a PW, provide sufficient functionality to carry RTP without an intervening transport layer. UDP MAY therefore be omitted from the protocol stack.

A PSN Convergence Layer is not needed, because all the tunnel protocols shown above are designed to operate directly over an IP

PSN.

As a special case, if the PW Demultiplexer is an MPLS label, the protocol architecture of [section 5.4.2](#) can be used instead of the protocol architecture of this section.

[5.4.2. PWE3 over an MPLS PSN](#)

The MPLS ethos places importance on wire efficiency. By using a control word, some components of the PWE3 protocol layers can be compressed to increase this efficiency.

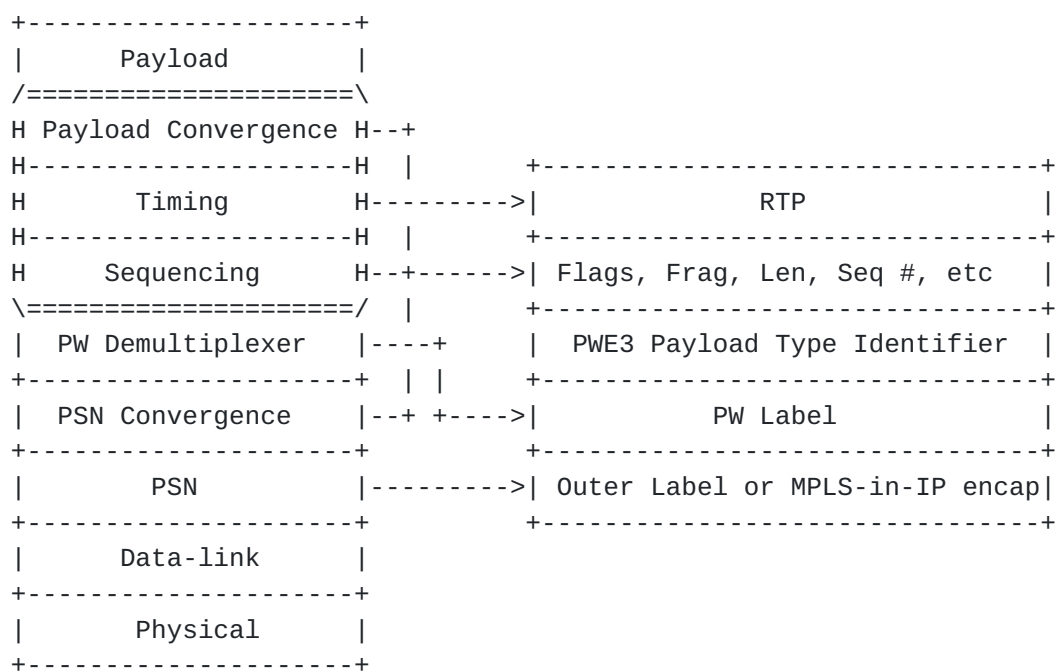


Figure 11: PWE3 over an MPLS PSN using a control word

Figure 11 shows the protocol layering for PWE3 over an MPLS PSN. An inner MPLS label is used to provide the PW demultiplexing function. A control word is used to carry most of the information needed by the PWE3 Encapsulation Layer and the PSN Convergence Layer in a compact format. The flags in the control word provide the necessary payload convergence. A sequence field provides support for both in-order payload delivery and (supported by a fragmentation control method) a PSN fragmentation service within the PSN Convergence Layer. Ethernet pads all frames to a minimum size of 64 bytes. The MPLS header does not include a length indicator. Therefore to allow PWE3 to be carried in MPLS to correctly pass over an Ethernet data-link, a length correction field is needed in the control word. Where the design of the control word would alias an IP packet, a PWE3 Payload Type

Identifier (PWE3 PID) should be interposed between the PW label and the control word (see 5.4.4). As with an IP PSN, where appropriate, timing is provided by RTP [[RFC3550](#)].

In some networks it may be necessary to carry PWE3 over MPLS over IP. In these circumstances, the PW is encapsulated for carriage over MPLS as described in this section, and then a method of carrying MPLS over an IP PSN (such as GRE [[RFC2784](#)], [[RFC2890](#)]) is applied to the resultant PW-PDU.

5.4.3. PW over MPLS Generic Control Word

To allow accurate packet inspection in an MPLS PSN, and/or to operate correctly over MPLS PSNs that have deployed equal-cost multiple-path load-balancing (ECMP), a PW packet MUST NOT alias an IP packet. IP packets are carried in MPLS label stacks without any protocol identifier. Historic values of the IP version number [[RFC791](#)] [[RFC1883](#)] are therefore used to distinguish between IP and non-IP MPLS payloads.

