

RTGWG
Internet-Draft
Intended status: Informational
Expires: December 31, 2012

S. Ning
Tata Communications
D. McDysan
Verizon
E. Osborne
Cisco
L. Yong
Huawei USA
C. Villamizar
Outer Cape Cod Network
Consulting
June 29, 2012

**Composite Link Framework in Multi Protocol Label Switching (MPLS)
draft-so-yong-rtgwg-cl-framework-06**

Abstract

This document specifies a framework for support of composite link in MPLS networks. A composite link consists of a group of homogenous or non-homogenous links that have the same forward adjacency and can be considered as a single TE link or an IP link in routing. A composite link relies on its component links to carry the traffic over the composite link. Applicability is described for a single pair of MPLS-capable nodes, a sequence of MPLS-capable nodes, or a set of layer networks connecting MPLS-capable nodes.

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 December 31, 2012.

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	4
1.1.	Architecture Summary	4
1.2.	Conventions used in this document	5
1.2.1.	Terminology	5
2.	Composite Link Key Characteristics	5
2.1.	Flow Identification	6
2.2.	Composite Link in Control Plane	8
2.3.	Composite Link in Data Plane	11
3.	Architecture Tradeoffs	11
3.1.	Scalability Motivations	12
3.2.	Reducing Routing Information and Exchange	12
3.3.	Reducing Signaling Load	13
3.3.1.	Reducing Signaling Load using LDP	14
3.3.2.	Reducing Signaling Load using Hierarchy	14
3.3.3.	Using Both LDP and RSVP-TE Hierarchy	14
3.4.	Reducing Forwarding State	14
3.5.	Avoiding Route Oscillation	15
4.	New Challenges	16
4.1.	Control Plane Challenges	16
4.1.1.	Delay and Jitter Sensitive Routing	17
4.1.2.	Local Control of Traffic Distribution	17
4.1.3.	Path Symmetry Requirements	17
4.1.4.	Requirements for Contained LSP	18
4.1.5.	Retaining Backwards Compatibility	19
4.2.	Data Plane Challenges	19
4.2.1.	Very Large LSP	20
4.2.2.	Very Large Microflows	20
4.2.3.	Traffic Ordering Constraints	20
4.2.4.	Accounting for IP and LDP Traffic	21
4.2.5.	IP and LDP Limitations	21
5.	Existing Mechanisms	22
5.1.	Link Bundling	22
5.2.	Classic Multipath	24

6.	Mechanisms Proposed in Other Documents	24
6.1.	Loss and Delay Measurement	24
6.2.	Link Bundle Extensions	25
6.3.	Fat PW and Entropy Labels	26
6.4.	Multipath Extensions	26
7.	Required Protocol Extensions and Mechanisms	27
7.1.	Brief Review of Requirements	27
7.2.	Required Document Coverage	28
7.2.1.	Component Link Grouping	28
7.2.2.	Delay and Jitter Extensions	29
7.2.3.	Path Selection and Admission Control	29
7.2.4.	Dynamic Multipath Balance	30
7.2.5.	Frequency of Load Balance	30
7.2.6.	Inter-Layer Communication	30
7.2.7.	Packet Ordering Requirements	31
7.2.8.	Minimally Disruption Load Balance	31
7.2.9.	Path Symmetry	31
7.2.10.	Performance, Scalability, and Stability	32
7.2.11.	IP and LDP Traffic	32
7.2.12.	LDP Extensions	32
7.2.13.	Pseudowire Extensions	33
7.2.14.	Multi-Domain Composite Link	33
7.3.	Open Issues Regarding Requirements	34
7.4.	Framework Requirement Coverage by Protocol	34
7.4.1.	OSPF-TE and ISIS-TE Protocol Extensions	35
7.4.2.	PW Protocol Extensions	35
7.4.3.	LDP Protocol Extensions	35
7.4.4.	RSVP-TE Protocol Extensions	35
7.4.5.	RSVP-TE Path Selection Changes	35
7.4.6.	RSVP-TE Admission Control and Preemption	35
7.4.7.	Flow Identification and Traffic Balance	35
8.	Security Considerations	35
9.	Acknowledgments	36
10.	References	36
10.1.	Normative References	36
10.2.	Informative References	37
	Authors' Addresses	40

1. Introduction

Composite Link functional requirements are specified in [\[I-D.ietf-rtgwg-cl-requirement\]](#). Composite Link use cases are described in [\[I-D.symmvo-rtgwg-cl-use-cases\]](#). This document specifies a framework to meet these requirements.

Classic multipath, including Ethernet Link Aggregation has been widely used in today's MPLS networks [\[RFC4385\]](#)[\[RFC4928\]](#). Classic multipath using non-Ethernet links are often advertised using MPLS Link bundling. A link bundle [\[RFC4201\]](#) bundles a group of homogeneous links as a TE link to make IGP-TE information exchange and RSVP-TE signaling more scalable. A composite link allows bundling non-homogenous links together as a single logical link. The motivations for using a composite link are described in [\[I-D.ietf-rtgwg-cl-requirement\]](#) and [\[I-D.symmvo-rtgwg-cl-use-cases\]](#).

This document describes a composite link framework in the context of MPLS networks using an IGP-TE and RSVP-TE MPLS control plane with GMPLS extensions [\[RFC3209\]](#)[\[RFC3630\]](#)[\[RFC3945\]](#)[\[RFC5305\]](#).

A composite link is a single logical link in MPLS network that contains multiple parallel component links between two MPLS LSR. Unlike a link bundle [\[RFC4201\]](#), the component links in a composite link can have different properties such as cost or capacity.

Specific protocol solutions are outside the scope of this document, however a framework for the extension of existing protocols is provided. Backwards compatibility is best achieved by extending existing protocols where practical rather than inventing new protocols. The focus is on examining where existing protocol mechanisms fall short with respect to [\[I-D.ietf-rtgwg-cl-requirement\]](#) and on extensions that will be required to accommodate functionality that is called for in [\[I-D.ietf-rtgwg-cl-requirement\]](#).

1.1. Architecture Summary

Networks aggregate information, both in the control plane and in the data plane, as a means to achieve scalability. A tradeoff exists between the needs of scalability and the needs to identify differing path and link characteristics and differing requirements among flows contained within further aggregated traffic flows. These tradeoffs are discussed in detail in [Section 3](#).

Some aspects of Composite Link requirements present challenges for which multiple solutions may exist. In [Section 4](#) various challenges and potential approaches are discussed.

A subset of the functionality called for in [\[I-D.ietf-rtgwg-cl-requirement\]](#) is available through MPLS Link Bundling [\[RFC4201\]](#). Link bundling and other existing standards applicable to Composite Link are covered in [Section 5](#).

The most straightforward means of supporting Composite Link requirements is to extend MPLS protocols and protocol semantics and in particular to extend link bundling. Extensions which have already been proposed in other documents which are applicable to Composite Link are discussed in [Section 6](#).

Goals of most new protocol work within IETF is to reuse existing protocol encapsulations and mechanisms where they meet requirements and extend existing mechanisms such that additional complexity is minimized while meeting requirements and such that backwards compatibility is preserved to the extent it is practical to do so. These goals are considered in proposing a framework for further protocol extensions and mechanisms in [Section 7](#).

[1.2.](#) 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](#) [\[RFC2119\]](#).

[1.2.1.](#) Terminology

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

The abbreviation IGP-TE is used as a shorthand indicating either OSPF-TE [\[RFC3630\]](#) or ISIS-TE [\[RFC5305\]](#).

[2.](#) Composite Link Key Characteristics

[\[I-D.ietf-rtgwg-cl-requirement\]](#) defines external behavior of Composite Links. The overall framework approach involves extending existing protocols in a backwards compatible manner and reusing ongoing work elsewhere in IETF where applicable, defining new protocols or semantics only where necessary. Given the requirements, and this approach of extending MPLS, Composite Link key characteristics can be described in greater detail than given requirements alone.

2.1. Flow Identification

Traffic mapping to component links is a data plane operation. Control over how the mapping is done may be directly dictated or constrained by the control plane or by the management plane. When unconstrained by the control plane or management plane, distribution of traffic is entirely a local matter. Regardless of constraints or lack of constraints, the traffic distribution is required to keep packets belonging to individual flows in sequence and meet QoS criteria specified per LSP by either signaling or management [[RFC2475](#)][[RFC3260](#)]. A key objective of the traffic distribution is to not overload any component link, and be able to perform local recovery when one of component link fails.

The network operator may have other objectives such as placing a bidirectional flow or LSP on the same component link in both direction, load balance over component links, composite link energy saving, and etc. These new requirements are described in [[I-D.ietf-rtgwg-cl-requirement](#)].

