

RTGWG
Internet-Draft
Intended status: Informational
Expires: January 12, 2014

C. Villamizar, Ed.
OCCNC, LLC
D. McDysan, Ed.
Verizon
S. Ning
Tata Communications
A. Malis
Verizon
L. Yong
Huawei USA
July 11, 2013

Requirements for Advanced Multipath in MPLS Networks
draft-ietf-rtgwg-cl-requirement-11

Abstract

This document provides a set of requirements for Advanced Multipath in MPLS Networks.

Advanced Multipath is a formalization of multipath techniques currently in use in IP and MPLS networks and a set of extensions to existing 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 January 12, 2014.

Copyright Notice

Copyright (c) 2013 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
1.1.	Requirements Language	3
2.	Definitions	3
3.	Functional Requirements	6
3.1.	Availability, Stability and Transient Response	6
3.2.	Component Links Provided by Lower Layer Networks	7
3.3.	Parallel Component Links with Different Characteristics	8
4.	Derived Requirements	11
5.	Management Requirements	12
6.	Acknowledgements	13
7.	IANA Considerations	13
8.	Security Considerations	13
9.	References	14
9.1.	Normative References	14
9.2.	Informative References	14
	Authors' Addresses	15

1. Introduction

There is often a need to provide large aggregates of bandwidth that are best provided using parallel links between routers or carrying traffic over multiple MPLS LSP. In core networks there is often no alternative since the aggregate capacities of core networks today far exceed the capacity of a single physical link or single packet processing element.

The presence of parallel links, with each link potentially comprised of multiple layers has resulted in additional requirements. Certain services may benefit from being restricted to a subset of the component links or a specific component link, where component link characteristics, such as latency, differ. Certain services require that an LSP be treated as atomic and avoid reordering. Other services will continue to require only that reordering not occur within a microflow as is current practice.

The purpose of this document is to clearly enumerate a set of requirements related to the protocols and mechanisms that provide MPLS based Advanced Multipath. The intent is to first provide a set of functional requirements that are as independent as possible of protocol specifications ([Section 3](#)). For certain functional requirements this document describes a set of derived protocol requirements ([Section 4](#)) and management requirements ([Section 5](#)).

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

Any statement which requires the solution to support some new functionality through use of [[RFC2119](#)] keywords, SHOULD be interpreted as follows. The implementation either MUST or SHOULD support the new functionality depending on the use of either MUST or SHOULD in the requirements statement. The implementation SHOULD in most or all cases allow any new functionality to be individually enabled or disabled through configuration. A service provider or other deployment MAY choose to enable or disable any feature in their network, subject to implementation limitations on sets of features which can be disabled.

2. Definitions

Multipath

The term multipath includes all techniques in which

1. Traffic can take more than one path from one node to a destination.
2. Individual packets take one path only. Packets are not subdivided and reassembled at the receiving end.
3. Packets are not resequenced at the receiving end.
4. The paths may be:
 - a. parallel links between two nodes, or
 - b. may be specific paths across a network to a destination node, or
 - c. may be links or paths to an intermediate node used to reach a common destination.

The paths need not have equal capacity. The paths may or may not have equal cost in a routing protocol.

Advanced Multipath

Advanced Multipath meets the requirements defined in this document. A key capability of advanced multipath is the support of non-homogeneous component links.

Composite Link

The term Composite Link had been a registered trademark of Avici Systems, but was abandoned in 2007. The term composite link is now defined by the ITU in [[ITU-T.G.800](#)]. The ITU definition includes multipath as defined here, plus inverse multiplexing which is explicitly excluded from the definition of multipath.

Inverse Multiplexing

Inverse multiplexing either transmits whole packets and resequences the packets at the receiving end or subdivides packets and reassembles the packets at the receiving end. Inverse multiplexing requires that all packets be handled by a common egress packet processing element and is therefore not useful for very high bandwidth applications.

Component Link

The ITU definition of composite link in [[ITU-T.G.800](#)] and the IETF definition of link bundling in [[RFC4201](#)] both refer to an individual link in the composite link or link bundle as a

component link. The term component link is applicable to all forms of multipath. The IEEE uses the term member rather than component link in Ethernet Link Aggregation [[IEEE-802.1AX](#)].