To disambiguate the PW from an IP flow the PW SHOULD employ either the generic PW control word shown in Figure 12, or a PWE3 PID. Note that an MPLS payload with bits 0..3 = 4 is an IPv4 packet and an MPLS payload with bits 0..3 = 6 is an IPv6 packet.

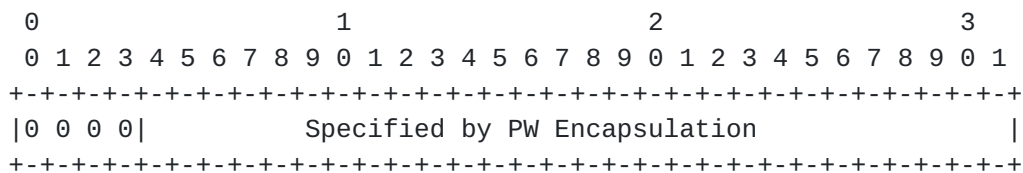


Figure 12: Generic PW Control Word

The PW set-up protocol determines whether a PW uses a control word. When a control word is used, it SHOULD have the following preferred form:

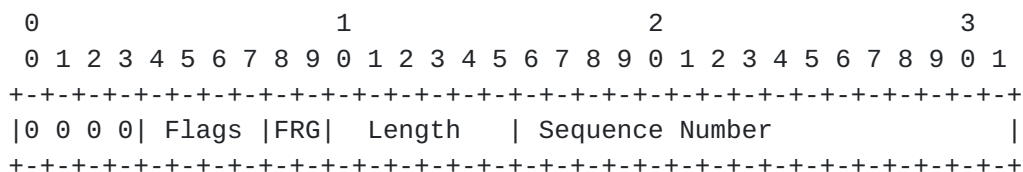


Figure 13: MPLS Preferred Control Word

The meaning of the fields of the MPLS Preferred Control Word (Figure 13) is as follows:

Flags (bits 4 to 7):

These bits are available for per payload signaling. Their definition is encapsulation specific.

FRG (bits 8 and 9):

These bits are used when fragmenting a PW payload. Their use is defined in [[FRAG](#)]. When the PW is of a type that will never need payload fragmentation, these bits may be used as general purpose flags.

Length (bits 10 to 15):

The length field is used to determine the size of a PW payload that might have been padded to the minimum Ethernet MAC frame size during its transit across the PSN. If the MPLS payload (defined as the CW + the PW payload + any additional PW headers) is less than 46 bytes, the length MUST be set to the length of the MPLS payload. If the MPLS payload is between 46 bytes and 63 bytes the implementation MAY either set to the length to the length of the MPLS payload, or it MAY set it to 0. If the length of the MPLS payload is greater than 63 bytes the length MUST be set to 0.

Sequence number (Bit 16 to 31):

If the sequence number is not used, it is set to zero by the sender and ignored by the receiver. Otherwise it specifies the sequence number of a packet. A circular list of sequence numbers is used. A sequence number takes a value from 1 to 65535 ($2^{16}-1$). If the payload is an OAM packet the sequence number MAY be used to mark the position in the sequence, in which case it has the same value as the last data PDU sent. The use of the sequence number is optional for OAM payloads.

5.4.4. PWE3 Payload Type Identifier

If technical considerations result in a PW control word that may alias an IP packet, the control word SHOULD be preceded by an PWE3 payload type identifier (PWE3 PID).

The PWE3 PID is defined as follows:

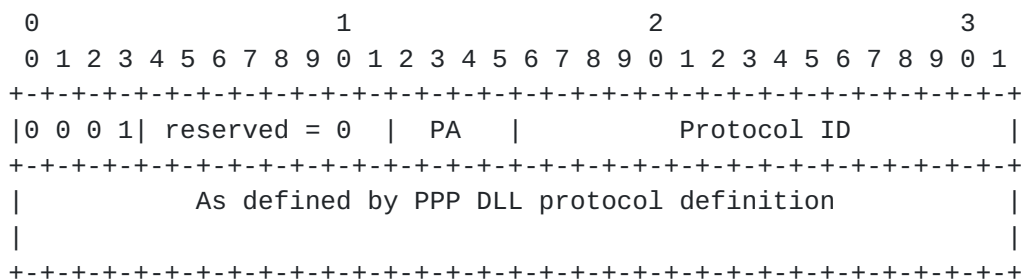


Figure 14: PWE3 PID

The meaning of the fields of the PWE3 PID (Figure 14) is as follows:

PA protocol authority for the user plane or the control plane
 protocol ID
 0 = PPP DLL
 1-15 = Reserved

Protocol ID
 Protocol ID following the format defined by the protocol
 authority identified in PA.

