

RTGWG
Internet-Draft
Intended status: Informational
Expires: February 13, 2013

S. Ning
Tata Communications
A. Malis
D. McDysan
Verizon
L. Yong
Huawei USA
C. Villamizar
Outer Cape Cod Network
Consulting
August 12, 2012

Composite Link Use Cases and Design Considerations
draft-ietf-rtgwg-cl-use-cases-01

Abstract

This document provides a set of use cases and design considerations for composite links.

Composite link is a formalization of multipath techniques currently in use in IP and MPLS networks and a set of extensions to multipath techniques.

Status of this Memo

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

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF). Note that other groups may also distribute working documents as Internet-Drafts. The list of current Internet-Drafts is at <http://datatracker.ietf.org/drafts/current/>.

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."

This Internet-Draft will expire on February 13, 2013.

Copyright Notice

Copyright (c) 2012 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

(<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document. Code Components extracted from this document must include Simplified BSD License text as described in Section 4.e of the Trust Legal Provisions and are provided without warranty as described in the Simplified BSD License.

Table of Contents

1.	Introduction	3
2.	Conventions used in this document	3
2.1.	Terminology	3
3.	Composite Link Foundation Use Cases	4
4.	Delay Sensitive Applications	7
5.	Large Volume of IP and LDP Traffic	7
6.	Composite Link and Packet Ordering	8
6.1.	MPLS-TP in network edges only	10
6.2.	Composite Link at core LSP ingress/egress	11
6.3.	MPLS-TP as a MPLS client	12
7.	IANA Considerations	12
8.	Security Considerations	12
9.	Acknowledgments	13
10.	References	13
10.1.	Normative References	13
10.2.	Informative References	13
Appendix A.	More Details on Existing Network Operator Practices and Protocol Usage	15
Appendix B.	Existing Multipath Standards and Techniques	17
B.1.	Common Multipath Load Splitting Techniques	18
B.2.	Simple and Adaptive Load Balancing Multipath	19
B.3.	Traffic Split over Parallel Links	20
B.4.	Traffic Split over Multiple Paths	20
Appendix C.	Characteristics of Transport in Core Networks	20
	Authors' Addresses	22

1. Introduction

Composite link requirements are specified in [\[I-D.ietf-rtgwg-cl-requirement\]](#). A composite link framework is defined in [\[I-D.ietf-rtgwg-cl-framework\]](#).

Multipath techniques have been widely used in IP networks for over two decades. The use of MPLS began more than a decade ago. Multipath has been widely used in IP/MPLS networks for over a decade with very little protocol support dedicated to effective use of multipath.

The state of the art in multipath prior to composite links is documented in [Appendix B](#).

Both Ethernet Link Aggregation [[IEEE-802.1AX](#)] and MPLS link bundling [[RFC4201](#)] have been widely used in today's MPLS networks. Composite link differs in the following characteristics.

1. A composite link allows bundling of non-homogenous links together as a single logical link.
2. A composite link provides more information in the TE-LSDB and supports more explicit control over placement of LSP.

2. 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 [RFC 2119](#) [[RFC2119](#)].

2.1. Terminology

Terminology defined in [\[I-D.ietf-rtgwg-cl-requirement\]](#) is used in this document.

In addition, the following terms are used:

classic multipath:

Classic multipath refers to the most common current practice in implementation and deployment of multipath (see [Appendix A](#)). The most common current practice makes use of a hash on the MPLS label stack and if IPv4 or IPv6 are indicated under the label stack, makes use of the IP source and destination addresses [[RFC4385](#)] [[RFC4928](#)].

classic link bundling:

Classic link bundling refers to the use of [[RFC4201](#)] where the "all ones" component is not used. Where the "all ones" component is used, link bundling behaves as classic multipath does. Classic link bundling selects a single component link on which to put any given LSP.

Among the important distinctions between classic multipath or classic link bundling and Composite Link are:

1. Classic multipath has no provision to retain order among flows within a subset of LSP. Classic link bundling retains order among all flows but as a result does a poor job of splitting load among components and therefore is rarely (if ever) deployed. Composite Link allows per LSP control of load split characteristics.
2. Classic multipath and classic link bundling do not provide a means to put some LSP on component links with lower delay. Composite Link does.
3. Classic multipath will provide a load balance for IP and LDP traffic. Classic link bundling will not. Neither classic multipath or classic link bundling will measure IP and LDP traffic and reduce the advertised "Available Bandwidth" as a result of that measurement. Composite Link better supports RSVP-TE used with significant traffic levels of native IP and native LDP.
4. Classic link bundling cannot support an LSP that is greater in capacity than any single component link. Classic multipath and Composite Link support this capability but will reorder traffic on such an LSP. Composite Link can retain order of an LSP that is carried within an LSP that is greater in capacity than any single component link if the contained LSP has such a requirement.

None of these techniques, classic multipath, classic link bundling, or Composite Link, will reorder traffic among IP microflows. None of these techniques will reorder traffic among PW, if a PWE3 Control Word is used [[RFC4385](#)].

3. Composite Link Foundation Use Cases

A simple composite link composed entirely of physical links is illustrated in Figure 1, where a composite link is configured between LSR1 and LSR2. This composite link has three component links.

Individual component links in a composite link may be supported by different transport technologies such as wavelength, Ethernet VLAN. Even if the transport technology implementing the component links is identical, the characteristics (e.g., bandwidth, latency) of the component links may differ.

