Routing Working Group Internet-Draft Intended status: Standards Track Expires: January 16, 2014

R. Shakir RT D. Vernals Vodafone A. Capello Telecom Italia July 15, 2013

Performance Engineered LSPs using the Segment Routing Data-Plane draft-shakir-rtgwg-sr-performance-engineered-lsps-00

Abstract

A number of applications and services running over IP/MPLS networks have strict requirements relating to their routing, or the performance of the path supporting their traffic flow, for instance, in terms of characteristics such as latency, loss, or bandwidth availability. Segment routing provides a means by which the dataplane of an IP/MPLS network can be programmed to support such "performance engineered" paths. This document describes an architecture for the use of such performance engineered label switched paths, and the control-plane functionality required to allow both distributed and centralised computation of acceptable forwarding paths.

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of <u>BCP 78</u> and <u>BCP 79</u>.

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 16, 2014.

Copyright Notice

Copyright (c) 2013 IETF Trust and the persons identified as the document authors. All rights reserved.

Shakir, et al. Expires January 16, 2014

[Page 1]

This document is subject to **BCP** 78 and the IETF Trust's Legal Provisions Relating to IETF Documents (<u>http://trustee.ietf.org/license-info</u>) 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. Motivation |
|---|
| $\overline{2}$. Conventions Used in This Document |
| $\overline{3}$. Data-plane Path Selection |
| 3.1. SID Selection for Non-Revertive Services |
| 3.2. SID Selection for Revertive Services |
| 3.2.1. Procedure for Link Protection of Adi-SIDs |
| 3.2.2. Procedure for Node Protection of Adj-SIDs |
| 3.2.3. Example of Revertive Adi-SID Protection |
| 3 3 Path Re-Ontimisation and Re-Routing |
| 4 Distributed Path Computation via Constrained Shortest-Path |
| Algorithms |
| 1 Path Selection based on Static TGP Path Attributes 13 |
| 4.1. Path Selection based on Performance Pelated ICP Path |
| 4.2. Facily Selection based on Ferrormance-Netated for Facily |
| $4.2.1 \qquad \text{Product remember for ICD Attributes Partaining to}$ |
| 4.2.1. Requirements for formance |
| Aujacency reflormance $\dots \dots \dots$ |
| <u>5</u> . Centralised Path Computation Using PCE <u>10</u> |
| 5.1. Use of Path Computation Element to Provide Inter-Area |
| SR LSPS \dots |
| 5.2. Providing Co-routed or Multi-Layer Aware LSPs using PCE . 1/ |
| 5.2.1. Co-Routed LSPs |
| 5.2.2. Resource Reservation and Admission Control through |
| a Stateful PCE |
| <u>5.2.3</u> . Multi-Layer Calculation through a Common PCE <u>18</u> |
| $\underline{6}$. Security Considerations |
| <u>7</u> . Acknowledgements |
| <u>8</u> . Normative References |
| Authors' Addresses |

1. Motivation

For numerous applications running over IP/MPLS networks, there is a requirement to provide paths that have guaranteed network performance. These resources guarantees may be in terms of sufficient bandwidth being available for a traffic flow, but also can be in terms of other characteristics (such as latency, packet delay variation, and packet loss). In addition to such characteristics of the underlying network, requirements related to path routing can exist to ensure that a path offers characteristics such as affinity to particular infrastructure, or disjointness to another service. For instance:

- o Where two services provided by an IP/MPLS network make up part of an active/backup or live/live service pair for a transported application it is required that the paths are wholly disjoint (shared risk, link, and node) to ensure that they do not fail simultaneously.
- o If the service a network provides supports an application that requires a particular end-to-end latency budget, then the service must be constrained to a path, or paths, meeting this budget and the path made unavailable if these characteristics cannot be met.
- o Where a service provided by the IP/MPLS network makes up part of another network's topology (e.g., an ATM PWE3 service provided within an IP/MPLS network may form a part of a wider client ATN network), then an affinity to particular links within the network (such as particular sub-sea cable systems) may be required. Where such a path is not available, it can be preferable to utilise protection within the client network rather than re-route the IP/ MPLS service.
- o Where services, such as paths for real-time voice and video trunking delivered over IP/MPLS services are routed according to paths with guaranteed resource availability (such as available bandwidth).
- o Where the service requires that both directions of the network path are co-routed, such as where an IP/MPLS network path is used to carry IEEE 1588 synchronisation traffic. In this case, two uni-directional LSPs must be routed in a co-ordinated manner, which may diverge from the shortest-path within the network.

