

Network Working Group
Internet Draft
Intended status: Standards Track
Expires: September 7, 2011

W. Lu
S. Kini
A. Csaszar
G. Enyedi
J. Tantsura
A. Tian
March 7, 2011

Transport of Fast Notification Messages
draft-lu-fn-transport-00.txt

Status of this Memo

This Internet-Draft is submitted in full conformance with the provisions of [BCP 78](#) and [BCP 79](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), 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/lid-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 September 7, 2011.

Copyright Notice

Copyright (c) 2011 IETF Trust and the persons identified as the document authors. All rights reserved.

This document is subject to [BCP 78](#) and the IETF Trust's Legal Provisions Relating to IETF Documents (<http://trustee.ietf.org/license-info>) in effect on the date of publication of this document. Please review these documents carefully, as they describe your rights and restrictions with respect to this document.

Internet-Draft

Transport of FN Messages

March 2011

Abstract

This document specifies a fast, dataplane-based low-overhead event notification protocol, called Fast Notification (FN) protocol. The draft discusses the design goals, the message container and options for delivering the notifications to all routers within a routing area.

Table of Contents

1.	Introduction.....	3
2.	Design Goals.....	3
3.	Transport Logic - Distribution of the Notifications.....	4
3.1.	Multicast Tree-based Transport.....	5
3.1.1.	Fault Tolerance of a Single Distribution Tree.....	5
3.1.2.	Pair of Redundant Trees.....	5
4.	Message Encoding.....	7
4.1.	Seamless Encapsulation.....	7
4.2.	Dedicated FN Message.....	7
4.2.1.	Authentication.....	9
4.2.1.1.	Areas-scoped and Link-scoped Authentication...10	10
4.2.1.2.	Simple Password Authentication.....10	10
4.2.1.3.	Cryptographic Authentication for FN.....10	10
5.	Security Considerations.....	12
6.	IANA Considerations.....	12
7.	References.....	12
7.1.	Normative References.....	12
7.2.	Informative References.....	12
8.	Acknowledgments.....	13
Appendix A.	Further Options for Transport Logic.....	14
A.1.	Unicast.....	14
A.1.1.	Method.....	14
A.1.2.	Sample Operation.....	15
A.2.	Gated Multicast through RPF Check.....	15
A.2.1.	Loop Prevention - RPF Check.....	16
A.2.2.	Operation.....	16
A.3.	Further Multicast Tree based Transport Options.....	18
A.3.1.	Source Specific Trees.....	18
A.3.2.	A Single Bidirectional Shared Tree.....	18
A.4.	Layer 2 Networks.....	19
Appendix B.	Computing maximally redundant trees.....	20
A.5.	Simple pair of maximally redundant trees in 2-connected	

networks.....	20
A.6. Non-2-connected networks.....	22
A.7. Finding maximally redundant trees in distributed environment	23

[1.](#) Introduction

Draft [[fn-framework](#)] describes the architectural framework to enable fast dissemination of a network event to routers in a limited domain. Existing use cases involve new approaches for IP Fast ReRoute such as [[ipfrr-fn](#)], and faster dissemination of link state information of routing protocols [[ospf-fn](#)] e.g. in order to speed up convergence.

Existing link state routing protocols such as OSPF and ISIS rely on a flooding mechanism to advertise information. These flooding mechanisms are implemented in the control plane, i.e. processed at the control plane in each hop before forwarding to the next hop. A new fast notification protocol should be implemented with low hop-by-hop processing overhead.

This draft proposes a fast notification (FN) protocol as a separate transport layer, whereby applications are proposed in separate drafts.

The functions offered by the FN transport layer include but are not limited to:

- o authentication
- o forwarding (lookup, policy, etc.)
- o topology setup (statically or dynamically configurable)

[2.](#) Design Goals

A low-overhead event notification mechanism which should be used to facilitate extremely quick dissemination of information in a limited area should have the following properties.

1. The mechanism should be fast in the sense that getting the notification packet to remote nodes through possible multiple hops should not require (considerably) more processing at each hop than

plain fast path packet forwarding. Latency due to involving the control plane in each hop is not preferable for some applications.

2. The signaling mechanism should be resilient or, at least, it should offer a delivery option which can disseminate event notifications to all interested nodes even in a network where a single link or a node is down.

3. The mechanism should be trustable; that is, it should provide means to verify the authenticity of the notifications.
4. No fate sharing with control plane, allowing easy graceful restart.
5. The mechanism should be applicable in nodes where control plane is separated from data plane (such as ForCES or OpenFlow forwarding elements).
6. The new protocol should not be dependent upon routing protocol flooding procedures, and should work with either IS-IS or OSPF.

The above reasons point towards a dataplane implementation with the following design goal:

7. The mechanism should involve simple and efficient processing to be feasible for implementation in the dataplane.

The last goal manifests itself in the following ways:

- o Origination of notification should very easy, e.g. creating a simple IP packet, the payload of which can be filled easily by an application.
- o When receiving the packet, it should be easy to recognize by dataplane linecards so that processing can commence after forwarding.
- o No complex operations should be required in order to extract the information from the packet.
- o Duplication of notification packets should be either strictly

bounded or handled without significant dataplane processing burden.