The composite link in Figure 1 may carry LSP traffic flows and control plane packets. Control plane packets may appear as IP packets or may be carried within a generic associated channel (G-Ach) [RFC5586]. A LSP may be established over the link by either RSVP-TE [RFC3209] or LDP [RFC5036] signaling protocols. All component links in a composite link are summarized in the same forwarding adjacency LSP (FA-LSP) routing advertisement [RFC3945]. The composite link is summarized as one TE-Link advertised into the IGP by the composite link end points. This information is used in path computation when a full MPLS control plane is in use. The individual component links or groups of component links may optionally be advertised into the IGP as sub-TLV of the composite link advertisement to indicate capacity available with various characteristics, such as a delay range.

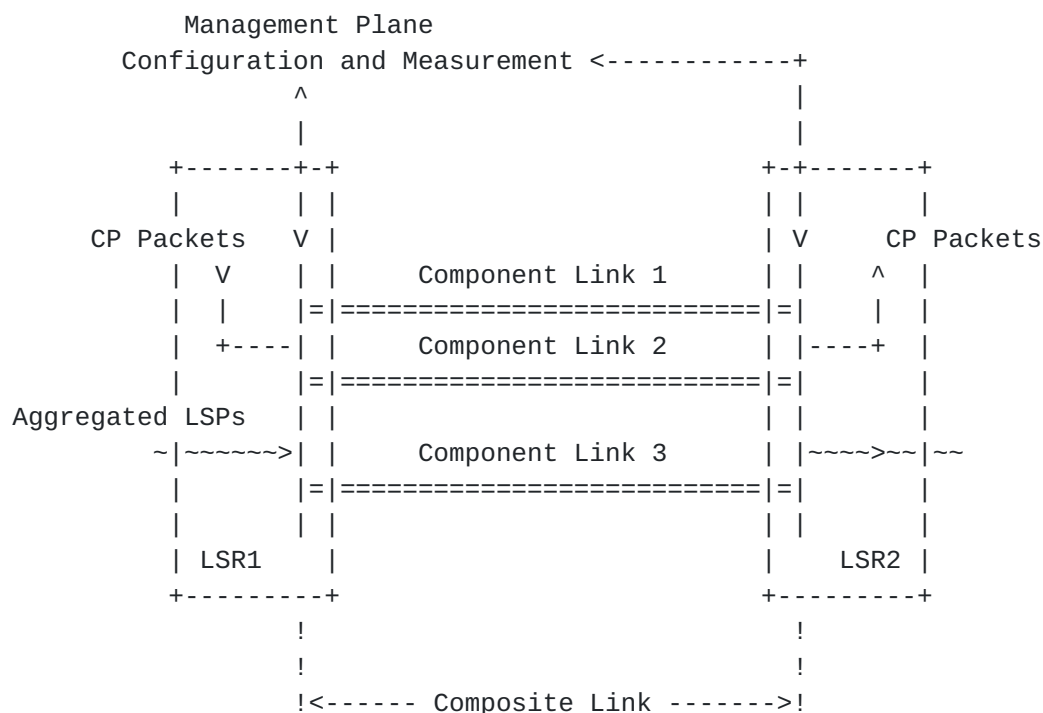


Figure 1: a composite link constructed with multiple physical links between two LSR

[I-D.ietf-rtgwg-cl-requirement] specifies that component links may themselves be composite links. Figure 2 shows three forms of component links which may be deployed in a network.

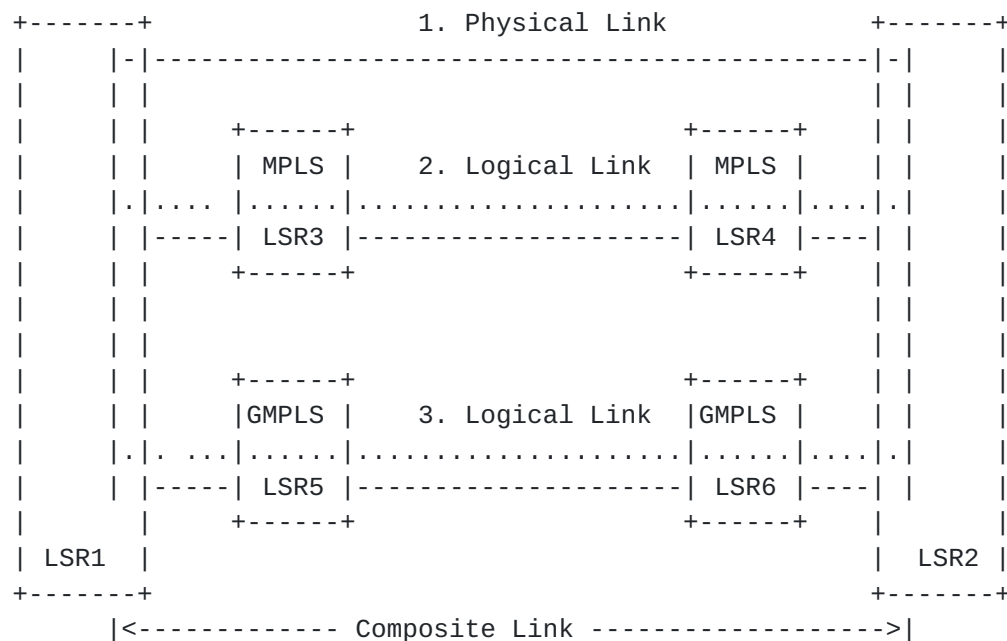


Figure 2: Illustration of Various Component Link Types

The three forms of component link shown in Figure 2 are:

1. The first component link is configured with direct physical media.
2. The second component link is a TE tunnel that traverses LSR3 and LSR4, where LSR3 and LSR4 are the nodes supporting MPLS, but supporting few or no GMPLS extensions.
3. The third component link is formed by lower layer network that has GMPLS enabled. In this case, LSR5 and LSR6 are not the nodes controlled by the MPLS but provide the connectivity for the component link.

