

Network Working Group
Internet-Draft
Expires: July 22, 2005

A. Atlas, Ed.
Avici Systems, Inc.
January 21, 2005

Basic Specification for IP Fast-Reroute: Loop-free Alternates
draft-ietf-rtgwg-ipfrr-spec-base-02

Status of this Memo

This document is an Internet-Draft and is subject to all provisions of [section 3 of RFC 3667](#). By submitting this Internet-Draft, each author represents that any applicable patent or other IPR claims of which he or she is aware have been or will be disclosed, and any of which he or she become aware will be disclosed, in accordance with [RFC 3668](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/1id-abstracts.txt>.

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

This Internet-Draft will expire on July 22, 2005.

Copyright Notice

Copyright (C) The Internet Society (2005).

Abstract

This document describes the use of loop-free alternates to provide local protection for IP unicast and/or LDP traffic in the event of a single failure, whether link, node or shared risk link group (SRLG). The goal of this technology is to reduce the micro-looping that and packet loss that happens while routers converge after a topology change due to a failure. When a topology change occurs, a router S determines for each prefix an alternate next-hop which can be used if the primary next-hop fails. An acceptable alternate next-hop must be

a loop-free alternate, which goes to a neighbor whose shortest path to the prefix does not go back through the router S.

Table of Contents

1.	Introduction	3
1.1	Failure Scenarios	4
2.	Alternate Next-Hop Calculation	6
2.1	Basic Loop-free Condition	7
2.2	Node-Protecting Alternate Next-Hops	7
2.3	Broadcast and NBMA Links	7
2.4	Interactions with ISIS Overload, RFC 3137 and Costed Out Links	8
2.5	Selection Procedure	9
3.	Using an Alternate	10
3.1	Terminating Use of Alternate	10
4.	Requirements on LDP Mode	12
5.	Routing Aspects	12
5.1	Multi-Homed Prefixes	12
5.2	OSPF External Routing	13
5.3	OSPF Virtual Links	14
5.4	BGP Next-Hop Synchronization	14
5.5	Multicast Considerations	14
6.	Security Considerations	14
7.	References	15
	Authors' Addresses	15
	Intellectual Property and Copyright Statements	17

1. Introduction

Applications for interactive multimedia services such as VoIP and pseudo-wires can be very sensitive to traffic loss, such as occurs when a link or router in the network fails. A router's convergence time is generally on the order of seconds; the application traffic may be sensitive to losses greater than 10s of milliseconds.

As discussed in [[FRAMEWORK](#)], minimizing traffic loss requires a mechanism for the router adjacent to a failure to rapidly invoke a repair path, which is minimally affected by any subsequent re-convergence. This specification describes such a mechanism which allows a router whose local link has failed to forward traffic to a pre-computed alternate until the router installs the new primary next-hops based upon the changed network topology. The terminology used in this specification is given in [[FRAMEWORK](#)].

When a local link fails, a router currently must signal the event to its neighbors via the IGP, recompute new primary next-hops for all affected prefixes, and only then install those new primary next-hops into the forwarding plane. Until the new primary next-hops are installed, traffic directed towards the affected prefixes is discarded. This process can take seconds.

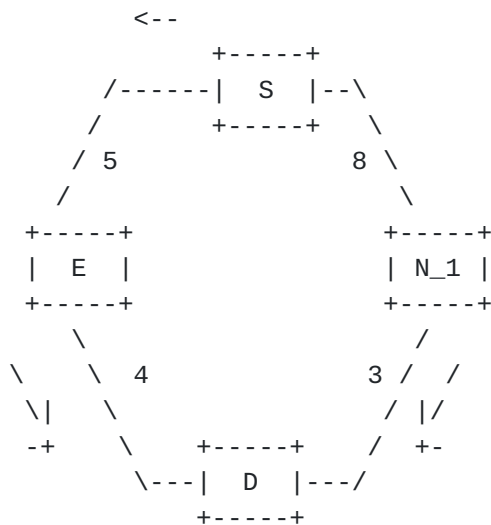


Figure 1: Basic Topology

The goal of IP Fast-Reroute is to reduce that traffic convergence time to 10s of milliseconds by using a pre-computed alternate next-hop, in the event that the currently selected primary next-hop fails, so that the alternate can be rapidly used when the failure is detected. A network with this feature experiences less traffic loss and less micro-looping of packets than a network without IPFRR.