These design goals present a trade-off. A proper balance needs to be found that offers good enough authentication and reliability while keeping processing complexity sufficiently low to be feasible for data plane implementation. This draft attempts to propose such a solution.

[3.](#) Transport Logic - Distribution of the Notifications

The distribution of a notification to multiple receivers can be implemented in various ways. The main body of this draft describes one such option that is especially advantageous: dual redundant

Lu et al.

Expires September 7, 2011

[Page 4]

Internet-Draft

Transport of FN Messages

March 2011

trees. This option is perfect to guarantee that each notification object reaches each node in the domain even with single link or single node failures.

[3.1.](#) Multicast Tree-based Transport

The optimal way of transporting an identical piece of information to several receivers at the same time is to use multicast distribution trees. A tree based transport solution is beneficial since multicast support is already implemented in all forwarding entities, so it is possible to leverage on existing implementations.

With multicast or tree based transport, the Fast Notification (FN) packet can be recognized by a pre-configured or well known destination IP address, denoted by MC-FN in the following, which is the group address of the FN service.

If the FN service is triggered to send out a notification, the notification will originate in a new IP packet, where the destination IP address is set to MC-FN.

[3.1.1.](#) Fault Tolerance of a Single Distribution Tree

The previous solutions and those listed in [Appendix A](#) use a single logical or explicit tree for disseminating a notification from one given source.

If an application requires that the FN transport is tolerant to

single link or node failures, a single distribution tree can be questionable, as a single tree is not redundant: a single failure may cut the tree into two halves, which means that until reconfiguration, a specific notification originating in one half of the tree will not reach nodes in the other part.

For those applications that disseminate failure notifications (as e.g. in [[ospf-fn](#)] or [[ipfrr-fn](#)]) in many failure conditions (e.g. link failure) a failure notification may not be necessary from each of the nodes first detecting the failure as the failure notification from one node may be sufficient to deduce the notification from the other node detecting the failure. This draft lets the applications choose their preferred transport options.

[3.1.2.](#) Pair of Redundant Trees

If an FPN application needs that the exact same data is distributed in the case of any single node or any single link failure, the FPN service should be run in "redundant tree mode".

A pair of "redundant trees" ensures that at each single node or link failure each node still reaches the common root of the trees through at least one of the trees. A redundant tree pair is a known prior-art graph-theoretical object that is possible to find on any 2-node connected network. Even better, it is even possible to find maximally redundant trees in networks where the 2-node connected criterion does not "fully" hold (e.g. there are a few cut vertices) [[Eny2009](#)].

Note that the referenced algorithm(s) build a pair of trees considering a specific root. The root can be selected in different ways, the only thing that is important that each node makes the same selection, consistently. For instance, the node with the highest or lowest router ID can be used.

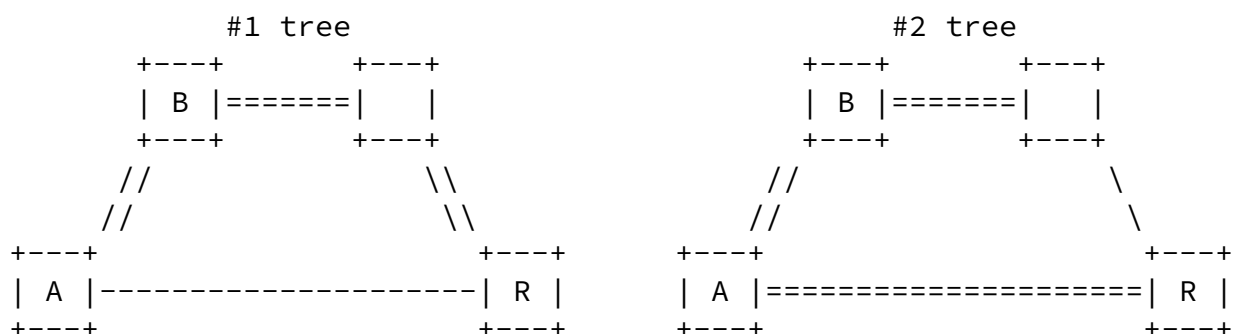




Figure 1 Example: a pair of redundant trees (double lines) of a common root R

There is one special constraint in building the redundant trees. A (maximally) redundant tree pair is needed, where in one of the trees the root has only one child in order to protect against the failure of the root itself. Algorithms presented in [Eny2009] produce such trees. The algorithm is also described in [Appendix B](#) in this specification.

In redundant-tree mode, each node multicasts the requested notification on both trees, if it is possible, but at least along one of the trees. Redundant trees require two multicast group addresses. MC-FN identifies one of the trees, and MC-FN-2 identifies the other tree.

Each node multicast forwards the received notification packet (on the same tree). The root node performs as every other node but in

addition it also multicast the notification on the other tree! I.e. it forwards a replica of the incoming notification in which it replaces the destination address identifying the other multicast distribution tree.

When the network remains connected and the root remains operable after a single failure, the root will be reached on at least one of the trees. Thus, since the root can reach every node along at least one of the trees, all the notifications will reach each node. However, when the root or the link to the root fails, that tree, in which the root has only one child, remains connected (the root is a leaf there), thus, all the nodes can be reached along that tree.