These requirements can be generalised into a need to support arbitrary constrained paths within an IP/MPLS network - with the constraints being both in terms of the path selected in the network (and its underlying characteristics) and the treatment of packets

forwarding onto this path during network events (such as during link failures, or re-convergence).

It is important to note that these routing requirements are inherently related to a particular service (e.g., customer service A must be disjoint to customer service B, or a customer service must only be available if paths with an end-to-end latency of less than 200 milliseconds are available between node X and node Y) - and apply to all traffic (as may be the case within a pair of PWE3 services) or a subset of traffic within the service (as may be the case with particular treatment of voice traffic within an L3VPN service). This results in the number of such routed paths in the network being dependent upon the number of services supported by the network (order of tens or hundreds of thousands), rather than to the number of edge devices within a network (typically of the order of hundreds, or thousands).

It is possible to utilise Segment Routing (SR) as described in [I-D.filsfils-rtgwg-segment-routing] to provide a means to support services with constrained path requirements via label-switched paths (LSPs) in an IP/MPLS network without requiring devices within the core of the network to maintain per-LSP state. Since in the limits described above, this number of LSPs is proportional to the number of services on the network utilising a stateless mechanism (such as SR) to provide such forwarding paths equates to avoiding per-service state on transit devices. The forwarding paths utilised to support such services are referred to as performance engineered LSPs within this document.

For example, considering the topology shown in Figure 1, if the ingress label edge router (LER) requires forwarding of traffic on a path to the egress LER which does not exceed 40 milliseconds, it must ensure that traffic traverses the iLER-A, A-B, and B-eLER links.

> lbl 1100 lbl 1200 lbl 1700 ---(10ms)--- A ---(5ms)--- B --(20ms)---\ / (30ms) lbl 1400 [iLER] [eLER] \ / ---(50ms)--- X -----(10ms) ----lbl 1500 lbl 1600

Figure 1

With segment routing, iLER can apply a label stack of {1200,1700} and next-hop of A to influence A to forward packets to B, and B to

forward to eLER, regardless of the shortest path selection due to the IGP metrics within the network topology.

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 [RFC2119].

3. Data-plane Path Selection

When considering services that are to be carried via performance engineered paths within a network, constraints are introduced relating to protection during failures of the explicit path selected through the network. Particularly, where a path is provided with particular link affinity, or as part of an active/active service pair where both primary and backup LSPs are routed within certain performance constraints, it is not desirable to perform dynamic rerouting or fast re-route protection for the traffic within the service, as a path violating the performance constraints may be introduced, where the alternate route made available continues to be compliant with the original routing criteria. In order to allow an operator to route services not requiring reversion from the primary path, an implementation MUST allow the advertisement of segment identifiers which are explicitly excluded from fast re-route, and reversion away from the primary path.

Where services are specified as being suitable for reversion, some consideration is required as to the means by which they are rerouted. For instance, where with some IP FRR mechanisms the protection path may follow the shortest-path tree from the point of local repair (PLR) to the packet destination (effectively making the LSP tail-end the merge point - MP), such re-routing behaviour is not desirable for all performance engineered paths. In such cases, it is more preferable to provide means by which (during protection events) traffic is routed back on to the original LSP path, as soon as possible, essentially minimising the divergence away from the path calculated to meet service constraints. An implementation SHOULD provide means by which an operator can influence MP selection to support such requirements.

