

Network Working Group
Request for Comments: 3623
Category: Standards Track

J. Moy
Sycamore Networks
P. Pillay-Esnault
Juniper Networks
A. Lindem
Redback Networks
November 2003

Graceful OSPF Restart

Status of this Memo

This document specifies an Internet standards track protocol for the Internet community, and requests discussion and suggestions for improvements. Please refer to the current edition of the "Internet Official Protocol Standards" (STD 1) for the standardization state and status of this protocol. Distribution of this memo is unlimited.

Copyright Notice

Copyright (C) The Internet Society (2003). All Rights Reserved.

Abstract

This memo documents an enhancement to the OSPF routing protocol, whereby an OSPF router can stay on the forwarding path even as its OSPF software is restarted. This is called "graceful restart" or "non-stop forwarding". A restarting router may not be capable of adjusting its forwarding in a timely manner when the network topology changes. In order to avoid the possible resulting routing loops, the procedure in this memo automatically reverts to a normal OSPF restart when such a topology change is detected, or when one or more of the restarting router's neighbors do not support the enhancements in this memo. Proper network operation during a graceful restart makes assumptions upon the operating environment of the restarting router; these assumptions are also documented.

Table of Contents

1. Overview	2
2. Operation of Restarting Router	3
2.1. Entering Graceful Restart.	4
2.2. When to Exit Graceful Restart.	5
2.3. Actions on Exiting Graceful Restart.	6
3. Operation of Helper Neighbor	7
3.1. Entering Helper Mode	7
3.2. Exiting Helper Mode.	8
4. Backward Compatibility	9
5. Unplanned Outages.	10
6. Interaction with Traffic Engineering	11
7. Possible Future Work	11
8. Intellectual Property Rights Notice.	11
9. References	11
9.1. Normative References	11
9.2. Informative References	11
A. Grace-LSA Format	13
B. Configurable Parameters.	15
Security Considerations.	16
Acknowledgements	16
Authors' Addresses	17
Full Copyright Statement	18

1. Overview

Today many Internet routers implement a separation of control and forwarding functions. Certain processors are dedicated to control and management tasks such as OSPF routing, while other processors perform the data forwarding tasks. This separation creates the possibility of maintaining a router's data forwarding capability while the router's control software is restarted/reloaded. We call such a possibility "graceful restart" or "non-stop forwarding".

The OSPF protocol presents a problem to graceful restart whereby, under normal operation, OSPF intentionally routes around a restarting router while it rebuilds its link-state database. OSPF avoids the restarting router to minimize the possibility of routing loops and/or black holes caused by lack of database synchronization. Avoidance is accomplished by having the router's neighbors reissue their LSAs, omitting links to the restarting router.

However, if (a) the network topology remains stable and (b) the restarting router is able to keep its forwarding table(s) across the restart, it would be safe to keep the restarting router on the forwarding path. This memo documents an enhancement to OSPF that makes such graceful restart possible, and automatically reverts back

to a standard OSPF restart for safety when network topology changes are detected.

In a nutshell, the OSPF enhancements for graceful restart are as follows:

- The router attempting a graceful restart originates link-local Opaque-LSAs, herein called Grace-LSAs, announcing its intention to perform a graceful restart within a specified amount of time or "grace period".
- During the grace period, its neighbors continue to announce the restarting router in their LSAs as if it were fully adjacent (i.e., OSPF neighbor state Full), but only if the network topology remains static (i.e., the contents of the LSAs in the link-state database having LS types 1-5,7 remain unchanged and periodic refreshes are allowed).

There are two roles being played by OSPF routers during graceful restart. First there is the router that is being restarted. The operation of this router during graceful restart, including how the router enters and exits graceful restart, is the subject of Section 2. Then there are the router's neighbors, which must cooperate in order for the restart to be graceful. During graceful restart, we say that the neighbors are running in "helper mode". Section 3 covers the responsibilities of a router running in helper mode, including entering and exiting helper mode.