Bits 4 to 11 inclusive are reserved for future use and must be zero.

6. PW Demultiplexer Layer and PSN Requirements

PWE3 places three service requirements on the protocol layers used to carry it across the PSN:

- o Multiplexing
- o Fragmentation
- o Length and Delivery

6.1 Multiplexing

The purpose of the PW Demultiplexer Layer is to allow multiple PWs to be carried in a single tunnel. This minimizes complexity and conserves resources.

Some types of native service are capable of grouping multiple circuits into a "trunk", e.g. multiple ATM VCs in a VP, multiple Ethernet VLANs on a physical media, or multiple DS0 services within a T1 or E1. A PW MAY interconnect two end-trunks. That trunk would have a single multiplexing identifier.

When a MPLS label is used as a PW Demultiplexer setting of the TTL value [[RFC3032](#)] in the PW label is application specific, however in a strict point to point application the TTL SHOULD be set to 2.

[6.2](#) Fragmentation

If the PSN provides a fragmentation and reassembly service of adequate performance, it MAY be used to obtain an effective MTU that is large enough to transport the PW PDUs. See [Section 5.3](#) for a full discussion of the PW fragmentation issues.

[6.3](#) Length and Delivery

PDU delivery to the egress PE is the function of the PSN Layer.

If the underlying PSN does not provide all the information necessary to determine the length of a PW-PDU, the Encapsulation Layer MUST provide it.

[6.4](#) PW-PDU Validation

It is a common practice to use an error detection mechanism such as a CRC or similar mechanism to assure end-to-end integrity of frames. The PW service-specific mechanisms MUST define whether the packet's checksum shall be preserved across the PW, or be removed from PE-bound PDUs and then be re-calculated for insertion in CE-bound data.

The former approach saves work, while the latter saves bandwidth. For a given implementation the choice may be dictated by hardware restrictions, which may not allow the preservation of the checksum.

For protocols such as ATM and FR, the scope of the checksum is restricted to a single link. This is because the circuit identifiers (e.g. FR DLCI or ATM VPI/VCI) have only local significance and are changed on each hop or span. If the circuit identifier (and thus checksum) were going to change as a part of the PW emulation, it would be more efficient to strip and re-calculate the checksum.

The service specific document for each protocol MUST describe the validation scheme to be used.

[6.5](#) Congestion Considerations

The PSN carrying the PW may be subject to congestion. The congestion characteristics will vary with the PSN type, the network architecture and configuration, and the loading of the PSN.

Where the traffic carried over the PW is known to be TCP friendly

(by, for example, packet inspection), packet discard in the PSN will trigger the necessary reduction in offered load, and no additional congestion avoidance action is necessary.

If the PW is operating over a PSN that provides enhanced delivery, the PEs SHOULD monitor packet loss to ensure that the service that was requested is actually being delivered. If it is not, then the PE SHOULD assume that the PSN is providing a best-effort service, and SHOULD use the best-effort service congestion avoidance measures described below.

If best-effort service is being used and the traffic is not known to be TCP friendly, the PEs SHOULD monitor packet loss to ensure that the packet loss rate is within acceptable parameters. Packet loss is considered acceptable if a TCP flow across the same network path and experiencing the same network conditions would achieve an average throughput, measured on a reasonable timescale, that is not less than the PW flow is achieving. This condition can be satisfied by implementing a rate-limiting measure in the NSP, or by shutting down one or more PWs. The choice of which approach to use depends upon the type of traffic being carried. Where congestion is avoided by shutting down a PW, a suitable mechanism MUST be provided to prevent it immediately returning to service, causing a series of congestion pulses.

The comparison to TCP cannot be specified exactly, but is intended as an "order-of-magnitude" comparison in timescale and throughput. The timescale on which TCP throughput is measured is the round-trip time of the connection. In essence, this requirement states that it is not acceptable to deploy an application (using PWE3 or any other transport protocol) on the best-effort Internet which consumes bandwidth arbitrarily and does not compete fairly with TCP within an order of magnitude. One method of determining an acceptable PW bandwidth is described in [[RFC3448](#)].

[7. Control Plane](#)

This section describes PWE3 control plane services.

[7.1 Set-up or Teardown of Pseudo-Wires](#)

A PW MUST be set up before an emulated service can be established, and MUST be torn down when an emulated service is no longer needed.

Set up or teardown of a PW can be triggered by an operator command, from the management plane of a PE, by signaling (i.e., set-up or teardown) of an AC, e.g., an ATM SVC, or by an auto-discovery

mechanism.