Client LSP

A client LSP is an LSP which has been set up over a server layer. In the context of this discussion, a client LSP is a LSP which has been set up over a multipath as opposed to an LSP representing the multipath itself or any LSP supporting a component links of that multipath.

Flow

A sequence of packets that should be transferred in order on one component link of a multipath.

Flow identification

The label stack and other information that uniquely identifies a flow. Other information in flow identification may include an IP header, pseudowire (PW) control word, Ethernet MAC address, etc. Note that a client LSP may contain one or more Flows or a client LSP may be equivalent to a Flow. Flow identification is used to locally select a component link, or a path through the network toward the destination.

Load Balance

Load split, load balance, or load distribution refers to subdividing traffic over a set of component links such that load is fairly evenly distributed over the set of component links and certain packet ordering requirements are met. Some existing techniques better achieve these objectives than others.

Performance Objective

Numerical values for performance measures, principally availability, latency, and delay variation. Performance objectives may be related to Service Level Agreements (SLA) as defined in [RFC2475](#) or may be strictly internal. Performance objectives may span links, edge-to-edge, or end-to-end. Performance objectives may span one provider or may span multiple providers.

A Component Link may be a point-to-point physical link (where a "physical link" includes one or more link layer plus a physical layer) or a logical link that preserves ordering in the steady state. A component link may have transient out of order events, but such events must not exceed the network's Performance Objectives. For example, a component link may be comprised of any supportable combination of link layers over a physical layer or over logical sub-layers, including those providing physical layer emulation.

The ingress and egress of a multipath may be midpoint LSRs with respect to a given client LSP. A midpoint LSR does not participate in the signaling of any clients of the client LSP. Therefore, in general, multipath endpoints cannot determine requirements of clients of a client LSP through participation in the signaling of the clients of the client LSP.

The term Advanced Multipath is intended to be used within the context of this document and the related documents, [[I-D.ietf-rtgwg-cl-use-cases](#)] and [[I-D.ietf-rtgwg-cl-framework](#)] and any other related document. Other advanced multipath techniques may in the future arise. If the capabilities defined in this document become commonplace, they would no longer be considered "advanced". Use of the term "advanced multipath" outside this document, if referring to the term as defined here, should indicate Advanced Multipath as defined by this document, citing the current document name. If using another definition of "advanced multipath", documents may optionally clarify that they are not using the term "advanced multipath" as defined by this document if clarification is deemed helpful.

3. Functional Requirements

The Functional Requirements in this section are grouped in subsections starting with the highest priority.

3.1. Availability, Stability and Transient Response

Limiting the period of unavailability in response to failures or transient events is extremely important as well as maintaining stability. The transient period between some service disrupting event and the convergence of the routing and/or signaling protocols MUST occur within a time frame specified by Performance Objective values.

FR#1 An advanced multipath MAY be announced in conjunction with detailed parameters about its component links, such as bandwidth and latency. The advanced multipath SHALL behave as a single IGP adjacency.

FR#2 The solution SHALL provide a means to summarize some routing advertisements regarding the characteristics of an advanced multipath such that the updated protocol mechanisms maintain convergence times within the timeframe needed to meet or no significantly exceed existing Performance Objective for convergence on the same network or convergence on a network with a similar topology.

- FR#3 The solution SHALL ensure that restoration operations happen within the timeframe needed to meet existing Performance Objective for restoration time on the same network or restoration time on a network with a similar topology.
- FR#4 The solution SHALL provide a mechanism to select a path for a flow across a network that contains a number of paths comprised of pairs of nodes connected by advanced multipath in such a way as to automatically distribute the load over the network nodes connected by advanced multipaths while meeting all of the other mandatory requirements stated above. The solution SHOULD work in a manner similar to that of current networks without any advanced multipath protocol enhancements when the characteristics of the individual component links are advertised.
- FR#5 If extensions to existing protocols are specified and/or new protocols are defined, then the solution SHOULD provide a means for a network operator to migrate an existing deployment in a minimally disruptive manner.
- FR#6 Any load balancing solutions MUST NOT oscillate. Some change in path MAY occur. The solution MUST ensure that path stability and traffic reordering continue to meet Performance Objective on the same network or on a network with a similar topology. Since oscillation may cause reordering, there MUST be means to control the frequency of changing the component link over which a flow is placed.
- FR#7 Management and diagnostic protocols MUST be able to operate over advanced multipaths.