2. Operation of Restarting Router

After the router restarts/reloads, it must change its OSPF processing somewhat until it re-establishes full adjacencies with all its former fully-adjacent neighbors. This time period, between the restart/reload and the reestablishment of adjacencies, is called "graceful restart". During graceful restart:

- 1) The restarting router does not originate LSAs with LS types 1-5,7. Instead, the restarting router wants the other routers in the OSPF domain to calculate routes using the LSAs that it originated prior to its restart. During this time, the restarting router does not modify or flush received self-originated LSAs, (see Section 13.4 of [1]). Instead they are accepted as valid. In particular, the grace-LSAs that the restarting router originated before the restart are left in place. Received self-originated LSAs will be dealt with when the router exits graceful restart (see Section 2.3).

- 2) The restarting router runs its OSPF routing calculations, as specified in Section 16 of [1]. This is necessary to return any OSPF virtual links to operation. However, the restarting router does **not** install OSPF routes into the system's forwarding table(s) and relies on the forwarding entries that it installed prior to the restart.
- 3) If the restarting router determines that it was the Designated Router on a given segment prior to the restart, it elects itself as the Designated Router again. The restarting router knows that it was the Designated Router if, while the associated interface is in Waiting state, a Hello packet is received from a neighbor listing the router as the Designated Router.

Otherwise, the restarting router operates the same as any other OSPF router. It discovers neighbors using OSPF's Hello protocol, elects Designated and Backup Designated Routers, performs the Database Exchange procedure to initially synchronize link-state databases with its neighbors, and maintains this synchronization through flooding.

The processes of entering graceful restart, and of exiting graceful restart (either successfully or not) are covered in the following sections.

2.1. Entering Graceful Restart

The router (call it Router X) is informed of the desire for its graceful restart when an appropriate command is issued by the network operator. The network operator may also specify the length of the grace period, or the necessary grace period may be calculated by the router's OSPF software. In order to avoid the restarting router's LSAs from aging out, the grace period should not exceed LSRefreshTime (1800 second) [1].

In preparation for the graceful restart, Router X must perform the following actions before its software is restarted/reloaded:

(Note that common OSPF shutdown procedures are **not** performed, since we want the other OSPF routers to act as if Router X remains in continuous service. For example, Router X does not flush its locally originated LSAs, since we want them to remain in other routers' link-state databases throughout the restart period.)

- 1) Router X must ensure that its forwarding table(s) is/are up-to-date and will remain in place across the restart.

- 2) The router may need to preserve the cryptographic sequence numbers being used on each interface in non-volatile storage. An alternative is to use the router's clock for cryptographic sequence number generation and ensure that the clock is preserved across restarts (either on the same or redundant route processors). If neither of these can be guaranteed, it can take up to RouterDeadInterval seconds after the restart before adjacencies can be reestablished and this would force the grace period to be lengthened greatly.

Router X then originates the grace-LSAs. These are link-local Opaque-LSAs (see Appendix A). Their LS Age field is set to 0, and the requested grace period (in seconds) is inserted into the body of the grace-LSA. The precise contents of the grace-LSA are described in Appendix A.

A grace-LSA is originated for each of the router's OSPF interfaces. If Router X wants to ensure that its neighbors receive the grace-LSAs, it should retransmit the grace-LSAs until they are acknowledged (i.e., perform standard OSPF reliable flooding of the grace-LSAs). If one or more fully adjacent neighbors do not receive grace-LSAs, they will more than likely cause premature termination of the graceful restart procedure (see Section 4).

After the grace-LSAs have been sent, the router should store the fact that it is performing graceful restart along with the length of the requested grace period in non-volatile storage. (Note to implementors: It may be easiest to simply store the absolute time of the end of the grace period). The OSPF software should then be restarted/reloaded. When the reloaded software starts executing the graceful restart, the protocol modifications in Section 2 are followed. (Note that prior to the restart, the router does not know whether its neighbors are going to cooperate as "helpers"; the mere reception of grace-LSAs does not imply acceptance of helper responsibilities. This memo assumes that the router would want to restart anyway, even if the restart is not going to be graceful).