There are cases where micro-looping is still a possibility since IPFRR coverage varies but in the worst possible situation a network with IPFRR is equivalent with respect traffic convergence to a network without IPFRR.

To clarify the behavior of IP Fast-Reroute, consider the simple topology in Figure 1. When router S computes its shortest path to router D, router S determines to use the link to router E as its primary next-hop. Without IP Fast-Reroute, that link is the only next-hop that router S computes to reach D. With IP Fast-Reroute, S also looks for an alternate next-hop to use. In this example, S would determine that it could send traffic destined to D by using the link to router N_1 and therefore S would install the link to N_1 as its alternate next-hop. At some later time, the link between router S and router E could fail. When that link fails, S and E will be the first to detect it. On detecting the failure, S will stop sending traffic destined for D towards E via the failed link, and instead send the traffic to S's pre-computed alternate next-hop, which is the link to N_1, until a new SPF is run and its results are installed. As with the primary next-hop, an alternate next-hop is computed for each destination. The process of computing an alternate next-hop does not alter the primary next-hop computed via a standard SPF.

If in the example of Figure 1, the link cost from N_1 to D increased to 30 from 3, then N_1 would not be a loop-free alternate, because the cost of the path from N_1 to D via S would be 17 while the cost from N_1 directly to D would be 30. In real networks, we may often face this situation. The existence of a suitable loop-free alternate next-hop is topology dependent.

A neighbor N can provide a loop-free alternate if and only if

$$\text{Distance_opt}(N, D) < \text{Distance_opt}(N, S) + \text{Distance_opt}(S, D)$$

Equation 1: Loop-Free Criterion

A sub-set of loop-free alternate are downstream paths which must meet the more restrictive condition of

$$\text{Distance_opt}(N, D) < \text{Distance_opt}(S, D)$$

Equation 2: Downstream Path Criterion

1.1 Failure Scenarios

The alternate next-hop can protect against a single link failure, a single node failure, one or more shared risk link group failure, or a

combination of these. Whenever a failure occurs that is more extensive than what the alternate was intended to protect, there is the possibility of looping traffic. The example where a node fails when the alternate provided only link protection is illustrated below. If unexpected simultaneous failures occur, then micro-looping may occur since the alternates are not pre-computed to avoid the set of failed links.

If only link protection is provided and the node fails, it is possible for traffic using the alternates to experience micro-looping. This issue is illustrated in Figure 2. If Link(S->E) fails, then the link-protecting alternate via N will work correctly. However, if router E fails, then both S and N will detect a failure and switch to their alternates. In this example, that would cause S to redirect the traffic to N and N to redirect the traffic to S and thus causing a forwarding loop. Such a scenario can arise because the key assumption, that all other routers in the network are forwarding based upon the shortest path, is violated because of a second simultaneous correlated failure - another link connected to the same primary neighbor. If there are not other protection mechanisms a node failure is still a concern when only using link protection.

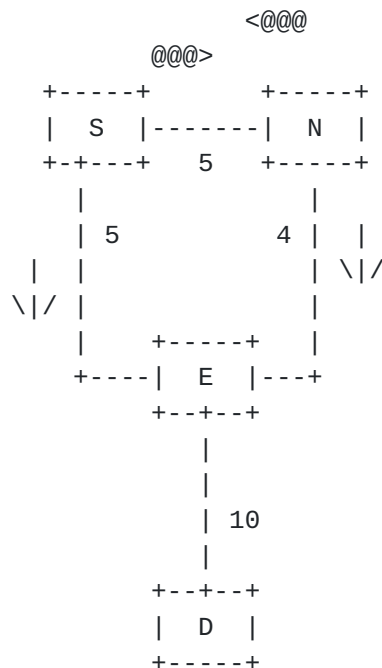


Figure 2: Link-Protecting Alternates Causing Loop on Node Failure

Micro-looping of traffic via the alternates caused when a more extensive failure than planned for can be prevented via selection of