Examples of means to identify a flow may in principle include:

1. an LSP identified by an MPLS label,
2. a sub-LSP [[I-D.kompella-mpls-rsvp-ecmp](#)] identified by an MPLS label,
3. a pseudowire (PW) [[RFC3985](#)] identified by an MPLS PW label,
4. a flow or group of flows within a pseudowire (PW) [[RFC6391](#)] identified by an MPLS flow label,
5. a flow or flow group in an LSP [[I-D.ietf-mpls-entropy-label](#)] identified by an MPLS entropy label,
6. all traffic between a pair of IP hosts, identified by an IP source and destination pair,
7. a specific connection between a pair of IP hosts, identified by an IP source and destination pair, protocol, and protocol port pair,
8. a layer-2 conversation within a pseudowire (PW), where the identification is PW payload type specific, such as Ethernet MAC addresses and VLAN tags within an Ethernet PW ([RFC4448](#)).

Although in principle a layer-2 conversation within a pseudowire (PW), may be identified by PW payload type specific information, in

practice this is impractical at LSP midpoints when PW are carried. The PW ingress may provide equivalent information in a PW flow label [[RFC6391](#)]. Therefore, in practice, item #8 above is covered by [[RFC6391](#)] and may be dropped from the list.

An LSR must at least be capable of identifying flows based on MPLS labels. Most MPLS LSP do not require that traffic carried by the LSP are carried in order. MPLS-TP is a recent exception. If it is assumed that no LSP require strict packet ordering of the LSP itself (only of flows within the LSP), then the entire label stack can be used as flow identification. If some LSP may require strict packet ordering but those LSP cannot be distinguished from others, then only the top label can be used as a flow identifier. If only the top label is used (for example, as specified by [[RFC4201](#)] when the "all-ones" component described in [[RFC4201](#)] is not used), then there may not be adequate flow granularity to accomplish well balanced traffic distribution and it will not be possible to carry LSP that are larger than any individual component link.

The number of flows can be extremely large. This may be the case when the entire label stack is used and is always the case when IP addresses are used in provider networks carrying Internet traffic. Current practice for native IP load balancing at the time of writing were documented in [[RFC2991](#)], [[RFC2992](#)]. These practices as described, make use of IP addresses. The common practices were extended to include the MPLS label stack and the common practice of looking at IP addresses within the MPLS payload. These extended practices are described in [[RFC4385](#)] and [[RFC4928](#)] due to their impact on pseudowires without a PWE3 Control Word. Additional detail on current multipath practices can be found in the appendices of [[I-D.symmvo-rtgwq-cl-use-cases](#)].

Using only the top label supports too coarse a traffic balance. Using the full label stack or IP addresses as flow identification provides a sufficiently fine traffic balance, but is capable of identifying such a high number of distinct flows, that a technique of grouping flows, such as hashing on the flow identification criteria, becomes essential to reduce the stored state, and is an essential scaling technique. Other means of grouping flows may be possible.

In summary:

1. Load balancing using only the MPLS label stack provides too coarse a granularity of load balance.
2. Tracking every flow is not scalable due to the extremely large number of flows in provider networks.

3. Existing techniques, IP source and destination hash in particular, have proven in over two decades of experience to be an excellent way of identifying groups of flows.
4. If a better way to identify groups of flows is discovered, then that method can be used.
5. IP address hashing is not required, but use of this technique is strongly encouraged given the technique's long history of successful deployment.

2.2. Composite Link in Control Plane

A composite Link is advertised as a single logical interface between two connected routers, which forms forwarding adjacency (FA) between the routers. The FA is advertised as a TE-link in a link state IGP, using either OSPF-TE or ISIS-TE. The IGP-TE advertised interface parameters for the composite link can be preconfigured by the network operator or be derived from its component links. Composite link advertisement requirements are specified in [\[I-D.ietf-rtgwg-cl-requirement\]](#).

In IGP-TE, a composite link is advertised as a single TE link between two connected routers. This is similar to a link bundle [\[RFC4201\]](#). Link bundle applies to a set of homogenous component links. Composite link allows homogenous and non-homogenous component links. Due to the similarity, and for backwards compatibility, extending link bundling is viewed as both simple and as the best approach.

In order for a route computation engine to calculate a proper path for a LSP, it is necessary for composite link to advertise the summarized available bandwidth as well as the maximum bandwidth that can be made available for single flow (or single LSP where no finer flow identification is available). If a composite link contains some non-homogeneous component links, the composite link also should advertise the summarized bandwidth and the maximum bandwidth for single flow per each homogeneous component link group.

Both LDP [\[RFC5036\]](#) and RSVP-TE [\[RFC3209\]](#) can be used to signal a LSP over a composite link. LDP cannot be extended to support traffic engineering capabilities [\[RFC3468\]](#).

When an LSP is signaled using RSVP-TE, the LSP MUST be placed on the component link that meets the LSP criteria indicated in the signaling message.

When an LSP is signaled using LDP, the LSP MUST be placed on the component link that meets the LSP criteria, if such a component link

is available. LDP does not support traffic engineering capabilities, imposing restrictions on LDP use of Composite Link. See [Section 4.2.5](#) for further details.

A composite link may contain non-homogeneous component links. The route computing engine may select one group of component links for a LSP. The routing protocol MUST make this grouping available in the TE-LSDB. The route computation used in RSVP-TE MUST be extended to include only the capacity of groups within a composite link which meet LSP criteria. The signaling protocol MUST be able to indicate either the criteria, or which groups may be used. A composite link MUST place the LSP on a component link or group which meets or exceeds the LSP criteria.

Composite link capacity is aggregated capacity. LSP capacity MAY be larger than individual component link capacity. Any aggregated LSP can determine a bounds on the largest microflow that could be carried and this constraint can be handled as follows.

1. If no information is available through signaling, management plane, or configuration, the largest microflow is bound by one of the following:
 - A. the largest single LSP if most traffic is RSVP-TE signaled and further aggregated,
 - B. the largest pseudowire if most traffic is carrying pseudowire payloads that are aggregated within RSVP-TE LSP,
 - C. or the largest source and sink interface if a large amount of IP or LDP traffic is contained within the aggregate.

If a very large amount of traffic being aggregated is IP or LDP, then the largest microflow is bound by the largest component link on which IP traffic can arrive. For example, if an LSR is acting as an LER and IP and LDP traffic is arriving on 10 Gb/s edge interfaces, then no microflow larger than 10 Gb/s will be present on the RSVP-TE LSP that aggregate traffic across the core, even if the core interfaces are 100 Gb/s interfaces.

2. The prior conditions provide a bound on the largest microflow when no signaling extensions indicate a bounds. If an LSP is aggregating smaller LSP for which the largest expected microflow carried by the smaller LSP is signaled, then the largest microflow expected in the containing LSP (the aggregate) is the maximum of the largest expected microflow for any contained LSP. For example, RSVP-TE LSP may be large but aggregate traffic for which the source or sink are all 1 Gb/s or smaller interfaces

(such as in mobile applications in which cell sites backhauls are no larger than 1 Gb/s). If this information is carried in the LSP originated at the cell sites, then further aggregates across a core may make use of this information.

3. The IGP must provide the bounds on the largest microflow that a composite link can accommodate, which is the maximum capacity on a component link that can be made available by moving other traffic. This information is needed by the ingress LER for path determination.
4. A means to signal an LSP whose capacity is larger than individual component link capacity is needed [[I-D.ietf-rtgwg-cl-requirement](#)] and also signal the largest microflow expected to be contained in the LSP. If a bounds on the largest microflow is not signaled there is no means to determine if an LSP which is larger than any component link can be subdivided into flows and therefore should be accepted by admission control.

When a bidirectional LSP request is signaled over a composite link, if the request indicates that the LSP must be placed on the same component link, the routers of the composite link MUST place the LSP traffic in both directions on a same component link. This is particularly challenging for aggregated capacity which makes use of the label stack for traffic distribution. The two requirements are mutually exclusive for any one LSP. No one LSP may be both larger than any individual component link and require symmetrical paths for every flow. Both requirements can be accommodated by the same composite link for different LSP, with any one LSP requiring no more than one of these two features.

Individual component link may fail independently. Upon component link failure, a composite link MUST support a minimally disruptive local repair, preempting any LSP which can no longer be supported. Available capacity in other component links MUST be used to carry impacted traffic. The available bandwidth after failure MUST be advertised immediately to avoid looped crankback.

When a composite link is not able to transport all flows, it preempts some flows based upon local management configuration and informs the control plane on these preempted flows. The composite link MUST support soft preemption [[RFC5712](#)]. This action ensures the remaining traffic is transported properly. FR#10 requires that the traffic be restored. FR#12 requires that any change be minimally disruptive. These two requirements are interpreted to include preemption among the types of changes that must be minimally disruptive.