For example, let us consider that in Figure 1 FN is used to disseminate failure information. If link A-B fails, the notifications originating from node B (e.g. reporting that the connectivity from B to A is lost) will reach R on tree #1. Notifications originating from A (e.g. reporting that the connectivity from A to B is lost) will

reach R on tree #2. From R, each node is reachable through one of the trees, so each node will be notified about both events.

4. Message Encoding

4.1. Seamless Encapsulation

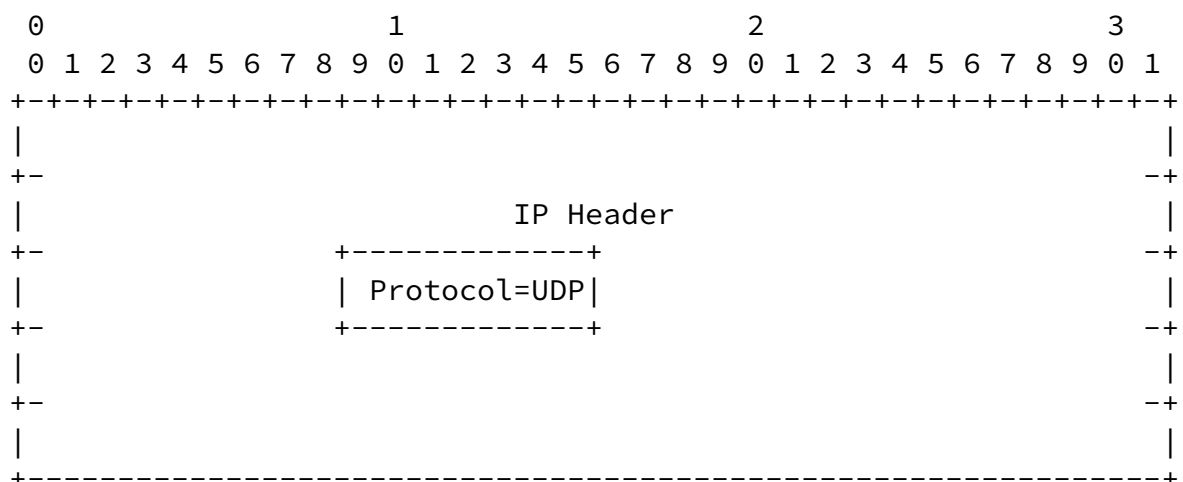
An application may define its own message to distribute quickly. In this case, only the special destination address (e.g. MC-FN) shows that the message was sent using the FN service.

In this case, the entire payload of the IP packet is determined by the application including sequence numbering and authentication. The IP packet's protocol field is also set by the application. The draft [\[ospf-fn\]](#) presents a use case, where the IP protocol field is set to "OSPF" and the payload is entirely handled by OSPF.

4.2. Dedicated FN Message

An alternative option is that the FN messages are distributed in UDP datagrams with well-known port values in the UDP header that need to be allocated by IANA.

The FN packet format inside a UDP datagram is the following:



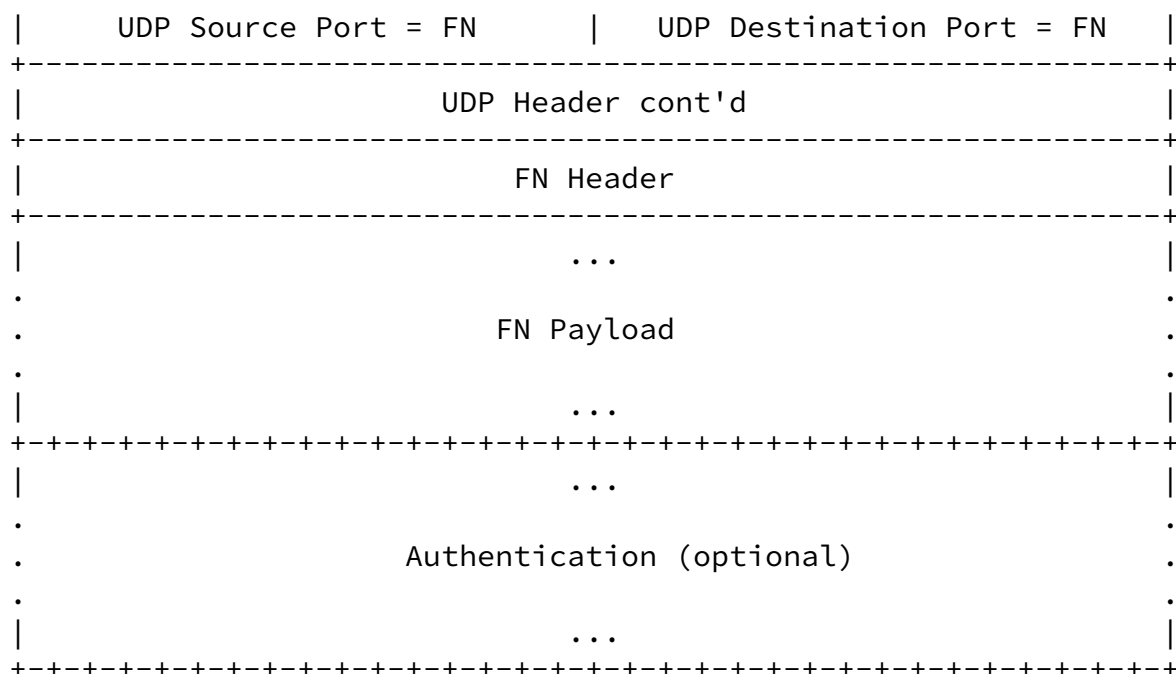


Figure 2 FN packet format as a UDP datagram

The encoding of the FN Header is as follows:

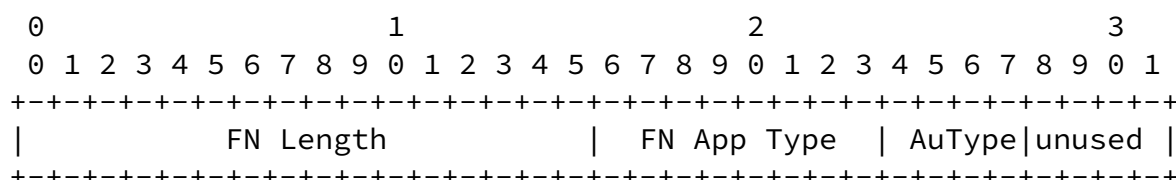


Figure 3 FN Header encoding

FN Length (16 bits)

The length of the FN message in bytes including the FN Header and the FN Payload. The authentication data optionally appended to the FN packet is not actually considered part of the FN message: the authentication data is not included in the FN Length field, although it is included in the length field of the packet's IP header.

FN App Type (8 bits)

Identifies the application which should be the receiver of the notification. A value for each application needs to be assigned by IANA.

AuType

Identifies the authentication procedure to be used for the packet. Authentication options are discussed in [Section 4.2.1](#). of the specification. Note

4.2.1. Authentication

Fast Notification intends to provide a trustable service option, so that receivers of FN packets are able to verify that the packet is sent by an authentic source. Simple password authentication and MD5 authentication is described in the following subsections.

If AuType is set to 0x0, then the FN packet is not carrying an Authentication field at the end of the packet. Note that even in this case the FN application in the payload may still use its own authentication mechanism.

If AuType is non null, an Authentication field must be appended after the FN message. The encoding of this field is as follows:

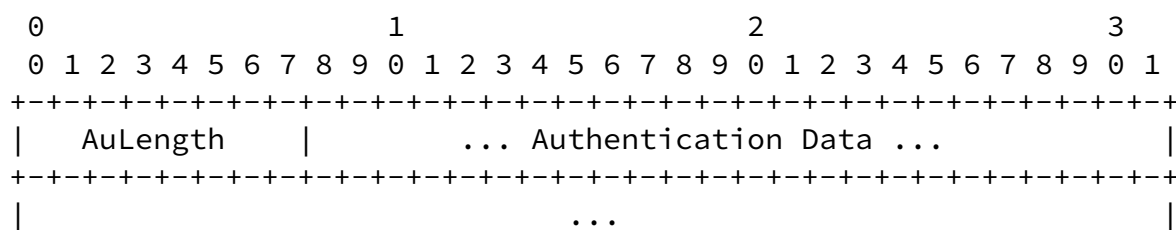


Figure 4 Authentication field in FN packets

AuLength

Describes the length of the entire Authentication field in bytes.

4.2.1.1. Areas-scoped and Link-scoped Authentication

Since FN is a solution to disseminate an event notification from one

source to a whole area of nodes, the simplest approach would be to use per-area authentication, i.e. a common password, a common pre-shared key among all nodes in the area as described in the following sub-sections.

Carriers may, however, prefer per-link authentication. In order not to lose the speed (simple per-hop processing, fast forwarding property) of FN, link-scoped authentication is suggested only if the dataplane supports it, i.e. if there is hardware support to verify and re-generate authentication hop-by-hop. In such cases, the operator may need to configure a common pre-shared key only on routers connected by the same link. It is even possible that there is no authentication on some links considered safe.

[4.2.1.2](#). Simple Password Authentication

Simple password authentication guards against routers inadvertently joining the routing area; each router must first be configured with a password before it can participate in Fast Notification.

The password is stored in the Authentication Data field. AuLength is set to the length of the password in bytes plus 1. Two AuType values for simple password authentication need to be allocated by IANA: one for area-scope and another for link-scoped.

With per-link authentication mode, the Authentication field must be stripped and regenerated hop-by-hop.

Simple password authentication, however, can be easily compromised as anyone with physical access to the network can read the password.

[4.2.1.3](#). Cryptographic Authentication for FN

Using this authentication type, a shared secret key is configured in all routers subscribed to the FN service. The key is used to generate/verify a "message digest" that is appended to the end of the FN packet. The message digest is a one-way function of the FN packet and the secret key. This authentication mechanism resembles the cryptographic authentication mechanism of [\[OSPF\]](#).

The packet signature is created by an MD5 hash performed on an object which is the concatenation of the FN message, including the FN header, and the pre-shared secret key. The resulting 16 byte MD5

message digest is appended to the FN message into the Authentication field as shown below.

The AuType in the FN header is set to indicate cryptographic authentication, the specific value is to be assigned by IANA both for area-scoped and for link-scoped versions.