only downstream paths as alternates. In Figure 2, S would be able to use N as an alternate, but N could not use S; therefore N would have no alternate and would discard the traffic, thus avoiding the micro-loop. A micro-loop due to the use of alternates can be avoided by using downstream paths because each router in the path to the destination must be closer to the destination (according to the topology prior to the failures). Although use of downstream paths ensures that the micro-looping via alternates does not occur, such a restriction can severely limit the coverage of alternates.

It may be desirable to find an alternate that can protect against other correlated failures (of which node failure is a specific instance). In the general case, these are handled by shared risk link groups (SRLGs) where any links in the network can belong to the SRLG. General SRLGs may add unacceptably to the computational complexity of finding a loop-free alternate.

However, a sub-category of SRLGs is of interest and can be applied only during the selection of an acceptable alternate. This sub-category is to express correlated failures of links that are connected to the same router. For example, if there are multiple logical sub-interfaces on the same physical interface, such as VLANs on an Ethernet interface, if multiple interfaces use the same physical port because of channelization, or if multiple interfaces share a correlated failure because they are on the same line-card. This sub-category of SRLGs will be referred to as local-SRLGs. A local-SRLG has all of its member links with one end connected to the same router. Thus, router S could select a loop-free alternate which does not use a link in the same local-SRLG as the primary next-hop. The local-SRLGs belonging to E can be protected against via node-protection; i.e. picking a loop-free node-protecting alternate.

2. Alternate Next-Hop Calculation

To support IP Fast-Reroute, a router must be able to determine if a next-hop will provide a loop-free alternate before the router installs that next-hop as an alternate. That next-hop must go to a loop-free neighbor.

To do this computation, a router could run an SPF from the perspective of each of its neighbors as well as from its own perspective. This provides the router with all the information necessary to test the equations given in this specification.

To determine SRLG protection, the set of SRLGs that include at least one link from the computing router could be determined. Then when the SPF is run from the perspective of a router's neighbor, the SRLGs traversed on each shortest path can be tracked.

2.1 Basic Loop-free Condition

Alternate next hops used by implementations following this specification MUST conform to at least the loop-freeness condition stated above in Equation 1. Further conditions may be applied when determining link-protecting and/or node-protecting alternate next-hops as described in Sections [Section 2.2](#) and [Section 2.3](#).

2.2 Node-Protecting Alternate Next-Hops

For an alternate next-hop to protect against node failure, the alternate next-hop MUST be loop-free with respect to the primary neighbor E and the destination.

An alternate will be node-protecting if it doesn't go through the same primary neighbor as the primary next-hop. This is the case if Equation 3 is true, where N is the neighbor providing a loop-free alternate.

$$\text{Distance_opt}(N, D) < \text{Distance_opt}(N, E) + \text{Distance_opt}(E, D)$$

Equation 3: Criteria for a Node-Protecting Loop-Free Alternate

If $\text{Distance_opt}(N, D) = \text{Distance_opt}(N, E) + \text{Distance_opt}(E, D)$, it is possible that the neighbor may have equal-cost paths and one of those could provide a loop-free node-protecting alternate. However, the decision as to which of equal-cost paths a router will use is a router-local decision. Therefore, a router MUST assume that an alternate next-hop does not offer node protection if Equation 3 is not met.

2.3 Broadcast and NBMA Links

The computation for link-protection is a bit more complicated for broadcast links. In an SPF computation, a broadcast link is represented as a pseudo-node with links of 0 cost exiting the pseudo-node. For an alternate to be considered link-protecting, it must be loop-free with regard to the pseudo-node. Consider the example in Figure 3.