2.3. Composite Link in Data Plane

The data plane must first identify groups of flows. Flow identification is covered in [Section 2.1](#). Having identified groups of flows the groups must be placed on individual component links. This second step is called traffic distribution or traffic placement. The two steps together are known as traffic balancing or load balancing.

Traffic distribution may be determined by or constrained by control plane or management plane. Traffic distribution may be changed due to component link status change, subject to constraints imposed by either the management plane or control plane. The distribution function is local to the routers in which a composite link belongs to and is not specified here.

When performing traffic placement, a composite link does not differentiate multicast traffic vs. unicast traffic.

In order to maintain scalability, existing data plane forwarding retains state associated with the top label only. The use of flow group identification is in a second step in the forwarding process. Data plane forwarding makes use of the top label to select a composite link, or a group of components within a composite link or for the case where an LSP is pinned (see [\[RFC4201\]](#)), a specific component link. For those LSP for which the LSP selects only the composite link or a group of components within a composite link, the load balancing makes use of the flow group identification.

The most common traffic placement techniques uses the a flow group identification as an index into a table. The table provides an indirection. The number of bits of hash is constrained to keep table size small. While this is not the best technique, it is the most common. Better techniques exist but they are outside the scope of this document and some are considered proprietary.

Requirements to limit frequency of load balancing can be adhered to by keeping track of when a flow group was last moved and imposing a minimum period before that flow group can be moved again. This is straightforward for a table approach. For other approaches it may be less straightforward but is achievable.

3. Architecture Tradeoffs

Scalability and stability are critical considerations in protocol design where protocols may be used in a large network such as today's service provider networks. Composite Link is applicable to networks

which are large enough to require that traffic be split over multiple paths. Scalability is a major consideration for networks that reach a capacity large enough to require Composite Link.

Some of the requirements of Composite Link could potentially have a negative impact on scalability. For example, Composite Link requires additional information to be carried in situations where component links differ in some significant way.

3.1. Scalability Motivations

In the interest of scalability information is aggregated in situations where information about a large amount of network capacity or a large amount of network demand provides is adequate to meet requirements. Routing information is aggregated to reduce the amount of information exchange related to routing and to simplify route computation (see [Section 3.2](#)).

In an MPLS network large routing changes can occur when a single fault occurs. For example, a single fault may impact a very large number of LSP traversing a given link. As new LSP are signaled to avoid the fault, resources are consumed elsewhere, and routing protocol announcements must flood the resource changes. If protection is in place, there is less urgency to converging quickly. If multiple faults occur that are not covered by shared risk groups (SRG), then some protection may fail, adding urgency to converging quickly even where protection was deployed.

Reducing the amount of information allows the exchange of information during a large routing change to be accomplished more quickly and simplifies route computation. Simplifying route computation improves convergence time after very significant network faults which cannot be handled by preprovisioned or precomputed protection mechanisms. Aggregating smaller LSP into larger LSP is a means to reduce path computation load and reduce RSVP-TE signaling (see [Section 3.3](#)).

Neglecting scaling issues can result in performance issues, such as slow convergence. Neglecting scaling in some cases can result in networks which perform so poorly as to become unstable.

3.2. Reducing Routing Information and Exchange

Link bundling at the very least provides a means of aggregating control plane information. Even where the all-ones component link supported by link bundling is not used, the amount of control information is reduced by the average number of component links in a bundle.

Fully deaggregating link bundle information would negate this benefit. If there is a need to deaggregate, such as to distinguish between groups of links within specified ranges of delay, then no more deaggregation than is necessary should be done.

For example, in supporting the requirement for heterogeneous component links, it makes little sense to fully deaggregate link bundles when adding support for groups of component links with common attributes within a link bundle can maintain most of the benefit of aggregation while adequately supporting the requirement to support heterogeneous component links.

Routing information exchange is also reduced by making sensible choices regarding the amount of change to link parameters that require link readvertisement. For example, if delay measurements include queuing delay, then a much more coarse granularity of delay measurement would be called for than if the delay does not include queuing and is dominated by geographic delay (speed of light delay).

3.3. Reducing Signaling Load

Aggregating traffic into very large hierarchical LSP in the core very substantially reduces the number of LSP that need to be signaled and the number of path computations any given LSR will be required to perform when a major network fault occurs.

In the extreme, applying MPLS to a very large network without hierarchy could exceed the 20 bit label space. For example, in a network with 4,000 nodes, with 2,000 on either side of a cutset, would have 4,000,000 LSP crossing the cutset. Even in a degree four cutset, an uneven distribution of LSP across the cutset, or the loss of one link would result in a need to exceed the size of the label space. Among provider networks, 4,000 access nodes is not at all large.

In less extreme cases, having each node terminate hundreds of LSP to achieve a full mesh creates a very large computational load. The time complexity of one CSPF computation is order($N \log N$), where L is proportional to N , and N and L are the number of nodes and number of links, respectively. If each node must perform order(N) computations when a fault occurs, then the computational load increases as order($N^2 \log N$) as the number of nodes increases. In practice at the time of writing, this imposes a limit of a few hundred nodes in a full mesh of MPLS LSP before the computational load is sufficient to result in unacceptable convergence times.

Two solutions are applied to reduce the amount of RSVP-TE signaling. Both involve subdividing the MPLS domain into a core and a set of

regions.

3.3.1. Reducing Signaling Load using LDP

LDP can be used for edge-to-edge LSP, using RSVP-TE to carry the LDP intra-core traffic and also optionally also using RSVP-TE to carry the LDP intra-region traffic within each region. LDP does not support traffic engineering, but does support multipoint-to-point (MPTP) LSP, which require less signaling than edge-to-edge RSVP-TE point-to-point (PTP) LSP. A drawback of this approach is the inability to use RSVP-TE protection (FRR or GMPLS protection) against failure of the border LSR sitting at a core/region boundary.

3.3.2. Reducing Signaling Load using Hierarchy

When the number of nodes grows too large, the amount of RSVP-TE signaling can be reduced using the MPLS PSC hierarchy [[RFC4206](#)]. A core within the hierarchy can divide the topology into M regions of on average N/M nodes. Within a region the computational load is reduced by more than M^2 . Within the core, the computational load generally becomes quite small since M is usually a fairly small number (a few tens of regions) and each region is generally attached to the core in typically only two or three places on average.

Using hierarchy improves scaling but has two consequences. First, hierarchy effectively forces the use of platform label space. When a containing LSP is rerouted, the labels assigned to the contained LSP cannot be changed but may arrive on a different interface. Second, hierarchy results in much larger LSP. These LSP today are larger than any single component link and therefore force the use of the all-ones component in link bundles.

3.3.3. Using Both LDP and RSVP-TE Hierarchy

It is also possible to use both LDP and RSVP-TE hierarchy. MPLS networks with a very large number of nodes may benefit from the use of both LDP and RSVP-TE hierarchy. The two techniques are certainly not mutually exclusive.

3.4. Reducing Forwarding State

Both LDP and MPLS hierarchy have the benefit of reducing the amount of forwarding state. Using the example from [Section 3.3](#), and using MPLS hierarchy, the worst case generally occurs at borders with the core.

For example, consider a network with approximately 1,000 nodes divided into 10 regions. At the edges, each node requires 1,000 LSP

to other edge nodes. The edge nodes also require 100 intra-region LSP. Within the core, if the core has only 3 attachments to each region the core LSR have less than 100 intra-core LSP. At the border cutset between the core and a given region, in this example there are 100 edge nodes with inter-region LSP crossing that cutset, destined to 900 other edge nodes. That yields forwarding state for on the order of 90,000 LSP at the border cutset. These same routers need only reroute well under 200 LSP when a multiple fault occurs, as long as only links are affected and a border LSR does not go down.

In the core, the forwarding state is greatly reduced. If inter-region LSP have different characteristics, it makes sense to make use of aggregates with different characteristics. Rather than exchange information about every inter-region LSP within the intra-core LSP it makes more sense to use multiple intra-core LSP between pairs of core nodes, each aggregating sets of inter-region LSP with common characteristics or common requirements.

3.5. Avoiding Route Oscillation

Networks can become unstable when a feedback loop exists such that moving traffic to a link causes a metric such as delay to increase, which then causes traffic to move elsewhere. For example, the original ARPANET routing used a delay based cost metric and proved prone to route oscillations [[DBP](#)].

Delay may be used as a constraint in routing for high priority traffic, where the movement of traffic cannot impact the delay. The safest way to measure delay is to make measurements based on traffic which is prioritized such that it is queued ahead of the traffic which will be affected. This is a reasonable measure of delay for high priority traffic for which constraints have been set which allow this type of traffic to consume only a fraction of link capacities with the remaining capacity available to lower priority traffic.

Any measurement of jitter (delay variation) that is used in route decision is likely to cause oscillation. Jitter that is caused by queuing effects and cannot be measured using a very high priority measurement traffic flow.