A composite link forms one logical link between connected LSR and is used to carry aggregated traffic [[I-D.ietf-rtgwg-cl-requirement](#)]. Composite link relies on its component links to carry the traffic over the composite link. The endpoints of the composite link maps incoming traffic into component links.

For example, LSR1 in Figure 1 distributes the set of traffic flows including control plane packets among the set of component links. LSR2 in Figure 1 receives the packets from its component links and sends them to MPLS forwarding engine with no attempt to reorder packets arriving on different component links. The traffic in the opposite direction, from LSR2 to LSR1, is distributed across the set of component links by the LSR2.

These three forms of component link are only example. Many other examples are possible. A component link may itself be a composite link. A segment of an LSP (single hop for that LSP) may be a composite link.

4. Delay Sensitive Applications

Most applications benefit from lower delay. Some types of applications are far more sensitive than others. For example, real time bidirectional applications such as voice communication or two way video conferencing are far more sensitive to delay than unidirectional streaming audio or video. Non-interactive bulk transfer is almost insensitive to delay if a large enough TCP window is used.

Some applications are sensitive to delay but unwilling to pay extra to insure lower delay. For example, many SIP end users are willing to accept the delay offered to best effort services as long as call quality is good most of the time.

Other applications are sensitive to delay and willing to pay extra to insure lower delay. For example, financial trading applications are extremely sensitive to delay and with a lot at stake are willing to go to great lengths to reduce delay.

Among the requirements of Composite Link are requirements to advertise capacity available within configured ranges of delay within a given composite link and the support the ability to place an LSP only on component links that meeting that LSP's delay requirements.

The Composite Link requirements to accommodate delay sensitive applications are analogous to diffserv requirements to accommodate applications requiring higher quality of service on the same infrastructure as applications with less demanding requirements. The ability to share capacity with less demanding applications, with best effort applications being the least demanding, can greatly reduce the cost of delivering service to the more demanding applications.

5. Large Volume of IP and LDP Traffic

IP and LDP do not support traffic engineering. Both make use of a shortest (lowest routing metric) path, with an option to use equal cost multipath (ECMP). Note that though ECMP is prohibited in LDP specifications, it is widely implemented. Where implemented for LDP, ECMP is generally disabled by default for standards compliance, but often enabled in LDP deployments.

Without traffic engineering capability, there must be sufficient capacity to accomodate the IP and LDP traffic. If not, persistent queuing delay and loss will occur. Unlike RSVP-TE, a subset of traffic cannot be routed using constraint based routing to avoid a congested portion of an infrastructure.

In existing networks which accomodate IP and/or LDP with RSVP-TE, either the IP and LDP can be carried over RSVP-TE, or where the traffic contribution of IP and LDP is small, IP and LDP can be carried native and the effect on RSVP-TE can be ignored. Ignoring the traffic contribution of IP is certainly valid on high capacity networks where native IP is used primarily for control and network management and customer IP is carried within RSVP-TE.

Where it is desireable to carry native IP and/or LDP and IP and/or LDP traffic volumes are not negligible, RSVP-TE needs improvement. The enhancement offered by Composite Link is an ability to measure the IP and LDP, filter the measurements, and reduce the capacity available to RSVP-TE to avoid congestion. The treatment given to the IP or LDP traffic is similar to the treatment when using the "auto-bandwidth" feature in some RSVP-TE implementations on that same traffic, and giving a higher priority (numerically lower setup priority and holding priority value) to the "auto-bandwidth" LSP. The difference is that the measurement is made at each hop and the reduction in advertised bandwidth is made more directly.

6. Composite Link and Packet Ordering

A strong motivation for Composite Link is the need to provide LSP capacity in IP backbones that exceeds the capacity of single wavelengths provided by transport equipment and exceeds the practical capacity limits acheivable through inverse multiplexing. [Appendix C](#) describes characteristics and limitations of transport systems today. [Section 2](#) defines the terms "classic multipath" and "classic link bundling" used in this section.

For purpose of discussion, consider two very large cities, city A and city Z. For example, in the US high traffic cities might be New York and Los Angeles and in Europe high traffic cities might be London and Amsterdam. Two other high volume cities, city B and city Y may share common provider core network infrastructure. Using the same examples, the city B and Y may Washington DC and San Francisco or Paris and Stockholm. In the US, the common infrastructure may span Denver, Chicago, Detroit, and Cleveland. Other major traffic contributors on either US coast include Boston, northern Virginia on the east coast, and Seattle, and San Diego on the west coast. The capacity of IP/MPLS links within the shared infrastructure, for

example city to city links in the Denver, Chicago, Detroit, and Cleveland path in the US example, have capacities for most of the 2000s decade that greatly exceeded single circuits available in transport networks.

For a case with four large traffic sources on either side of the shared infrastructure, up to sixteen core city to core city traffic flows in excess of transport circuit capacity may be accommodated on the shared infrastructure.

Today the most common IP/MPLS core network design makes use of very large links which consist of many smaller component links, but use classic multipath techniques rather than classic link bundling or Composite Link. A component link typically corresponds to the largest circuit that the transport system is capable of providing (or the largest cost effective circuit). IP source and destination address hashing is used to distribute flows across the set of component links as described in [Appendix B.3](#).