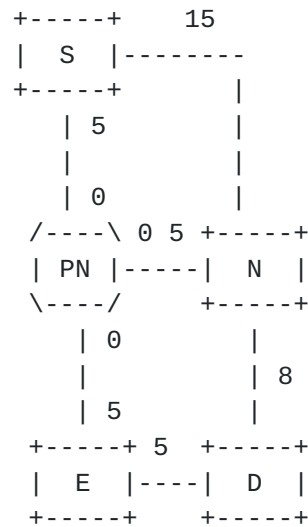


Figure 3: Loop-Free Alternate that is Link-Protecting

In Figure 3, N offers a loop-free alternate which is link-protecting. If the primary next-hop uses a broadcast link, then an alternate must be loop-free with respect to that link's pseudo-node to provide link protection. This requirement is described in Equation 4 below.

$$D_{\text{opt}}(N, D) < D_{\text{opt}}(N, \text{pseudo}) + D_{\text{opt}}(\text{pseudo}, D)$$

Equation 4: Loop-Free Link-Protecting Criterion for Broadcast Links

Because the shortest path from the pseudo-node goes through E, if a loop-free alternate from a neighbor N is node-protecting, the alternate will also be link-protecting unless the router S can only reach the neighbor N via the same pseudo-node. This can occur because S will direct traffic away from the shortest path to use an alternate. Therefore link protection must be considered during the alternate selection.

2.4 Interactions with ISIS Overload, [RFC 3137](#) and Costed Out Links

As described in [[RFC3137](#)], there are cases where it is desirable not to have a router used as a transit node. For those cases, it is also desirable not to have the router used on an alternate path.

For computing an alternate, a router MUST not consider diverting from the SPF tree along a link whose cost or reverse cost is LSInfinity (for OSPF) or the maximum cost (for ISIS) or whose next-hop router has the overload bit set (for ISIS).

In the case of OSPF, if all links from router S to a neighbor N_i have a reverse cost of LSInfinity, then router S MUST NOT consider

using N_i as an alternate.

Similarly in the case of ISIS, if N_i has the overload bit set, then S MUST NOT consider using N_i as an alternate.

This preserves the desired behavior of diverting traffic away from a router which is following [[RFC3137](#)] and it also preserves the desired behavior when an operator sets the cost of a link to LSInfinity for maintenance which is not permitting traffic across that link unless there is no other path.

If a link or router which is costed out was the only possible alternate to protect traffic from a particular router S to a particular destination, then there will be no alternate provided for protection.

[2.5](#) Selection Procedure

A router supporting this specification SHOULD select a loop-free alternate next-hop for each primary next-hop used for a given prefix. A router MAY decide to not use an available loop-free alternate next-hop. A reason for such a decision might be that the loop-free alternate next-hop does not provide protection for the failure scenario of interest.

The alternate selection should maximize the coverage of the failure cases.

S SHOULD select a loop-free node-protecting alternate next-hop, if one is available. If S has a choice between a loop-free link-protecting node-protecting alternate and a loop-free node-protecting alternate which is not link-protecting, S SHOULD select a loop-free node-protecting alternate which is also link-protecting. This can occur as explained in [Section 2.3](#). If S has multiple primary next-hops, then S SHOULD select as a loop-free alternate either one of the other primary next-hops or a loop-free node-protecting alternate. If no loop-free node-protecting alternate is available, then S MAY select a loop-free link-protecting alternate.

Each next-hop can be categorized as to the type of alternate it can provide to a particular destination D from router S for a particular primary next-hop which goes to a neighbor E. A next-hop may provide one of the following types of paths:

Primary Path - This is the primary next-hop.

Loop-Free Node-Protecting Alternate - This next-hop satisfies Equation 1 and Equation 3. The path avoids S, S's primary neighbor E, and the link from S to E.

Loop-Free Link-Protecting Alternate - This next-hop satisfies Equation 1 but not Equation 3. If the primary next-hop uses a broadcast link, then this next-hop satisfies Equation 4.