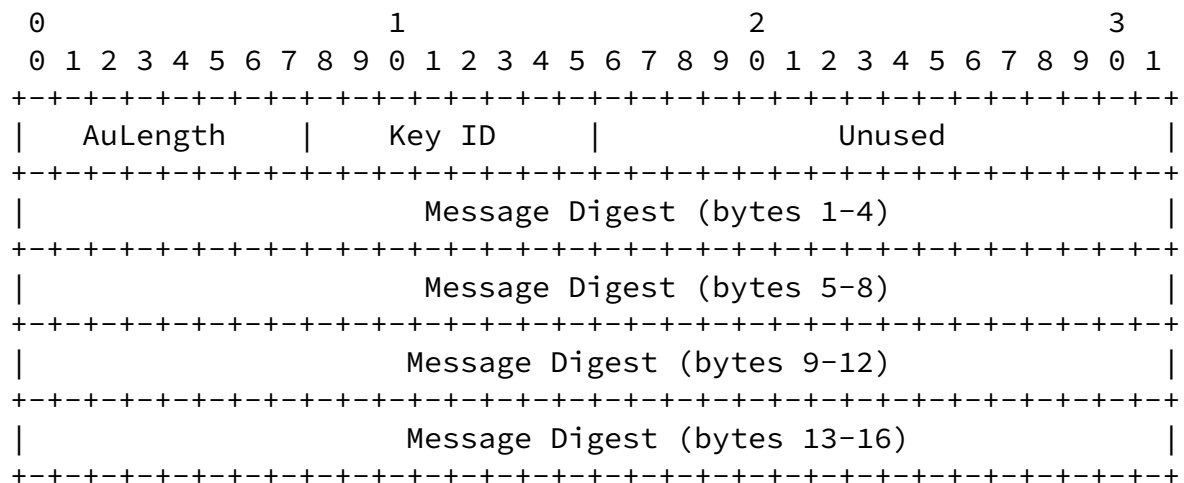


Figure 5 Authentication field in FN packets with MD5 cryptographic authentication.

AuLength

AuLength is set to 20 bytes.

Key ID

This field identifies the algorithm and secret key used to create the message digest appended to the FN packet. This field allows that multiple pre-shared keys may exist in parallel.

Message Digest

The 16 byte long MD5 hash performed on an object which is the concatenation of the FN message, including the FN header, and the pre-shared secret key identified by Key ID.

When receiving an FN message, if the FN header indicates MD5 authentication, then the last 20 bytes of the FN message are set aside. The recipient data plane element calculates a new MD5 digest of the remainder of the FN message to which it appends its own known secret key identified by Key ID. The calculated and received digests are compared. In case of mismatch, the FN message is discarded.

In per-link authentication mode, the Authentication field must be stripped and regenerated hop-by-hop using the key of the outgoing link.

[5. Security Considerations](#)

This draft has described basic optional procedures for authentication. The mechanism, however, does not protect against replay attacks.

If applications of FN require protection against replay attacks, then these applications should provide their own specific sequence numbering within the FN payload. Recipient applications should accept FN messages only if the included sequence number is what they expected.

Since the message digest of cryptographic authentication also covers the payload, even if an attacker knew how to construct the new sequence number, it would not be able to generate a correct message digest without the pre shared key. This way, a sequence number in the payload combined with FN's cryptographic authentication offers sufficient protection against replay attacks.

[6. IANA Considerations](#)

An IP protocol value needs to be assigned by IANA for FN. IANA also needs to maintain values for FN App Type as applications are being proposed.

Multicast addresses used for the distribution trees are either allocated by IANA or they can be a configuration parameter within the local domain.

[7. References](#)

[7.1. Normative References](#)

[BIDIR-PIM] M. Handley et al., "Bidirectional Protocol Independent Multicast (BIDIR-PIM)", [RFC 5015](#), IETF, 2007

[7.2. Informative References](#)

[Eny2009] Gabor Enyedi, Gabor Retvari, Andras Csaszar, "On Finding Maximally Redundant Trees in Strictly Linear Time", IEEE Symposium on Computers and Communications (ISCC), 2009.

[fn-framework] W. Lu, A. Tian, S. Kini, "Fast Notification

Framework", [draft-lu-fast-notification-framework-00](#)

[ipfrr-fn] A. Csaszar, G. Enyedi, S. Kini, "IP Fast Re-Route with Fast Notification", [draft-csaszar-ipfrr-fn-00](#)

Lu et al.

Expires September 7, 2011

[Page 12]

Internet-Draft

Transport of FN Messages

March 2011

[OSPF] J. Moy, "OSPF Version 2", IETF [RFC 2328](#), 1998

[ospf-fn] S. Kini, W. Lu, A. Tian, "OSPF Fast Notifications", [draft-kini-ospf-fast-notification-00](#)

8. Acknowledgments

The authors owe thanks to Acee Lindem and Joel Harpern for their review and comments.

This document was prepared using 2-Word-v2.0.template.dot.

[Appendix A.](#) Further Options for Transport Logic

The options described in this appendix represent alternative solutions to the redundant tree based approach described in [Section 3.1.2](#).