Existing scaling techniques used in MPLS networks apply to MPLS networks which support Advanced Multipaths. Scalability and stability are covered in more detail in [\[I-D.ietf-rtgwg-cl-framework\]](#).

3.2. Component Links Provided by Lower Layer Networks

A component link may be supported by a lower layer network. For example, the lower layer may be a circuit switched network or another MPLS network (e.g., MPLS-TP)). The lower layer network may change the latency (and/or other performance parameters) seen by the client layer. Currently, there is no protocol for the lower layer network to inform the higher layer network of a change in a performance parameter. Communication of the latency performance parameter is a very important requirement. Communication of other performance parameters (e.g., delay variation) is desirable.

- FR#8 The solution SHALL specify a protocol means to allow a lower layer server network to communicate latency to the higher layer client network.
- FR#9 The precision of latency reporting SHOULD be configurable. A reasonable default SHOULD be provided. Implementations SHOULD support precision of at least 10% of the one way latencies for latency of 1 ms or more.
- FR#10 The solution SHALL provide a means to limit the latency to meet a Performance Objective target on a per flow basis or group of flow basis, where flows or groups of flows are identifiable in the forwarding plane and are signaled using in the control plane or set up using the management plane.

The Performance Objectives differ across the services, and some services have different Performance Objectives for different QoS classes, for example, one QoS class may have a much larger latency bound than another. Overload can occur which would violate a Performance Objective parameter (e.g., loss) and some remedy to handle this case for an advanced multipath is required.

- FR#11 If the total demand offered by traffic flows exceeds the capacity of the advanced multipath, the solution SHOULD define a means to cause some traffic flows or groups of flows to move to some other point in the network that is not congested. These "preempted flows" may not be restored if there is no uncongested path in the network.

The intent is to measure the predominant latency in uncongested service provider networks, where geographic delay dominates and is on the order of milliseconds or more. The argument for including queuing delay is that it reflects the delay experienced by applications. The argument against including queuing delay is that it if used in routing decisions it can result in routing instability. This tradeoff is discussed in detail in [\[I-D.ietf-rtgwg-cl-framework\]](#).

3.3. Parallel Component Links with Different Characteristics

As one means to provide high availability, network operators deploy a topology in the MPLS network using lower layer networks that have a certain degree of diversity at the lower layer(s). Many techniques have been developed to balance the distribution of flows across component links that connect the same pair of nodes. When the path for a flow can be chosen from a set of candidate nodes connected via advanced multipaths, other techniques have been developed. Refer to

the Appendices in [[I-D.ietf-rtgwg-cl-use-cases](#)] for a description of existing techniques and a set of references.

- FR#12 The solution SHALL measure traffic flows or groups of traffic flows and dynamically select the component link on which to place this traffic in order to balance the load so that no component link in the advanced multipath between a pair of nodes is overloaded.
- FR#13 When a traffic flow is moved from one component link to another in the same advanced multipath between a set of nodes (or sites), it MUST be done so in a minimally disruptive manner.
- FR#14 Load balancing MAY be used during sustained low traffic periods to reduce the number of active component links for the purpose of power reduction.
- FR#15 The solution SHALL provide a means to identify flows whose rearrangement frequency needs to be bounded by a configured value and MUST provide a means to bound the rearrangement frequency for these flows.
- FR#16 The solution SHALL provide a means that communicates whether the flows within an client LSP can be split across multiple component links. The solution SHOULD provide a means to indicate the flow identification field(s) which can be used along the flow path which can be used to perform this function.
- FR#17 The solution SHALL provide a means to indicate that a traffic flow will traverse a component link with the minimum latency value.
- FR#18 The solution SHALL provide a means to indicate that a traffic flow will traverse a component link with a maximum acceptable latency value as specified by protocol.
- FR#19 The solution SHALL provide a means to indicate that a traffic flow will traverse a component link with a maximum acceptable delay variation value as specified by protocol.
- FR#20 The solution SHALL provide a means local to a node that automatically distributes flows across the component links in the advanced multipath such that Performance Objectives are met as described in prior requirements.