2.2. When to Exit Graceful Restart

A Router X exits graceful restart when any of the following occurs:

- 1) Router X has reestablished all its adjacencies. Router X can determine this by examining the router-LSAs that it last originated before the restart (called the "pre-restart router-LSA"), and, on those segments where the router is the Designated Router, the pre-restart network-LSAs. These LSAs will have been received from the helping neighbors, and need not have been stored in non-volatile storage across the

restart. All previous adjacencies will be listed as type-1 and type-2 links in the router-LSA, and as neighbors in the body of the network-LSA.

- 2) Router X receives an LSA that is inconsistent with its pre-restart router-LSA. For example, X receives a router-LSA originated by router Y that does not contain a link to X, even though X's pre-start router-LSA did contain a link to Y. This indicates that either a) Y does not support graceful restart, b) Y never received the grace-LSA or c) Y has terminated its helper mode for some reason (Section 3.2). A special case of LSA inconsistency is when Router X establishes an adjacency with router Y and doesn't receive an instance of its own pre-restart router LSA.

- 3) The grace period expires.

2.3. Actions on Exiting Graceful Restart

Upon exiting "graceful restart", the restarting router reverts back to completely normal OSPF operation, reoriginating LSAs based on the router's current state and updating its forwarding table(s) based on the current contents of the link-state database. In particular, the following actions should be performed when exiting, either successfully or unsuccessfully, graceful restart:

- 1) The router should reoriginate its router-LSAs for all attached areas in order to make sure they have the correct contents.
- 2) The router should reoriginate network-LSAs on all segments where it is the Designated Router.
- 3) The router reruns its OSPF routing calculations (Section 16 of [1]), this time installing the results into the system forwarding table, and originating summary-LSAs, Type-7 LSAs and AS-external-LSAs as necessary.
- 4) Any remnant entries in the system forwarding table that were installed before the restart, but that are no longer valid, should be removed.
- 5) Any received self-originated LSAs that are no longer valid should be flushed.
- 6) Any grace-LSAs that the router originated should be flushed.

3. Operation of Helper Neighbor

The helper relationship is per network segment. As a "helper neighbor" on a segment S for a restarting router X, router Y has several duties. It monitors the network for topology changes, and as long as there are none, continues to advertise its LSAs as if X had remained in continuous OSPF operation. This means that Y's LSAs continue to list an adjacency to X over network segment S, regardless of the adjacency's current synchronization state. This logic affects the contents of both router-LSAs and network-LSAs, and also depends on the type of network segment S (see Sections 12.4.1.1 through 12.4.1.5 and Section 12.4.2 of [1]). When helping over a virtual link, the helper must also continue to set bit V in its router-LSA for the virtual link's transit area (Section 12.4.1 of [1]).

Also, if X was the Designated Router on network segment S when the helping relationship began, Y maintains X as the Designated Router until the helping relationship is terminated.

3.1. Entering Helper Mode

When a router Y receives a grace-LSA from router X, it enters helper mode for X on the associated network segment, as long as all the following checks pass:

- 1) Y currently has a full adjacency with X (neighbor state Full) over the associated network segment. On broadcast, NBMA and Point-to-MultiPoint segments, the neighbor relationship with X is identified by the IP interface address in the body of the grace-LSA (see Appendix A). On all other segment types, X is identified by the grace-LSA's Advertising Router field.
- 2) There have been no changes in content to the link-state database (LS types 1-5,7) since router X restarted. This is determined as follows:
 - Router Y examines the link-state retransmission list for X over the associated network segment.
 - If there are any LSAs with LS types 1-5,7 on the list, then they all must be periodic refreshes.
 - If there are instead LSAs on the list whose contents have changed (see Section 3.3 of [7]), Y must refuse to enter helper mode.

Router Y may optionally disallow graceful restart with Router X on other network segments. Determining whether

changed LSAs have been successfully flooded to router Y on other network segments is feasible but beyond the scope of this document.