Classic multipath can handle large LSP up to the total capacity of the multipath (within limits, see [Appendix B.2](#)). A disadvantage of classic multipath is the reordering among traffic within a given core city to core city LSP. While there is no reordering within any microflow and therefore no customer visible issue, MPLS-TP cannot be used across an infrastructure where classic multipath is in use, except within pseudowires.

These capacity issues force the use of classic multipath today. Classic multipath excludes a direct use of MPLS-TP. The desire for OAM, offered by MPLS-TP, is in conflict with the use of classic multipath. There are a number of alternatives that satisfy both requirements. Some alternatives are described below.

MPLS-TP in network edges only

A simple approach which requires no change to the core is to disallow MPLS-TP across the core unless carried within a pseudowire (PW). MPLS-TP may be used within edge domains where classic multipath is not used. PW may be signaled end to end using single segment PW (SS-PW), or stitched across domains using multisegment PW (MS-PW). The PW and anything carried within the PW may use OAM as long as fat-PW [[RFC6391](#)] load splitting is not used by the PW.

Composite Link at core LSP ingress/egress

The interior of the core network may use classic link bundling, with the limitation that no LSP can exceed the capacity of a

single circuit. Larger non-MPLS-TP LSP can be configured using multiple ingress to egress component MPLS-TP LSP. This can be accomplished using existing IP source and destination address hashing configured at LSP ingress and egress, or using Composite Link configured at ingress and egress. Each component LSP, if constrained to be no larger than the capacity of a single circuit, can make use of MPLS-TP and offer OAM for all top level LSP across the core.

MPLS-TP as a MPLS client

A third approach involves modifying the behavior of LSR in the interior of the network core, such that MPLS-TP can be used on a subset of LSP, where the capacity of any one LSP within that MPLS-TP subset of LSP is not larger than the capacity of a single circuit. This requirement is accommodated through a combination of signaling to indicate LSP for which traffic splitting needs to be constrained, the ability to constrain the depth of the label stack over which traffic splitting can be applied on a per LSP basis, and the ability to constrain the use of IP addresses below the label stack for traffic splitting also on a per LSP basis.

The above list of alternatives allow packet ordering within an LSP to be maintained in some circumstances and allow very large LSP capacities. Each of these alternatives are discussed further in the following subsections.

6.1. MPLS-TP in network edges only

Classic MPLS link bundling is defined in [[RFC4201](#)] and has existed since early in the 2000s decade. Classic MPLS link bundling place any given LSP entirely on a single component link. Classic MPLS link bundling is not in widespread use as the means to accomodate large link capacities in core networks due to the simplicity and better multiplexing gain, and therefore lower network cost of classic multipath.

If MPLS-TP OAM capability in the IP/MPLS network core LSP is not required, then there is no need to change existing network designs which use classic multipath and both label stack and IP source and destination address based hashing as a basis for load splitting.

If MPLS-TP is needed for a subset of LSP, then those LSP can be carried within pseudowires. The pseudowires adds a thin layer of encapsulation and therefore a small overhead. If only a subset of LSP need MPLS-TP OAM, then some LSP must make use of the pseudowires and other LSP avoid them. A straightforward way to accomplish this is with administrative attributes [[RFC3209](#)].

6.2. Composite Link at core LSP ingress/egress

Composite Link can be configured only for large LSP that are made of smaller MPLS-TP component LSP. This approach is capable of supporting MPLS-TP OAM over the entire set of component link LSP and therefore the entire set of top level LSP traversing the core.

There are two primary disadvantage of this approach. One is the number of top level LSP traversing the core can be dramatically increased. The other disadvantage is the loss of multiplexing gain that results from use of classic link bundling within the interior of the core network.

If component LSP use MPLS-TP, then no component LSP can exceed the capacity of a single circuit. For a given composite LSP there can either be a number of equal capacity component LSP or some number of full capacity component links plus one LSP carrying the excess. For example, a 350 Gb/s composite LSP over a 100 Gb/s infrastructure may use five 70 Gb/s component LSP or three 100 Gb/s LSP plus one 50 Gb/s LSP. Classic MPLS link bundling is needed to support MPLS-TP and suffers from a bin packing problem even if LSP traffic is completely predictable, which it never is in practice.

The common means of setting composite link bandwidth parameters uses long term statistical measures. For example, many providers base their LSP bandwidth parameters on the 95th percentile of carried traffic as measured over a one week period. It is common to add 10-30% to the 95th percentile value measured over the prior week and adjust bandwidth parameters of LSP weekly. It is also possible to measure traffic flow at the LSR and adjust bandwidth parameters somewhat more dynamically. This is less common in deployments and where deployed, make use of filtering to track very long term trends in traffic levels. In either case, short term variation of traffic levels relative to signaled LSP capacity are common. Allowing a large overallocation of LSP bandwidth parameters (ie: adding 30% or more) avoids overutilization of any given LSP, but increases unused network capacity and increases network cost. Allowing a small overallocation of LSP bandwidth parameters (ie: 10-20% or less) results in both underutilization and overutilization but statistically results in a total utilization within the core that is under capacity most or all of the time.

The classic multipath solution accommodates the situation in which some composite LSP are underutilizing their signaled capacity and others are overutilizing their capacity with the need for far less unused network capacity to accommodate variation in actual traffic levels. If the actual traffic levels of LSP can be described by a probability distribution, the variation of the sum of LSP is less

than the variation of any given LSP for all but a constant traffic level (where the variation of the sum and the components are both zero).

There are two situations which can motivate the use of this approach. This design is favored if the provider values MPLS-TP OAM across the core more than efficiency (or is unaware of the efficiency issue). This design can also make sense if transport equipment or very low cost core LSR are available which support only classic link bundling and regardless of loss of multiplexing gain, are more cost effective at carrying transit traffic than using equipment which supports IP source and destination address hashing.