Unavailable - This may be because the path goes through S to reach D, because the link is costed out, etc.

An alternate path may also provide none, some or complete SRLG protection as well as node and link or link protection. For instance, a link may belong to two SRLGs G1 and G2. The alternate path might avoid other links in G1 but not G2, in which case the alternate would only provide partial SRLG protection.

3. Using an Alternate

If an alternate next-hop is available, the router SHOULD redirect traffic to the alternate next-hop when the primary next-hop has failed.

When a local interface failure is detected, traffic that was destined to go out the failed interface must be redirected to the appropriate alternate next-hops. Other failure detection mechanisms which detect the loss of a link or a node may also be used to trigger redirection of traffic to the appropriate alternate next-hops. The mechanisms available for failure detection are discussed in [[FRAMEWORK](#)] and are outside the scope of this specification.

The alternate next-hop MUST be used only for traffic types which are routed according to the shortest path. Multicast traffic is specifically out of scope for this specification.

3.1 Terminating Use of Alternate

A router MUST limit the amount of time an alternate next-hop is used after the primary next-hop has become unavailable. This ensures that the router will start using the new primary next-hops. It ensures that all possible transient conditions are removed and the network converges according to the deployed routing protocol.

It is desirable to avoid micro-forwarding loops involving S. An example illustrating the problem is given in Figure 4. If the link from S to E fails, S will use N1 as an alternate and S will compute

N2 as the new primary next-hop to reach D. If S starts using N2 as soon as S can compute and install its new primary, it is probable that N2 will not have yet installed its new primary next-hop. This would cause traffic to loop and be dropped until N2 has installed the new topology. This can be avoided by S delaying its installation and leaving traffic on the alternate next-hop.

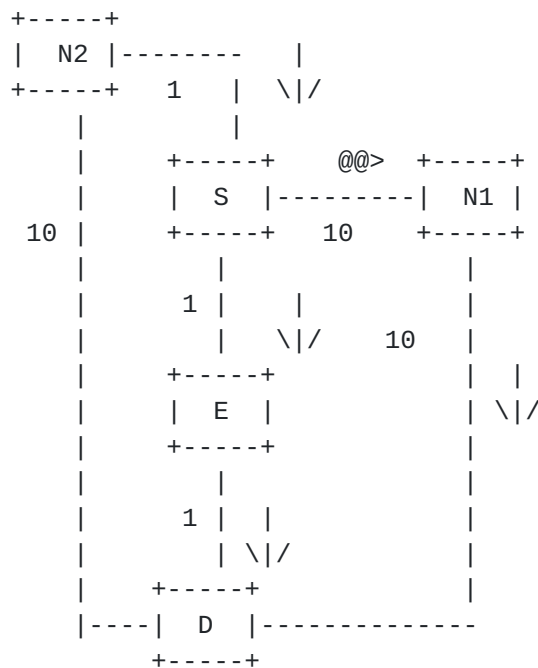


Figure 4: Example where Continued Use of Alternate is Desirable

This is an example of a case where the new primary is not a loop-free alternate before the failure and therefore may have been forwarding traffic through S. This will occur when the path via a previously upstream node is shorter than the the path via a loop-free alternate neighbor. In these cases, it is useful to give sufficient time to ensure that the new primary neighbor and other nodes on the new primary path have switched to the new route.

If the newly selected primary was loop-free before the failure, then it is safe to switch to that new primary immediately; the new primary wasn't dependent on the failure and therefore its path will not have changed.

Given that there is an alternate providing appropriate protection and while the assumption of a single failure holds, it is safe to delay the installation of the new primaries; this will not create forwarding loops because the alternate's path to the destination is known to not go via S or the failed element and will therefore not be affected by the failure.

An implementation SHOULD continue to use the alternate next-hops for packet forwarding even after the new routing information is available based on the new network topology. The use of the alternate next-hops for packet forwarding SHOULD terminate:

- a. if the new primary next-hop was loop-free prior to the topology change, or
- b. if a configured hold-down, which represents a worst-case bound on the length of the network convergence transition, has expired, or
- c. if notification of an unrelated topological change in the network is received.