- FR#21 The solution SHALL provide a means to distribute flows from a single client LSP across multiple component links to handle at least the case where the traffic carried in an client LSP exceeds that of any component link in the advanced multipath. As defined in [Section 2](#), a flow is a sequence of packets that should be transferred on one component link and should be transferred in order.
- FR#22 The solution SHOULD support the use case where an advanced multipath itself is a component link for a higher order advanced multipath. For example, an advanced multipath comprised of MPLS-TP bi-directional tunnels viewed as logical links could then be used as a component link in yet another advanced multipath that connects MPLS routers.
- FR#23 The solution MUST support an optional means for client LSP signaling to bind a client LSP to a particular component link within an advanced multipath. If this option is not exercised, then a client LSP that is bound to an advanced multipath may be bound to any component link matching all other signaled requirements, and different directions of a bidirectional client LSP can be bound to different component links.
- FR#24 The solution MUST support a means to indicate that both directions of co-routed bidirectional client LSP MUST be bound to the same component link.

A minimally disruptive change implies that as little disruption as is practical occurs. Such a change can be achieved with zero packet loss. A delay discontinuity may occur, which is considered to be a minimally disruptive event for most services if this type of event is sufficiently rare. A delay discontinuity is an example of a minimally disruptive behavior corresponding to current techniques.

A delay discontinuity is an isolated event which may greatly exceed the normal delay variation (jitter). A delay discontinuity has the following effect. When a flow is moved from a current link to a target link with lower latency, reordering can occur. When a flow is moved from a current link to a target link with a higher latency, a time gap can occur. Some flows (e.g., timing distribution, PW circuit emulation) are quite sensitive to these effects. A delay discontinuity can also cause a jitter buffer underrun or overrun affecting user experience in real time voice services (causing an audible click). These sensitivities may be specified in a Performance Objective.

As with any load balancing change, a change initiated for the purpose

of power reduction may be minimally disruptive. Typically the disruption is limited to a change in delay characteristics and the potential for a very brief period with traffic reordering. The network operator when configuring a network for power reduction should weigh the benefit of power reduction against the disadvantage of a minimal disruption.

4. Derived Requirements

This section takes the next step and derives high-level requirements on protocol specification from the functional requirements.

- DR#1 The solution SHOULD attempt to extend existing protocols wherever possible, developing a new protocol only if this adds a significant set of capabilities.
- DR#2 A solution SHOULD extend LDP capabilities to meet functional requirements (without using TE methods as decided in [\[RFC3468\]](#)).
- DR#3 Coexistence of LDP and RSVP-TE signaled LSPs MUST be supported on an advanced multipath. Other functional requirements should be supported as independently of signaling protocol as possible.
- DR#4 When the nodes connected via an advanced multipath are in the same MPLS network topology, the solution MAY define extensions to the IGP.
- DR#5 When the nodes are connected via an advanced multipath are in different MPLS network topologies, the solution SHALL NOT rely on extensions to the IGP.
- DR#6 The solution SHOULD support advanced multipath IGP advertisement that results in convergence time better than that of advertising the individual component links. The solution SHALL be designed so that it represents the range of capabilities of the individual component links such that functional requirements are met, and also minimizes the frequency of advertisement updates which may cause IGP convergence to occur.

Examples of advertisement update triggering events to be considered include: client LSP establishment/release, changes in component link characteristics (e.g., latency, up/down state), and/or bandwidth utilization.

DR#7 When a worst case failure scenario occurs, the number of RSVP-TE client LSPs to be resigned will cause a period of unavailability as perceived by users. The resignaling time of the solution MUST support protocol mechanisms meeting existing provider Performance Objective for the duration of unavailability without significantly relaxing those existing Performance Objectives for the same network or for networks with similar topology. For example, the processing load due to IGP readvertisement MUST NOT increase significantly and the resignaling time of the solution MUST NOT increase significantly as compared with current methods.

