

PWE3
Internet-Draft
Intended status: Standards Track
Expires: September 3, 2009

S. Bryant, Ed.
C. Filsfils
Cisco Systems
U. Drafz
Deutsche Telekom
V. Kompella
J. Regan
Alcatel-Lucent
S. Amante
Level 3 Communications
March 2, 2009

Flow Aware Transport of MPLS Pseudowires
draft-bryant-filsfils-fat-pw-03

Status of this Memo

This Internet-Draft is submitted to IETF in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

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/lid-abstracts.txt>.

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>.

This Internet-Draft will expire on September 3, 2009.

Copyright Notice

Copyright (c) 2009 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents in effect on the date of publication of this document (<http://trustee.ietf.org/license-info>). Please review these documents carefully, as they describe your rights

Internet-Draft

FAT-PW

March 2009

and restrictions with respect to this document.

Abstract

Where the payload carried over a pseudowire carries a number of identifiable flows it can in some circumstances be desirable to carry those flows over the equal cost multiple paths (ECMPs) that exist in the packet switched network. Most forwarding engines are able to hash based on label stacks and use this to balance flows over ECMPs. This draft describes a method of identifying the flows, or flow groups, to the label switched routers by including an additional label in the label stack.

Requirements Language

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](#) [1].

Internet-Draft

FAT-PW

March 2009

Table of Contents

1.	Introduction	4
1.1.	ECMP in Label Switched Routers	5
1.2.	Flow Label	5
2.	Native Service Processing Function	6
3.	Pseudowire Forwarder	6
3.1.	Encapsulation	6
4.	Signaling the Presence of the Flow Label	7
4.1.	Structure of Flow Label TLV	8
5.	OAM	9
6.	Applicability	10
6.1.	ECMP	11
6.2.	Link Aggregation Groups	12
6.3.	The Single Large Flow Case	12
7.	Applicability to MPLS	13
8.	Security Considerations	14
9.	IANA Considerations	14
10.	Congestion Considerations	14
11.	Acknowledgements	15
12.	References	15
12.1.	Normative References	15
12.2.	Informative References	16
	Authors' Addresses	16

1. Introduction

A pseudowire [[11](#)] is normally transported over one single network path, even if multiple Equal Cost Multiple Paths (ECMP) exit between the ingress and egress PEs[2] [[3](#)]. This is required to preserve the characteristics of the emulated service (e.g. to avoid misordering SAToP pseudowire's [[4](#)]). The use of a single path to preserve order remains the default mode of operation of a pseudowire (PW). The new capability proposed in this document is an OPTIONAL mode which may be used when the use of ECMP paths for is known to be beneficial (and not harmful) to the operation of the PW.

Some pseudowires are used to transport large volumes of IP traffic between routers at two locations. One example of this is the use of an Ethernet pseudowire to create a virtual direct link between a pair of routers. Such pseudowire's may carry from hundred's of Mbps to Gbps of traffic. Such pseudowire's do not require strict ordering to be preserved between packets of the pseudowire. They only require ordering to be preserved within the context of each individual transported IP flow. Some operators have requested the ability to explicitly configure such a pseudowire to leverage the availability of multiple ECMP paths. This allows for better capacity planning as the statistical multiplexing of a larger number of smaller flows is more efficient than with a smaller set of larger flows. Although Ethernet is used as an example above, the mechanisms described in this draft are general mechanisms that may be applied to any pseudowire type in which there are identifiable flows, and in which there is no requirement to preserve the order between those flows.

Typically, forwarding hardware can deduce that an IP payload is being directly carried by an MPLS label stack, and is capable of looking at some fields in packets to construct hash buckets for conversations or flows. However, an intermediate node has no information on the type pseudowire being carried in the packet. This limits the forwarder at the intermediate node to only being able to make an ECMP choice based on a hash of the label stack. In the case of a pseudowire emulating a high bandwidth trunk, the granularity obtained by hashing the default label stack is inadequate for satisfactory load-balancing. The ingress node, however, is in the special position of being able to look at the un-encapsulated packet and spread flows amongst an available ECMP paths, or even Loop-Free Alternates I [\[12\]](#) . This draft proposes a method to introduce granularity on the hashing of traffic running over pseudowires by introducing an additional label, chosen by the ingress node, and placed at the bottom of the label stack.