3.1. SID Selection for Non-Revertive Services

Following the calculation of a route meeting a particular set of constraints, the ingress service edge device should select the data path through determining the relevant SIDs to be pushed to the packet. During the translation of a route to a set of SIDs the ingress device MUST consider the service requirements in order to ensure that protection within the network does not result in path constraints being violated where this is unacceptable. Where a service is specified to be non-revertive, the iLER MUST utilise only segment identifiers which are explicitly identified to be nonrevertive. Since protection requirements vary on per-service basis, the control of reversion to other paths during failures SHOULD be specified on a per-LSP basis on the ingress router.

It should be noted that where strict path constraints are required,

this results in the number of SIDs applied to a packet being proportion to the number of links traversed through the topology (through disallowing use of Node-SIDs), which may have implications for certain hardware implementations.

A network operator MAY create forwarding adjacencies consisting of multiple SIDs utilising the Explicit Route Object encoding of the Adjacency-SID specified in

[I-D.previdi-isis-segment-routing-extensions] for IS-IS and [I-D.psenak-ospf-segment-routing-extensions] for OSPF, such that the total SID-stack to be imposed is minimised. Where such forwarding adjacency LSPs (FA-LSPs), or other paths consisting of multiple segments are re-advertised, the path characteristics of the underlying links MUST be included, such that other constraints used for calculation (such as shared risk link groups) can be considered by the calculating iLER. In order that the routing of non-segment routed traffic and devices not supporting SR is not influenced by the advertisement of such adjacencies:

- o It MUST be possible for an operator to specify policies as to whether such forwarding-adjacencies are utilised for non-Segment Routed traffic within the IP/MPLS network.
- o The advertisement of FA-LSPs for the use of performance engineered LSPs SHOULD NOT negatively impact the scaling and performance of devices running vanilla SPF calculations of the network topology (in order to avoid introducing additional computational overhead to legacy devices within the IGP domain within which such LSPs are to be introduced).

3.2. SID Selection for Revertive Services

3.2.1. Procedure for Link Protection of Adj-SIDs

By default, Adj-SID values refer to an particular individual or set of physical or logical adjacencies between two devices. It is therefore linked specifically to a specific path between two nodes , and hence (by default) does not have a viable alternate route. Where a revertive Adj-SID is advertised, specified through the 'B' flag of the Adj-SID advertisement described in

[I-D.previdi-isis-segment-routing-extensions] and

[I-D.psenak-ospf-segment-routing-extensions] the advertising LSR MUST calculate a backup path for this adjacency.

By default an LSR SHOULD calculate a link-protecting tunnel to the node to which the adjacency is received on - this can be achieved through mechanisms such as Loop-Free Alternates [RFC5286]. During failure of a path advertised with a revertive Adj-SID, the LSR

detecting the adjacency failure should act as the point of local repair (PLR) and SHOULD pop the adjacency segment (as per the default Adj-SID action). In order to reach the merge point (MP), it is possible for the PLR to utilise either:

- o A set of SIDs relating to the loop free alternate path to reach the MP - in this case, it should be noted that such a set of SIDs may relate to multiple node and/or adjacency SIDs, where a Remote or Directed LFA is required to reach the MP.
- o The adjacency segments relating to the calculated path between the PLR and the MP. Utilising Adj-SIDs requires the PLR to perform no calculation of the path between its neighbours and the MP, however, may result in a less survivable service, in cases where simultaneous failures result in the backup SR-LSP specified by the set of Adj-SIDs becoming unavailable.

In cases where particular policies should be enforced for the protection path for an Adj-SID, an implementation SHOULD utilise a set of Adj-SIDs that indicate the links to be traversed between the PLR and the MP, based on characteristics of these adjacencies (e.g., the maximum total link bandwidth path). Where such Adj-SID based backup path selection is utilised, the path selected SHOULD be influenced by operator policy in a similar manner as the LFA selection considered in [I-D.ietf-rtgwg-lfa-manageability].

It should be noted that since the selection of the protecting set of SIDs is calculated on a per-Adj-SID basis, no particular backup path selection can be performed by a transit LSR on a per-service basis. Therefore, where revertive SIDs are utilised an operator SHOULD recognise that during protection events, no path characteristics, or resource constraints can be met whilst re-routing results in the service diverging from the specified explicit path.