It is left for WG discussion and further evaluation to decide whether any of these options should potentially be preferred instead of redundant trees.

[A.1.](#) Unicast

This method addresses the need in a unique way. It has the following properties:

- o Extremely simple, without the need of any data plane change or cooperation;
- o Short turnaround time (i.e. ready for next hit);
- o 100% link break coverage (may not work in certain node failure cases);
- o Little change to OSPF (need encapsulation for IS-IS).

[A.1.1.](#) Method

The method is simple in design, easy to implement and quick to deploy. It requires no topology changes or specific configurations. It adds little overhead to the overall system.

The method is literally sending the event message to every router in the domain in form of IP packet. This appears burdensome to the sending router which has to duplicate the packet sending effort many

times. Practical experience has shown, however, that the amount of effort is not a big concern in reasonable sized networks.

Normal flooding (regular or fast) process requires a router to duplicate the packet to all flooding eligible interfaces. All routers have to be fast-flooding-aware. This implies new code to every router in control plane and/or data plane.

The method uses a different approach. It takes advantage of the given routing/forwarding table in each router in the IP domain. The originating router of the flooding information simply sends multiple copies of the packet to each and every router in the domain. These packets are forwarded to the destination routers at data plane speed,

just like the way the regular IP data traffic is handled. No special handling in any other routers is needed.

This small delay on the sender can be minimized by pre-downloading the link-broken message packets to the data plane. Since the data plane already has the list of all routers which are part of the IGP routing table, the data plane can dispatch the packet directly.

By essence, the flooding in this method is tree based, just like a multicast tree. The key is that no special tree is generated for this purpose; the normal routing table which is an SPF tree (SPT) plays a role of the flooding tree. This logic guarantees that the flooding follows the shortest path and no flooding loop is created.

[A.1.2.](#) Sample Operation

1.1. depicts a scenario where router A wants to flood its message to all other routers in the domain using the unicast flooding method.

Instead of sending one packet to each of its neighbor, and let the neighbor to flood the packet further, router A directly send the same packet to each router in the domain, one at a time. In this sample network, router A literally sends out 5 packets.

```
A---B---C---D
 \
  --E---F
```


1. Packet(A->B);
2. Packet(A->C);
3. Packet(A->D);
4. Packet(A->E);
5. Packet(A->F).

Figure 1 Multiple Unicast Packets

The unicast flooding procedure is solely controlled by the sending router. No action is needed from other routers other than their normal forwarding functionalities. This method is extremely simple and useful for quick prototyping and deployment.

[A.2.](#) Gated Multicast through RPF Check

This method fulfills the purpose with the following characters:

Lu et al.	Expires September 7, 2011	[Page 15]
-----------	---------------------------	-----------

Internet-Draft	Transport of FN Messages	March 2011
----------------	--------------------------	------------

1. No need to build the multicast tree. It is the same as the SPT computed by the IGP routing process;
2. The flooding loop is prevented through RPF Check.

The method has all the benefit of multicast flooding. It, however, does not require running multicast protocol to setup the multicast tree. The unicast shortest path tree is leveraged as a multicast tree.

[A.2.1.](#) Loop Prevention - RPF Check

In this mechanism, the distribution tree is not explicitly built; rather each node will first do a Reverse Path Forwarding (RPF) check before it floods the notification to other links.

A special multicast address is defined and is subject to IANA approval. This address is used to qualify the notification packet for fast flooding. When a notification packet arrives, the receiving node will perform an IP unicast routing table lookup for the originator IP address of the notification and find the outgoing interface. Only when the arriving interface of the notification is the same as the outgoing interface leading towards the originator IP address, the notification will be flooded to other interfaces.

IP Multicast forwarding with RPF check is available on most of the routing/switch platforms. To support flooding with RPF check, a special IP multicast group must be used. A bi-directional IP multicast forwarding entry is created that consists of all interfaces with in the flooding scope, typically an IGP area.

[A.2.2.](#) Operation

The Gated flooding operation is illustrated in 1.1. .

All Routers, IGP Process:

```
if (SPT ready) {  
    duplicate the SPT as Bidir_Multicast_tree;  
    download the multicast_tree to data plane;  
}  
add FNF_multicast_group_addr;
```

Sender of the FNF notification:

```
if (breakage detected) {  
    pack the notification in a packet;  
    send the packet to the FNF_multicast_group_addr;  
}
```

Receiver of the FNF notification:

```
if (notification received) {  
    if (RPC_interface == incoming_interface) {  
        multicast the notification to all other interfaces;  
    }  
}
```

```

        forward the notification to IGP for processing;
    }

```

Figure 2 Gated flooding operation

1.1. shows a sample operation on a four-router mesh network. The left figure is the topology. The right figure is the shortest path tree rooted at A.

Router A initiates the flooding. But the downstream routers B, C, and D will drop all messages except the ones that come from their shortest path parent node. For example, A's message to C via B is dropped by C, because C knows that its reverse path forwarding (RPF) nexthop is A.