It may be possible to find links with constrained queuing delay or jitter using a theoretical maximum or a probability based bound on queuing delay or jitter at a given priority based on the types and amounts of traffic accepted and combining that theoretical limit with a measured delay at very high priority.

Instability can occur due to poor performance and interaction with protocol timers. In this way a computational scaling problem can

become a stability problem when a network becomes sufficiently large. For this reason, [[I-D.ietf-rtgwg-cl-requirement](#)] has a number of requirements focusing on minimally impacting scalability.

4. New Challenges

New technical challenges are posed by [[I-D.ietf-rtgwg-cl-requirement](#)] in both the control plane and data plane.

Among the more difficult challenges are the following.

1. requirements related delay or jitter (see [Section 4.1.1](#)),
2. the combination of ingress control over LSP placement and retaining an ability to move traffic as demands dictate can pose challenges and such requirements can even be conflicting (see `target="sect.local-control" />`),
3. path symmetry requires extensions and is particularly challenging for very large LSP (see [Section 4.1.3](#)),
4. accommodating a very wide range of requirements among contained LSP can lead to inefficiency if the most stringent requirements are reflected in aggregates, or reduce scalability if a large number of aggregates are used to provide a too fine a reflection of the requirements in the contained LSP (see [Section 4.1.4](#)),
5. backwards compatibility is somewhat limited due to the need to accommodate legacy multipath interfaces which provide too little information regarding their configured default behavior, and legacy LSP which provide too little information regarding their requirements (see [Section 4.1.5](#)),
6. data plane challenges include those of accommodating very large LSP, large microflows, traffic ordering constraints imposed by a subset of LSP, and accounting for IP and LDP traffic (see [Section 4.2](#)).

4.1. Control Plane Challenges

Some of the control plane requirements are particularly challenging. Handling large flows which aggregate smaller flows must be accomplished with minimal impact on scalability. Potentially conflicting are requirements for jitter and requirements for stability. Potentially conflicting are the requirements for ingress control of a large number of parameters, and the requirements for local control needed to achieve traffic balance across a composite

link. These challenges and potential solutions are discussed in the following sections.

4.1.1. Delay and Jitter Sensitive Routing

Delay and jitter sensitive routing are called for in [\[I-D.ietf-rtgwg-cl-requirement\]](#) in requirements FR#2, FR#7, FR#8, FR#9, FR#15, FR#16, FR#17, FR#18. Requirement FR#17 is particularly problematic, calling for constraints on jitter.

A tradeoff exists between scaling benefits of aggregating information, and potential benefits of using a finer granularity in delay reporting. To maintain the scaling benefit, measured link delay for any given composite link SHOULD be aggregated into a small number of delay ranges. IGP-TE extensions MUST be provided which advertise the available capacities for each of the selected ranges.

For path selection of delay sensitive LSP, the ingress SHOULD bias link metrics based on available capacity and select a low cost path which meets LSP total path delay criteria. To communicate the requirements of an LSP, the ERO MUST be extended to indicate the per link constraints. To communicate the type of resource used, the RRO SHOULD be extended to carry an identification of the group that is used to carry the LSP at each link bundle hop.

4.1.2. Local Control of Traffic Distribution

Many requirements in [\[I-D.ietf-rtgwg-cl-requirement\]](#) suggest that a node immediately adjacent to a component link should have a high degree of control over how traffic is distributed, as long as network performance objectives are met. Particularly relevant are FR#18 and FR#19.

The requirements to allow local control are potentially in conflict with requirement FR#21 which gives full control of component link select to the LSP ingress. While supporting this capability is mandatory, use of this feature is optional per LSP.

A given network deployment will have to consider this pair of conflicting requirements and make appropriate use of local control of traffic placement and ingress control of traffic placement to best meet network requirements.

4.1.3. Path Symmetry Requirements

Requirement FR#21 in [\[I-D.ietf-rtgwg-cl-requirement\]](#) includes a provision to bind both directions of a bidirectional LSP to the same component. This is easily achieved if the LSP is directly signaled

across a composite link. This is not as easily achieved if a set of LSP with this requirement are signaled over a large hierarchical LSP which is in turn carried over a composite link. The basis for load distribution in such as case is the label stack. The labels in either direction are completely independent.

This could be accommodated if the ingress, egress, and all midpoints of the hierarchical LSP make use of an entropy label in the distribution, and use only that entropy label. A solution for this problem may add complexity with very little benefit. There is little or no true benefit of using symmetrical paths rather than component links of identical characteristics.

Traffic symmetry and large LSP capacity are a second pair of conflicting requirements. Any given LSP can meet one of these two requirements but not both. A given network deployment will have to make appropriate use of each of these features to best meet network requirements.

4.1.4. Requirements for Contained LSP

[I-D.ietf-rtgwg-cl-requirement] calls for new LSP constraints. These constraints include frequency of load balancing rearrangement, delay and jitter, packet ordering constraints, and path symmetry.

When LSP are contained within hierarchical LSP, there is no signaling available at midpoint LSR which identifies the contained LSP let alone providing the set of requirements unique to each contained LSP. Defining extensions to provide this information would severely impact scalability and defeat the purpose of aggregating control information and forwarding information into hierarchical LSP. For the same scalability reasons, not aggregating at all is not a viable option for large networks where scalability and stability problems may occur as a result.

As pointed out in [Section 4.1.3](#), the benefits of supporting symmetric paths among LSP contained within hierarchical LSP may not be sufficient to justify the complexity of supporting this capability.

A scalable solution which accommodates multiple sets of LSP between given pairs of LSR is to provide multiple hierarchical LSP for each given pair of LSR, each hierarchical LSP aggregating LSP with common requirements and a common pair of endpoints. This is a network design technique available to the network operator rather than a protocol extension. This technique can accommodate multiple sets of delay and jitter parameters, multiple sets of frequency of load balancing parameters, multiple sets of packet ordering constraints, etc.

4.1.5. Retaining Backwards Compatibility

Backwards compatibility and support for incremental deployment requires considering the impact of legacy LSR in the role of LSP ingress, and considering the impact of legacy LSR advertising ordinary links, advertising Ethernet LAG as ordinary links, and advertising link bundles.

Legacy LSR in the role of LSP ingress cannot signal requirements which are not supported by their control plane software. The additional capabilities supported by other LSR has no impact on these LSR. These LSR however, being unaware of extensions, may try to make use of scarce resources which support specific requirements such as low delay. To a limited extent it may be possible for a network operator to avoid this issue using existing mechanisms such as link administrative attributes and attribute affinities [[RFC3209](#)].

Legacy LSR advertising ordinary links will not advertise attributes needed by some LSP. For example, there is no way to determine the delay or jitter characteristics of such a link. Legacy LSR advertising Ethernet LAG pose additional problems. There is no way to determine that packet ordering constraints would be violated for LSP with strict packet ordering constraints, or that frequency of load balancing rearrangement constraints might be violated.

Legacy LSR advertising link bundles have no way to advertise the configured default behavior of the link bundle. Some link bundles may be configured to place each LSP on a single component link and therefore may not be able to accommodate an LSP which requires bandwidth in excess of the size of a component link. Some link bundles may be configured to spread all LSP over the all-ones component. For LSR using the all-ones component link, there is no documented procedure for correctly setting the "Maximum LSP Bandwidth". There is currently no way to indicate the largest microflow that could be supported by a link bundle using the all-ones component link.

Having received the RRO, it is possible for an ingress to look for the all-ones component to identify such link bundles after having signaled at least one LSP. Whether any LSR collects this information on legacy LSR and makes use of it to set defaults, is an implementation choice.

4.2. Data Plane Challenges

Flow identification is briefly discussed in [Section 2.1](#). Traffic distribution is briefly discussed in [Section 2.3](#). This section discusses issues specific to particular requirements specified in

[[I-D.ietf-rtgwg-cl-requirement](#)].

4.2.1. Very Large LSP

Very large LSP may exceed the capacity of any single component of a composite link. In some cases contained LSP may exceed the capacity of any single component. These LSP may the use of the equivalent of the all-ones component of a link bundle, or may use a subset of components which meet the LSP requirements.

Very large LSP can be accommodated as long as they can be subdivided (see [Section 4.2.2](#)). A very large LSP cannot have a requirement for symmetric paths unless complex protocol extensions are proposed (see [Section 2.2](#) and [Section 4.1.3](#)).

4.2.2. Very Large Microflows

Within a very large LSP there may be very large microflows. A very large microflow is a very large flows which cannot be further subdivided. Flows which cannot be subdivided must be no larger than the capacity of any single component.

Current signaling provides no way to specify the largest microflow that can be supported on a given link bundle in routing advertisements. Extensions which address this are discussed in [Section 6.4](#). Absent extensions of this type, traffic containing microflows that are too large for a given composite link may be present. There is no data plane solution for this problem that would not require reordering traffic at the composite link egress.