6.3. MPLS-TP as a MPLS client

Accommodating MPLS-TP as a MPLS client requires a small change to forwarding behavior and is therefore most applicable to major network overbuilds or new deployments. The change to forwarding is an ability to limit the depth of MPLS labels used in hashing on the label stack on a per LSP basis. Some existing hardware, particularly microprogrammed hardware, may be able to accommodate this forwarding change. Providing support in new hardware is not difficult, a much smaller change than, for example, changes required to disable PHP in an environment where LSP hierarchy is used.

The advantage of this approach is an ability to accommodate MPLS-TP as a client LSP but retain the high multiplexing gain and therefore efficiency and low network cost of a pure MPLS deployment. The disadvantage is the need for a small change in forwarding.

7. IANA Considerations

This memo includes no request to IANA.

8. Security Considerations

This document is a use cases document. Existing protocols are referenced such as MPLS. Existing techniques such as MPLS link bundling and multipath techniques are referenced. These protocols and techniques are documented elsewhere and contain security considerations which are unchanged by this document.

This document also describes use cases for Composite Link, which is a work-in-progress. Composite Link requirements are defined in [\[I-D.ietf-rtgwg-cl-requirement\]](#). [\[I-D.ietf-rtgwg-cl-framework\]](#) defines a framework for Composite Link. Composite Link bears many

similarities to MPLS link bundling and multipath techniques used with MPLS. Additional security considerations, if any, beyond those already identified for MPLS, MPLS link bundling and multipath techniques, will be documented in the framework document if specific to the overall framework of Composite Link, or in protocol extensions if specific to a given protocol extension defined later to support Composite Link.

9. Acknowledgments

Authors would like to thank [no one so far] for their reviews and great suggestions.

In the interest of full disclosure of affiliation and in the interest of acknowledging sponsorship, past affiliations of authors are noted. Much of the work done by Ning So occurred while Ning was at Verizon. Much of the work done by Curtis Villamizar occurred while at Infinera. Infinera continues to sponsor this work on a consulting basis.

10. References

10.1. Normative References

[RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", [BCP 14](#), [RFC 2119](#), March 1997.

10.2. Informative References

- [I-D.ietf-rtgwg-cl-framework]
Ning, S., McDysan, D., Osborne, E., Yong, L., and C. Villamizar, "Composite Link Framework in Multi Protocol Label Switching (MPLS)", [draft-ietf-rtgwg-cl-framework-00](#) (work in progress), August 2012.
- [I-D.ietf-rtgwg-cl-requirement]
Villamizar, C., McDysan, D., Ning, S., Malis, A., and L. Yong, "Requirements for MPLS Over a Composite Link", [draft-ietf-rtgwg-cl-requirement-07](#) (work in progress), June 2012.
- [IEEE-802.1AX]
IEEE Standards Association, "IEEE Std 802.1AX-2008 IEEE Standard for Local and Metropolitan Area Networks - Link Aggregation", 2006, <<http://standards.ieee.org/getieee802/download/802.1AX-2008.pdf>>.

[ITU-T.G.694.2]

ITU-T, "Spectral grids for WDM applications: CWDM wavelength grid", 2003,
<<http://www.itu.int/rec/T-REC-G.694.2-200312-I>>.

[ITU-T.G.800]

ITU-T, "Unified functional architecture of transport networks", 2007,
<<http://www.itu.int/rec/T-REC-G.800-200709-I>>.

[ITU-T.Y.1540]

ITU-T, "Internet protocol data communication service - IP packet transfer and availability performance parameters", 2007, <<http://www.itu.int/rec/T-REC-Y.1540/en>>.

[ITU-T.Y.1541]

ITU-T, "Network performance objectives for IP-based services", 2006, <<http://www.itu.int/rec/T-REC-Y.1541/en>>.

[RFC1717] Sklower, K., Lloyd, B., McGregor, G., and D. Carr, "The PPP Multilink Protocol (MP)", [RFC 1717](#), November 1994.

[RFC2475] Blake, S., Black, D., Carlson, M., Davies, E., Wang, Z., and W. Weiss, "An Architecture for Differentiated Services", [RFC 2475](#), December 1998.

[RFC2597] Heinanen, J., Baker, F., Weiss, W., and J. Wroclawski, "Assured Forwarding PHB Group", [RFC 2597](#), June 1999.

[RFC2615] Malis, A. and W. Simpson, "PPP over SONET/SDH", [RFC 2615](#), June 1999.

[RFC2991] Thaler, D. and C. Hopps, "Multipath Issues in Unicast and Multicast Next-Hop Selection", [RFC 2991](#), November 2000.

[RFC2992] Hopps, C., "Analysis of an Equal-Cost Multi-Path Algorithm", [RFC 2992](#), November 2000.

[RFC3209] Awduche, D., Berger, L., Gan, D., Li, T., Srinivasan, V., and G. Swallow, "RSVP-TE: Extensions to RSVP for LSP Tunnels", [RFC 3209](#), December 2001.

[RFC3260] Grossman, D., "New Terminology and Clarifications for Diffserv", [RFC 3260](#), April 2002.

[RFC3809] Nagarajan, A., "Generic Requirements for Provider Provisioned Virtual Private Networks (PPVPN)", [RFC 3809](#), June 2004.