During the set-up process, the PEs need to exchange some information (e.g. learn each other's capabilities). The tunnel signaling protocol MAY be extended to provide mechanisms to enable the PEs to exchange all necessary information on behalf of the PW.

Manual configuration of PWs can be considered a special kind of signaling, and is allowed.

7.2 Status Monitoring

Some native services have mechanisms for status monitoring. For example, ATM supports OAM for this purpose. For such services, the corresponding emulated services MUST specify how to perform status monitoring.

7.3 Notification of Pseudo-wire Status Changes

7.3.1. Pseudo-wire Up/Down Notification

If a native service REQUIRES bi-directional connectivity, the corresponding emulated service can only be signaled as being up when the associated PWs, and PSN tunnels if any, are functional in both directions.

Because the two CEs of an emulated service are not adjacent, a failure may occur at a place such that one or both physical links between the CEs and PEs remain up. For example, in Figure 2, if the physical link between CE1 and PE1 fails, the physical link between CE2 and PE2 will not be affected and will remain up. Unless CE2 is notified about the remote failure, it will continue to send traffic over the emulated service to CE1. Such traffic will be discarded at PE1. Some native services have failure notification so that when the services fail, both CEs will be notified. For such native services, the corresponding PWE3 service MUST provide a failure notification mechanism.

Similarly, if a native service has notification mechanisms so that when a network failure is fixed, all the affected services will change status from "Down" to "Up", the corresponding emulated service MUST provide a similar mechanism for doing so.

These mechanisms may already be built into the tunneling protocol. For example, the L2TP control protocol [[RFC2661](#)] [[L2TPv3](#)] has this capability and LDP has the ability to withdraw the corresponding MPLS label.

7.3.2. Misconnection and Payload Type Mismatch

With PWE3, misconnection and payload type mismatch can occur. If a misconnection occurs it can breach the integrity of the system. If a payload mismatch occurs it can disrupt the customer network. In both instances, there are security and operational concerns.

The services of the underlying tunneling mechanism, and its associated control protocol, can be used to mitigate this. As part of the PW set-up a PW-TYPE identifier is exchanged. This is then used by the forwarder and the NSP to verify the compatibility of the ACs.

7.3.3. Packet Loss, Corruption, and Out-of-order Delivery

A PW can incur packet loss, corruption, and out-of-order delivery on the PSN path between the PEs. This can impact the working condition of an emulated service. For some payload types, packet loss, corruption, and out-of-order delivery can be mapped to either a bit error burst, or loss of carrier on the PW. If a native service has some mechanism to deal with bit error, the corresponding PWE3 service should provide a similar mechanism.

7.3.4. Other Status Notification

A PWE3 approach MAY provide a mechanism for other status notification, if any are needed.

7.3.5. Collective Status Notification

Status of a group of emulated services may be affected identically by a single network incident. For example, when the physical link (or sub-network) between a CE and a PE fails, all the emulated services that go through that link (or sub-network) will fail. It is likely that there exists a group of emulated services that all terminate at a remote CE. There may also be multiple such CEs affected by the failure. Therefore, it is desirable that a single notification message be used to notify failure of the whole group of emulated services.

A PWE3 approach MAY provide some mechanism for notifying status changes of a group of emulated circuits. One possible method is to associate each emulated service with a group ID when the PW for that emulated service is set up. Multiple emulated services can then be grouped by associating them with the same group ID. In status notification, that group ID can be used to refer all the emulated services in that group. The group ID mechanism should be a mechanism provided by the underlying tunnel signaling protocol.

7.4 Keep-alive

If a native service has a keep-alive mechanism, the corresponding emulated service **MUST** provide a mechanism to propagate this across the PW. An approach following the principle of minimum intervention would be to transparently transport keep-alive messages over the PW. However, to accurately reproduce the semantics of the native mechanism, some PWs **MAY REQUIRE** an alternative approach, such as piggy-backing on the PW signaling mechanism.

7.5 Handling Control Messages of the Native Services

Some native services use control messages for circuit maintenance. These control messages **MAY** be in-band, e.g. Ethernet flow control, ATM performance management, or TDM tone signaling, or they **MAY** be out-of-band, e.g. the signaling VC of an ATM VP, or TDM CCS signaling.

From the principle of minimum intervention, it is desirable that the PEs participate as little as possible in the signaling and maintenance of the native services. This principle **SHOULD NOT**, however, override the need to satisfactorily emulate the native service.

If control messages are passed through, it may be desirable to send them using either a higher priority or a reliable channel provided by the PW Demultiplexer layer. See PWE3 Channel Types.

8. Management and Monitoring

This section describes the management and monitoring architecture for PWE3.