4. Requirements on LDP Mode

Since LDP traffic will follow the path specified by the IGP, it is also possible for the LDP traffic to follow the loop-free alternates indicated by the IGP. To do so, it is necessary for LDP to have the appropriate labels available for the alternate so that the appropriate out-segments can be installed in the forwarding plane before the failure occurs.

This means that a Label Switched Router (LSR) running LDP must distribute its labels for the FECs it can provide to all its neighbors, regardless of whether or not they are upstream. Additionally, LDP must be acting in liberal label retention mode so that the labels which correspond to neighbors that aren't currently the primary neighbor are stored. Similarly, LDP should be in downstream unsolicited mode, so that the labels for the FEC are distributed other than along the SPT.

If these requirements are met, then LDP can use the loop-free alternates without requiring any targeted sessions or signaling extensions for this purpose.

5. Routing Aspects

5.1 Multi-Homed Prefixes

An SPF-like computation is run for each topology, which corresponds to a particular OSPF area or ISIS level. The IGP needs to determine loop-free alternates to multi-homed routes. Multi-homed routes occur for routes obtained from outside the routing domain by multiple routers, for subnets on links where the subnet is announced from multiple ends of the link, and for routes advertised by multiple routers to provide resiliency.

Figure 5 demonstrates such a topology. In this example, the shortest path to reach the prefix p is via E. The prefix p will have the link to E as its primary next-hop. If the alternate next-hop for the prefix p is simply inherited from the router advertising it on the shortest path to p, then the prefix p's alternate next-hop would be the link to C. This would provide link protection, but not the node protection that is possible via A.

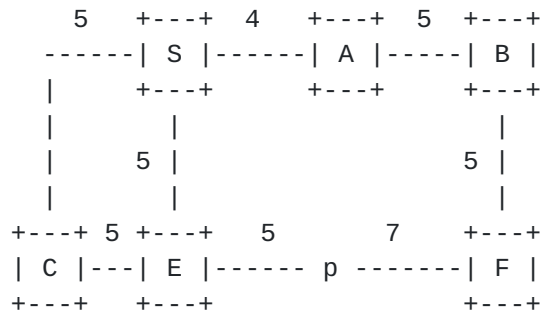


Figure 5: Multi-homed prefix

To determine the best protection possible, the prefix p can be treated in the SPF computations as a node with uni-directional links to it from those routers that have advertised the prefix. Such a node need never have its links explored, as it has no out-going links.

If there exist multiple multi-homed prefixes exist that share the same connectivity and the difference in metrics to those routers, then a single node can be used to represent the set. For instance, if in Figure 5 there were another prefix X that was connected to E with a metric of 1 and to F with a metric of 3, then that prefix X could use the same alternate next-hop as was computed for prefix p.

A router SHOULD compute the alternate next-hop for an IGP multi-homed prefix by considering alternate paths via all routers that have announced that prefix.

5.2 OSPF External Routing

An additional complication comes from forwarding addresses, where an ASBR uses a forwarding address to indicate to all routers in the Autonomous System to use the specified address instead of going through the ASBR. When a forwarding address has been indicated, all routers in the topology calculate the shortest path to the link specified in the external LSA. In this case, the alternate next-hop should be computed by selecting among the alternate paths to the forwarding link(s) instead of among alternate paths to the ASBR.

5.3 OSPF Virtual Links

OSPF virtual links are used to connect two disjoint backbone areas using a transit area. A virtual link is configured at the border routers of the disjoint area. If router S is itself an ABR or one of the endpoints of the disjoint area, then router S must resolve its paths to the destination on the other side of the disjoint area by using the summary links in the transit area and using the closest ABR summarizing them into the transit area. This means that the data path may diverge from the virtual neighbor's control path. An ABR's primary and alternate next-hops are calculated by S on the transit area.