In addition to providing an indication of the flow structure for use

in ECMP forwarding decisions, the mechanism described in the document may also be used to select flows for distribution over an 802.1ad link aggregation group that has been used in an MPLS network.

[1.1.](#) ECMP in Label Switched Routers

Label switched routers commonly hash the label stack or some elements of the label stack as a method of discriminating between flows, in order to distribute those flows over the available equal cost multiple paths that exist in the network. Since the label at the bottom of stack is usually the label most closely associated with the flow, this normally provides the greatest entropy, and hence is usually included in the hash. This draft describes a method of adding an additional label at the bottom of stack in order to facilitate the load balancing of the flows within a pseudowire over the available ECMPs. A similar design for general MPLS use has also been proposed [\[13\]](#), however that is outside the scope of this draft.

An alternative method of load balancing by creating a number of pseudowires and distributing the flows amongst them was considered, but was rejected because:

- o It did not introduce as much entropy as the load balance label method.
- o It required additional pseudowires to be set up and maintained.

[1.2.](#) Flow Label

An additional label is interposed between the pseudowire label and the control word, or if the control word is not present, between the pseudowire label and the pseudowire payload. This additional label is called the Flow label. Indivisible flows within the pseudowire MUST be mapped to the same Flow label by the ingress PE. The flow label stimulates the correct ECMP load balancing behaviour in the PSN. On receipt of the pseudowire packet at the egress PE (which knows this additional label is present) the flow label is discarded without processing.

Note that the flow label MUST NOT be an MPLS reserved label (values in the range 0..15) [[5](#)], but is otherwise unconstrained by the protocol.

Considerations of the TTL value are described in the Security section of this document. In the case of a pseudowire there are no lower restrictions on the label value since the TTL is never the top label. The designers of the generalized solution [[13](#)].

[2.](#) Native Service Processing Function

The Native Service Processing (NSP) function is a component of a PE that has knowledge of the structure of the emulated service and is able to take action on the service outside the scope of the pseudowire. In this case it is required that the NSP in the ingress PE identify flows, or groups of flows within the service, and indicate the flow (group) identity of each packet as it is passed to the pseudowire forwarder. Since this is an NSP function, by definition, the method used to identify a flow is outside the scope of the pseudowire design. Similarly, since the NSP is internal to the PE, the method of flow indication to the pseudowire forwarder is outside the scope of this document

3. Pseudowire Forwarder

The pseudowire forwarder must be provided with a method of mapping flows to load balanced paths.

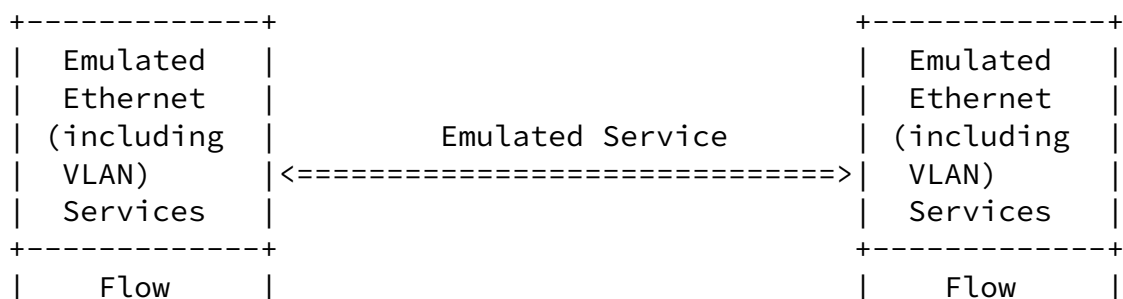
The forwarder must generate a label for the flow or group of flows. How the load balance label values are determined is outside the scope of this document, however the load balance label allocated to a flow MUST NOT be an MPLS reserved label and SHOULD remain constant for the life of the flow. It is recommended that the method chosen to generate the load balancing labels introduces a high degree of entropy in their values, to maximise the entropy presented to the ECMP path selection mechanism in the LSRs in the PSN, and hence distribute the flows as evenly as possible over the available PSN ECMP paths. The forwarder at the ingress PE prepends the pseudowire control word (if applicable), and then pushes the flow label, followed by the pseudowire label.