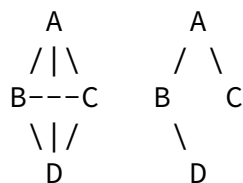


Figure 3 Loop Prevention through the RPF check

[A.3.](#) Further Multicast Tree based Transport Options

[A.3.1.](#) Source Specific Trees

One implementation option is to rely on source specific multicast. This means, that even though there is only a single multicast group address (MC-FN) allocated to the FN service, the FIB of each router is configured with forwarding information for as many trees as many FN sources (nodes) there are in the routing area, i.e. to each (S_i,MC-FN) pair.

[A.3.2.](#) A Single Bidirectional Shared Tree

In the previous solution each source specific tree is a spanning tree. It is possible to reduce the complexity of managing and

configuring n spanning trees in the area by using bidirectional shared trees. By building a bidirectional shared tree, all nodes on the tree can send and receive traffic using that single tree. Each sent packet from any source is multicasted on the tree to all other receivers.

The tree must be consistently computed at all routers. For this, the following rules may be given:

The tree can be computed as a shortest path tree rooted at e.g. the highest router-id. When multiple paths are available, the neighbouring node in the graph e.g. with highest router-id can be picked. When multiple paths are available through multiple interfaces to a neighbouring node, e.g. a numbered interface may be preferred over an unnumbered interface. A higher IP address may be preferred among numbered interfaces and a higher ifIndex may be preferred among unnumbered interfaces.

Note, however, that the important point is that the rules are consistent among nodes. That is, a router may pick the lower router IDs if it is ensured that ALL routers will do the same to ensure consistency.

Multicast forwarding state is installed using such a tree as a bi-directional tree. Each router on the tree can send packets to all other routers on that tree.

Note that the multicast spanning tree can be built using [[BIDIR-PIM](#)] so that each router within an area subscribes to the same multicast group address. Using BIDIR-PIM in such a way will eventually build a multicast spanning tree among all routers within the area. (BIDIR-PIM

is normally used to build a shared, bidirectional multicast tree among multiple sources and receivers.)

[A.4](#). Layer 2 Networks

Layer 2 (e.g. Ethernet) networks offer further options for distributing the notification (e.g. using spanning trees offered by STP). Definition of these is being considered and will be included in a future revision of this draft.

[Appendix B.](#)

Computing maximally redundant trees

Here we describe a possible, not optimal, way of computing maximally redundant trees. First, we suppose that the network is 2-connected and that we have a central processor computing the trees, then we lift these assumptions.

[A.5.](#) Simple pair of maximally redundant trees in 2-connected networks

Finding a simple pair of maximally redundant trees in a 2-connected network is quite simple. We call a node "ready", if it was already added to the trees. Initially, the only ready node is the common root (node r in the sequel).

When we have at least one node x in the network, which is not ready, find two node-disjoint paths from x either to r or to two distinct ready nodes. Since the network is 2-connected, there are always two node-disjoint paths from x to r . It is possible that one or both of these paths reaches another ready node sooner than r , in which case we have the two node-disjoint paths to distinct nodes. Combining the undirected links of these paths makes up an *ear*.

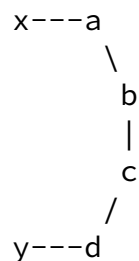


Figure 6 An *ear* connected to node x and y (x and y are ready)

Let x and y be the two ready endpoints of an ear, and first suppose that they are different nodes and none of them is r . Note that both x and y are in the two trees (since they are "ready") and if x is an ancestor of y in the first tree (x is on the path from y to r), then x cannot be the ancestor of y along the second tree at the same time. Thus, it is safe to connect the nodes of the freshly found ear to x in the first tree and to y in the second tree, if either x is an ancestor of y in the first tree, or y is an ancestor of x in the second tree. Considering the example in Figure 6, this means that links $d-c-b-a-x$ should be added to the first tree, and $a-b-c-d-y$ should be added to the second one.

In the case, when either $x=r$ or $y=r$ or when neither x is an ancestor of y nor y is an ancestor of x in any of the trees, the endpoints are

not firmly bound to one of the trees, it is only important to put the

links to one endpoint in one of the trees and put the links towards the other endpoint to the other tree. In our example this means that either d-c-b-a-x or a-b-c-d-y could be added to the first tree. Naturally, then the other endpoint must be selected for the second tree.

In order to protect against the failure of the root r , we need to construct (maximally) redundant trees, where there is only one edge entering to the root on one of the trees. This makes the root a leaf in that tree. To achieve this, we should add the ear to the second tree through r only if both endpoints are r . Moreover, we need to select an ear with different endpoints when it is possible.

Finding an ear is relatively simple and can be done in different ways. Probably the simplest way is to find a ready node q (q is not the root) with a non-ready neighbor w , (virtually) remove q from the topology, and to find a path from w to r ; since the network is 2-connected, such a path either reaches r , or reach another ready node. Moreover, when only r is ready such a node q does not exist, so we select one of r 's neighbors as w , and remove not r but the link between them.

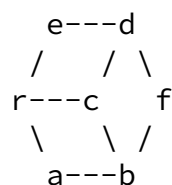


Figure 7 A 2-connected network

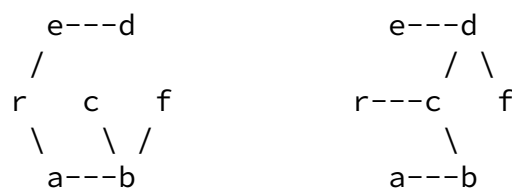


Figure 8 The two maximally redundant trees found in the network depicted in Figure 7.