- [RFC3945] Mannie, E., "Generalized Multi-Protocol Label Switching (GMPLS) Architecture", [RFC 3945](#), October 2004.
- [RFC4201] Kompella, K., Rekhter, Y., and L. Berger, "Link Bundling in MPLS Traffic Engineering (TE)", [RFC 4201](#), October 2005.
- [RFC4301] Kent, S. and K. Seo, "Security Architecture for the Internet Protocol", [RFC 4301](#), December 2005.
- [RFC4385] 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.
- [RFC4928] Swallow, G., Bryant, S., and L. Andersson, "Avoiding Equal Cost Multipath Treatment in MPLS Networks", [BCP 128](#), [RFC 4928](#), June 2007.
- [RFC5036] Andersson, L., Minei, I., and B. Thomas, "LDP Specification", [RFC 5036](#), October 2007.
- [RFC5586] Bocci, M., Vigoureux, M., and S. Bryant, "MPLS Generic Associated Channel", [RFC 5586](#), June 2009.
- [RFC6391] Bryant, S., Filsfils, C., Drafz, U., Kompella, V., Regan, J., and S. Amante, "Flow-Aware Transport of Pseudowires over an MPLS Packet Switched Network", [RFC 6391](#), November 2011.

Appendix A. More Details on Existing Network Operator Practices and Protocol Usage

Often, network operators have a contractual Service Level Agreement (SLA) with customers for services that are comprised of numerical values for performance measures, principally availability, latency, delay variation. Additionally, network operators may have Service Level Specification (SLS) that is for internal use by the operator. See [[ITU-T.Y.1540](#)], [[ITU-T.Y.1541](#)], [RFC3809, Section 4.9](#) [[RFC3809](#)] for examples of the form of such SLA and SLS specifications. In this document we use the term Network Performance Objective (NPO) as defined in section 5 of [[ITU-T.Y.1541](#)] since the SLA and SLS measures have network operator and service specific implications. Note that the numerical NPO values of Y.1540 and Y.1541 span multiple networks and may be looser than network operator SLA or SLS objectives. Applications and acceptable user experience have an important relationship to these performance parameters.

Consider latency as an example. In some cases, minimizing latency

relates directly to the best customer experience (e.g., in TCP closer is faster). In other cases, user experience is relatively insensitive to latency, up to a specific limit at which point user perception of quality degrades significantly (e.g., interactive human voice and multimedia conferencing). A number of NPOs have a bound on point-point latency, and as long as this bound is met, the NPO is met -- decreasing the latency is not necessary. In some NPOs, if the specified latency is not met, the user considers the service as unavailable. An unprotected LSP can be manually provisioned on a set of to meet this type of NPO, but this lowers availability since an alternate route that meets the latency NPO cannot be determined.

Historically, when an IP/MPLS network was operated over a lower layer circuit switched network (e.g., SONET rings), a change in latency caused by the lower layer network (e.g., due to a maintenance action or failure) this was not known to the MPLS network. This resulted in latency affecting end user experience, sometimes violating NPOs or resulting in user complaints.

A response to this problem was to provision IP/MPLS networks over unprotected circuits and set the metric and/or TE-metric proportional to latency. This resulted in traffic being directed over the least latency path, even if this was not needed to meet an NPO or meet user experience objectives. This results in reduced flexibility and increased cost for network operators. Using lower layer networks to provide restoration and grooming is expected to be more efficient, but the inability to communicate performance parameters, in particular latency, from the lower layer network to the higher layer network is an important problem to be solved before this can be done.

Latency NPOs for point-to-point services are often tied closely to geographic locations, while latency for multipoint services may be based upon a worst case within a region.

Section 7 of [\[ITU-T.Y.1540\]](#) defines availability for an IP service in terms of loss exceeding a threshold for a period on the order of 5 minutes. However, the timeframes for restoration (i.e., as implemented by pre-determined protection, convergence of routing protocols and/or signaling) for services range from on the order of 100 ms or less (e.g., for VPWS to emulate classical SDH/SONET protection switching), to several minutes (e.g., to allow BGP to reconverge for L3VPN) and may differ among the set of customers within a single service.

The presence of only three Traffic Class (TC) bits (previously known as EXP bits) in the MPLS shim header is limiting when a network operator needs to support QoS classes for multiple services (e.g., L2VPN VPWS, VPLS, L3VPN and Internet), each of which has a set of QoS

classes that need to be supported. In some cases one bit is used to indicate conformance to some ingress traffic classification, leaving only two bits for indicating the service QoS classes. The approach that has been taken is to aggregate these QoS classes into similar sets on LER-LSR and LSR-LSR links.

Labeled LSPs and use of link layer encapsulation have been standardized in order to provide a means to meet these needs.

The IP DSCP cannot be used for flow identification since [RFC 4301 Section 5.5](#) [[RFC4301](#)] requires Diffserv transparency, and in general network operators do not rely on the DSCP of Internet packets. In addition, the use of IP DSCP for flow identification is incompatible with Assured Forwarding services [[RFC2597](#)] or any other service which may use more than one DSCP code point to carry traffic for a given microflow.

A label is pushed onto Internet packets when they are carried along with L2/L3VPN packets on the same link or lower layer network provides a mean to distinguish between the QoS class for these packets.

Operating an MPLS-TE network involves a different paradigm from operating an IGP metric-based LDP signaled MPLS network. The multipoint-to-point LDP signaled MPLS LSPs occur automatically, and balancing across parallel links occurs if the IGP metrics are set "equally" (with equality a locally definable relation).

Traffic is typically comprised of a few large (some very large) flows and many small flows. In some cases, separate LSPs are established for very large flows. This can occur even if the IP header information is inspected by a LSR, for example an IPsec tunnel that carries a large amount of traffic. An important example of large flows is that of a L2/L3 VPN customer who has an access line bandwidth comparable to a client-client composite link bandwidth -- there could be flows that are on the order of the access line bandwidth.