3.2.2. Procedure for Node Protection of Adj-SIDs

An LSR MAY provide node-protection for an Adj-SID if such a nodeprotecting path exists within the network topology.

In the case where such a path is available, the LSR acting as the PLR must be capable of programming its forwarding plane based on the tuple of the top two labels of the SID stack. Where this can be achieved, a backup action to push the corresponding set of SIDs to reach the next-next-hop node (indicated by the advertising entity of the second entry in the label stack) during the failure of the primary Adj-SID. The PLR must therefore program an action to pop the first two labels of the ingress packet, and subsequently push the SIDs relating to this path.

Again, it is possible for a node to utilise either a set of Node or Adjacency SIDs to reach the next-next-hop node (MP). Where Node-SIDs are utilised, means to determine a loop-free path MUST be used to determine the set of Node-SIDs required. Where Adj-SIDs are utilised for such functionality, no such calculation is required. Where certain policies are to be enforced for the protecting path, an implementation SHOULD allow the use of Adj-SIDs to determine the path utilised, and the selection of these SIDs SHOULD be influenced by operator policy.

3.2.3. Example of Revertive Adj-SID Protection





In Figure 2 a source X sends traffic utilising Adj-SIDs to a destination Y utilising through pushing the relevant Adj-SIDs to traverse A-B, B-C, C-D. A therefore receives a packet with {10,20,30} segments applied. B's FIB is programmed with a pop() operation for Adj-SID 20, along with a next-hop of the B-C link.

To provide a link-protecting backup, if E has been calculated as a link-protecting LFA for segment 20, then B programs a backup action of push(9003) for the ingress label of 20, with an egress interface of the B-E link. When E receives this labelled packet, it swaps label 9003 for label 9003 (as per standard behaviour for a Node-SID) , and forwards directly to C, who receives the packet with the label 30 exposed, and hence acts as the MP.

Clearly, an alternative is that B programs a backup action to push label 90 to this stack (with a next-hop of E) to ensure that the adjacency between E-C is utilised, rather than E's IGP shortest-path to C.

To provide node protection for the Adj-SID, then additional complexity is introduced. If B receives a packet destined for the Adj-SID indicating link B-C, B can examine the following label within the stack (in this case, label 30) which is advertised by D within the topology. Since there is a node-protecting LFA via E to D, B may therefore pop() the subsequent Adj-SID and push the Node-SID of D (9004), whilst forwarding the packet via the B-E link.

Alternatively, it is possible for B to utilise an explicitly derived path to reach D (namely, the Adj-SID of the E-D link, with a next-hop of E), to reach the MP. In this case, no calculation as to the routing behaviour of E is required to determine this protection path (trading computational complexity for increasing the length of the protection SID stack). Through this behaviour, a packet can be forwarded during the complete failure of C. It should be noted that this requires two look-ups on the PLR rather than the single look-up required in the link protecting case.

3.3. Path Re-Optimisation and Re-Routing

Where performance-engineered SR paths are selected by a head-end, the calculation of the path is based on the information that is available to the computing entity (be it centralised or distributed) at the time of calculation. In order to ensure that such paths are rerouted onto more optimal paths where available, an ingress LER MUST perform periodic re-optimisation whereby the path selected for a service is recalculated. The period between such re-optimisations SHOULD be configurable by a network operator.

Unlike other network technologies which can be utilised to specify explicit paths within an IP/MPLS network, the mid-point network elements are unaware of the LSPs that traverse them. There is therefore a requirement for an SR head-end to determine when specific segments are no longer valid to be utilised for a service to be routed. In order to ensure that traffic is forwarded onto paths that remain valid, a head-end device MUST trigger re-calculation of explicit paths within the network when it receives an IGP update relating to the segment utilised within a particular service topology. In order to provide means to tolerate short-lived failures (particularly where such services are revertive), it SHOULD be possible to delay such recalculation on a per-service basis. Such triggered re-optimisation MUST be performed for IGP updates that withdraw segments from the topology and MAY be triggered based on updates to other attributes within the network. Where updates are triggered on information that may rapidly change within the IGP (e.g., information relating to bandwidth reservation or utilisation) an iLER device SHOULD provide means to limit the period between reoptimisations, or provide thresholds over which re-optimisation is triggered.