8.1 Status and Statistics

The PE should report the status of the interface and tabulate statistics that help monitor the state of the network, and to help with measurement of service level agreements (SLAs). Typical counters include:

- o Counts of PW-PDUs sent and received, with and without errors.
- o Counts of sequenced PW-PDUs lost.
- o Counts of service PDUs sent and received over the PSN, with and without errors (non-TDM).
- o Service-specific interface counts.

- o One way delay and delay variation.

These counters would be contained in a PW-specific MIB, and they should not replicate existing MIB counters.

8.2 PW SNMP MIB Architecture

This section describes the general architecture for SNMP MIBs used to manage PW services and the underlying PSN. The intent here is to provide a clear picture of how all of the pertinent MIBs fit together to form a cohesive management framework for deploying PWE3 services.

8.2.1. MIB Layering

The SNMP MIBs created for PWE3 should fit the architecture shown in Figure 15.

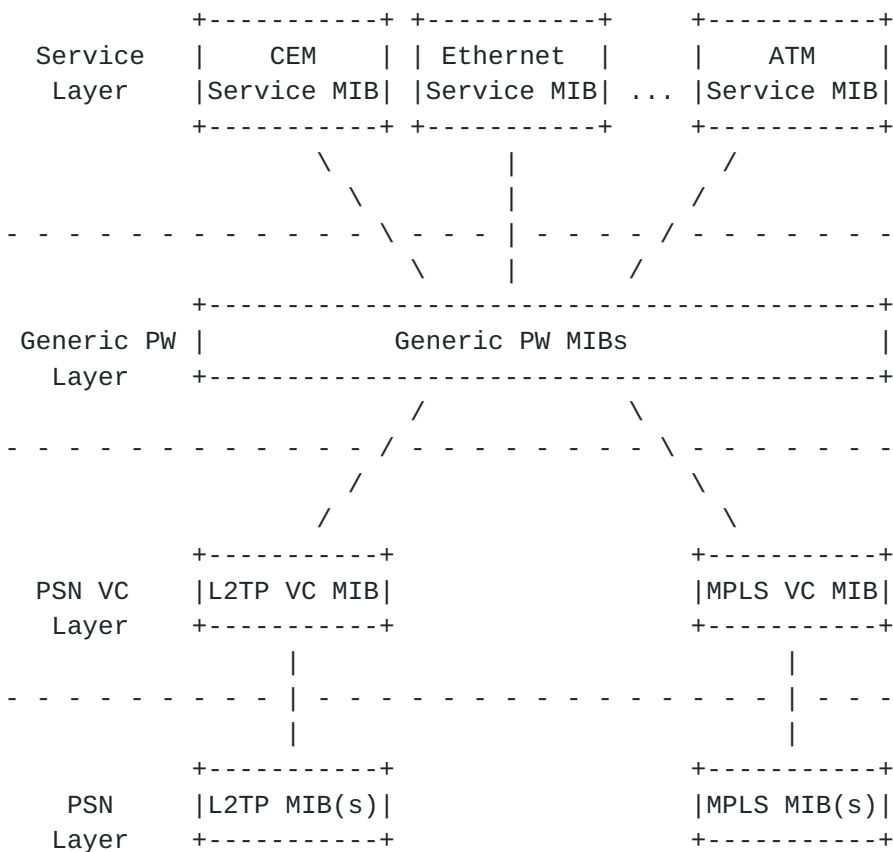


Figure 15: Relationship of SNMP MIBs

Figure 16 shows an example for a SONET PW carried over MPLS.