Some techniques are susceptible to statistical collisions where an algorithm to distribute traffic is unable to disambiguate traffic among two or more very large microflow where their sum is in excess of the capacity of any single component. Hash based algorithms which use too small a hash space are particularly susceptible and require a change in hash seed in the event that this were to occur. A change in hash seed is highly disruptive, causing traffic reordering among all traffic flows over which the hash function is applied.

4.2.3. Traffic Ordering Constraints

Some LSP have strict traffic ordering constraints. Most notable among these are MPLS-TP LSP. In the absence of aggregation into hierarchical LSP, those LSP with strict traffic ordering constraints can be placed on individual component links if there is a means of identifying which LSP have such a constraint. If LSP with strict traffic ordering constraints are aggregated in hierarchical LSP, the hierarchical LSP capacity may exceed the capacity of any single

component link. In such a case the load balancing for the containing may be constrained to look only at the top label and the first contained label. This and related issues are discussed further in [Section 6.4](#).

[4.2.4](#). Accounting for IP and LDP Traffic

Networks which carry RSVP-TE signaled MPLS traffic generally carry low volumes of native IP traffic, often only carrying control traffic as native IP. There is no architectural guarantee of this, it is just how network operators have made use of the protocols.

[I-D.ietf-rtgwg-cl-requirement] requires that native IP and native LDP be accommodated. In some networks, a subset of services may be carried as native IP or carried as native LDP. Today this may be accommodated by the network operator estimating the contribution of IP and LDP and configuring a lower set of available bandwidth figures on the RSVP-TE advertisements.

The only improvement that Composite Link can offer is that of measuring the IP and LDP traffic levels and automatically reducing the available bandwidth figures on the RSVP-TE advertisements. The measurements would have to be significantly filtered. This is similar to a feature in existing LSR, commonly known as "autobandwidth" with a key difference. In the "autobandwidth" feature, the bandwidth request of an RSVP-TE signaled LSP is adjusted in response to traffic measurements. In this case the IP or LDP traffic measurements are used to reduce the link bandwidth directly, without first encapsulating in an RSVP-TE LSP.

This may be a subtle and perhaps even a meaningless distinction if Composite Link is used to form a Sub-Path Maintenance Element (SPME). A SPME is in practice essentially an unsignaled single hop LSP with PHP enabled [[RFC5921](#)]. A Composite Link SPME looks very much like classic multipath, where there is no signaling, only management plane configuration creating the multipath entity (of which Ethernet Link Aggregation is a subset).

[4.2.5](#). IP and LDP Limitations

IP does not offer traffic engineering. LDP cannot be extended to offer traffic engineering [[RFC3468](#)]. Therefore there is no traffic engineered fallback to an alternate path for IP and LDP traffic if resources are not adequate for the IP and/or LDP traffic alone on a given link in the primary path. The only option for IP and LDP would be to declare the link down. Declaring a link down due to resource exhaustion would reduce traffic to zero and eliminate the resource exhaustion. This would cause oscillations and is therefore not a

viable solution.

Congestion caused by IP or LDP traffic loads is a pathologic case that can occur if IP and/or LDP are carried natively and there is a high volume of IP or LDP traffic. This situation can be avoided by carrying IP and LDP within RSVP-TE LSP.

It is also not possible to route LDP traffic differently for different FEC. LDP traffic engineering is specifically disallowed by [\[RFC3468\]](#). It may be possible to support multi-topology IGP extensions to accommodate more than one set of criteria. If so, the additional IGP could be bound to the forwarding criteria, and the LDP FEC bound to a specific IGP instance, inheriting the forwarding criteria. Alternately, one IGP instance can be used and the LDP SPF can make use of the constraints, such as delay and jitter, for a given LDP FEC. [Note: WG needs to discuss this and decide first whether to solve this at all and then if so, how.]

5. Existing Mechanisms

In MPLS the one mechanisms which support explicit signaling of multiple parallel links is Link Bundling [\[RFC4201\]](#). The set of techniques known as "classis multipath" support no explicit signaling, except in two cases. In Ethernet Link Aggregation the Link Aggregation Control Protocol (LACP) coordinates the addition or removal of members from an Ethernet Link Aggregation Group (LAG). The use of the "all-ones" component of a link bundle indicates use of classis multipath, however the ability to determine if a link bundle makes use of classis multipath is not yet supported.

5.1. Link Bundling

Link bundling supports advertisement of a set of homogenous links as a single route advertisement. Link bundling supports placement of an LSP on any single component link, or supports placement of an LSP on the all-ones component link. Not all link bundling implementations support the all-ones component link. There is no way for an ingress LSR to tell which potential midpoint LSR support this feature and use it by default and which do not. Based on [\[RFC4201\]](#) it is unclear how to advertise a link bundle for which the all-ones component link is available and used by default. Common practice is to violate the specification and set the Maximum LSP Bandwidth to the Available Bandwidth. There is no means to determine the largest microflow that could be supported by a link bundle that is using the all-ones component link.

[\[RFC6107\]](#) extends the procedures for hierarchical LSP but also

extends link bundles. An LSP can be explicitly signaled to indicate that it is an LSP to be used as a component of a link bundle. Prior to that the common practice was to simply not advertise the component link LSP into the IGP, since only the ingress and egress of the link bundle needed to be aware of their existence, which they would be aware of due to the RSVP-TE signaling used in setting up the component LSP.

While link bundling can be the basis for composite links, a significant number of small extension needs to be added.

1. To support link bundles of heterogeneous links, a means of advertising the capacity available within a group of homogeneous needs to be provided.
2. Attributes need to be defined to support the following parameters for the link bundle or for a group of homogeneous links.
 - A. delay range
 - B. jitter (delay variation) range
 - C. group metric
 - D. all-ones component capable
 - E. capable of dynamically balancing load
 - F. largest supportable microflow
 - G. abilities to support strict packet ordering requirements within contained LSP
3. For each of the prior extended attributes, the constraint based routing path selection needs to be extended to reflect new constraints based on the extended attributes.
4. For each of the prior extended attributes, LSP admission control needs to be extended to reflect new constraints based on the extended attributes.
5. Dynamic load balance must be provided for flows within a given set of links with common attributes such that NPO are not violated including frequency of load balance adjustment for any given flow.

5.2. Classic Multipath

Classic multipath is defined in [[I-D.symmvo-rtgwg-cl-use-cases](#)].

Classic multipath refers to the most common current practice in implementation and deployment of multipath. 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 multipath provides a highly scalable means of load balancing. Adaptive multipath has proven value in assuring an even loading on component link and an ability to adapt to change in offered load that occurs over periods of hundreds of milliseconds or more. Classic multipath scalability is due to the ability to effectively work with an extremely large number of flows (IP host pairs) using relatively little resources (a data structure accessed using a hash result as a key or using ranges of hash results).

Classic multipath meets a small subset of Composite Link requirements. Due to scalability of the approach, classic multipath seems to be an excellent candidate for extension to meet the full set of Composite Link forwarding requirements.

Additional detail can be found in [[I-D.symmvo-rtgwg-cl-use-cases](#)].

6. Mechanisms Proposed in Other Documents

A number of documents which at the time of writing are works in progress address parts of the requirements of Composite Link, or assist in making some of the goals achievable.

6.1. Loss and Delay Measurement

Procedures for measuring loss and delay are provided in [[RFC6374](#)]. These are OAM based measurements. This work could be the basis of delay measurements and delay variation measurement used for metrics called for in [[I-D.ietf-rtgwg-cl-requirement](#)].

Currently there are two additional Internet-Drafts that address delay and delay variation metrics.

[draft-wang-ccamp-latency-te-metric](#)

[[I-D.wang-ccamp-latency-te-metric](#)] is designed specifically to meet this requirement. OSPF-TE and ISIS-TE extensions are defined to indicate link delay and delay variance. The RSVP-TE ERO is extended to include service level requirements. A latency

accumulation object is defined to provide a means of verification of the service level requirements. This draft is intended to proceed in the CCAMP WG. It is currently an individual submission. The 03 version of this draft expired in September 2012.

[draft-giacalone-ospf-te-express-path](#)

This document proposes to extend OSPF-TE only. Extensions support delay, delay variance, loss, residual bandwidth, and available bandwidth. No extensions to RSVP-TE are proposed. This draft is intended to proceed in the CCAMP WG. It is currently an individual submission. The 02 version will expire in March 2012.

A possible course of action may be to combine these two drafts. The delay variance, loss, residual bandwidth, and available bandwidth extensions are particularly prone to network instability. The question as to whether queuing delay and delay variation should be considered, and if so for which diffserv Per-Hop Service Class (PSC) is not addressed.