Now, a simple example is in order. Consider the network depicted in Figure 7, and suppose that the common root is node r . We have only r

in the trees, so we select one of its neighbors, let it be *a*, remove the link between them, and select a path (let it be the shortest one) from *a* to *r*; this path is *a-b-c-r*, so the ear is *r-a-b-c-r*. Since both endpoints of the ear are *r*, selecting the right tree is not important, e.g., we can add *c-b-a-r* to the first tree, and *a-b-c-r* to the second one (Figure 8). This way, *r*, *a*, *b* and *c* form the set of "ready" nodes. From the ready set, *c* and *d* are not the root and have non-ready neighbors. Let us select, e.g., *c*. The shortest path from *d* to *r* when *c* is removed is *d-e-r*, so we have ear *c-d-e-r*, we add *d-e-r* to the first tree and *e-d-c* to the second one (recall that we do not want to create a new neighbor for *r* in the second tree). Finally, the last non-ready node is *f*, and the ear is *b-f-d*. Since neither is *b* an ancestor of *d* nor is *d* an ancestor of *b* in any of the trees, we can connect *f* to the trees in both ways. E.g., add *f-b* to the first tree, and *f-d* to the second one.

[A.6.](#) Non-2-connected networks

When, however, the network is not 2-connected, it is not always possible to find a pair of node-disjoint paths from any node *x* to root *r*, which makes our previous algorithm unable to find the trees. However, while the network is connected, it is made up by 2-connected components bordered by "cut-vertices" (naturally, some of these components may contain only one node). A node is a cut-vertex, if removing that node splits the network into two.

A simple algorithm to find the components and the cut-vertices can be to (virtually) remove each vertex one by one, and check connectivity with BFS or DFS. Moreover, nodes *a* and *b* are in the same 2-connected component, if *a* remains reachable from *b* after removing any single node. Note that linear time algorithms do exist that find both the 2-connected components and the cut-vertices.

Now, we can build up redundant trees in each component. In components containing *r*, the root of such trees must be *r*. Otherwise, in the remaining components the root must be the last node in the component along a path to the root. Recall, that this must be a cut-vertex, so it is the same for each path emanating from that component.

At this point, we are ready, if there is no cut-edge in the network. However, if some 2-connected components are connected by a cut-edge, we must add that edge to both of the trees.

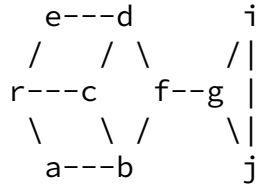


Figure 9 Non-2-connected network

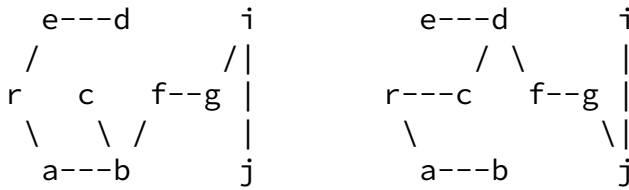


Figure 10 The two maximally redundant trees found in the network depicted Figure 9.

As an example consider the network depicted in Figure 9. Observe that now we have two 2-connected components, one contains r, a, b, c, d, e, f and the other contains g, i, j . Moreover, these components have no common node, they are connected with a cut-edge.

Finding the trees in the component containing r is simple, these trees are the same as previously. Moreover, the other component is a cycle, so it will be covered by a single ear. Finally we must add link $f-g$ to both of the trees, to get the trees depicted in Figure 10.

[A.7.](#) Finding maximally redundant trees in distributed environment

If we need to compute exactly the same maximally redundant trees at each of the routers, consistency needs to be ensured by tie-breaking mechanisms. Observe that the previous algorithm has multiple choices when it selects how to connect nodes to the trees when only r is ready, how to select ready node q and non-ready node w for a later ear and when neither of the endpoints is an ancestor of the other one.

All the previous problems can easily be handled. E.g., the first ear should be connected in such a way, that the neighbor of r with the lowest ID must be directly connected to r in the first tree.

Moreover, later we should choose ready router with non-ready neighbor as q and its non-ready neighbor with the lowest ID as w . Finally, when neither of the endpoint is an ancestor of the other one, connect the ear to the endpoint with the lower ID in the first tree.

Authors' Addresses

Wenhu Lu
Ericsson
300 Holger Way, San Jose, CA 95134
Email: wenhu.lu@ericsson.com

Sriganesh Kini
Ericsson
300 Holger Way, San Jose, CA 95134
Email: sriganesh.kini@ericsson.com

Andras Csaszar
Ericsson
Irinyi utca 4-10, Budapest, Hungary, 1117
Email: Andras.Csaszar@ericsson.com

Gabor Sandor Enyedi
Ericsson
Irinyi utca 4-10, Budapest, Hungary, 1117
Email: Gabor.Sandor.Enyedi@ericsson.com

Jeff Tantsura
Ericsson
300 Holger Way, San Jose, CA 95134
Email: jeff.tantsura@ericsson.com

Albert Tian
Ericsson
300 Holger Way, San Jose, CA 95134
Email: albert.tian@ericsson.com