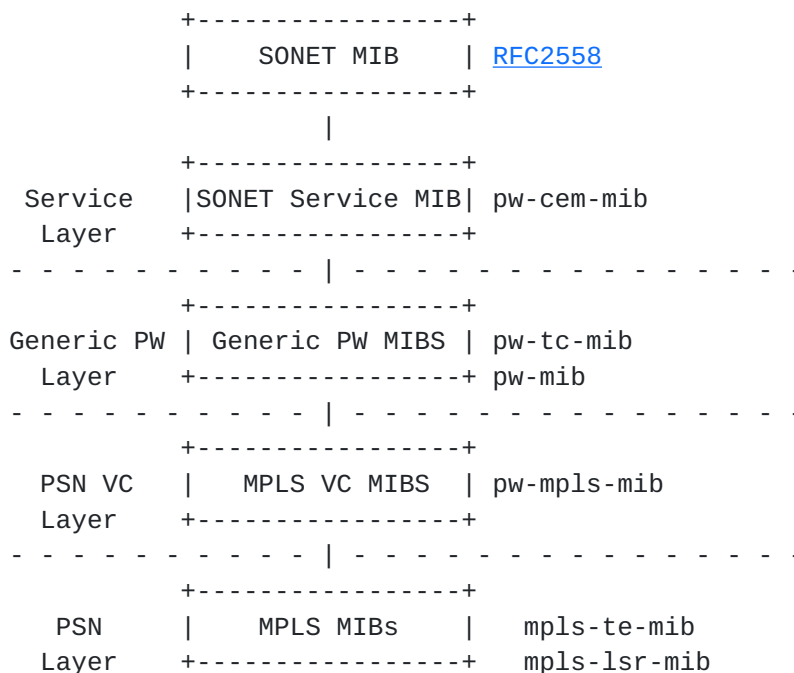


Figure 16: Service-specific Example for MIBs

Note that there is a separate MIB for each emulated service as well as one for each underlying PSN. These MIBs MAY be used in various combinations as needed.

[8.2.2.](#) Service Layer MIBs

The first layer is referred to as the Service Layer. It contains MIBs for PWE3 services such as Ethernet, ATM, circuits and Frame Relay. This layer contains those corresponding MIBs used to mate or adapt those emulated services to the underlying services. This working group should not produce any MIBs for managing the general service; rather, it should produce just those MIBs that are used to interface or adapt the emulated service onto the PWE3 management framework. For example, the standard SONET MIB [[RFC2558](#)] is designed and maintained by another working group. Also, the SONET MIB is designed to manage the native service without PW emulation. Since the PWE3 working group is chartered to produce the corresponding adaptation MIB, in this case, it would produce the PW-CEM-MIB [[PWMPLSMIB](#)] that would be used to adapt SONET services to the underlying PSN that carries the PWE3 service.

8.2.3. Generic PW MIBs

The second layer is referred to as the Generic PW Layer. This layer is composed of two MIBs: the PWE-TC-MIB [[PWTCMIB](#)] and the PWE-MIB [[PWMIB](#)]. These MIBs are responsible for providing general PWE3 counters and service models used for monitoring and configuration of PWE3 services over any supported PSN service. That is, this MIB provides a general model of PWE3 abstraction for management purposes. This MIB is used to interconnect the Service Layer MIBs to the PSN VC Layer MIBs. The latter will be described in the next section. This layer also provides the PW-TC-MIB [[PWTCMIB](#)]. This MIB contains common SMI textual conventions [[RFC1902](#)] that MAY be used by any PW MIB.

8.2.4. PSN VC Layer MIBs

The third layer in the PWE3 management architecture is referred to as the PSN VC layer. This layer is comprised of MIBs that are specifically designed to interface general PWE3 services (VCs) onto those underlying PSN services. In general this means that the MIB provides a means with which an operator can map the PW service onto the native PSN service. For example, in the case of MPLS, it is required that the general VC service be layered onto MPLS LSPs or Traffic Engineered (TE) Tunnels [[RFC3031](#)]. In this case, the PW-MPLS-MIB [[PWMPLSMIB](#)] was created to adapt the general PWE3 circuit services onto MPLS. Like the Service Layer described above the PWE3 working group should produce these MIBs.

8.2.5. PSN Layer MIBs

The fourth and final layer in the PWE3 management architecture is referred to as the PSN layer. This layer is comprised of those MIBs that control the PSN service-specific services. For example, in the case of the MPLS [[RFC3031](#)] PSN service, the MPLS-LSR-MIB [[LSRMIB](#)] and the MPLS-TE-MIB [[TEMIB](#)] are used to interface the general PWE3 VC services onto native MPLS LSPs and/or TE tunnels to carry the emulated services. In addition, the MPLS-LDP-MIB [[LDPMIB](#)] MAY be used to reveal the MPLS labels that are distributed over the MPLS PSN in order to maintain the PW service. The MIBs in this layer are produced by other working groups that design and specify the native PSN services. These MIBs should contain the appropriate mechanisms for monitoring and configuring the PSN service such that the emulated PWE3 service will function correctly.

8.3 Connection Verification and Traceroute

A connection verification mechanism should be supported by PWs. Connection verification as well as other alarm mechanisms can alert the operator that a PW has lost its remote connection. The opaque nature of a PW means that it is not possible to specify a generic connection verification or traceroute mechanism that passes this status to the CEs over the PW. If connection verification status of the PW is needed by the CE, it MUST be mapped to the native connection status method.

For troubleshooting purposes, it is sometimes desirable to know the exact functional path of a PW between PEs. This is provided by the traceroute service of the underlying PSN. The opaque nature of the PW means that this traceroute information is only available within the provider network, e.g., at the PEs.