A virtual link MUST NOT be used as an alternate next-hop.

5.4 BGP Next-Hop Synchronization

Typically BGP prefixes are advertised with AS exit routers router-id, and AS exit routers are reached by means of IGP routes. BGP resolves its advertised next-hop to the immediate next-hop by potential recursive lookups in the routing database. IP Fast-Reroute computes the alternate next-hops to all IGP destinations, which include alternate next-hops to the AS exit router's router-id. BGP simply inherits the alternate next-hop from IGP. The BGP decision process is unaltered; BGP continues to use the IGP optimal distance to find the nearest exit router. MBGP routes do not need to copy the alternate next hops.

It is possible to provide ASBR protection if BGP selected a set of IGP next-hops and allowed the IGP to determine the primary and alternate next-hops as if the BGP route were a multi-homed prefix. This is for future study.

5.5 Multicast Considerations

Multicast traffic is out of scope for this specification of IP Fast-Reroute. The alternate next-hops SHOULD not be used for multi-cast RPF checks.

6. Security Considerations

This document does not introduce any new security issues. The mechanisms described in this document depend upon the network topology distributed via an IGP, such as OSPF or ISIS. It is dependent upon the security associated with those protocols.

7 References

- [FRAMEWORK] Shand, M., "IP Fast Reroute Framework",
 [draft-ietf-rtgwg-ipfrr-framework-02.txt](#) (work in
 progress), October 2004.
- [RFC3036] Andersson, L., Doolan, P., Feldman, N., Fredette, A. and
 B. Thomas, "LDP Specification", [RFC 3036](#), January 2001.
- [RFC3137] Retana, A., Nguyen, L., White, R., Zinin, A. and D.
 McPherson, "OSPF Stub Router Advertisement", [RFC 3137](#),
 June 2001.
- [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.

Authors' Addresses

Alia K. Atlas (editor)
Avici Systems, Inc.
101 Billerica Avenue
N. Billerica, MA 01862
USA

Phone: +1 978 964 2070
EMail: aatlas@avici.com

Raveendra Torvi
Avici Systems, Inc.
101 Billerica Avenue
N. Billerica, MA 01862
USA

Phone: +1 978 964 2026
EMail: rtorvi@avici.com

Gagan Choudhury
AT&T
200 Laurel Avenue, Room D5-3C21
Middletown, NJ 07748
USA

Phone: +1 732 420-3721
EMail: gchoudhury@att.com

Christian Martin
Verizon
1880 Campus Commons Drive
Reston, VA 20191
USA

Brent Imhoff
LightCore
14567 North Outer Forty Rd.
Chesterfield, MO 63017
USA

Phone: +1 314 880 1851
EMail: brent@lightcore.net

Don Fedyk
Nortel Networks
600 Technology Park
Billerica, MA 01821
USA

Phone: +1 978 288 3041
EMail: dwfedyk@nortelnetworks.com

Intellectual Property Statement

The IETF takes no position regarding the validity or scope of any Intellectual Property Rights or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; nor does it represent that it has made any independent effort to identify any such rights. Information on the procedures with respect to rights in RFC documents can be found in [BCP 78](#) and [BCP 79](#).

Copies of IPR disclosures made to the IETF Secretariat and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the IETF on-line IPR repository at <http://www.ietf.org/ipr>.

The IETF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights that may cover technology that may be required to implement this standard. Please address the information to the IETF at ietf-ipr@ietf.org.

Disclaimer of Validity

This document and the information contained herein are provided on an "AS IS" basis and THE CONTRIBUTOR, THE ORGANIZATION HE/SHE REPRESENTS OR IS SPONSORED BY (IF ANY), THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIM ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

Copyright Statement

Copyright (C) The Internet Society (2005). This document is subject to the rights, licenses and restrictions contained in [BCP 78](#), and except as set forth therein, the authors retain all their rights.

Acknowledgment

Funding for the RFC Editor function is currently provided by the Internet Society.