In addition to re-optimisation based on failures, an iLER SHOULD provide means by which per-service OAM measuring performance or liveliness characteristics of a particular path can trigger a path to be withdrawn from use, and/or re-optimisation of the SID selection for the path. Such per-service OAM is critical within multi-area environments where it cannot be guaranteed that a head-end device

will have all routing information propagated to it - in such deployments, an implementation MUST support per-service OAM. In all environments, per-service OAM can be utilised to ensure that a service can be withdrawn more quickly than IGP-updates relating to segment failures can be propagated, or a head-end is able to react to "grey" failure events, where data-plane traffic forwarding has failed, but no IGP update is generated.

<u>4</u>. Distributed Path Computation via Constrained Shortest-Path Algorithms

4.1. Path Selection based on Static IGP Path Attributes

To determine the relevant set of segment identifiers to be utilised for a service, existing constrained shortest-path (CSPF) functions such as those used for other route selection mechanisms can be utilised. This can provide route selection based on IGP traffic engineering metrics (such as those specified for OSPF in <u>RFC3630</u>, or IS-IS in <u>RFC5305</u>), or GMPLS IGP extensions such as shared risk link groups (SRLGs) carried in attributes such as those that are described for OSPF in <u>RFC5307</u> and IS-IS in <u>RFC4203</u>. An implementation providing distributed CSPF to provide performance-engineered SR paths SHOULD support path selection through consideration of such traffic engineering IGP attributes.

Where adjacency segments are created for use as forwarding adjacency LSPs, or as a means to provide compression of SID stacks, an implementation MUST include the relevant IGP traffic engineering attributes indicating the characteristics of the underlying Adj-SIDs within IGP attributes relating to such segments.

<u>4.2</u>. Path Selection based on Performance-Related IGP Path Attributes

In addition to considering such static attributes of links within the IGP topology, distributed path computation can be triggered based on performance monitoring information propagated into IGP attributes such as those described in [I-D.giacalone-ospf-te-express-path] and [I-D.ietf-isis-te-metric-extensions]. Consideration of such attributes allows paths to be calculated based on the underlying loss and delay characteristics of a network path. Through monitoring updates to these attributes advertised through IGP update messages, re-routes based on changes in the performance characteristics of a path can be achieved. An iLER supporting performance engineered LSPs utilising the SR dataplane SHOULD allow consideration of these attributes when performing ERO calculation, and SHOULD provide means to trigger re-routes based on changes in their values.

In many implementations of MPLS Traffic Engineering, a mechanism referred to as "auto-bandwidth" is implemented. In this case, the traffic forwarded via a particular label switched path signalled by RSVP-TE is monitored, and the utilisation observed over a set period of time utilised as the bandwidth requested when a service is periodically re-optimised. Whilst a SR-based implementation cannot provide control-plane resource reservation based on this approach, through monitoring these attributes, three forms of bandwidth-aware routing can be achieved:

- o Least-Fill When selecting an particular path for a service to be routed by, where the service has affinity to an individual link within an ECMP, a typical means to ensure balancing of traffic between the different candidate links is to route the service via the link within the ECMP that is least utilised. During the computation of a Segment ERO, an ingress LSR SHOULD provide means by which such services can select the least utilised link from an set of ECMP candidate links through consideration of the Available Bandwidth sub-TLV within such IGP extensions.
- o Reaction to bandwidth utilisation within the network to re-route services based on load of links. In this case, through monitoring a set of particular (potentially high-bandwidth) services against the bandwidth utilisation of the links that they follow, it is possible to re-optimise the routing of services such that traffic is re-routed away from links experiencing congestion in a reactive manner.
- o Reaction to the bandwidth consumed per-service for instance, in cases where traffic is routed via a network with mixed maximum link bandwidths (e.g., some paths may have a maximum of 2.5Gbps where others have a maximum of 10Gbps) it is advantageous for a head-end device to split traffic flows into multiple sub-elements, with some diverging from the SPT. In this case, no knowledge of the utilisation of the network is required, however, the maximum available bandwidth of adjacencies within the SPT combined with explicit routed LSPs can be utilised to achieve traffic balance across the network.