The forwarder at the egress PE uses the pseudowire label to identify the pseudowire. From the context associated with the pseudowire label, the egress PE can determine whether a flow label is present. If a flow label is present, the label is discarded.

All other pseudowire forwarding operations are unmodified by the inclusion of the flow label.

3.1. Encapsulation

The PWE3 Protocol Stack Reference Model modified to include flow label is shown in Figure 1 below



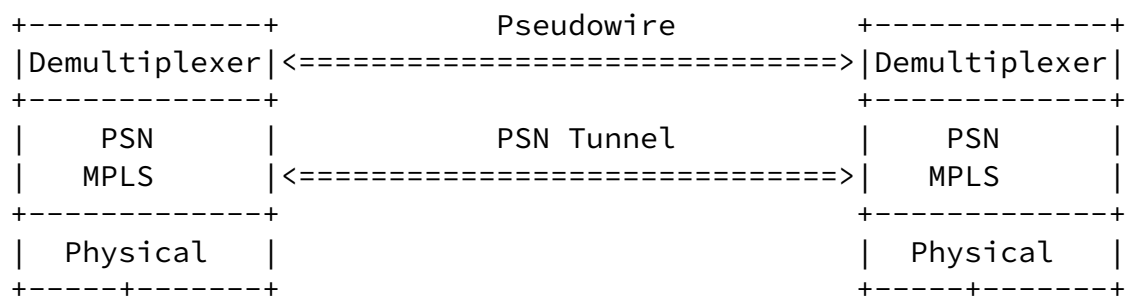


Figure 1: PWE3 Protocol Stack Reference Model

The encapsulation of a pseudowire with a flow label is shown in Figure 2 below

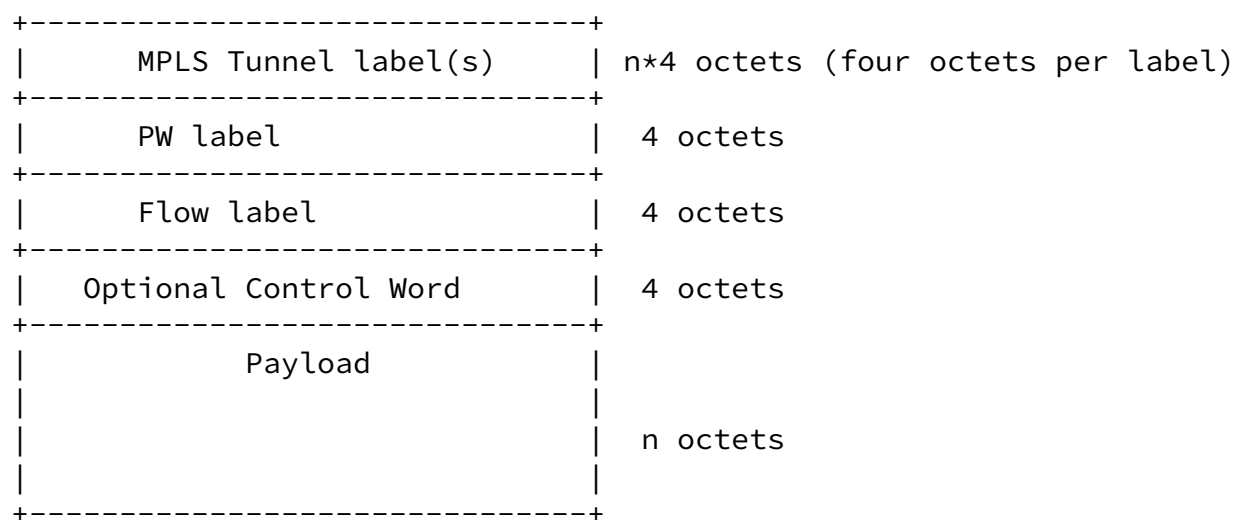


Figure 2: Encapsulation of a pseudowire with a pseudowire load balancing label

4. Signaling the Presence of the Flow Label

When using the signalling procedures in [6], there is a Pseudowire Interface Parameter Sub-TLV type used to synchronize the flow label states between the ingress and egress PEs.