9. IANA considerations

Sections [5.4.3](#) and [5.4.4](#) discuss the issue of aliasing between PW and IP packets on an MPLS PSN. This aliasing is resolved by using two historic IP version numbers to indicate that the payload is an MPLS preferred control word, or a PWE3 PID. The IP version number registry needs to be updated to allocate IP version number 0 (currently reserved) to MPLS preferred control word, and IP version number 1 (currently unassigned) to PWE3 PID.

10. Security Considerations

PWE3 provides no means of protecting the integrity, confidentiality or delivery of the native data units. The use of PWE3 can therefore expose a particular environment to additional security threats. Assumptions that might be appropriate when all communicating systems are interconnected via a point to point or circuit-switched network may no longer hold when they are interconnected using an emulated wire carried over some types of PSN. It is outside the scope of this specification, to fully analyze and review the risks of PWE3, particularly as these risks will depend on the PSN. An example should make the concern clear. A number of IETF standards employ relatively weak security mechanisms when communicating nodes are expected to be connected to the same local area network. The Virtual Router Redundancy Protocol [[RFC2338](#)] is one instance. The relatively weak

security mechanisms represent a greater vulnerability in an emulated Ethernet connected via a PW.

Exploitation of vulnerabilities from within the PSN may be directed to the PW Tunnel end-point so that PW Demultiplexer and PSN tunnel services are disrupted. Controlling PSN access to the PW Tunnel end-point is one way to protect against this. By restricting PW Tunnel end-point access to legitimate remote PE sources of traffic, the PE may reject traffic that would interfere with the PW Demultiplexing and PSN tunnel services.

Protection mechanisms MUST also address the spoofing of tunneled PW data. The validation of traffic addressed to the PW Demultiplexer end-point is paramount in ensuring integrity of PW encapsulation. Security protocols such as IPSec [[RFC2401](#)] MAY be used by the PW Demultiplexer Layer in order to maintain the integrity of the PW by authenticating data between the PW Demultiplexer End-points.

IPSec MAY provide authentication, integrity, non-repudiation, and confidentiality of data transferred between two PEs. It cannot provide the equivalent services to the native service.

Based on the type of data being transferred, the PW MAY indicate to the PW Demultiplexer Layer that enhanced security services are required. The PW Demultiplexer Layer MAY define multiple protection profiles based on the requirements of the PW emulated service. CE-to-CE signaling and control events emulated by the PW and some data types may require additional protection mechanisms. Alternatively, the PW Demultiplexer Layer may use peer authentication for every PSN packet to prevent spoofed native data units from being sent to the destination CE.

The unlimited transformation capability of the NSP may be perceived as a security risk. In practise the type of operation that the NSP may perform will be limited to those that have been implemented in the data path. The access controls that are in place in the PE to protect and validate its configuration will be sufficient to ensure that the NSP performs as expected.

Acknowledgments

We thank: Sasha Vainshtein for his work on Native Service Processing and advice on bit-stream over PW services. Thomas K. Johnson for his work on the background and motivation for PWs.

We also thank: Ron Bonica, Stephen Casner, Durai Chinnaiiah, Jayakumar Jayakumar, Ghassem Koleyni, Danny McPherson, Eric Rosen, John Rutenmiller, Scott Wainner and David Zelig for their comments and contributions.

Normative References

Internet-drafts are works in progress available from
<<http://www.ietf.org/internet-drafts/>>

- [FRAG] Malis and Townsley, "PWE3 Fragmentation and Reassembly", <[draft-ietf-pwe3-fragmentation-02.txt](#)>, work in progress, June 2003.
- [L2TPv3] Layer Two Tunneling Protocol (Version 3)'L2TPv3', J Lau, et. al. <[draft-ietf-l2tpext-l2tp-base-08.txt](#)>, work in progress, June 2003.
- [RFC791] [RFC-791](#): DARPA Internet Program, Protocol Specification, ISI, September 1981.
- [RFC1883] [RFC-1883](#): Internet Protocol, Version 6 (IPv6), S. Deering, et al, December 1995
- [RFC1902] [RFC-1902](#): Structure of Management Information for Version 2 of the Simple Network Management Protocol (SNMPv2), Case et al, January 1996.
- [RFC2119] [RFC-2119](#), [BCP-14](#): Key words for use in RFCs to Indicate Requirement Levels, S. Bradner.
- [RFC2401] [RFC-2401](#): Security Architecture for the Internet Protocol. S. Kent, R. Atkinson.
- [RFC2474] [RFC-2474](#): Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers, K. Nichols, et. al.
- [RFC2558] K. Tesink, "Definitions of Managed Objects for the SONET/SDH Interface Type", [RFC2558](#), March 1999.
- [RFC2661] [RFC-2661](#): Layer Two Tunneling Protocol "L2TP". W. Townsley, et. al.
- [RFC2784] [RFC-2784](#): Generic Routing Encapsulation (GRE). D. Farinacci et al.