Note to co-authors: The ccamp-latency-te-metric draft refers to [\[I-D.ietf-rtgwg-cl-requirement\]](#) and is well matched to those requirements, including stability. The ospf-te-express-path draft refers to the "Alto Protocol" ([draft-ietf-alto-protocol](#)) and therefore may not be intended for RSVP-TE use. The authors of the two drafts may be able to resolve this. It may be best to drop ospf-te-express-path from this framework document.

6.2. Link Bundle Extensions

A set of link bundling extensions are defined in [\[I-D.ietf-mpls-explicit-resource-control-bundle\]](#). This document provides extensions to the ERO and RRO to explicitly control the labels and resources within a bundle used by an LSP.

The extensions in this document could be further extended to support indicating a group of component links in the ERO or RRO, where the group is given an interface identification like the bundle itself. The extensions could also be further extended to support specification of the all-ones component link in the ERO or RRO.

[\[I-D.ietf-mpls-explicit-resource-control-bundle\]](#) does not provide a means to advertise the link bundle components. It is not certain how the ingress LSR would determine the set of link bundle component links available for a given link bundle.

[\[I-D.ospf-cc-stlv\]](#) provides a baseline draft for extending link

bundling to advertise components. A new component TVL (C-TLV) is proposed, which must reference a Composite Link Link TLV. [\[I-D.ospf-cc-stlv\]](#) is intended for the OSPF WG and submitted for the "Experimental" track. The 00 version expired in February 2012.

6.3. Fat PW and Entropy Labels

Two documents provide a means to add entropy for the purpose of improving load balance. MPLS encapsulation can bury information that is needed to identify microflows. These two documents allow a pseudowire ingress and LSP ingress respectively to add a label solely for the purpose of providing a finer granularity of microflow groups.

[RFC6391] allows pseudowires which carry a large volume of traffic, where microflows can be identified to be load balanced across multiple members of an Ethernet LAG or an MPLS link bundle. This is accomplished by adding a flow label below the pseudowire label in the MPLS label stack. For this to be effective the link bundle load balance must make use of the label stack up to and including this flow label.

[I-D.ietf-mpls-entropy-label] provides a means for a LER to put an additional label known as an entropy label on the MPLS label stack. As defined, only the LER can add the entropy label.

Core LSR acting as LER for aggregated LSP can add entropy labels based on deep packet inspection and place an entropy label indicator (ELI) and entropy label (EL) just below the label being acted on. This would be helpful in situations where the label stack depth to which load distribution can operate is limited by implementation or is limited for other reasons such as carrying both MPLS-TP and MPLS with entropy labels within the same hierarchical LSP.

6.4. Multipath Extensions

The multipath extensions drafts address one aspect of Composite Link. These drafts deal with the issue of accommodating LSP which have strict packet ordering constraints in a network containing multipath. MPLS-TP has become the one important instance of LSP with strict packet ordering constraints and has driven this work.

[I-D.villamizar-mpls-tp-multipath] outlines requirements and gives a number of options for dealing with the apparent incompatibility of MPLS-TP and multipath. A preferred option is described.

[I-D.villamizar-mpls-tp-multipath-te-extn] provides protocol extensions needed to implement the preferred option described in [\[I-D.villamizar-mpls-tp-multipath\]](#).

Other issues pertaining to multipath are also addressed. Means to advertise the largest microflow supportable are defined. Means to indicate the largest expected microflow within an LSP are defined. Issues related to hierarchy are addressed.

7. Required Protocol Extensions and Mechanisms

Prior sections have reviewed key characteristics, architecture tradeoffs, new challenges, existing mechanisms, and relevant mechanisms proposed in existing new documents.

This section first summarizes and groups requirements. A set of documents coverage groupings are proposed with existing works-in-progress noted where applicable. The set of extensions are then grouped by protocol affected as a convenience to implementors.

7.1. Brief Review of Requirements

The following list provides a categorization of requirements specified in [[I-D.ietf-rtgwg-cl-requirement](#)] along with a short phrase indication what topic the requirement covers.

routing information aggregation

FR#1 (routing summarization), FR#20 (composite link may be a component of another composite link)

restoration speed

FR#2 (restoration speed meeting NPO), FR#12 (minimally disruptive load rebalance), DR#6 (fast convergence), DR#7 (fast worst case failure convergence)

load distribution, stability, minimal disruption

FR#3 (automatic load distribution), FR#5 (must not oscillate), FR#11 (dynamic placement of flows), FR#12 (minimally disruptive load rebalance), FR#13 (bounded rearrangement frequency), FR#18 (flow placement must satisfy NPO), FR#19 (flow identification finer than per top level LSP), MR#6 (operator initiated flow rebalance)

backward compatibility and migration

FR#4 (smooth incremental deployment), FR#6 (management and diagnostics must continue to function), DR#1 (extend existing protocols), DR#2 (extend LDP, no LDP TE)

delay and delay variation

FR#7 (expose lower layer measured delay), FR#8 (precision of latency reporting), FR#9 (limit latency on per LSP basis), FR#15 (minimum delay path), FR#16 (bounded delay path), FR#17 (bounded jitter path)

admission control, preemption, traffic engineering

FR#10 (admission control, preemption), FR#14 (packet ordering), FR#21 (ingress specification of path), FR#22 (path symmetry), DR#3 (IP and LDP traffic), MR#3 (management specification of path)

single vs multiple domain

DR#4 (IGP extensions allowed within single domain), DR#5 (IGP extensions disallowed in multiple domain case)

general network management

MR#1 (polling, configuration, and notification), MR#2 (activation and de-activation)

path determination, connectivity verification

MR#4 (path trace), MR#5 (connectivity verification)

The above list is not intended as a substitute for [\[I-D.ietf-rtgwg-cl-requirement\]](#), but rather as a concise grouping and reminder or requirements to serve as a means of more easily determining requirements coverage of a set of protocol documents.

[7.2.](#) Required Document Coverage

The primary areas where additional protocol extensions and mechanisms are required include the topics described in the following subsections.

There are candidate documents for a subset of the topics below. This grouping of topics does not require that each topic be addressed by a separate document. In some cases, a document may cover multiple topics, or a specific topic may be addressed as applicable in multiple documents.

[7.2.1.](#) Component Link Grouping

An extension to link bundling is needed to specify a group of components with common attributes. This can be a TLV defined within the link bundle that carries the same encapsulations as the link bundle. Two interface indices would be needed for each group.

- a. An index is needed that if included in an ERO would indicate the need to place the LSP on any one component within the group.
- b. A second index is needed that if included in an ERO would indicate the need to balance flows within the LSP across all components of the group. This is equivalent to the "all-ones" component for the entire bundle.

[I-D.ospf-cc-stlv] can be extended to include multipath treatment capabilities. An ISIS solution is also needed. An extension of RSVP-TE signaling is needed to indicate multipath treatment preferences.

If a component group is allowed to support all of the parameters of a link bundle, then a group TE metric would be accommodated. This can be supported with the component TLV (C-TLV) defined in [\[I-D.ospf-cc-stlv\]](#).

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) is the "routing information aggregation" set of requirements. The "restoration speed", "backward compatibility and migration", and "general network management" requirements must also be considered.

[7.2.2.](#) Delay and Jitter Extensions

A extension is needed in the IGP-TE advertisement to support delay and delay variation for links, link bundles, and forwarding adjacencies. Whatever mechanism is described must take precautions that insure that route oscillations cannot occur. [\[I-D.wang-ccamp-latency-te-metric\]](#) may be a good starting point.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) is the "delay and delay variation" set of requirements. The "restoration speed", "backward compatibility and migration", and "general network management" requirements must also be considered.

[7.2.3.](#) Path Selection and Admission Control

Path selection and admission control changes must be documented in each document that proposes a protocol extension that advertises a new capability or parameter that must be supported by changes in path selection and admission control.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) are the "load distribution, stability, minimal disruption" and "admission control, preemption, traffic engineering"

sets of requirements. The "restoration speed" and "path determination, connectivity verification" requirements must also be considered. The "backward compatibility and migration", and "general network management" requirements must also be considered.

7.2.4. Dynamic Multipath Balance

FR#11 explicitly calls for dynamic load balancing similar to existing adaptive multipath. In implementations where flow identification uses a coarse granularity, the adjustments would have to be equally coarse, in the worst case moving entire LSP. The impact of flow identification granularity and potential adaptive multipath approaches may need to be documented in greater detail than provided here.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) are the "restoration speed" and the "load distribution, stability, minimal disruption" sets of requirements. The "path determination, connectivity verification" requirements must also be considered. The "backward compatibility and migration", and "general network management" requirements must also be considered.

7.2.5. Frequency of Load Balance

IGP-TE and RSVP-TE extensions are needed to support frequency of load balancing rearrangement called for in FR#13, and FR#15-FR#17. Constraints are not defined in RSVP-TE, but could be modeled after administrative attribute affinities in [RFC3209](#) and elsewhere.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) is the "load distribution, stability, minimal disruption" set of requirements. The "path determination, connectivity verification" must also be considered. The "backward compatibility and migration" and "general network management" requirements must also be considered.