- 3) The grace period has not yet expired. This means that the LS age of the grace-LSA is less than the grace period specified in the body of the grace-LSA (Appendix A).
- 4) Local policy allows Y to act as the helper for X. Examples of configured policies might be a) never act as helper, b) never allow the grace period to exceed a Time T, c) only help on software reloads/upgrades, or d) never act as a helper for specific routers (specified by OSPF Router ID).
- 5) Router Y is not in the process of graceful restart.

There is one exception to the above requirements. If Y was already helping X on the associated network segment, the new grace-LSA should be accepted and the grace period should be updated accordingly.

Note that Router Y may be helping X on some network segments, and not on others. However, that circumstance will probably lead to the premature termination of X's graceful restart, as Y will not continue to advertise adjacencies on the segments where it is not helping (see Section 2.2).

Alternately, Router Y may choose to enter helper mode when a grace-LSA is received and the above checks pass for all adjacencies with Router X. This implementation alternative of aggregating the adjacencies with respect to helper mode is compatible with implementations considering each adjacency independently.

A single router is allowed to simultaneously serve as a helper for multiple restarting neighbors.

3.2. Exiting Helper Mode

Router Y ceases to perform the helper function for its neighbor Router X on a given segment when one of the following events occurs:

- 1) The grace-LSA originated by X on the segment is flushed. This indicates the successful termination of graceful restart.
- 2) The grace-LSA's grace period expires.
- 3) A change in link-state database contents indicates a network topology change, which forces termination of a graceful restart. Specifically, if router Y installs a new LSA in its

database with LS types 1-5,7 and having the following two properties, it should cease helping X. The two properties of the LSA are:

- a) the contents of the LSA have changed; this includes LSAs with no previous link-state database instance and the flushing of LSAs from the database, but excludes periodic LSA refreshes (see Section 3.3 of [7]), and
- b) the LSA would have been flooded to X, had Y and X been fully adjacent. As an example of the second property, if Y installs a changed AS-external-LSA, it should not terminate a helping relationship with a neighbor belonging to a stub area, as that neighbor would not see the AS-external-LSA in any case. An implementation MAY provide a configuration option to disable link-state database options from terminating graceful restart. Such an option will, however, increase the risk of transient routing loops and black holes.

When Router Y exits helper mode for X on a given network segment, it reoriginates its LSAs based on the current state of its adjacency to Router X over the segment. In detail, Y takes the following actions:

- a) Y recalculates the Designated Router for the segment,
- b) Y reoriginates its router-LSA for the segment's OSPF area,
- c) if Y is Designated Router for the segment, it reoriginates the network-LSA for the segment and
- d) if the segment was a virtual link, Y reoriginates its router-LSA for the virtual link's transit area.

If Router Y aggregated adjacencies with Router X when entering helper mode (as described in section 3.1), it must also exit helper mode for all adjacencies with Router X when any one of the exit events occurs for an adjacency with Router X.

4. Backward Compatibility

Backward-compatibility with unmodified OSPF routers is an automatic consequence of the functionality documented above. If one or more neighbors of a router requesting graceful restart are unmodified, or if they do not receive the grace-LSA, the graceful restart reverts to a normal OSPF restart.

The unmodified routers will start routing around the restarted router X as it performs initial database synchronization by reissuing their LSAs with links to X omitted. These LSAs will be interpreted by helper neighbors as a topology change, and by X as an LSA inconsistency, in either case, reverting to normal OSPF operation.

5. Unplanned Outages

The graceful restart mechanisms in this memo can be used for unplanned outages. (Examples of unplanned outages include the crash of a router's control software, an unexpected switchover to a redundant control processor, etc). However, implementors and network operators should note that attempting graceful restart from an unplanned outage may not be a good idea, owing to the router's inability to properly prepare for the restart (see Section 2.1). In particular, it seems unlikely that a router could guarantee the sanity of its forwarding table(s) across an unplanned restart. In any event, implementors providing the option to recover gracefully from unplanned outages must allow a network operator to turn the option off.