The absence of a flow label (FL) TLV by either party indicates that

the PE concerned is unable to recognize this TLV and the sender of the FL TLV MUST send a new label mapping without the FL TLV. This preserves backwards compatibility with existing PEs that do not understand the FL TLV or that cannot, do not wish to, process the flow label.

A PE that wishes to use a flow label sends an FL TLV with the F bit set. A PE that can correctly process a flow label and is willing to receive on, but does not wish to send a flow label sends an FL TLV with the F bit clear. A PE that sends an FL TLV with the F bit set and receives an FL TLV with or without the F bit set **MUST** include the flow label between the pseudowire label and the control word (or is the control word is not present between the pseudowire label and the pseudowire payload).

If PWE3 signalling [6] is not in use for a pseudowire, then whether the flow label is used MUST be identically provisioned in both PEs at the pseudowire endpoints. If there is no provisioning support for this option, the default behaviour is not to include the flow label.

Note that what is signalled is the desire to include the flow label in the label stack. The value of the label is a local matter for the ingress PE, and the label value itself is not signalled.

4.1. Structure of Flow Label TLV

The structure of the flow label TLV is shown in Figure 3.

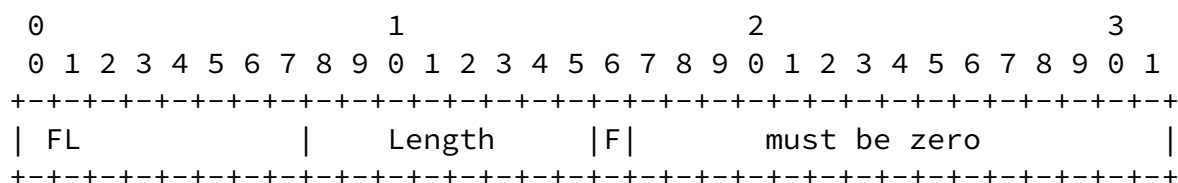


Figure 3: Multiple VC TLV

Where:

- o FL is the flow label TLV identifier assigned by IANA.
- o Length is the length of the TLV in octets and is 4.
- o When F=1 a flow label will be pushed. When F=0 a flow label will not be pushed.

5. OAM

The following OAM considerations apply to this method of load balancing.

Where the OAM is only to be used to perform a basic test that the pseudowires have been configured at the PEs, VCCV [7] messages may be sent using any load balance pseudowire path, i.e. using any value for the flow label.

Where it is required to verify that a pseudowire is fully functional for all flows, VCCV [7] connection verification message MUST be sent over each ECMP path to the pseudowire egress PE. This problem is difficult to solve and scales poorly. We believe that this problem is addressed by the following two methods:

1. If a failure occurs within the PSN, this failure will normally be detected by the PSN's IGP (link/node failure, link or BFD or IGP hello detection), and the IGP convergence will naturally modify the ECMP set of network paths between the Ingress and Egress PE's. Hence the PW is only impacted during the normal IGP convergence time.
2. If the failure is related to the individual corruption of an LFIB entry in a router, then only the network path using that specific entry is impacted. If the PW is load balanced over multiple network paths, then this failure can only be detected if, by chance, the transported OAM flow is mapped onto the impacted network path, or all paths are tested. This type of error may be better solved by other means such as LSP self test [14].

To troubleshoot the MPLS PSN, including multiple paths, the techniques described in [8] and [9] can be used.

Where the pseudowire OAM is carried out of band (VCCV Type 2) it is necessary to insert an "MPLS Router Alert Label" in the label stack. The resultant label stack is as follows:

Internet-Draft

FAT-PW

March 2009

MPLS Tunnel label(s)	n*4 octets (four octets per label)

Router Alert label	4 octets

PW label	4 octets

Flow label	4 octets

Optional Control Word	4 octets

Payload	
	n octets

Figure 4: Use of Router Alert LABEL

6. Applicability

A node within the PSN is not able to perform deep-packet-inspection (DPI) of the PW as the PW technology is not self-describing: the structure of the PW payload is only known to the ingress and egress PE devices. The method proposed in this document provides a statistical mitigation of the problem of load balance in those cases where a PE is able to discern flows embedded in the traffic received on the attachment circuit.