Through utilising the uni-directional residual and available bandwidth TLVs described in the aforementioned performance attributes, the current utilisation (or available bandwidth remaining on a link) can be considered within a path calculation. An SR implementation providing performance engineered LSPs SHOULD provide means by which residual or available bandwidth can utilised as a means to calculate an ERO, and trigger subsequent re-routing. Where re-routes are triggered based on available bandwidth an iLER MUST provide means by which the time between re-optimisations can be limited, and SHOULD provide means by which such recalculations can be jittered, such that periodic re-optimisation is not performed simultaneously for all LSPs on a particular iLER.

Requirements for IGP Attributes Pertaining to Adjacency 4.2.1. Performance

The use of extended IGP attributes to determine underlying path characteristics for the selection of performance engineered paths requires some considerations to ensure that the routing information

utilised is sufficient and timely - and to balance this accuracy against resource utilisation of the systems within the IGP.

Where dynamically measured performance statistics are advertised into the IGP - such as latency measurement, or bandwidth utilisation there is a requirement to ensure that a head-end performing rerouting of an LSP calculates performance in a manner which balances:

- o Accuracy of the resource consideration at the time of routing ensuring that the current performance of the adjacency meets the path selection criteria. This requirement may lead to frequent updates of performance information into the IGP - hence, in order to minimise the impact to the overall network system, a receiving implementation in such a network SHOULD provide means by which such updates do not result in a recalculation of (the complete, or a subset of) the network's topology. In some networks, especially those with legacy systems, it is not possible to make such changes to all elements within the IGP, and therefore an advertising implementation MUST provide means by which the flooding of bandwidth information can be limited to cases where particular (operator specified) thresholds in performance are exceeded.
- o Consideration of the medium-term performance of the network link for instance, where residual bandwidth-based path selection is to be performed, it is of advantage to consider both the instantaneous bandwidth utilisation, along with the measured average over a previous time period such that the longer-term performance guarantees can be considered during route selection. Whilst such criteria does not provide strict admission control for services, it provides means by which further accuracy can be added to calculations based on instantaneous measures. To this end, an advertising implementation SHOULD provide a moving average performance measure when advertising real-time performance information within the IGP, where such attributes are not available.

In some cases, it may be advantageous for distributed path selection to consider per-forwarding class performance - in these cases an advertising implementation MAY provide performance measures on a perconfigured forwarding class basis for a particular adjacency.

<u>5</u>. Centralised Path Computation using PCE

In addition to the utilisation of distributed computation, existing PCE mechanisms can be provided to provide centralised path computation for performance engineered services. Such PCE-based computation have utility both for providing inter-area or multi-layer aware information, alongside providing globally aware service functions.

It is envisaged that the interface between the PCE and head-end LSR utilises interfaces which may:

- Exploit the Path Computation Element Protocol described in <u>RFC5440</u>.
- Utilise other real-time protocols providing interaction between forwarding elements and centralised routing entities such as those described in [I-D.amante-i2rs-topology-use-cases].

Where an implementation provides support for performance engineered LSPs it SHOULD provide means by which a remote path calculation entity can be utilised to provide both explicit route object consisting of IP addresses that can be translated into SIDs by the iLER and SHOULD support receiving a set of SID values directly from the PCE.

5.1. Use of Path Computation Element to Provide Inter-Area SR LSPs

In a multi-area network deployment, where there is restricted information propagated into stub areas, an iLER within the stub area does not have full visibility of the Adj-SIDs required to build a particular network path. Such a LER is therefore unable to determine (based on distributed computation) which SIDs should be utilised for a path to a remote node. Whilst one solution to providing such visibility is to implement a single-area IGP, or propagate all topology information to all areas, non-engineering constraints can prevent such implementations. Through utilising a PCE which has information relating to the SIDs within the network, any iLER may be provided with the relevant SIDs create a particular path through the network.