In contrast to the procedure for planned restart/reloads that was described in Section 2.1, a router attempting graceful restart after an unplanned outage must originate grace-LSAs *after* its control software resumes operation. The following points must be observed during this grace-LSA origination.

- o The grace-LSAs must be originated and be sent *before* the restarted router sends any OSPF Hello Packets. On broadcast networks, this LSA must be flooded to the AllSPFRouters multicast address (224.0.0.5) since the restarting router is not aware of its previous DR state.
- o The grace-LSAs are encapsulated in Link State Update Packets and sent out to all interfaces, even though the restarted router has no adjacencies and no knowledge of previous adjacencies.
- o To improve the probability that grace-LSAs will be delivered, an implementation may send them multiple times (see for example the Robustness Variable in [8]).
- o The restart reason in the grace-LSAs must be set to 0 (unknown) or 3 (switch to redundant control processor). This enables the neighbors to decide whether they want to help the router through an unplanned restart.

6. Interaction with Traffic Engineering

The operation of the Traffic Engineering Extensions to OSPF [4] during OSPF Graceful Restart is specified in [6].

7. Possible Future Work

Devise a less conservative algorithm for graceful restart helper termination that provides a comparable level of black hole and routing loop avoidance.

8. Intellectual Property Rights Notice

The IETF takes no position regarding the validity or scope of any intellectual property 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; neither does it represent that it has made any effort to identify any such rights. Information on the IETF's procedures with respect to rights in standards-track and standards-related documentation can be found in BCP-11. Copies of claims of rights made available for publication 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 implementors or users of this specification can be obtained from the IETF Secretariat.

The IETF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights which may cover technology that may be required to practice this standard. Please address the information to the IETF Executive Director.

9. References

9.1. Normative References

- [1] Moy, J., "OSPF Version 2", STD 54, RFC 2328, April 1998.
- [2] Coltun, R., "The OSPF Opaque LSA Option", RFC 2370, July 1998.

9.2. Informative References

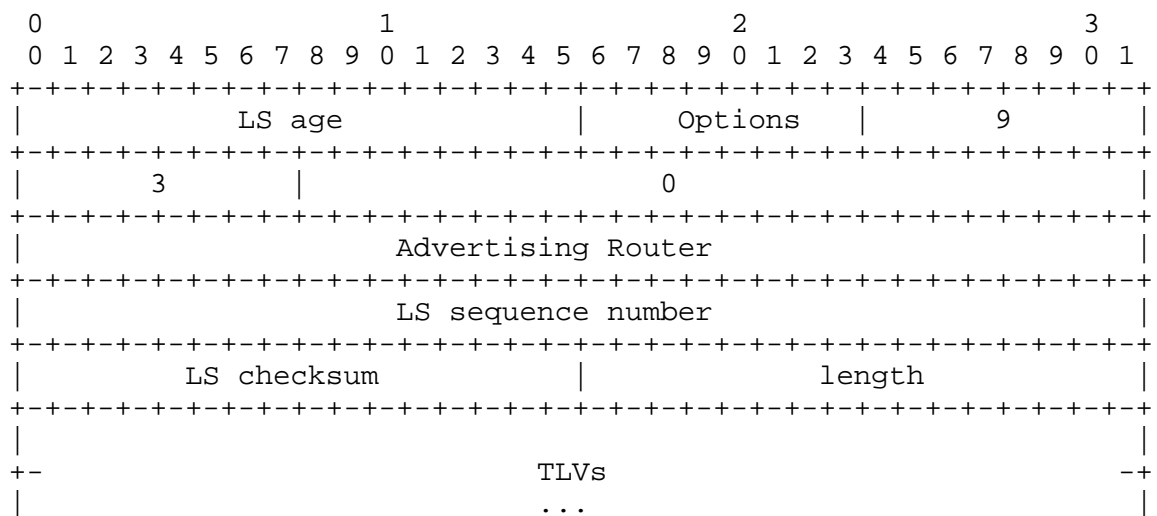
- [3] Murphy, S., Badger, M. and B. Wellington, "OSPF with Digital Signatures", RFC 2154, June 1997.
- [4] Katz, D., Kompella, K. and D. Yeung, "Traffic Engineering (TE) Extensions to OSPF Version 2", RFC 3630, September 2003.