The methods describe in this document are transparent to the PSN and as such do not require any new capability from the PSN.

The requirement to load-balance over multiple PSN paths occurs when the ratio between the PW access speed and the PSN's core link bandwidth is large (e.g. $\geq 10\%$). ATM and FR are unlikely to meet this property. Ethernet does and this is the reason why this document focuses on Ethernet. Applications for other high-access-bandwidth PW's (fiber-channel) may be defined in the future.

This design applies to MPLS pseudowires where it is meaningful to deconstruct the packets presented to the ingress PE into flows. The mechanism described in this document promotes the distribution of flows within the pseudowire over different network paths. This in turn means that whilst packets within a flow are delivered in order (subject to normal IP delivery perturbations due to topology variation), order is not maintained amongst packets of different flows. It is not proposed to associate a different sequence number

with each flow. If sequence number support is required this mechanism is not applicable.

Where it is known that the traffic carried by the Ethernet pseudowire is IP the method of identifying the flows are well known and can be applied. Such methods typically include hashing on the source and destination addresses, the protocol ID and higher-layer flow-dependent fields such as TCP/UDP ports, L2TPv3 Session ID's etc.

Where it is known that the traffic carried by the Ethernet pseudowire is non-IP, techniques used for link bundling between Ethernet switches may be reused. In this case however the latency distribution would be larger than is found in the link bundle case. The acceptability of the increased latency is for further study. Of particular importance the Ethernet control frames SHOULD always be mapped to the same PSN path to ensure in-order delivery.

[6.1.](#) ECMP

ECMP in packet switched networks is statistical in nature. The mapping of flows to a particular path does not take into account the bandwidth of the flow being mapped or the current bandwidth usage of the members of the ECMP set. This simplification works well when the distribution of flows is evenly spread over the ECMP set and there are a large number of flows that have low bandwidth relative to the paths. A random allocation of a flow to a path provides a good approximation to an even spread provided polarization effects are avoided. The method proposed in this document has the same statistical properties as an IP PSN.

ECMP is a load-sharing mechanism that is based on sharing the load over a number of layer 3 paths through the PSN. Often however

multiple links exist between a pair of LSRs that are considered by the IGP to be a single link. These are known as link bundles. The mechanism described in this document can also be used to distribute the flows within a pseudowire over the members of the link bundle by using the flow label value to identify candidate flows. How that mapping takes place is outside the scope of this specification. Similar considerations apply to link aggregation groups.

In the ECMP case and the link bundling case the NSP may attempt to take bandwidth into consideration when allocating groups of flows to a common path. That is permitted, but it must be borne in mind that the semantics of a label stack entry (LSE) as defined by [\[5\]](#) cannot be modified, the value of the flow label cannot be modified at any point on the LSP, and the interpretation of bit patterns in or values of the flow label by an LSR are undefined.

A different type of load balancing is the desire to carry a pseudowire over a set of PSN links in which the bandwidth of members of the link set is less than the bandwidth of the pseudowire. This problem is addressed in [\[15\]](#). Such a mechanism can be considered complementary to this mechanism.

[6.2.](#) Link Aggregation Groups

Link Aggregation (LAG) is used to bond together several physical circuits between two adjacent nodes so they appear to higher-layer protocols as a single, higher bandwidth "virtual" pipe. These may co-exist in various parts of a given network. An advantage of LAGs is that they reduce the number of routing and signaling protocol adjacencies between devices, reducing control plane processing overhead. As with ECMP key problem related to LAG is, due to inefficiencies in LAG load-distribution algorithms, a particular component- link may experience congestion, and the mechanism proposed here may be able to assist in producing a more uniform flow distribution.

The same considerations requiring a flow to go over a single member of an ECMP path set apply to a member of a LAG.