5. Management Requirements

- MR#1 Management Plane MUST support polling of the status and configuration of an advanced multipath and its individual advanced multipath and support notification of status change.
- MR#2 Management Plane MUST be able to activate or de-activate any component link in an advanced multipath in order to facilitate operation maintenance tasks. The routers at each end of an advanced multipath MUST redistribute traffic to move traffic from a de-activated link to other component links based on the traffic flow TE criteria.
- MR#3 Management Plane MUST be able to configure a client LSP over an advanced multipath and be able to select a component link for the client LSP.
- MR#4 Management Plane MUST be able to trace which component link a client LSP is assigned to and monitor individual component link and advanced multipath performance.
- MR#5 Management Plane MUST be able to verify connectivity over each individual component link within an advanced multipath.
- MR#6 Component link fault notification MUST be sent to the management plane.
- MR#7 Advanced multipath fault notification MUST be sent to the management plane and MUST be distributed via link state message in the IGP.
- MR#8 Management Plane SHOULD provide the means for an operator to initiate an optimization process.

MR#9 An operator initiated optimization MUST be performed in a minimally disruptive manner as described in [Section 3.3](#).

6. Acknowledgements

Frederic Jouray of France Telecom and Yuji Kamite of NTT Communications Corporation co-authored a version of this document.

A rewrite of this document occurred after the IETF77 meeting. Dimitri Papadimitriou, Lou Berger, Tony Li, the former WG chairs John Scuder and Alex Zinin, the current WG chair Alia Atlas, and others provided valuable guidance prior to and at the IETF77 RTGWG meeting.

Tony Li and John Drake have made numerous valuable comments on the RTGWG mailing list that are reflected in versions following the IETF77 meeting.

Iftekhhar Hussain and Kireeti Kompella made comments on the RTGWG mailing list after IETF82 that identified a new requirement. Iftekhar Hussain made numerous valuable comments on the RTGWG mailing list that resulted in improvements to document clarity.

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.

Tom Yu and Francis Dupont provided the SecDir and GenArt reviews respectively. Both reviews provided useful comments. Lou Berger provided the RtgDir review which resulted in substantial clarification of terminology and document wording, particularly in the Abstract, Introduction, and Definitions sections.

7. IANA Considerations

This memo includes no request to IANA.

8. Security Considerations

The security considerations for MPLS/GMPLS and for MPLS-TP are documented in [[RFC5920](#)] and [[RFC6941](#)]. This document does not impact the security of MPLS, GMPLS, or MPLS-TP.

The additional information that this document requires does not provide significant additional value to an attacker beyond the information already typically available from attacking a routing or signaling protocol. If the requirements of this document are met by extending an existing routing or signaling protocol, the security considerations of the protocol being extended apply. If the requirements of this document are met by specifying a new protocol, the security considerations of that new protocol should include an evaluation of what level of protection is required by the additional information specified in this document, such as data origin authentication.

9. References

9.1. Normative References

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

9.2. 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-01](#) (work in progress), August 2012.
- [I-D.ietf-rtgwg-cl-use-cases]
Ning, S., Malis, A., McDysan, D., Yong, L., and C. Villamizar, "Composite Link Use Cases and Design Considerations", [draft-ietf-rtgwg-cl-use-cases-01](#) (work in progress), August 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.800]
ITU-T, "Unified functional architecture of transport networks", 2007, <<http://www.itu.int/rec/T-REC-G/recommendation.asp?parent=T-REC-G.800>>.
- [RFC3468] Andersson, L. and G. Swallow, "The Multiprotocol Label Switching (MPLS) Working Group decision on MPLS signaling protocols", [RFC 3468](#), February 2003.

- [RFC4201] Kompella, K., Rekhter, Y., and L. Berger, "Link Bundling in MPLS Traffic Engineering (TE)", [RFC 4201](#), October 2005.
- [RFC5920] Fang, L., "Security Framework for MPLS and GMPLS Networks", [RFC 5920](#), July 2010.
- [RFC6941] Fang, L., Niven-Jenkins, B., Mansfield, S., and R. Graveman, "MPLS Transport Profile (MPLS-TP) Security Framework", [RFC 6941](#), April 2013.

Authors' Addresses

Curtis Villamizar (editor)
OCCNC, LLC

Email: curtis@occnc.com

Dave McDysan (editor)
Verizon
22001 Loudoun County PKWY
Ashburn, VA 20147
USA

Email: dave.mcdysan@verizon.com

So Ning
Tata Communications

Email: ning.so@tatacommunications.com

Andrew Malis
Verizon
60 Sylvan Road
Waltham, MA 02451
USA

Phone: +1 781-466-2362

Email: andrew.g.malis@verizon.com

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

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