[Appendix B](#). Existing Multipath Standards and Techniques

Today the requirement to handle large aggregations of traffic, much larger than a single component link, can be handled by a number of techniques which we will collectively call multipath. Multipath applied to parallel links between the same set of nodes includes Ethernet Link Aggregation [[IEEE-802.1AX](#)], link bundling [[RFC4201](#)], or other aggregation techniques some of which may be vendor specific. Multipath applied to diverse paths rather than parallel links

includes Equal Cost MultiPath (ECMP) as applied to OSPF, ISIS, or even BGP, and equal cost LSP, as described in [Appendix B.4](#). Various multipath techniques have strengths and weaknesses.

the term Composite Link is more general than terms such as Link Aggregation which is generally considered to be specific to Ethernet and its use here is consistent with the broad definition in [\[ITU-T.G.800\]](#). The term multipath excludes inverse multiplexing and refers to techniques which only solve the problem of large aggregations of traffic, without addressing the other requirements outlined in this document, particularly those described in [Section 4](#) and [Section 5](#).

[B.1](#). Common Multipath Load Splitting Techniques

Identical load balancing techniques are used for multipath both over parallel links and over diverse paths.

Large aggregates of IP traffic do not provide explicit signaling to indicate the expected traffic loads. Large aggregates of MPLS traffic are carried in MPLS tunnels supported by MPLS LSP. LSP which are signaled using RSVP-TE extensions do provide explicit signaling which includes the expected traffic load for the aggregate. LSP which are signaled using LDP do not provide an expected traffic load.

MPLS LSP may contain other MPLS LSP arranged hierarchically. When an MPLS LSR serves as a midpoint LSR in an LSP carrying other LSP as payload, there is no signaling associated with these inner LSP. Therefore even when using RSVP-TE signaling there may be insufficient information provided by signaling to adequately distribute load based solely on signaling.

Generally a set of label stack entries that is unique across the ordered set of label numbers in the label stack can safely be assumed to contain a group of flows. The reordering of traffic can therefore be considered to be acceptable unless reordering occurs within traffic containing a common unique set of label stack entries. Existing load splitting techniques take advantage of this property in addition to looking beyond the bottom of the label stack and determining if the payload is IPv4 or IPv6 to load balance traffic accordingly.

MPLS-TP OAM violates the assumption that it is safe to reorder traffic within an LSP. If MPLS-TP OAM is to be accommodated, then existing multipath techniques must be modified. Such modifications are outside the scope of this document.

For example, a large aggregate of IP traffic may be subdivided into a

large number of groups of flows using a hash on the IP source and destination addresses. This is as described in [RFC2475] and clarified in [RFC3260]. For MPLS traffic carrying IP, a similar hash can be performed on the set of labels in the label stack. These techniques are both examples of means to subdivide traffic into groups of flows for the purpose of load balancing traffic across aggregated link capacity. The means of identifying a set of flows should not be confused with the definition of a flow.

Discussion of whether a hash based approach provides a sufficiently even load balance using any particular hashing algorithm or method of distributing traffic across a set of component links is outside of the scope of this document.

The current load balancing techniques are referenced in [RFC4385] and [RFC4928]. The use of three hash based approaches are described in [RFC2991] and [RFC2992]. A mechanism to identify flows within PW is described in [RFC6391]. The use of hash based approaches is mentioned as an example of an existing set of techniques to distribute traffic over a set of component links. Other techniques are not precluded.

B.2. Simple and Adaptive Load Balancing Multipath

Simple multipath generally relies on the mathematical probability that given a very large number of small microflows, these microflows will tend to be distributed evenly across a hash space. Early very simple multipath implementations assumed that all component links are of equal capacity and perform a modulo operation across the hashed value. An alternate simple multipath technique uses a table generally with a power of two size, and distributes the table entries proportionally among component links according to the capacity of each component link.

Simple load balancing works well if there are a very large number of small microflows (i.e., microflow rate is much less than component link capacity). However, the case where there are even a few large microflows is not handled well by simple load balancing.

An adaptive load balancing multipath technique is one where the traffic bound to each component link is measured and the load split is adjusted accordingly. As long as the adjustment is done within a single network element, then no protocol extensions are required and there are no interoperability issues.

Note that if the load balancing algorithm and/or its parameters is adjusted, then packets in some flows may be briefly delivered out of sequence, however in practice such adjustments can be made very

infrequent.

B.3. Traffic Split over Parallel Links

The load splitting techniques defined in [Appendix B.1](#) and [Appendix B.2](#) are both used in splitting traffic over parallel links between the same pair of nodes. The best known technique, though far from being the first, is Ethernet Link Aggregation [[IEEE-802.1AX](#)]. This same technique had been applied much earlier using OSPF or ISIS Equal Cost MultiPath (ECMP) over parallel links between the same nodes. Multilink PPP [[RFC1717](#)] uses a technique that provides inverse multiplexing, however a number of vendors had provided proprietary extensions to PPP over SONET/SDH [[RFC2615](#)] that predated Ethernet Link Aggregation but are no longer used.

Link bundling [[RFC4201](#)] provides yet another means of handling parallel LSP. [RFC4201](#) explicitly allow a special value of all ones to indicate a split across all members of the bundle. This "all ones" component link is signaled in the MPLS RESV to indicate that the link bundle is making use of classic multipath techniques.

B.4. Traffic Split over Multiple Paths

OSPF or ISIS Equal Cost MultiPath (ECMP) is a well known form of traffic split over multiple paths that may traverse intermediate nodes. ECMP is often incorrectly equated to only this case, and multipath over multiple diverse paths is often incorrectly equated to ECMP.