In such a deployment, the PCE element MUST have a live view of the IGP topology for all areas. This allows knowledge of the SIDs that are to be utilised to the head-end. It is envisaged that such information be provided to the PCE through interfaces such as BGP-LS as described in [I-D.ietf-idr-ls-distribution], with extensions to encode Segment Routing IGP attributes within the information propagated.

Since it is not only during signalling that visibility into the IGP topology is required, a PCE supporting such SR LSPs MUST offer functionality to inform the ingress LER supporting the SR LSP of a change to the underlying path of the LSP. For non-revertive LSPs, an iLER SHOULD offer a mechanism by which a secondary (backup) path can be requested for a service, which can be switched to based on local failure detection mechanisms (such as in-band OAM) to allow fast restoration of a service independent on interaction with the PCE at the time of failure.

5.2. Providing Co-routed or Multi-Layer Aware LSPs using PCE

5.2.1. Co-Routed LSPs

For a number of use-cases, there is a requirement for a path (be it a complete service, or a subset of traffic within a service) to be routed according to the route of another service within the network. For example:

- o Where there is a requirement diversity between a pair of services within a network - particularly where there services are instantiated on different iLER devices. In such cases, global visibility is required in order to jointly route two the services in a manner such that they are diverse (SRLG, link, and node) to one another (in addition to meeting any other performance constraints required of them).
- o Where two services form a bi-directional service within an IP/MPLS network. In this case, some services (e.g., those supporting BFD for end-to-end monitoring) may have a requirement for a return path from the eLER to the iLER which requires similar performance characteristics to the path from the iLER to eLER. Other services have tighter coupling requirements, such that the forwarding path used from iLER to eLER is symmetrical (i.e., utilises the exact same links) to the return path.

In both cases, there is a requirement for association between a set of LSPs, which span multiple head-end LERs. An implementation supporting such co-routing requirements MUST support the use of a stateful PCE, such as that described in [I-D.ietf-pce-stateful-pce] to provide calculated paths for such services.

In cases where an LSR initiates an path computation request, it MUST be possible to communicate associated paths, and relationship between the calculated path and the associated paths (in terms of co-routing or disjointness) to the PCE device. Both PCE and LSRs SHOULD provide means by which the associated paths can be specified in terms of an arbitrary service identifier (e.g., the service provider's "circuit

identifier").

Since associated LSPs may also have performance requirements, it MUST be possible for an LSR to communicate performance constraints along with associations. Where a bi-directional service is specified, path computation SHOULD support the communication of common, or differing, performance requirements for each LSP within the path.

Such PCE functions MUST provide means by which the protection requirements for a particular service can be specified - both in terms of the expected reversion behaviour for the instantiated SR-LSPs and requirements for any supporting paths (e.g., path protection).

5.2.2. Resource Reservation and Admission Control through a Stateful PCE

Implementing explicit path routing via the segment routing dataplane, rather than alternatives such as RSVP-TE, results in the inability of transit LSR devices to provide admission control (since it is unaware of the existing flows and their resource reservations). For some premium applications - such as carrying broadcast traffic the reservation of bandwidth (and its guarantee via the relevant data-plane queueing configuration) continues to be a requirement.

A stateful PCE can be utilised to perform admission control into one or more forwarding classes - allowing reservation to be achieved for these services. Where reservations are required, an iLER MUST provide means by which a bandwidth reservation in (one or more) classes can be requested of the PCE. LERs supporting such requests MUST provide means by which the resources requested for a particular SR-LSP can be statically configured by an operator, and SHOULD provide means by which dynamic observation of traffic forwarded via an LSP can be used in the subsequent requests (in order to achieve PCE-controlled auto-bandwidth functionality).

5.2.3. Multi-Layer Calculation through a Common PCE

A further case in which such PCE-based computation can be utilised is to provide correlation between layers of network infrastructure. For instance, where a common network model is available to a PCE across an optical and IP/MPLS network infrastructure, SRLG diverse paths can be provided without the requirement to encode this information (or communicate it) between the layers of the network. Through utilising a PCE able to support such multi-layer optimisations SRLG-disjoint paths can be computed and provided to iLERs for use as performance engineered paths.