[6.3.](#) The Single Large Flow Case

Clearly the operator should make sure that the service offered using PW technology and the method described in this document does not exceed the maximum planned link capacity unless it can be guaranteed that it conforms to the Internet traffic profile of a very large number of small flows.

If the payload on a PW is made of a single inner flow (i.e. an encrypted connection between two routers), or the flow identifiers are too deeply buried in the packet then the functionality described in this document does not give any benefits, though neither does it cause harm relative to the existing situation. The most common case where a single flow dominated the traffic on a PW is when it is used to transport enterprise traffic. Enterprise traffic may well consist of a large single TCP flows , or encrypted flows that cannot be handled by the methods described in this document.

An operator has six options under these circumstances:

1. The operator can do nothing and the system will work as it does without the flow label.
2. The operator can make the customer aware that the service offering has a restriction on flow bandwidth and police flows to

that restriction. This would allow customers offering multiple flows to use a larger fraction their access bandwidth, whilst preventing an single flow from consuming a fraction of internal link bandwidth that the operator considered excessive.

3. The operator could configure the ingress PE to assign a constant flow label to all high bandwidth flows so that only one path was affected by these flows,
4. The operator could configure the ingress PE to assign a random flow label to all high bandwidth flows so as to minimise the disruption to the network as a cost of out of order traffic to the user.
5. The operator could configure the ingress to assign a label of special significance to all high bandwidth flows so that some other action (not specified in this document) could be taken on the flow.

The issues described above are mitigated by the following two factors:

- o Firstly, the customer of a high-bandwidth PW service has an incentive to get the best transport service because an inefficient use of the PSN leads to jitter and eventually to loss to the PW's payload.
- o Secondly, the customer is usually able to tailor their applications to generate many flows in the PSN. A well-known example is massive data transport between servers which use many parallel TCP sessions. This same technique can be used by any transport protocol: multiple UDP ports, multiple L2TPv3 Session ID's, multiple GRE keys may be used to decompose a large flow into smaller components. This approach may be applied to IPsec where multiple SPI's may be allocated to the same security association.

[7.](#) Applicability to MPLS

A further application of this technique would be to create a basis for hash diversity without having to peek below the label stack for IP traffic carried over LDP LSPs. Work on the generalization of this to MPLS has been described in [draft-kompella-mpls-entropy-label](#). This is can be regarded as a complementary but distinct approach since although similar consideration may apply to the identification of flows and the allocation of flow label values, the flow labels are imposed by different network components and the associated signalling mechanisms are different.

[8.](#) Security Considerations

The pseudowire generic security considerations described in [\[11\]](#) and the security considerations applicable to a specific pseudowire type (for example, in the case of an Ethernet pseudowire [\[10\]](#) apply.

The ingress PE SHOULD take steps to ensure that the load-balance label is not used as a covert channel.

It is useful to give consideration to the choice of TTL value in the flow label LSE. Since the flow label is the bottom of stack and even

when PHP is employed will on arrival at the egress PE be prepended by the PW label, the flow label TTL MAY be set to a value of 1. This will prevent the packet being inadvertently forwarded based on the value of the flow label. Note that this may be a departure from considerations that apply to the general MPLS case.

[9.](#) IANA Considerations

IANA is requested to allocate the next available values from the IETF Consensus range in the Pseudowire Interface Parameters Sub-TLV type Registry as a Flow Label indicator.

Parameter	Length	Description
TBD	4	Load Balancing Label

[10.](#) Congestion Considerations

The congestion considerations applicable to pseudowires as described in [[11](#)] and any additional congestion considerations developed at the time of publication apply to this design.

The ability to explicitly configure a PW to leverage the availability of multiple ECMP paths is beneficial to capacity planning as, all other parameters being constant, the statistical multiplexing of a larger number of smaller flows is more efficient than with a smaller number of larger flows.

Note that if the classification into flows is only performed on IP packets the behaviour of those flows in the face of congestion will be as already defined by the IETF for packets of that type and no additional congestion processing is required.

Where flows that are not IP are classified pseudowire congestion avoidance must be applied to each non-IP load balance group.