- [5] Murphy, P., "The OSPF Not-So-Stubby Area (NSSA) Option", RFC 3101, January 2003.
- [6] Kompella, K., et al., "Routing Extensions in Support of Generalized MPLS", Work in Progress.
- [7] Moy, J., "Extending OSPF to Support Demand Circuits", RFC 1793, April 1995.
- [8] Cain, B., Deering, S., Kouvelas, I., Fenner, B. and A. Thyagarajan, "Internet Group Management Protocol, Version 3", RFC 3376, October 2002.

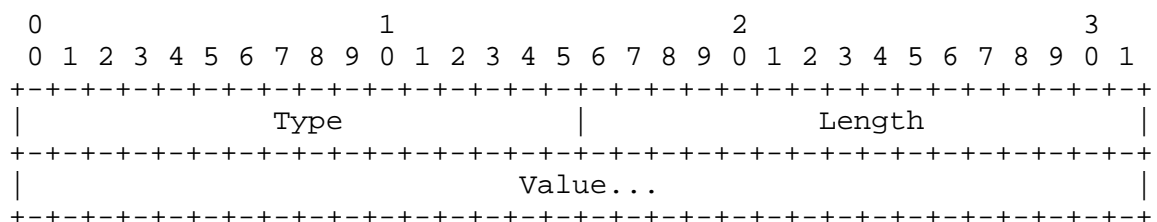
A. Grace-LSA Format

The grace-LSA is a link-local scoped Opaque-LSA [2], having an Opaque Type of 3 and an Opaque ID equal to 0. Grace-LSAs are originated by a router that wishes to execute a graceful restart of its OSPF software. A grace-LSA requests that the router's neighbors aid in its graceful restart by continuing to advertise the router as fully adjacent during a specified grace period.

Each grace-LSA has an LS age field set to 0 when the LSA is first originated; the current value of the LS age then indicates how long ago the restarting router made its request. The body of the LSA is TLV-encoded. The TLV-encoded information includes the length of the grace period, the reason for the graceful restart and, when the grace-LSA is associated with a broadcast, NBMA or Point-to-MultiPoint network segment, the IP interface address of the restarting router.



The format of the TLVs within the body of a grace-LSA is the same as the format used by the Traffic Engineering Extensions to OSPF [4]. The LSA payload consists of one or more nested Type/Length/Value (TLV) triplets. The format of each TLV is:



The Length field defines the length of the value portion in octets (thus a TLV with no value portion would have a length of zero). The TLV is padded to four-octet alignment; padding is not included in the length field (so a three octet value would have a length of three, but the total size of the TLV would be eight octets). Nested TLVs are also 32-bit aligned. For example, a one byte value would have the length field set to 1, and three bytes of padding would be added to the end of the value portion of the TLV. Unrecognized types are ignored.

The following is the list of TLVs that can appear in the body of a grace-LSA:

- o Grace Period (Type=1, length=4). The number of seconds that the router's neighbors should continue to advertise the router as fully adjacent, regardless of the state of database synchronization between the router and its neighbors. Since this time period began when grace-LSA's LS age was equal to 0, the grace period terminates when either:
 - a) the LS age of the grace-LSA exceeds the value of a Grace Period or
 - b) the grace-LSA is flushed. See Section 3.2 for other conditions that terminate graceful restart.

This TLV must always appear in a grace-LSA.