7.2.6. Inter-Layer Communication

Lower layer to upper layer communication called for in FR#7 and FR#20. This is addressed for a subset of parameters related to packet ordering in [[I-D.villamizar-mpls-tp-multipath](#)] where layers are MPLS. Remaining parameters, specifically delay and delay variation, need to be addressed. Passing information from a lower non-MPLS layer to an MPLS layer needs to be addressed, though this may largely be generic advice encouraging a coupling of MPLS to lower layer management plane or control plane interfaces. This topic can be addressed in each document proposing a protocol extension, where applicable.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) is the "restoration speed" set of requirements. The "backward compatibility and migration" and "general network management" requirements must also be considered.

[7.2.7.](#) Packet Ordering Requirements

A document is needed to define extensions supporting various packet ordering requirements, ranging from requirements to preserve microflow ordering only, to requirements to preserve full LSP ordering (as in MPLS-TP). This is covered by [\[I-D.villamizar-mpls-tp-multipath\]](#) and [\[I-D.villamizar-mpls-tp-multipath-te-extn\]](#).

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) are the "admission control, preemption, traffic engineering" and the "path determination, connectivity verification" sets of requirements. The "backward compatibility and migration" and "general network management" requirements must also be considered.

[7.2.8.](#) Minimally Disruption Load Balance

The behavior of hash methods used in classic multipath needs to be described in terms of FR#12 which calls for minimally disruptive load adjustments. For example, reseeding the hash violates FR#12. Using modulo operations is significantly disruptive if a link comes or goes down, as pointed out in [\[RFC2992\]](#). In addition, backwards compatibility with older hardware needs to be accommodated.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) is the "load distribution, stability, minimal disruption" set of requirements.

[7.2.9.](#) Path Symmetry

Protocol extensions are needed to support dynamic load balance as called for to meet FR#22 (path symmetry) and to meet FR#11 (dynamic placement of flows). Currently path symmetry can only be supported in link bundling if the path is pinned. When a flow is moved both ingress and egress must make the move as close to simultaneously as possible to satisfy FR#22 and FR#12 (minimally disruptive load rebalance). If a group of flows are identified using a hash, then the hash must be identical on the pair of LSR at the endpoint, using the same hash seed and with one side swapping source and destination. If the label stack is used, then either the entire label stack must be a special case flow identification, since the set of labels in either direction are not correlated, or the two LSR must conspire to use the same flow identifier. For example, using a common entropy

label value, and using only the entropy label in the flow identification would satisfy this requirement.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) are the "load distribution, stability, minimal disruption" and the "admission control, preemption, traffic engineering" sets of requirements. The "backward compatibility and migration" and "general network management" requirements must also be considered. Path symmetry simplifies support for the "path determination, connectivity verification" set of requirements, but with significant complexity added elsewhere.

[7.2.10](#). Performance, Scalability, and Stability

A separate document providing analysis of performance, scalability, and stability impacts of changes may be needed. The topic of traffic adjustment oscillation must also be covered. If sufficient coverage is provided in each document covering a protocol extension, a separate document would not be needed.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) is the "restoration speed" set of requirements. This is not a simple topic and not a topic that is well served by scattering it over multiple documents, therefore it may be best to put this in a separate document and put citations in documents called for in [Section 7.2.1](#), [Section 7.2.2](#), [Section 7.2.3](#), [Section 7.2.9](#), [Section 7.2.11](#), [Section 7.2.12](#), [Section 7.2.13](#), and [Section 7.2.14](#). Citation may also be helpful in [Section 7.2.4](#), and [Section 7.2.5](#).

[7.2.11](#). IP and LDP Traffic

A document is needed to define the use of measurements native IP and native LDP traffic levels to reduce link advertised bandwidth amounts.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) are the "load distribution, stability, minimal disruption" and the "admission control, preemption, traffic engineering" set of requirements. The "path determination, connectivity verification" must also be considered. The "backward compatibility and migration" and "general network management" requirements must also be considered.

[7.2.12](#). LDP Extensions

Extending LDP is called for in DR#2. LDP can be extended to couple FEC admission control to local resource availability without providing LDP traffic engineering capability. Other LDP extensions

such as signaling a bound on microflow size and LDP LSP requirements would provide useful information without providing LDP traffic engineering capability.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) is the "admission control, preemption, traffic engineering" set of requirements. The "backward compatibility and migration" and "general network management" requirements must also be considered.

[7.2.13.](#) Pseudowire Extensions

PW extensions such as signaling a bound on microflow size and PW requirements would provide useful information.

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) is the "admission control, preemption, traffic engineering" set of requirements. The "backward compatibility and migration" and "general network management" requirements must also be considered.

[7.2.14.](#) Multi-Domain Composite Link

DR#5 calls for Composite Link to span multiple network topologies. Component LSP may already span multiple network topologies, though most often in practice these are LDP signaled. Component LSP which are RSVP-TE signaled may also span multiple network topologies using at least three existing methods (per domain [[RFC5152](#)], BRPC [[RFC5441](#)], PCE [[RFC4655](#)]). When such component links are combined in a Composite Link, the Composite Link spans multiple network topologies. It is not clear in which document this needs to be described or whether this description in the framework is sufficient. The authors and/or the WG may need to discuss this. DR#5 mandates that IGP-TE extension cannot be used. This would disallow the use of [[RFC5316](#)] or [[RFC5392](#)] in conjunction with [[RFC5151](#)].

The primary focus of this document, among the sets of requirements listed in [Section 7.1](#) are "single vs multiple domain" and "admission control, preemption, traffic engineering". The "routing information aggregation" and "load distribution, stability, minimal disruption" requirements need attention due to their use of the IGP in single domain Composite Link. Other requirements such as "delay and delay variation", can more easily be accommodated by carrying metrics within BGP. The "path determination, connectivity verification" requirements need attention due to requirements to restrict disclosure of topology information across domains in multi-domain deployments. The "backward compatibility and migration" and "general network management" requirements must also be considered.

7.3. Open Issues Regarding Requirements

Note to co-authors: This section needs to be reduced to an empty section and then removed.

The following topics in the requirements document are not addressed. Since they are explicitly mentioned in the requirements document some mention of how they are supported is needed, even if to say nother needed to be done. If we conclude any particular topic is irrelevant, maybe the topic should be removed from the requirement document. At that point we could add the management requirements that have come up and were missed.

1. L3VPN [RFC 4364](#), [RFC 4797](#), L2VPN [RFC 4664](#), VPWS, VPLS [RFC 4761](#), [RFC 4762](#) and VPMS VPMS Framework ([draft-ietf-l2vpn-vpms-frmwk-requirements](#)). It is not clear what additional Composite Link requirements these references imply, if any. If no additional requirements are implied, then these references are considered to be informational only.
2. Migration may not be adequately covered in [Section 4.1.5](#). It might also be necessary to say more here on performance, scalability, and stability as it related to migration. Comments on this from co-authors or the WG?
3. We may need a performance section in this document to specifically address #DR6 (fast convergence), and #DR7 (fast worst case failure convergence), though we do already have scalability discussion. The performance section would have to say "no worse than before, except were there was no alternative to make it very slightly worse" (in a bit more detail than that). It would also have to better define the nature of the performance criteria.

7.4. Framework Requirement Coverage by Protocol

As an aid to implementors, this section summarizes requirement coverage listed in [Section 7.2](#) by protocol or LSR functionality affected.

Some documentation may be purely informational, proposing no changes and proposing usage at most. This includes [Section 7.2.3](#), [Section 7.2.8](#), [Section 7.2.10](#), and [Section 7.2.14](#).

[Section 7.2.9](#) may require a new protocol.

7.4.1. OSPF-TE and ISIS-TE Protocol Extensions

Many of the changes listed in [Section 7.2](#) require IGP-TE changes, though most are small extensions to provide additional information. This set includes [Section 7.2.1](#), [Section 7.2.2](#), [Section 7.2.5](#), [Section 7.2.6](#), and [Section 7.2.7](#). An adjustment to existing advertised parameters is suggested in [Section 7.2.11](#).

7.4.2. PW Protocol Extensions

The only suggestion of pseudowire (PW) extensions is in [Section 7.2.13](#).

7.4.3. LDP Protocol Extensions

Potential LDP extensions are described in [Section 7.2.12](#).

7.4.4. RSVP-TE Protocol Extensions

RSVP-TE protocol extensions are called for in [Section 7.2.1](#), [Section 7.2.5](#), [Section 7.2.7](#), and [Section 7.2.9](#).