[11.](#) Acknowledgements

The authors wish to thank Joerg Kuechemann, Wilfried Maas, Luca Martini, Mark Townsley, Kireeti Kompella and Lucy Yong for valuable

comments on this document.

12. References

12.1. Normative References

- [1] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.
- [2] Bryant, S., Swallow, G., Martini, L., and D. McPherson, "Pseudowire Emulation Edge-to-Edge (PWE3) Control Word for Use over an MPLS PSN", [RFC 4385](#), February 2006.
- [3] Swallow, G., Bryant, S., and L. Andersson, "Avoiding Equal Cost Multipath Treatment in MPLS Networks", [BCP 128](#), [RFC 4928](#), June 2007.
- [4] Vainshtein, A. and YJ. Stein, "Structure-Agnostic Time Division Multiplexing (TDM) over Packet (SAToP)", [RFC 4553](#), June 2006.
- [5] Rosen, E., Tappan, D., Fedorkow, G., Rekhter, Y., Farinacci, D., Li, T., and A. Conta, "MPLS Label Stack Encoding", [RFC 3032](#), January 2001.
- [6] Martini, L., Rosen, E., El-Aawar, N., Smith, T., and G. Heron, "Pseudowire Setup and Maintenance Using the Label Distribution Protocol (LDP)", [RFC 4447](#), April 2006.
- [7] Nadeau, T. and C. Pignataro, "Pseudowire Virtual Circuit Connectivity Verification (VCCV): A Control Channel for Pseudowires", [RFC 5085](#), December 2007.
- [8] Allan, D. and T. Nadeau, "A Framework for Multi-Protocol Label Switching (MPLS) Operations and Management (OAM)", [RFC 4378](#), February 2006.
- [9] Kompella, K. and G. Swallow, "Detecting Multi-Protocol Label Switched (MPLS) Data Plane Failures", [RFC 4379](#), February 2006.
- [10] Martini, L., Rosen, E., El-Aawar, N., and G. Heron, "Encapsulation Methods for Transport of Ethernet over MPLS Networks", [RFC 4448](#), April 2006.

12.2. Informative References

- [11] Bryant, S. and P. Pate, "Pseudo Wire Emulation Edge-to-Edge (PWE3) Architecture", [RFC 3985](#), March 2005.
- [12] Zinin, A., Torvi, R., Choudhury, G., Martin, C., Imhoff, B., and D. Fedyk, "Basic Specification for IP Fast-Reroute: Loop-free Alternates", [draft-ietf-rtgwg-ipfrr-spec-base-12](#) (work in progress), March 2008.
- [13] Kompella, K. and S. Amante, "The Use of Entropy Labels in MPLS Forwarding", [draft-kompella-mpls-entropy-label-00](#) (work in progress), July 2008.
- [14] Swallow, G., "Label Switching Router Self-Test", [draft-ietf-mpls-lsr-self-test-07](#) (work in progress), May 2007.
- [15] Stein, Y., Mendelsohn, I., and R. Insler, "PW Bonding", [draft-stein-pwe3-pwbonding-01](#) (work in progress), November 2008.

Authors' Addresses

Stewart Bryant (editor)
Cisco Systems
250 Longwater Ave
Reading RG2 6GB
United Kingdom

Phone: +44-208-824-8828
Email: stbryant@cisco.com

Clarence Filsfils
Cisco Systems
Brussels
Belgium

Email: cfilsfil@cisco.com

Internet-Draft

FAT-PW

March 2009

Ulrich Drafz
Deutsche Telekom
Muenster,
Germany

Phone:
Fax:
Email: Ulrich.Drafz@t-com.net
URI:

Vach Kompella
Alcatel-Lucent

Phone:
Fax:
Email: Alcatel-Lucent vach.kompella@alcatel-lucent.com
URI:

Joe Regan
Alcatel-Lucent

Phone:
Fax:
Email: joe.regan@alcatel-lucent.comRegan
URI:

Shane Amante
Level 3 Communications

Phone:
Fax:
Email: shane@castlepoint.net
URI:

Bryant, et al.

Expires September 3, 2009

[Page 17]