Many implementations are able to create more than one LSP between a pair of nodes, where these LSP are routed diversely to better make use of available capacity. The load on these LSP can be distributed proportionally to the reserved bandwidth of the LSP. These multiple LSP may be advertised as a single PSC FA and any LSP making use of the FA may be split over these multiple LSP.

Link bundling [[RFC4201](#)] component links may themselves be LSP. When this technique is used, any LSP which specifies the link bundle may be split across the multiple paths of the LSP that comprise the bundle.

Appendix C. Characteristics of Transport in Core Networks

The characteristics of primary interest are the capacity of a single circuit and the use of wave division multiplexing (WDM) to provide a large number of parallel circuits.

Wave division multiplexing (WDM) supports multiple independent channels (independent ignoring crosstalk noise) at slightly different wavelengths of light, multiplexed onto a single fiber. Typical in the early 2000s was 40 wavelengths of 10 Gb/s capacity per wavelength. These wavelengths are in the C-band range, which is about 1530-1565 nm, though some work has been done using the L-band 1565-1625 nm.

The C-band has been carved up using a 100 GHz spacing from 191.7 THz to 196.1 THz by [[ITU-T.G.694.2](#)]. This yields 44 channels. If the outermost channels are not used, due to poorer transmission characteristics, then typically 40 are used. For practical reasons, a 50 GHz or 25 GHz spacing is used by more recent equipment, yielding 80 or 160 channels in practice.

The early optical modulation techniques used within a single channel yielded 2.5Gb/s and 10 Gb/s capacity per channel. As modulation techniques have improved 40 Gb/s and 100 Gb/s per channel have been achieved.

The 40 channels of 10 Gb/s common in the mid 2000s yields a total of 400 Gb/s. Tighter spacing and better modulations are yielding up to 8 Tb/s or more in more recent systems.

Over the optical is an electrical encoding. In the 1990s this was typically Synchronous Optical Networking (SONET) or Synchronous Digital Hierarchy (SDH), with a maximum defined circuit capacity of 40 Gb/s (OC-768), though the 10 Gb/s OC-192 is more common. More recently the low level electrical encoding has been Optical Transport Network (OTN) defined by ITU-T. OTN currently defines circuit capacities up to a nominal 100 Gb/s (ODU4). Both SONET/SDH and OTN make use of time division multiplexing (TDM) where the a higher capacity circuit such as a 100 Gb/s ODU4 in OTN may be subdivided into lower fixed capacity circuits such as ten 10 Gb/s ODU2.

In the 1990s, all IP and later IP/MPLS networks either used a fraction of maximum circuit capacity, or at most the full circuit capacity toward the end of the decade, when full circuit capacity was 2.5 Gb/s or 10 Gb/s. Beyond 2000, the TDM circuit multiplexing capability of SONET/SDH or OTN was rarely used.

Early in the 2000s both transport equipment and core LSR offered 40 Gb/s SONET OC-768. However 10 Gb/s transport equipment was predominantly deployed throughout the decade, partially because LSR 10GbE ports were far more cost effective than either OC-192 or OC-768 and became practical in the second half of the decade.

Entering the 2010 decade, LSR 40GbE and 100GbE are expected to become

widely available and cost effective. Slightly preceeding this transport equipment making use of 40 Gb/s and 100 Gb/s modulations are becoming available. This transport equipment is capable of carrying 40 Gb/s ODU3 and 100 Gb/s ODU4 circuits.

Early in the 2000s decade IP/MPLS core networks were making use of single 10 Gb/s circuits. Capacity grew quickly in the first half of the decade but more IP/MPLS core networks had only a small number of IP/MPLS links requiring 4-8 parallel 10 Gb/s circuits. However, the use of multipath was necessary, was deemed the simplest and most cost effective alternative, and became thoroughly entrenched. By the end of the 2000s decade nearly all major IP/MPLS core service provider networks and a few content provider networks had IP/MPLS links which exceeded 100 Gb/s, long before 40GbE was available and 40 Gb/s transport in widespread use.

It is less clear when IP/MPLS LSP exceeded 10 Gb/s, 40 Gb/s, and 100 Gb/s. By 2010, many service providers have LSP in excess of 100 Gb/s, but few are willing to disclose how many LSP have reached this capacity.

At the time of writing 40GbE and 100GbE LSR products are being evaluated by service providers and content providers and are in use in network trials. The cost of components required to deliver 100 GbE products remains high making these products less cost effective. This is expected to change within years.

The important point is that IP/MPLS core network links have long ago exceeded 100 Gb/s and a small number of IP/MPLS LSP exceed 100 Gb/s. By the time 100 Gb/s circuits are widely deployed, IP/MPLS core network links are likely to exceed 1 Tb/s and many IP/MPLS LSP capacities are likely to exceed 100 Gb/s. Therefore multipath techniques are likely here to stay.

Authors' Addresses

So Ning
Tata Communications

Email: ning.so@tatacommunications.com

Andrew Malis
Verizon
117 West St.
Waltham, MA 02451

Phone: +1 781-466-2362
Email: andrew.g.malis@verizon.com

Dave McDysan
Verizon
22001 Loudoun County PKWY
Ashburn, VA 20147

Email: dave.mcdysan@verizon.com

Lucy Yong
Huawei USA
5340 Legacy Dr.
Plano, TX 75025

Phone: +1 469-277-5837
Email: lucy.yong@huawei.com

Curtis Villamizar
Outer Cape Cod Network Consulting

Email: curtis@occnc.com