Implementations providing PCE-based calculation with multi-layer awareness SHOULD provide means by which an arbitrary SRLG identifier can be provided to the PCE to allow path calculation.

<u>6</u>. Security Considerations

твс

7. Acknowledgements

The authors would like to thank Clarence Filsfils, Pierre Francois, Hannes Gredler, and Siva Sivabalan for their feedback and suggestions pertaining to this document.

<u>8</u>. Normative References

[I-D.amante-i2rs-topology-use-cases] Amante, S., Medved, J., Previdi, S., and T. Nadeau, "Topology API Use Cases", draft-amante-i2rs-topology-use-cases-00 (work in progress), February 2013. [I-D.filsfils-rtgwg-segment-routing] Filsfils, C., Previdi, S., Bashandy, A., Decraene, B., Litkowski, S., Horneffer, M., Milojevic, I., Shakir, R., Ytti, S., Henderickx, W., Tantsura, J., and E. Crabbe, "Segment Routing Architecture", draft-filsfils-rtgwg-segment-routing-00 (work in progress), June 2013. [I-D.giacalone-ospf-te-express-path] Giacalone, S., Ward, D., Drake, J., Atlas, A., and S. Previdi, "OSPF Traffic Engineering (TE) Express Path", draft-giacalone-ospf-te-express-path-02 (work in progress), September 2011. [I-D.ietf-idr-ls-distribution] Gredler, H., Medved, J., Previdi, S., Farrel, A., and S. Ray, "North-Bound Distribution of Link-State and TE Information using BGP", <u>draft-ietf-idr-ls-distribution-03</u> (work in progress), May 2013. [I-D.ietf-isis-te-metric-extensions] Previdi, S., Giacalone, S., Ward, D., Drake, J., Atlas, A., and C. Filsfils, "IS-IS Traffic Engineering (TE) Metric Extensions", draft-ietf-isis-te-metric-extensions-00 (work in progress), June 2013. [I-D.ietf-pce-stateful-pce] Crabbe, E., Medved, J., Minei, I., and R. Varga, "PCEP Extensions for Stateful PCE", draft-ietf-pce-stateful-pce-05 (work in progress), July 2013. [I-D.ietf-rtgwg-lfa-manageability] Litkowski, S., Decraene, B., Filsfils, C., and K. Raza, "Operational management of Loop Free Alternates",

[I-D.previdi-isis-segment-routing-extensions]

May 2013.

draft-ietf-rtgwg-lfa-manageability-00 (work in progress),

Previdi, S., Filsfils, C., Bashandy, A., Gredler, H., and S. Litkowski, "IS-IS Extensions for Segment Routing", draft-previdi-isis-segment-routing-extensions-02 (work in progress), July 2013.

- [I-D.psenak-ospf-segment-routing-extensions] Psenak, P., Previdi, S., Filsfils, C., Gredler, H., and R. Shakir, "OSPF Extensions for Segment Routing", draft-psenak-ospf-segment-routing-extensions-02 (work in progress), July 2013.
- [RFC2119] Bradner, S., "Key words for use in RFCs to Indicate Requirement Levels", <u>BCP 14</u>, <u>RFC 2119</u>, March 1997.
- [RFC5286] Atlas, A. and A. Zinin, "Basic Specification for IP Fast Reroute: Loop-Free Alternates", <u>RFC 5286</u>, September 2008.

Authors' Addresses

Rob Shakir BT pp. C3L, BT Centre 81, Newgate Street London EC1A 7AJ UK Email: rob.shakir@bt.com URI: http://www.bt.com/

Danny Vernals Vodafone Melbourne Street Leeds LS2 7PS UK

Email: danny.vernals@vodafone.com URI: http://www.vodafone.com/

Alessandro Capello Telecom Italia Via Reiss Romoli, 274 Torino 10148 Italy

Email: alessandro.capello@telecomitalia.it URI: http://www.telecomitalia.com/