7.4.5. RSVP-TE Path Selection Changes

[Section 7.2.3](#) calls for path selection to be addressed in individual documents that require change. These changes would include those proposed in [Section 7.2.1](#), [Section 7.2.2](#), [Section 7.2.5](#), and [Section 7.2.7](#).

7.4.6. RSVP-TE Admission Control and Preemption

When a change is needed to path selection, a corresponding change is needed in admission control. The same set of sections applies: [Section 7.2.1](#), [Section 7.2.2](#), [Section 7.2.5](#), and [Section 7.2.7](#). Some resource changes such as a link delay change might trigger preemption. The rules of preemption remain unchanged, still based on holding priority.

7.4.7. Flow Identification and Traffic Balance

The following describe either the state of the art in flow identification and traffic balance or propose changes: [Section 7.2.4](#), [Section 7.2.5](#), [Section 7.2.7](#), and [Section 7.2.8](#).

8. Security Considerations

The security considerations for MPLS/GMPLS and for MPLS-TP are

documented in [[RFC5920](#)] and [[I-D.ietf-mpls-tp-security-framework](#)].

The types protocol extensions proposed in this framework document provide additional information about links, forwarding adjacencies, and LSP requirements. The protocol semantics changes described in this framework document propose additional LSP constraints applied at path computation time and at LSP admission at midpoints LSR. The additional information and constraints provide no additional security considerations beyond the security considerations already documented in [[RFC5920](#)] and [[I-D.ietf-mpls-tp-security-framework](#)].

9. Acknowledgments

Authors would like to thank Adrian Farrel, Fred Jounay, Yuji Kamite for his extensive comments and suggestions regarding early versions of this document, Ron Bonica, Nabil Bitar, Eric Gray, Lou Berger, and Kireeti Kompella for their reviews of early versions and great suggestions.

Authors would like to thank Iftekhar Hussain for review and suggestions regarding recent versions of this document.

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.
- [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.
- [RFC3630] Katz, D., Kompella, K., and D. Yeung, "Traffic Engineering (TE) Extensions to OSPF Version 2", [RFC 3630](#), September 2003.
- [RFC4201] Kompella, K., Rekhter, Y., and L. Berger, "Link Bundling in MPLS Traffic Engineering (TE)", [RFC 4201](#), October 2005.

- [RFC4206] Kompella, K. and Y. Rekhter, "Label Switched Paths (LSP) Hierarchy with Generalized Multi-Protocol Label Switching (GMPLS) Traffic Engineering (TE)", [RFC 4206](#), October 2005.
- [RFC5036] Andersson, L., Minei, I., and B. Thomas, "LDP Specification", [RFC 5036](#), October 2007.
- [RFC5305] Li, T. and H. Smit, "IS-IS Extensions for Traffic Engineering", [RFC 5305](#), October 2008.
- [RFC5712] Meyer, M. and JP. Vasseur, "MPLS Traffic Engineering Soft Preemption", [RFC 5712](#), January 2010.
- [RFC6107] Shiimoto, K. and A. Farrel, "Procedures for Dynamically Signaled Hierarchical Label Switched Paths", [RFC 6107](#), February 2011.
- [RFC6374] Frost, D. and S. Bryant, "Packet Loss and Delay Measurement for MPLS Networks", [RFC 6374](#), September 2011.
- [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.

10.2. Informative References

- [DBP] Bertsekas, D., "Dynamic Behavior of Shortest Path Routing Algorithms for Communication Networks", IEEE Trans. Auto. Control 1982.
- [I-D.ietf-mpls-entropy-label]
Drake, J., Kompella, K., Yong, L., Amante, S., and W. Henderickx, "The Use of Entropy Labels in MPLS Forwarding", [draft-ietf-mpls-entropy-label-01](#) (work in progress), October 2011.
- [I-D.ietf-mpls-explicit-resource-control-bundle]
Zamfir, A., Ali, Z., and P. Dimitri, "Component Link Recording and Resource Control for TE Links", [draft-ietf-mpls-explicit-resource-control-bundle-10](#) (work in progress), April 2011.
- [I-D.ietf-mpls-tp-security-framework]
Niven-Jenkins, B., Fang, L., Graveman, R., and S. Mansfield, "MPLS-TP Security Framework", [draft-ietf-mpls-tp-security-framework-02](#) (work in progress), October 2011.

[I-D.ietf-rtgwg-cl-requirement]

Malis, A., Villamizar, C., McDysan, D., Yong, L., and N. So, "Requirements for MPLS Over a Composite Link", [draft-ietf-rtgwg-cl-requirement-05](#) (work in progress), January 2012.

[I-D.kompella-mpls-rsvp-ecmp]

Kompella, K., "Multi-path Label Switched Paths Signaled Using RSVP-TE", [draft-kompella-mpls-rsvp-ecmp-01](#) (work in progress), October 2011.

[I-D.ospf-cc-stlv]

Osborne, E., "Component and Composite Link Membership in OSPF", [draft-ospf-cc-stlv-00](#) (work in progress), August 2011.

[I-D.symmvo-rtgwg-cl-use-cases]

Malis, A., Villamizar, C., McDysan, D., Yong, L., and N. So, "Composite Link Use Cases and Design Considerations", [draft-symmvo-rtgwg-cl-use-cases-00](#) (work in progress), February 2012.

[I-D.villamizar-mpls-tp-multipath]

Villamizar, C., "Use of Multipath with MPLS-TP and MPLS", [draft-villamizar-mpls-tp-multipath-01](#) (work in progress), March 2011.

[I-D.villamizar-mpls-tp-multipath-te-extn]

Villamizar, C., "Multipath Extensions for MPLS Traffic Engineering", [draft-villamizar-mpls-tp-multipath-te-extn-00](#) (work in progress), July 2011.

[I-D.wang-ccamp-latency-te-metric]

Fu, X., Betts, M., Wang, Q., McDysan, D., and A. Malis, "GMPLS extensions to communicate latency as a traffic engineering performance metric", [draft-wang-ccamp-latency-te-metric-03](#) (work in progress), March 2011.

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

[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.

- [RFC3260] Grossman, D., "New Terminology and Clarifications for Diffserv", [RFC 3260](#), April 2002.
- [RFC3468] Andersson, L. and G. Swallow, "The Multiprotocol Label Switching (MPLS) Working Group decision on MPLS signaling protocols", [RFC 3468](#), February 2003.
- [RFC3945] Mannie, E., "Generalized Multi-Protocol Label Switching (GMPLS) Architecture", [RFC 3945](#), October 2004.
- [RFC3985] Bryant, S. and P. Pate, "Pseudo Wire Emulation Edge-to-Edge (PWE3) Architecture", [RFC 3985](#), March 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.
- [RFC4655] Farrel, A., Vasseur, J., and J. Ash, "A Path Computation Element (PCE)-Based Architecture", [RFC 4655](#), August 2006.
- [RFC4928] Swallow, G., Bryant, S., and L. Andersson, "Avoiding Equal Cost Multipath Treatment in MPLS Networks", [BCP 128](#), [RFC 4928](#), June 2007.
- [RFC5151] Farrel, A., Ayyangar, A., and JP. Vasseur, "Inter-Domain MPLS and GMPLS Traffic Engineering -- Resource Reservation Protocol-Traffic Engineering (RSVP-TE) Extensions", [RFC 5151](#), February 2008.
- [RFC5152] Vasseur, JP., Ayyangar, A., and R. Zhang, "A Per-Domain Path Computation Method for Establishing Inter-Domain Traffic Engineering (TE) Label Switched Paths (LSPs)", [RFC 5152](#), February 2008.
- [RFC5316] Chen, M., Zhang, R., and X. Duan, "ISIS Extensions in Support of Inter-Autonomous System (AS) MPLS and GMPLS Traffic Engineering", [RFC 5316](#), December 2008.
- [RFC5392] Chen, M., Zhang, R., and X. Duan, "OSPF Extensions in Support of Inter-Autonomous System (AS) MPLS and GMPLS Traffic Engineering", [RFC 5392](#), January 2009.
- [RFC5441] Vasseur, JP., Zhang, R., Bitar, N., and JL. Le Roux, "A Backward-Recursive PCE-Based Computation (BRPC) Procedure to Compute Shortest Constrained Inter-Domain Traffic Engineering Label Switched Paths", [RFC 5441](#), April 2009.

[RFC5920] Fang, L., "Security Framework for MPLS and GMPLS Networks", [RFC 5920](#), July 2010.

[RFC5921] Bocci, M., Bryant, S., Frost, D., Levrau, L., and L. Berger, "A Framework for MPLS in Transport Networks", [RFC 5921](#), July 2010.

Authors' Addresses

So Ning
Tata Communications

Email: ning.so@tatacommunications.com

Dave McDysan
Verizon
22001 Loudoun County PKWY
Ashburn, VA 20147

Email: dave.mcdysan@verizon.com

Eric Osborne
Cisco

Email: eosborne@cisco.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