- o Graceful restart reason (Type=2, length=1). Encodes the reason for the router restart as one of the following: 0 (unknown), 1 (software restart), 2 (software reload/upgrade) or 3 (switch to redundant control processor). This TLV must always appear in a grace-LSA.
- o IP interface address (Type=3, length=4). The router's IP interface address on the subnet associated with the grace-LSA. Required on broadcast, NBMA and Point-to-MultiPoint segments, where the helper uses the IP interface address to identify the restarting router (see Section 3.1).

DoNotAge is never set in a grace-LSA, even if the grace-LSA is flooded over a demand circuit [7]. This is because the grace-LSA's LS age field is used to calculate the duration of the grace period.

Grace-LSAs have link-local scope because they only need to be seen by the router's direct neighbors.

Additional Grace-LSA TLVs must be described in an Internet Draft and will be subject to the expert review of the OSPF Working Group.

B. Configurable Parameters

OSPF graceful restart parameters are suggested below. Section B.1 contains a minimum subset of parameters that should be supported. B.2 includes some additional configuration parameters that an implementation may choose to support.

B.1. Global Parameters (Minimum subset)

RestartSupport

The router's level of support for OSPF graceful restart. Allowable values are none, planned restart only, and planned/unplanned.

RestartInterval

The graceful restart interval in seconds. The range is from 1 to 1800 seconds, with a suggested default of 120 seconds.

B.2. Global Parameters (Optional)

RestartHelperSupport

The router's support for acting as an OSPF restart helper. Allowable values are none, planned restart only, and planned/unplanned.

RestartHelperStrictLSAChecking

Indicates whether or not an OSPF restart helper should terminate graceful restart when there is a change to an LSA that would be flooded to the restarting router or when there is a changed LSA on the restarting router's retransmission list when graceful restart is initiated. The suggested default is enabled.

Security Considerations

One of the ways to attack a link-state protocol such as OSPF is to inject false LSAs into, or corrupt existing LSAs in, the link-state database. Injecting a false grace-LSA would allow an attacker to spoof a router that, in reality, has been withdrawn from service. The standard way to prevent such corruption of the link-state database is to secure OSPF protocol exchanges using the cryptographic authentication specified in [1]. An even stronger way of securing link-state database contents has been proposed in [3].

When cryptographic authentication [1] is used on the restarting router the preservation of received sequence numbers in non-volatile storage is not mandatory. There is a risk that a replayed Hello packet could cause neighbor state for a deceased neighbor to be created. However, the risk is no greater than during normal operation.

Acknowledgments

The authors wish to thank John Drake, Vishwas Manral, Kent Wong, and Don Goodspeed for their helpful comments. We also wish to thank Alex Zinin and Bill Fenner for their thorough review.

Authors' Addresses

J. Moy
Sycamore Networks, Inc.
150 Apollo Drive
Chelmsford, MA 01824

Phone: (978) 367-2505
Fax: (978) 256-4203
EMail: jmoy@sycamorenet.com

Padma Pillay-Esnault
Juniper Networks
1194 N, Mathilda Avenue
Sunnyvale, CA 94089-1206

EMail: padma@juniper.net

Acee Lindem
Redback Networks
102 Carric Bend Court
Cary, NC 27519

EMail: acee@redback.com

Full Copyright Statement

Copyright (C) The Internet Society (2003). All Rights Reserved.

This document and translations of it may be copied and furnished to others, and derivative works that comment on or otherwise explain it or assist in its implementation may be prepared, copied, published and distributed, in whole or in part, without restriction of any kind, provided that the above copyright notice and this paragraph are included on all such copies and derivative works. However, this document itself may not be modified in any way, such as by removing the copyright notice or references to the Internet Society or other Internet organizations, except as needed for the purpose of developing Internet standards in which case the procedures for copyrights defined in the Internet Standards process must be followed, or as required to translate it into languages other than English.

The limited permissions granted above are perpetual and will not be revoked by the Internet Society or its successors or assignees.

This document and the information contained herein is provided on an "AS IS" basis and THE INTERNET SOCIETY AND THE INTERNET ENGINEERING TASK FORCE DISCLAIMS 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.

Acknowledgement

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