- [RFC2890] [RFC-2890](#): Key and Sequence Number Extensions to GRE.
G. Dommety.
- [RFC3031] [RFC3031](#): Multiprotocol Label Switching Architecture,
E. Rosen, January 2001.
- [RFC3032] [RFC3032](#): MPLS Label Stack Encoding, E. Rosen,
January 2001.
- [RFC3550] [RFC-3550](#): RTP: A Transport Protocol for Real-Time
Applications. H. Schulzrinne et. al.

Informative References

Internet-drafts are works in progress available from
<<http://www.ietf.org/internet-drafts/>>

- [DVB] EN 300 744 Digital Video Broadcasting (DVB); Framing
structure, channel coding and modulation for digital
terrestrial television (DVB-T), European
Telecommunications Standards Institute (ETSI)
- [LDPMIB] Cucchiara, J., Sjostrand, H., and Luciani, J.,
"Definitions of Managed Objects for the Multiprotocol
Label Switching, Label Distribution Protocol (LDP)",
<[draft-ietf-mpls-ldp-mib-11.txt](#)>, work in progress,
June 2003.
- [LSRMIB] Srinivasan et al, "MPLS Label Switch Router Management
Information Base Using SMIV2",
<[draft-ietf-mpls-lsr-mib-10.txt](#)>, work in progress,
June 2003.
- [PPPoL2TP] PPP Tunneling Using Layer Two Tunneling Protocol,
J Lau et al. <[draft-ietf-l2tpext-l2tp-ppp-02.txt](#)>,
work in progress, June 2002.
- [PWMIB] Zelig et al, "Pseudo Wire (PW) Management Information
Base Using SMIV2", <[draft-ietf-pwe3-pw-mib-01.txt](#)>,
work in progress, June 2003.
- [PWTCMIB] Nadeau et al, "Definitions for Textual Conventions and
OBJECT-IDENTITIES for Pseudo-Wires Management"
<[draft-ietf-pwe3-pw-tc-mib-01.txt](#)>, work in progress,
June 2003.
- [PWMLSMIB] Danenberg et al, "SONET/SDH Circuit Emulation Service

Over MPLS (CEM) Management Information Base Using SMIV2", <[draft-ietf-pwe3-cep-mib-01.txt](#)>, work in progress, October 2002.

- [RFC1191] [RFC-1191](#): Path MTU discovery. J.C. Mogul, S.E. Deering.
- [RFC1958] [RFC-1958](#): Architectural Principles of the Internet, B. Carpenter et al.
- [RFC1981] [RFC-1981](#): Path MTU Discovery for IP version 6. J. McCann, S. Deering, J. Mogul.
- [RFC2022] [RFC-2022](#): Support for Multicast over UNI 3.0/3.1 based ATM Networks, G. Armitage.
- [RFC2338] [RFC-2338](#): Virtual Router Redundancy Protocol, S. Knight, M. Shand et. al.
- [RFC3022] [RFC-3022](#): Traditional IP Network Address Translator (Traditional NAT). P Srisuresh et al.
- [RFC3448] [RFC3448](#): TCP Friendly Rate Control (TFRC): Protocol Specification, M. Handley et al. January 2003.
- [TEMIB] Srinivasan et al, "Traffic Engineering Management Information Base Using SMIV2", <[draft-ietf-mpls-te-mib-10.txt](#)>, work in progress, June 2003.
- [XIAO] Xiao et al, "Requirements for Pseudo-Wire Emulation Edge-to-Edge (PWE3)", ([draft-ietf-pwe3-requirements-06.txt](#)), X Xiao et al. work in progress, June 2002.

Editors' Addresses

Stewart Bryant
Cisco Systems,
250, Longwater,
Green Park,
Reading, RG2 6GB,
United Kingdom.

Email: stbryant@cisco.com

Prayson Pate
Overture Networks, Inc.
507 Airport Boulevard
Morrisville, NC, USA 27560 Email: prayson.pate@overturenetworks.com

Full copyright statement

Copyright (C) The Internet Society (2002).
All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the Internet Society or other Internet organizations, except as needed for the purpose of developing Internet standards in which case the procedures for copyrights defined in the Internet Standards process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the Internet Society or its successors or assigns.

This document and the information contained herein is provided on an "AS IS" basis and THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIMS ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

